

```
In [6]: #Title & Objective - Exploratory Data Analysis (EDA) Extract insights using visual and statistical exploration.
#Dataset (source, rows x cols) - titanic.csv 893*12
#Imports & settings
import pandas as pd
```

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import missingno as mso
from scipy import stats
from scipy.stats import chi2_contingency, ttest_ind
import statsmodels.api as sm
from statsmodels.stats.outliers_influence import variance_inflation_factor
import os

# Load data & quick plots
# Change file path as needed
data_path = "C:\\Users\\David\\OneDrive\\Documents\\Levante Lab\\tasks\\task1\\titanic_tasks.csv"
df = pd.read_csv(data_path)
```

```
display(df$info)
display(df$describe(include='all'),7)
```

```
Shape: (891, 12)
```

	PassengerId	Survived	Pclass		Name	Sex	Age	SibS
0	1	0	3		Braund, Mr. Owen Harris	male	22.0	
1	2	1	1		Cummings, Mrs. John Bradley Florence Briggs Toth	female	38.0	
2	3	1	3		Heikkinen, Mrs. Laina	female	26.0	
3	4	1	1		Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	
4	5	0	3		Allen, Mr. William Henry	male	35.0	

```
<class 'pandas.core.frame.DataFrame'>
```

```
Rankings: 891 entries, 8 to 10
```

```

0 PassengerId 891 non-null int64
1 Survived 891 non-null int64
2 Pclass 891 non-null int64
3 Name 891 non-null object
4 Sex 891 non-null object
5 Age 714 non-null float64
6 SibSp 891 non-null int64
7 Parch 891 non-null int64
8 Ticket 891 non-null object
9 Fare 891 non-null float64
10 Cabin 204 non-null object
11 Embarked 889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 18.7+ MB
None

```

```

      count unique      top freq      mean      std      min      25%      50%      75%      max
PassengerId 891.0 NaN NaN NaN 446.0 257.353842 1.0 223.5 446.0 668.5 891.0
Survived 891.0 NaN NaN NaN 0.363838 0.486592 0.0 0.0 0.0 1.0 1.0
Pclass 891.0 NaN NaN NaN 2.308642 0.836071 1.0 2.0 3.0 3.0 3.0
Name 891 891 Braund, Mr. Owen Harris 1 NaN NaN NaN NaN NaN NaN NaN NaN
Sex 891 2 male 577 NaN NaN NaN NaN NaN NaN NaN NaN
Age 714.0 NaN NaN NaN 29.691118 14.524947 0.42 20.0 28.0 38.0 80.0
SibSp 891.0 NaN NaN NaN 0.523008 1.102743 0.0 0.0 0.0 1.0 8.0
Parch 891.0 NaN NaN NaN 0.381594 0.806057 0.0 0.0 0.0 0.0 6.0
Ticket 891 681 347082 7 NaN NaN NaN NaN NaN NaN NaN NaN
Fare 891.0 NaN NaN NaN 32.264208 49.693429 0.0 7.9104 14.4542 31.0 512.3292
Cabin 204 147 896 898 4 NaN NaN NaN NaN NaN NaN NaN NaN
Embarked 889 3 S 644 NaN NaN NaN NaN NaN NaN NaN NaN

```

```

In [12]: #Clean column names & duplicates
# Remove 1 column names
df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_')

# Drop exact duplicates (if any)
df.drop_duplicates(inplace=True)
print('Duplicate rows (drop count)')
df = df.drop_duplicates()
print('Duplicate rows: 0')

```

```

In [14]: #missing values overview
missing = df.isnull().sum().sort_values(ascending=False)
missing_count = df.isnull().sum()['Parch'] * 100 / df.values.shape[0]
pd.concat([missing, missing_percent], axis=1, keys=['missing_count', 'missing_percent'])

```

```

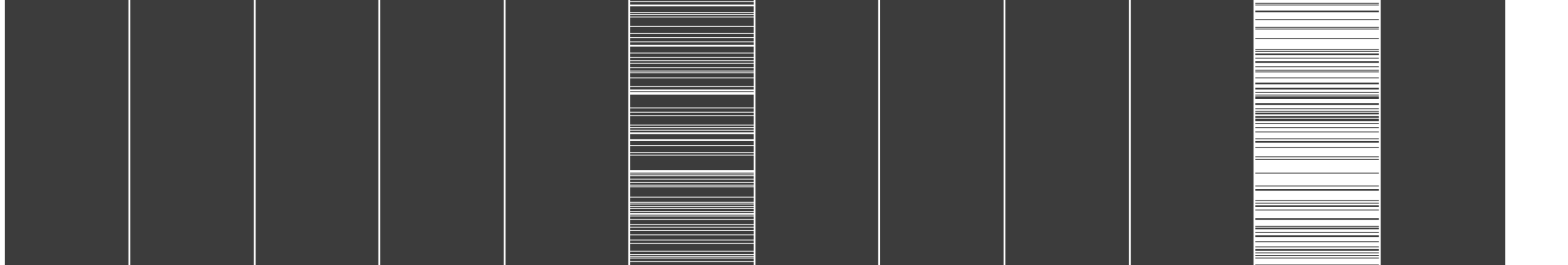
Out[14]:
      missing_count  missing_percent
cabin            687      77.164377
age             177     19.865320
embarked         2      0.224667
passengerid      0      0.000000
survived          0      0.000000
pclass           0      0.000000
name             0      0.000000
sex              0      0.000000
parch            0      0.000000
ticket           0      0.000000
fare             0      0.000000

```

```

plt.savefig('C:\\Users\\{diyaashree}\\OneDrive\\Documents\\elevate lab
plt.show()

```



