

AI and ML internship task 2

```
In [ ]: #loading the libraries

In [3]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
from sklearn.linear_model import LinearRegression

In [ ]: # loading the dataset(train.csv)

In [7]: data=pd.read_csv('train.csv')

In [9]: data

Out[9]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S
...	...	...	...	...	...	...	...	...	...	...	...	...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W/C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

```
In [11]: data.head()

Out[11]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
In [22]: # Generate summary statistics for numerical columns
summary_stats = data.describe().T # Transpose for readability
summary_stats

Out[22]:
```

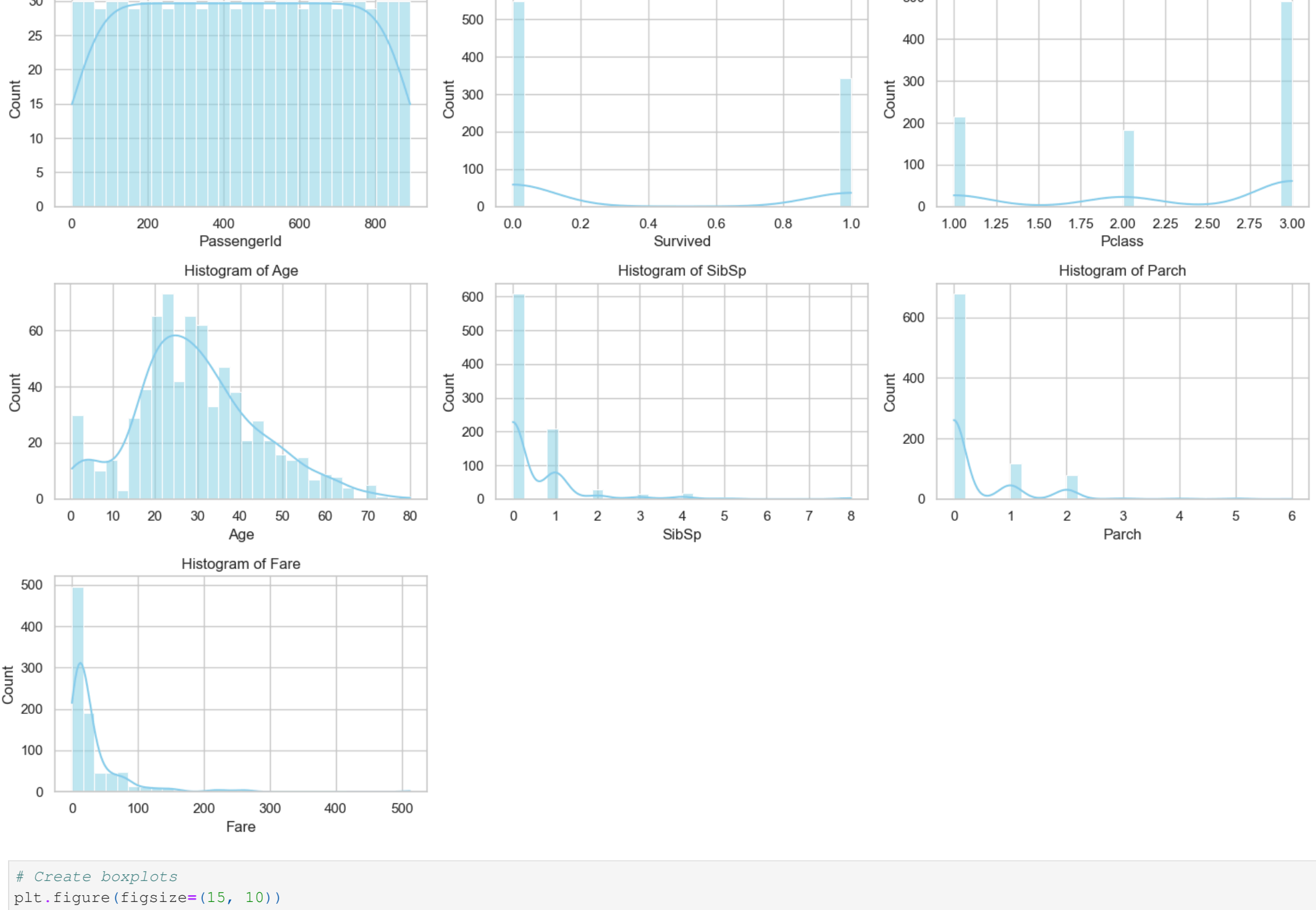
	count	mean	std	min	25%	50%	75%	max
PassengerId	891.0	446.000000	257.353842	1.00	223.5000	446.0000	668.5	891.0000
Survived	891.0	0.383838	0.486592	0.00	0.0000	0.0000	1.0	1.0000
Pclass	891.0	2.308642	0.836071	1.00	2.0000	3.0000	3.0	3.0000
Age	714.0	29.699118	14.526497	0.42	20.1250	28.0000	38.0	80.0000
SibSp	891.0	0.523008	1.102743	0.00	0.0000	0.0000	1.0	8.0000
Parch	891.0	0.381594	0.806057	0.00	0.0000	0.0000	0.0	6.0000
Fare	891.0	32.204208	49.693429	0.00	7.9104	14.4542	31.0	512.3292

```
In [15]: data.info()

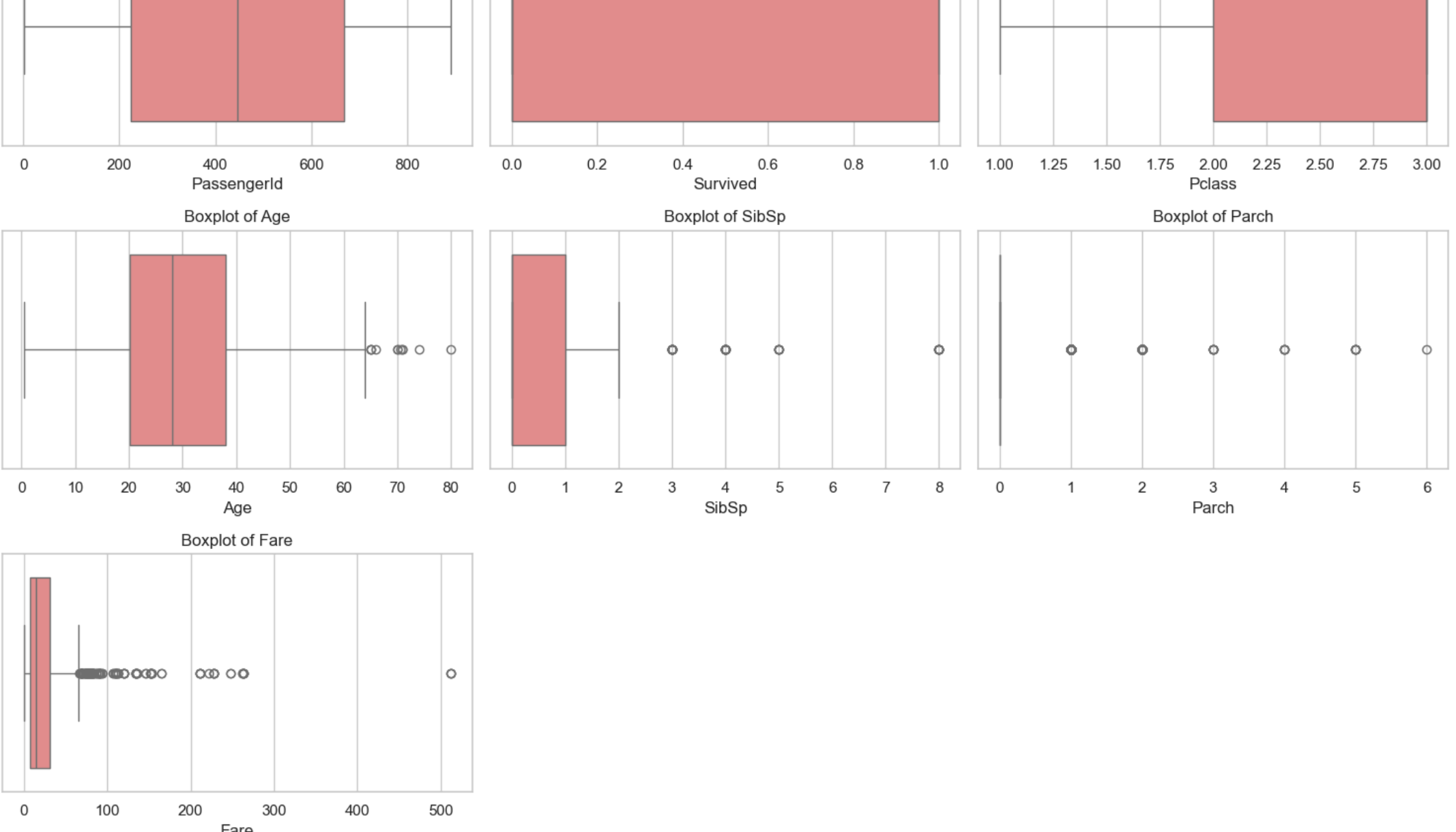
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   PassengerId  891 non-null    int64
 1   Survived     891 non-null    int64
 2   Pclass       891 non-null    int64
 3   Name         891 non-null    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: 83.7+ KB
```

Exploratory Data Analysis (EDA)