```

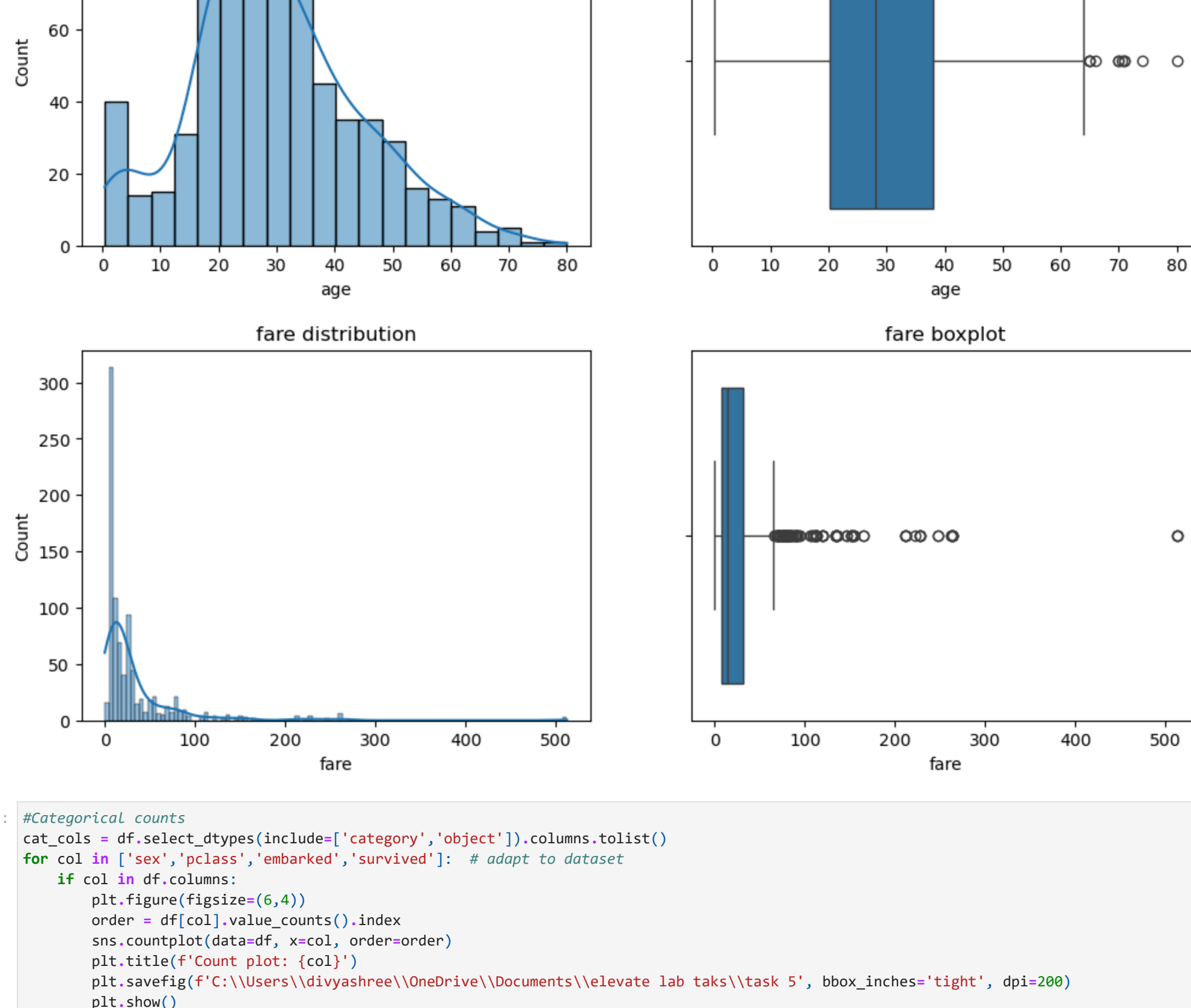
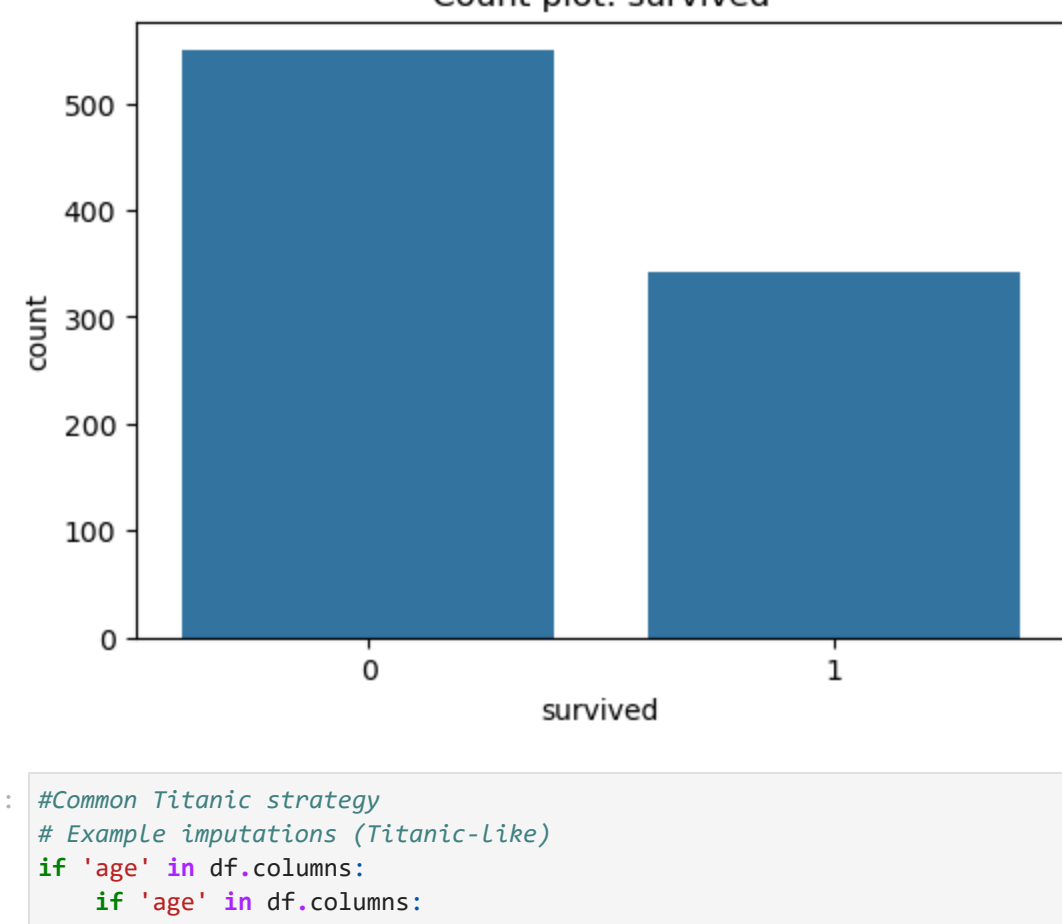
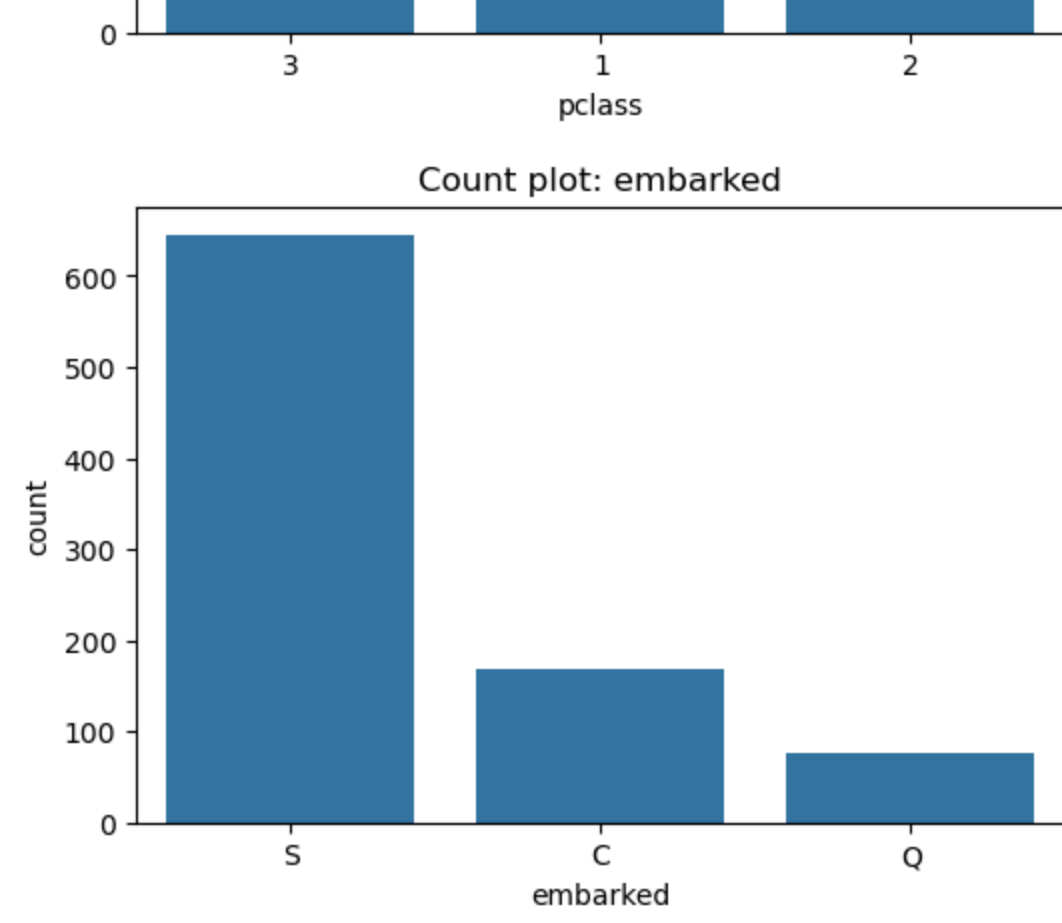
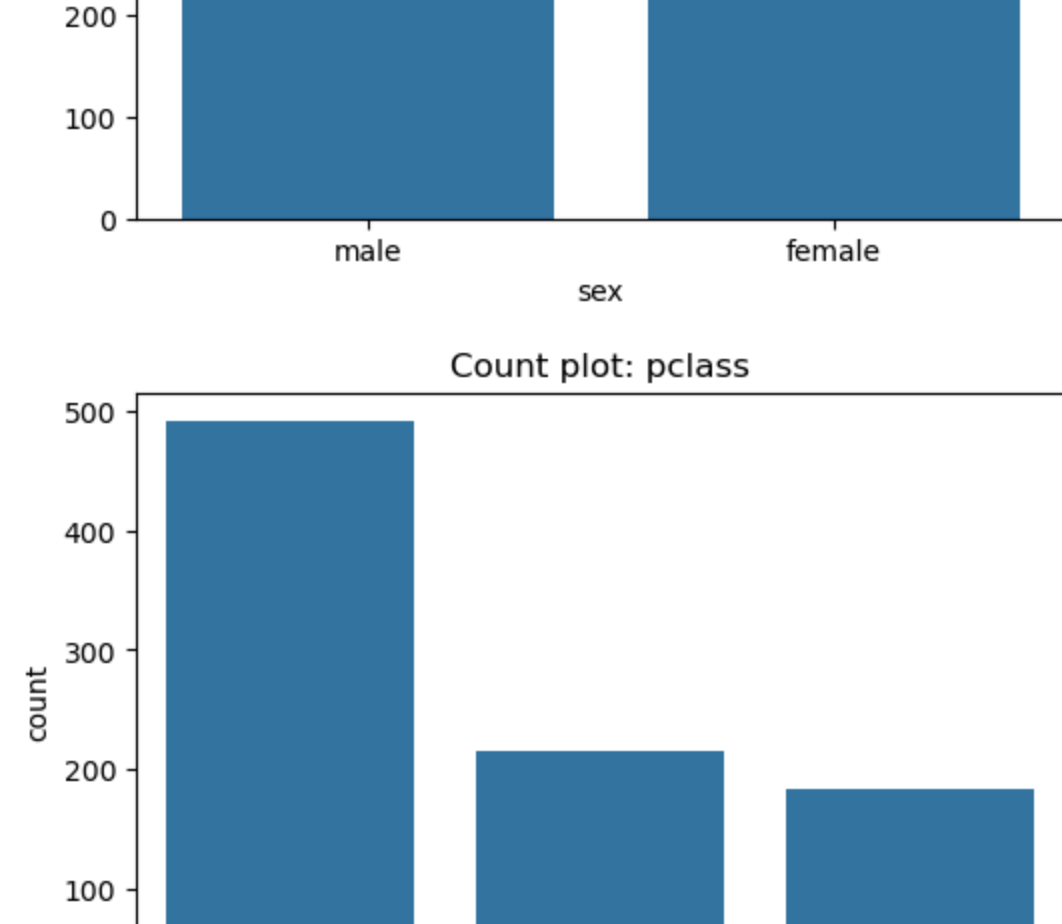
In [26]: #Convert dtyypes (dates, categories)
          # Convert categorical-looking columns to category dtype
          For col in df.select_dtypes(include='object').columns:
              if df[col].nunique() > 30: # heuristic
                  df[col] = df[col].astype('category')

In [30]: #Numeric Histograms v. histograms
          num_cols = df.select_dtypes(include=[np.number]).columns.tolist()
          for col in ['age', 'year']: # adapt to columns present
              if col in num_cols:
                  plt.figure(figsize=(12,4))
                  plt.subplot(1,2,1)
                  sns.histplot(df[col].dropna(), kde=True)
                  plt.title(f'{col} distribution')
  
```

```

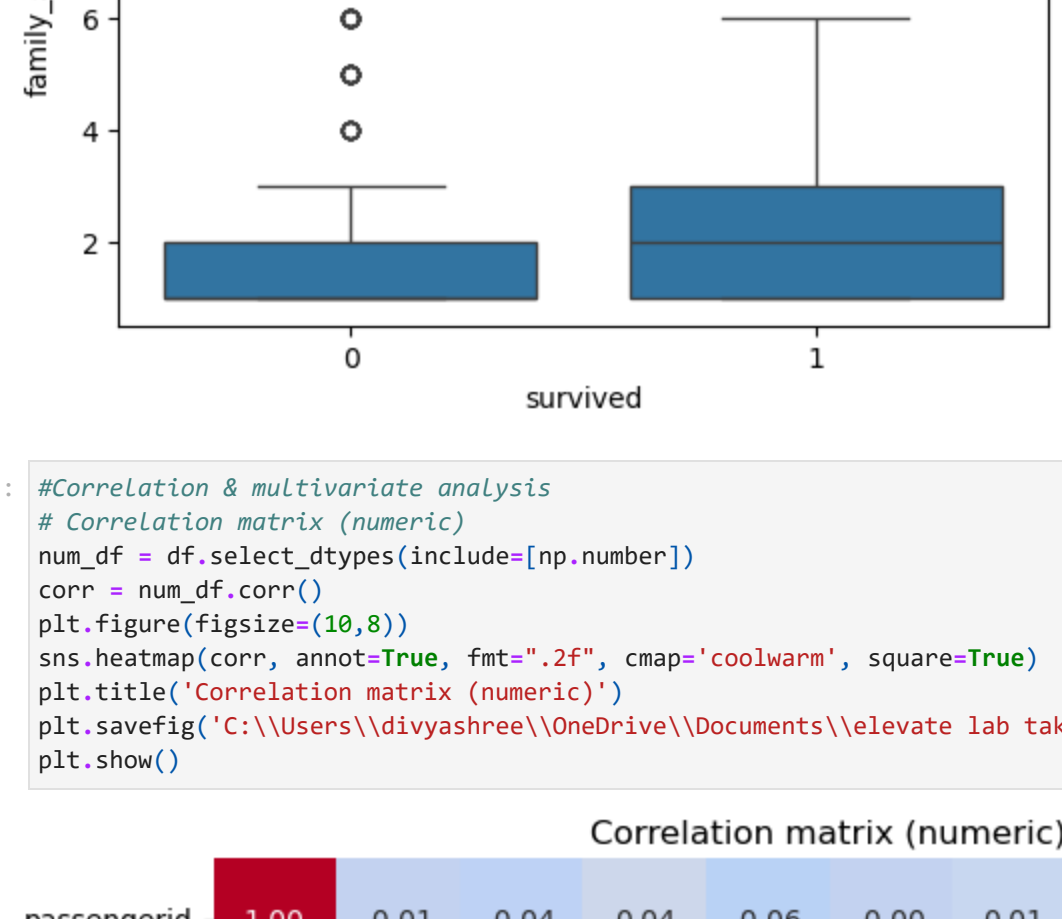
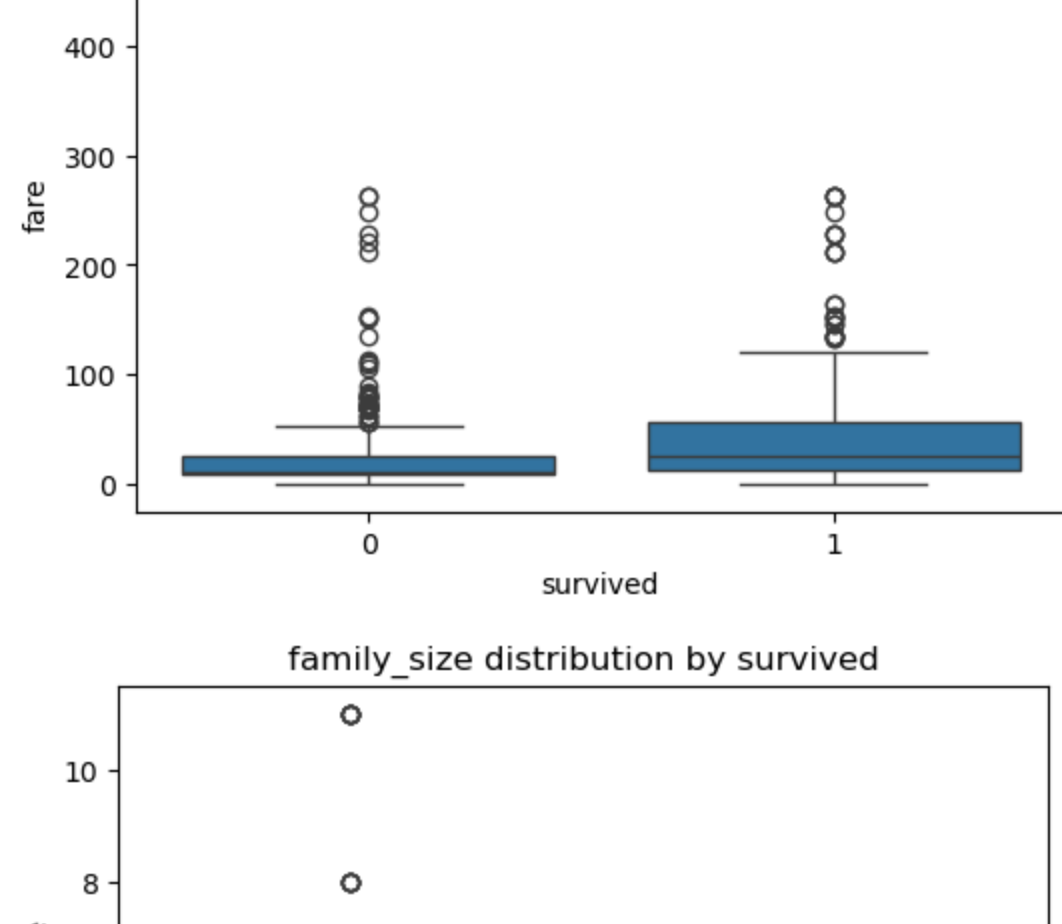
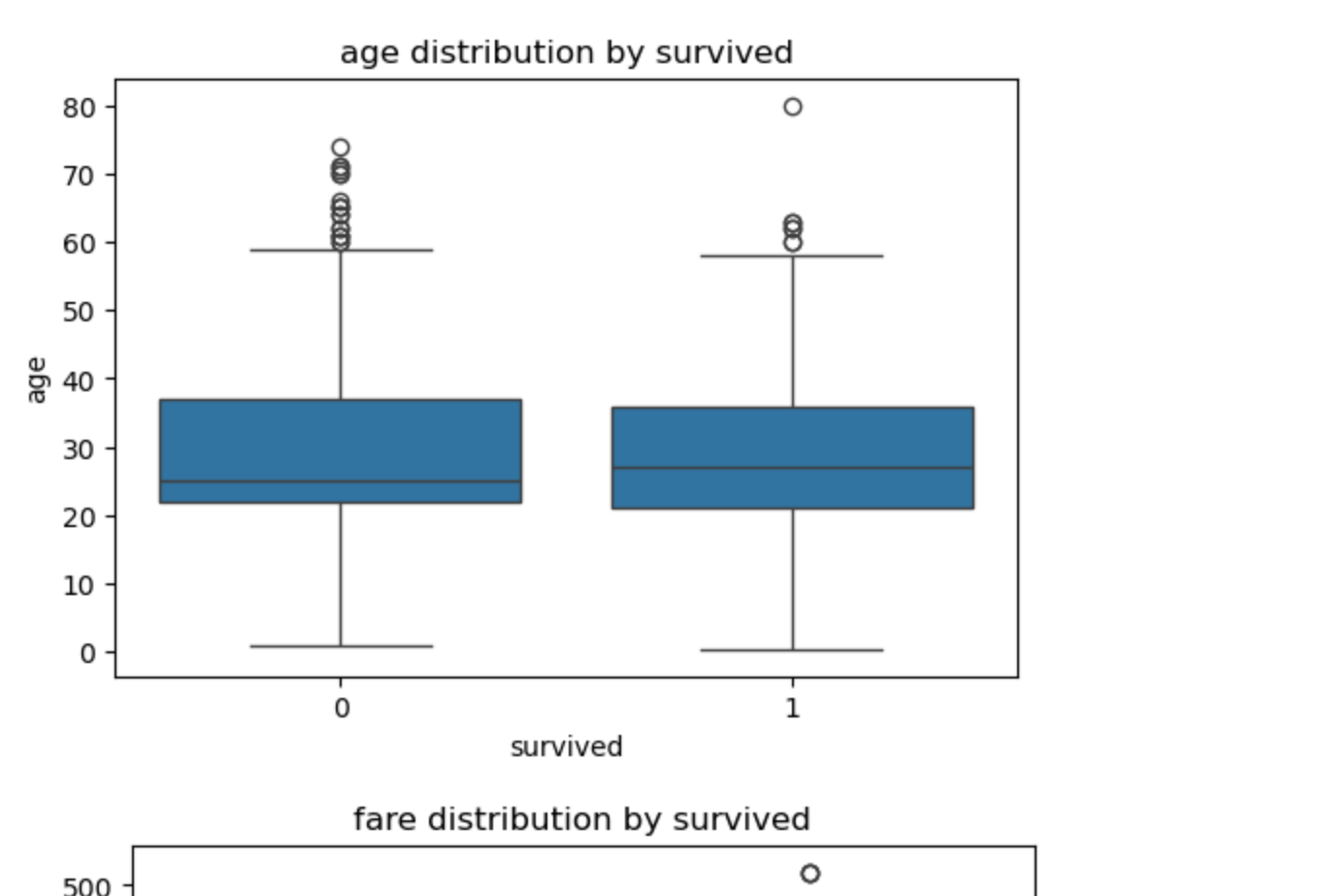
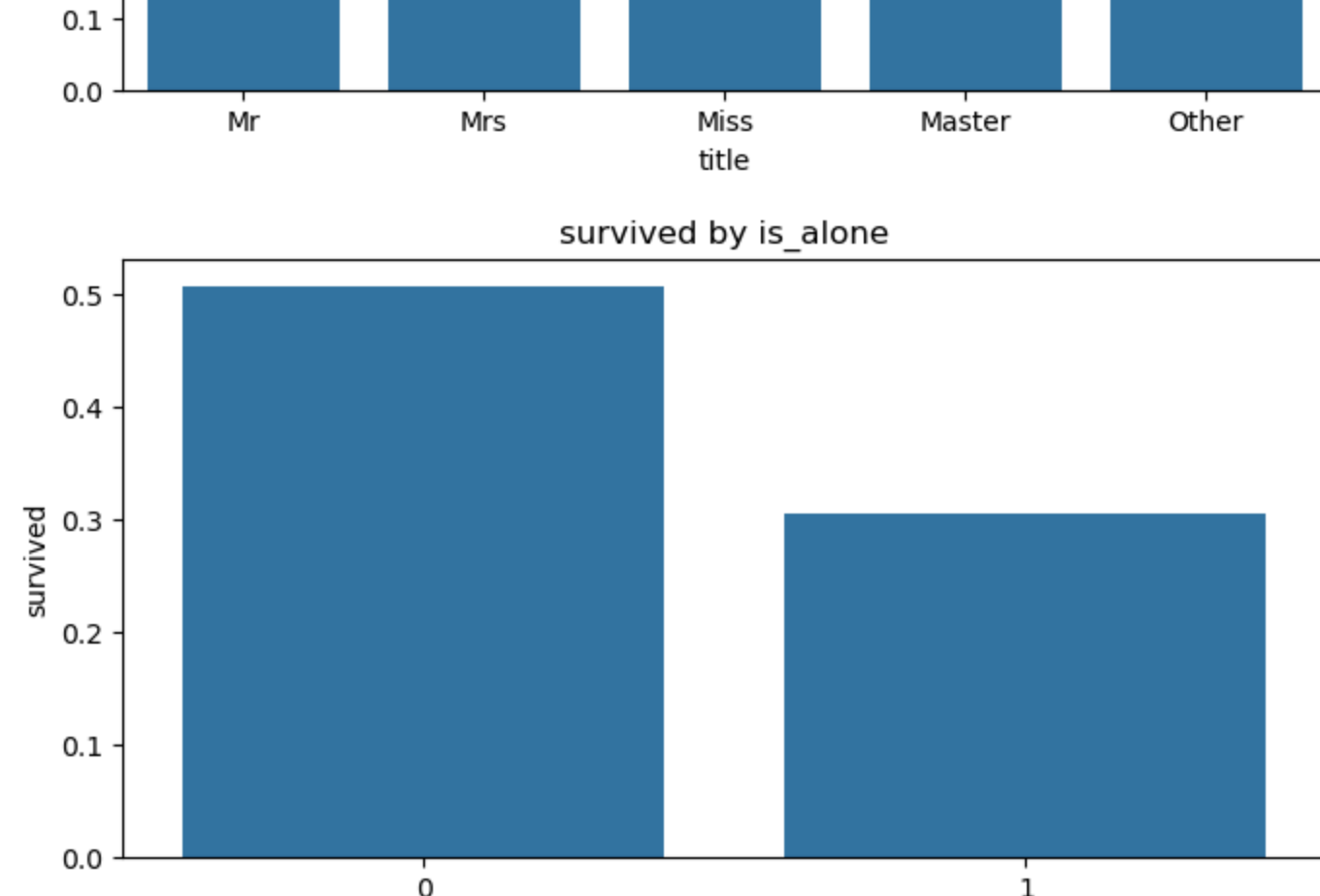