```
In [24]: # Create histograms
sns.set(style="whitegrid")
numeric_cols = data.select_dtypes(include=['float64', 'int64']).columns
plt.figure(figsize=(15, 10))
for i, col in enumerate(numeric_cols):
    plt.subplot(3, 3, i + 1)
    sns.histplot(data[col], kde=True, bins=30, color='skyblue')
plt.title(f'Histogram of {col}')
plt.tight_layout()
plt.show()
```



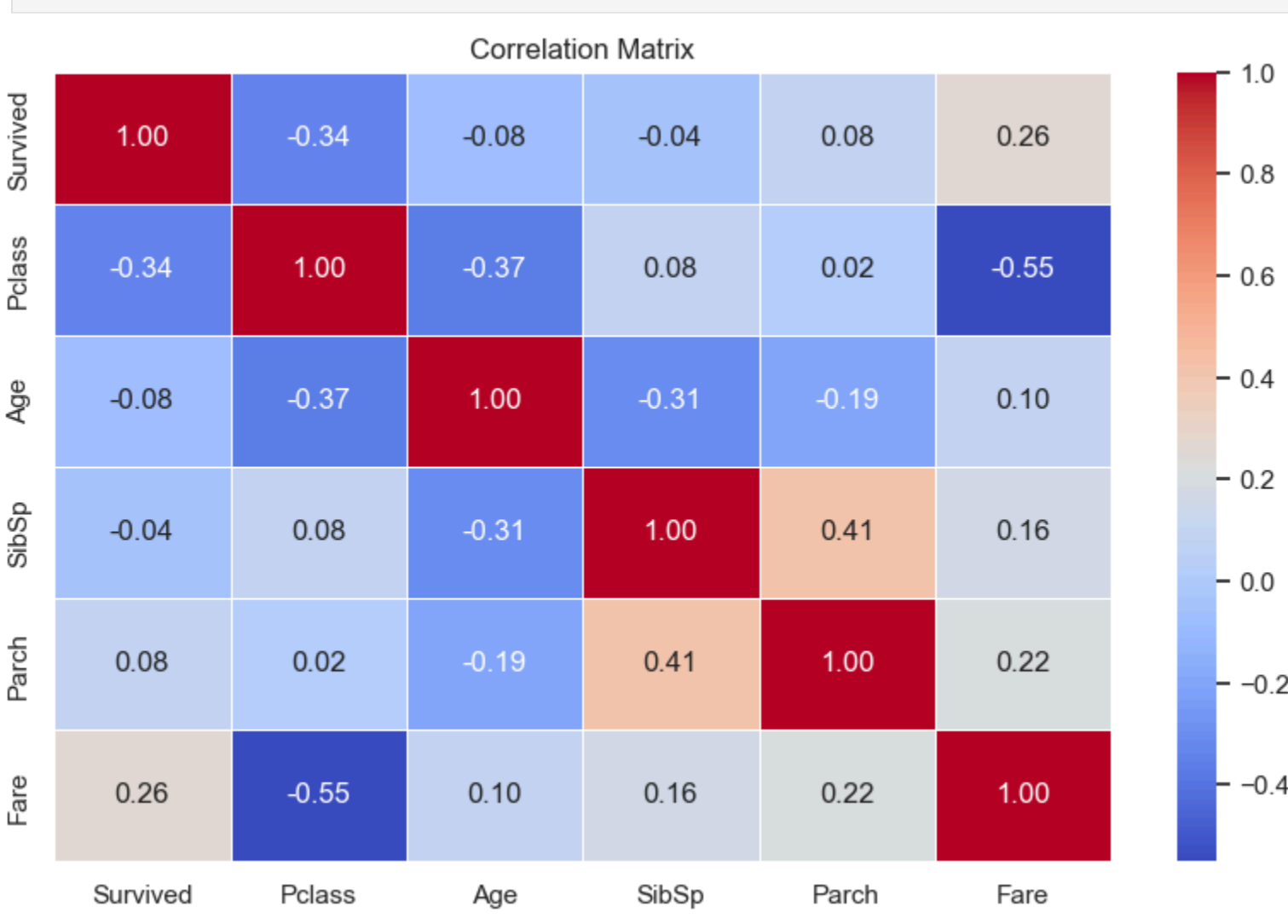
```
In [26]: # Create boxplots
plt.figure(figsize=(15, 10))
for i, col in enumerate(numeric_cols):
    plt.subplot(3, 3, i + 1)
    sns.boxplot(x=data[col], color='lightcoral')
plt.title(f'Boxplot of {col}')
plt.tight_layout()
plt.show()
```