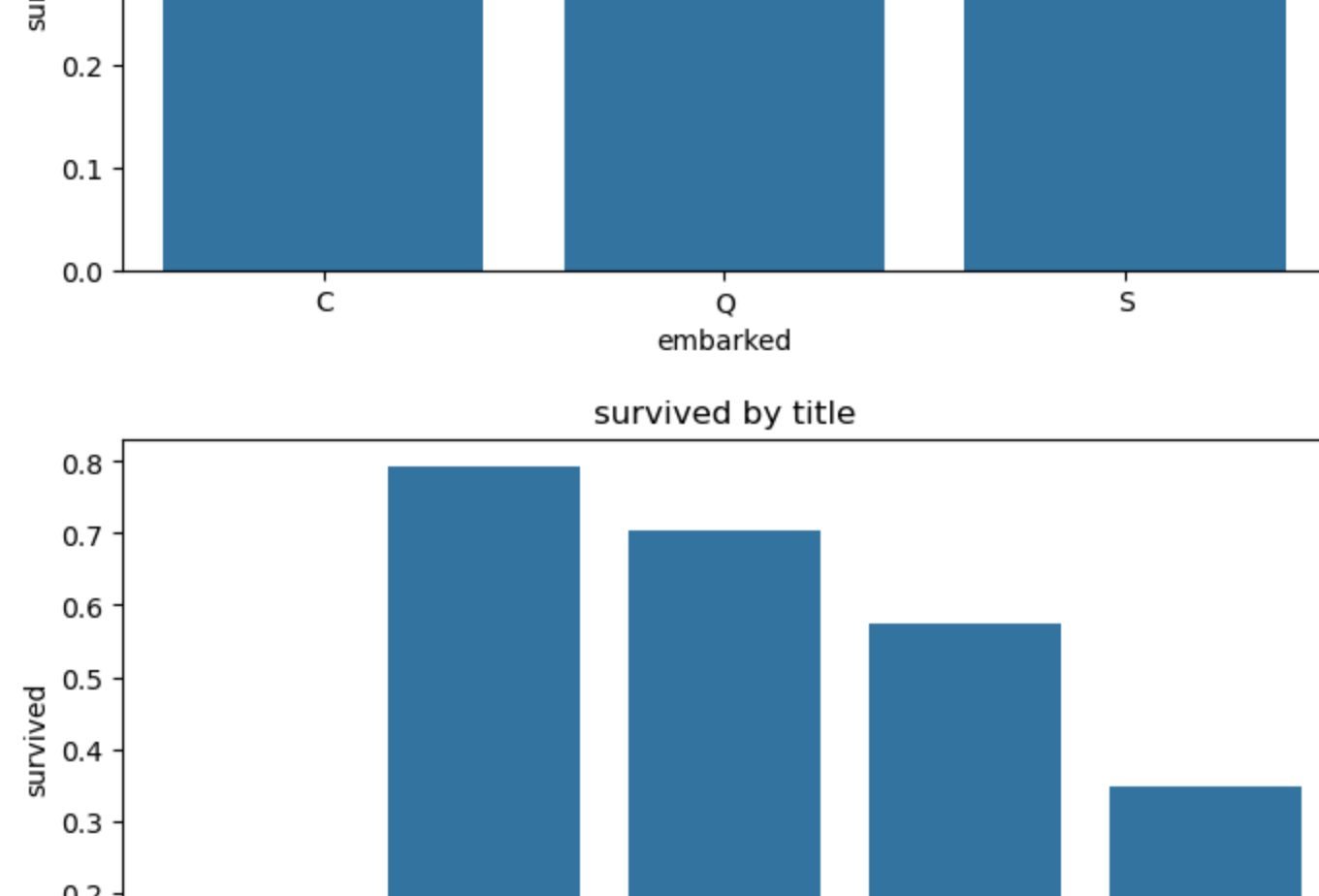
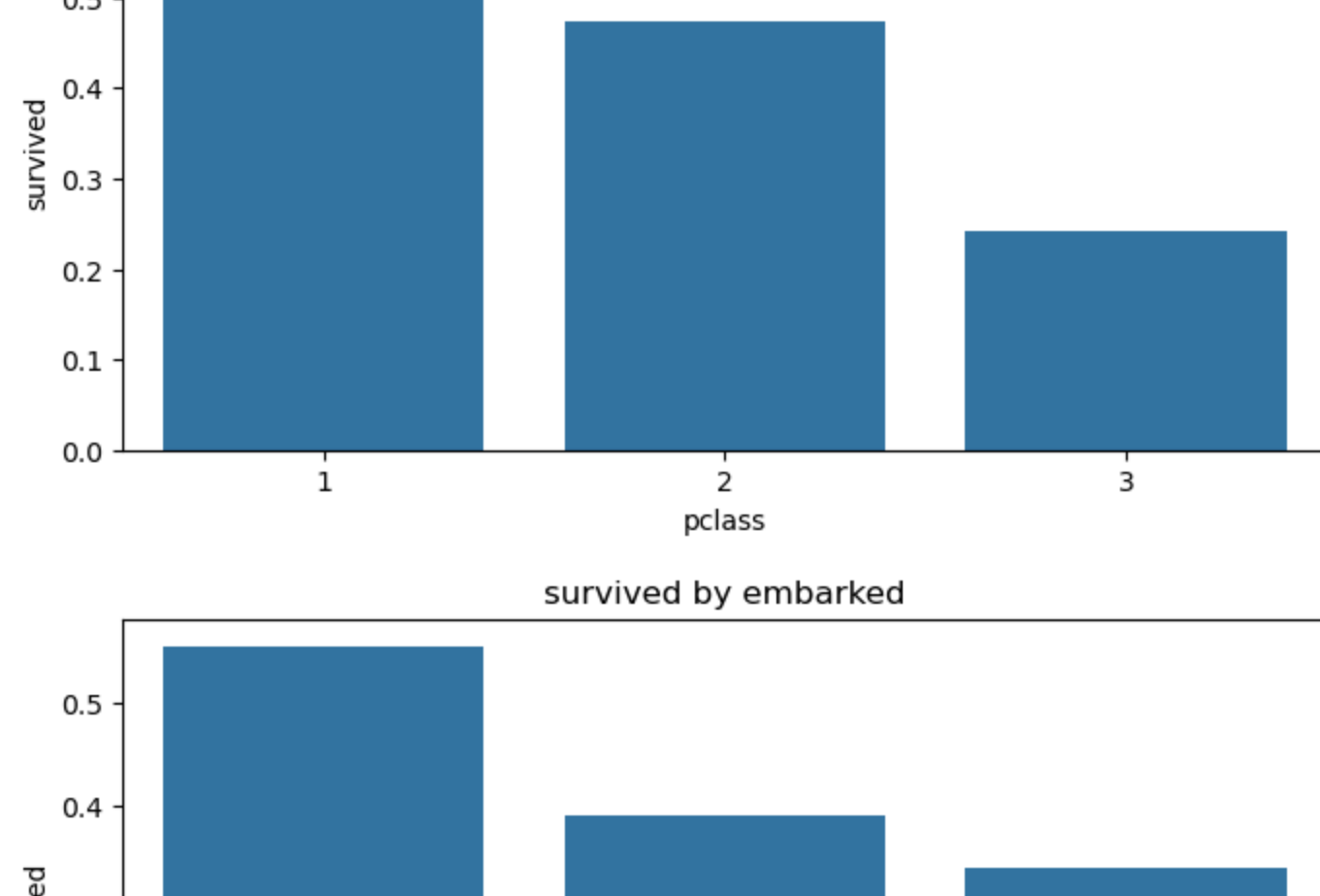
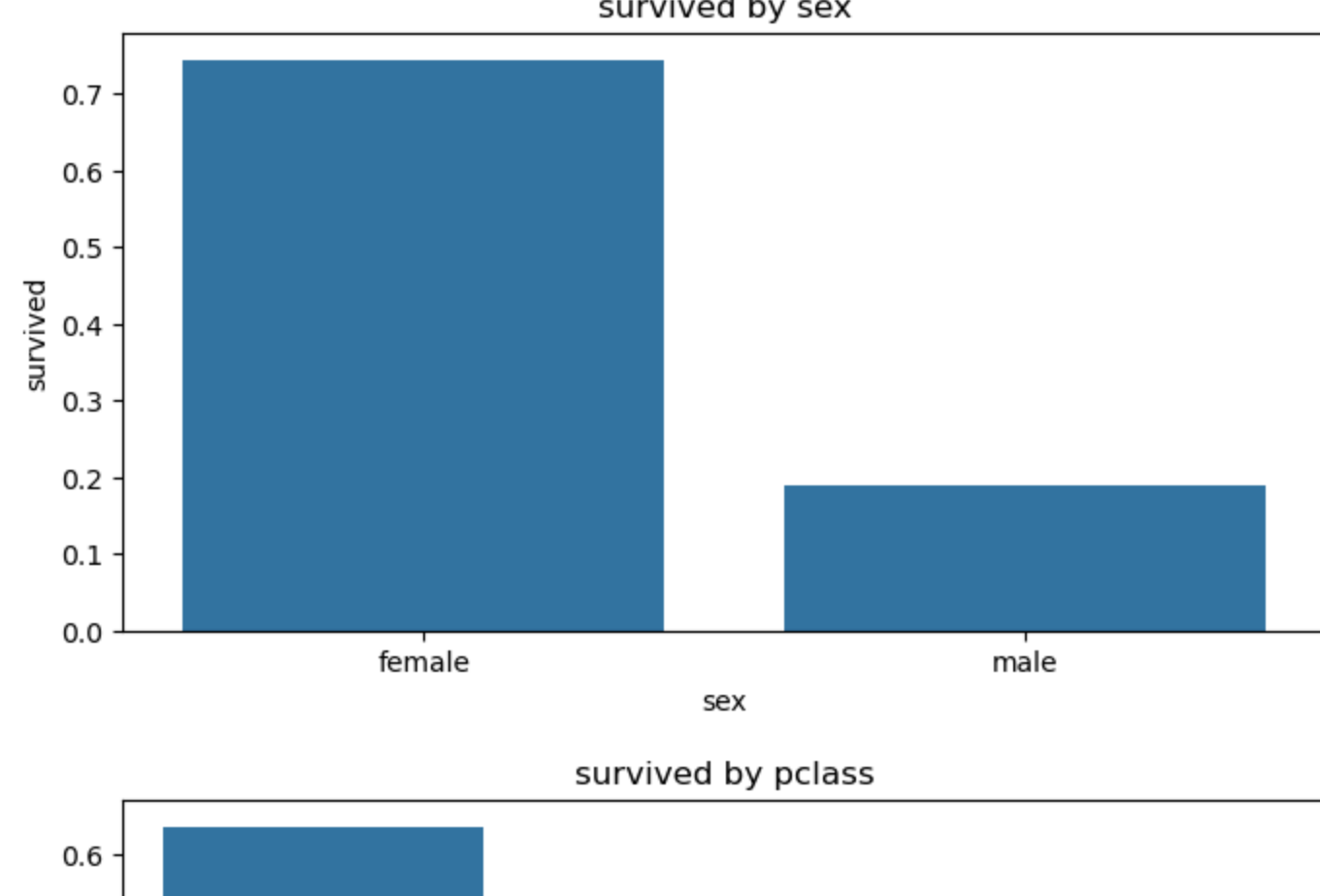
sns.boxplot(x=df[col])
plt.title('col boxplot')
plt.savefig('C:/Users/yelasyahree/OneDrive/Documents/valerate lab tasks/task 5', bbox_inches='tight', dpi=200)
plt.show()

```