```
In [28]: # Pairplot for selected relevant features
selected_features = ['Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare']
sns.pairplot(data[selected_features], hue='Survived', palette='coolwarm')
plt.suptitle('Pairplot of Selected Features', y=1.02)
plt.show()
```



```
In [30]: # Correlation matrix
plt.figure(figsize=(10, 6))
corr_matrix = data[selected_features].corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
plt.title('Correlation Matrix')
plt.show()
```



```
In [32]: #Identify patterns, trends, or anomalies in the data.
```

Missing Data:

Age: ~20% missing.

Cabin: ~77% missing — heavily incomplete, may need to drop or engineer.

Embarked: 2 missing values.

Trends & Patterns:

Survival Rate:

Overall survival rate is about 38%.

The Survived column is binary, with a slight imbalance (more did not survive).

Class Differences:

Pclass is skewed towards 3rd class.

First-class passengers likely had a higher survival rate (to confirm with group analysis).

Age Distribution:

Skewed right; most passengers were between 20–40.

Some passengers were infants (youngest is 0.42 years).

Fare Distribution:

Strong right-skew (outliers above 100, withmax:512).

Median fare is ~14.45, butmostpassengerspaidlessthan50.

SibSp & Parch:

Most people traveled alone (SibSp=0, Parch=0), but a few had large families aboard.

Anomalies:

Fare outliers: A few passengers paid extremely high fares.

Age includes infants: Worth analyzing infant survival separately.

Cabin: Highly sparse — may indicate only certain classes had cabin data.

```
In [ ]:
```