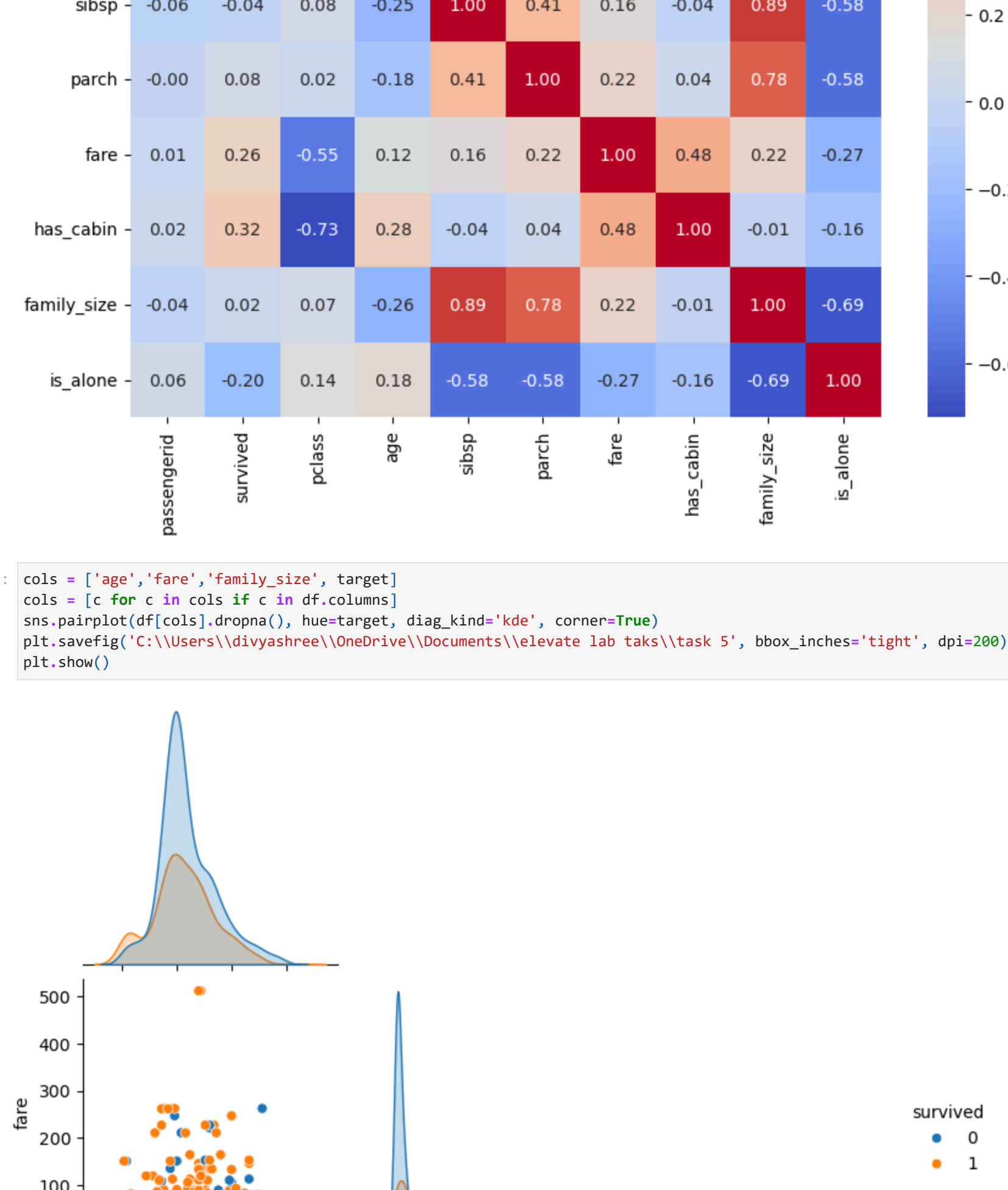
A bar chart titled 'Count plot: sex' showing the distribution of sex. The y-axis is labeled 'count' and ranges from 0 to 600. There are two bars: one for 'male' with a count of approximately 580, and one for 'female' with a count of approximately 310.

```
plt.title(f'[Target] by [col]')
plt.savefig(f'C:\Users\judyahrahe\OneDrive\Documents\eleivate lab tasks\task 5', bbox_inches='tight', dpi=200)
plt.show()

# Numeric vs target
for col in ['age', 'sex', 'family_size']:
    if col in df.columns:
        plt.figure(figsize=(6,4))
        sns.boxplot(x=target, y=col, data=df)
        plt.title(f'[col] distribution by [target]')
        plt.savefig(f'C:\Users\judyahrahe\OneDrive\Documents\eleivate lab tasks\task 5', bbox_inches='tight', dpi=200)
        plt.show()
```



	survived	pclass	age
survived	-0.01	1.00	-0.34
pclass	-0.04	-0.34	1.00
age	-0.04	-0.06	-0.41



```
In [52]: #Multicollinearity detection (VIF)
# Prepare numeric matrix for VIF (drop MA)
X = df.select_dtypes(include=[np.number]).drop(columns=[target], errors='ignore').dropna()
X_const = sm.add_constant(X)
vif_const = pd.DataFrame(X)
```

```
C:\Users\drydenhew\anaconda3\lib\site-packages\statsmodels\regression\linear_model.py:178: RuntimeWarning: divide by zero
return 1 - self.sr_self_centered_tss
C:\Users\drydenhew\anaconda3\lib\site-packages\statsmodels\outliers_influence.py:197: RuntimeWarning: divide by zero
vif = 1. / (1. - r_squared)
```

```
2      pclass 2.779021
7      has_cabin 2.168982
9      is_alone 2.097761
6      fare 1.618458
3      age 1.313374
1      passengerid 1.008039
0      count 0.000000
```

```
[54]: #Statistical tests
if ('sex', target).issubset(df.columns):
    ct = pd.crosstab(df['sex'], df[target])
    chi2, p, exp, expected = chi2_contingency(ct)
    print('Chi-square p-value (sex vs target):', p)

if 'age' in df.columns and target in df.columns:
    grp0 = df[df[target]=='0']['age'].dropna()
    grp1 = df[df[target]=='1']['age'].dropna()
    tstat, pval = ttest_ind(grp0, grp1, non_pooled=True)
    print('t-test p-value (age vs target):', pval)
```

```

t-test p-value (age by target): 0.40754851805863080942

In [58]: #finding skew & outliers
# Check skew
for col in ['fare', 'age']:
    if col in df.columns:
        print(col, 'skew:', df[col].dropna().skew())

# Log transform if skewed
if 'fare' in df.columns:
    df['fare_log'] = np.log1p(df['fare'])

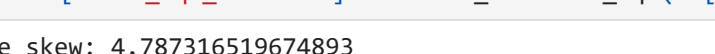
```

```
plt.title('log-transformed iqr')
plt.savefig('C:\Users\ldigheer\OneDrive\Documents\levente lab tasks\Task 5', bbox_inches='tight', dpi=200)
plt.show()

def remove_outliers(iqr(series, k=1.5)):
    q1 = series.quantile(0.25)
    q3 = series.quantile(0.75)
    iqr = q3 - q1
    lower = q1 - k * iqr
    upper = q3 + k * iqr
    return series.where(series.between(lower, upper))
```

```
df['fare_log_removed'] = remove_outliers_log(df['fare'])
```

fare skew: 4.787316519674893  
age skew: 0.5340834483875482



Log-transformed fare
285

```
In [64]: df.to_csv("C:\\Users\\diydashkov\\OneDrive\\Documents\\elevate lab tasks\\Task 5\\cleaned dataset.csv", index=False)
# Save a few key figures (already saved above with plt.savefig). Confirm files:
import glob
print(sorted(glob.glob("C:\\Users\\diydashkov\\OneDrive\\Documents\\elevate lab tasks\\Task 5\\figures"))[:18])
```



