

```
In [1]: import pandas as pd
```

```
In [2]: df=pd.read_csv("titanic (1).csv")
```

```
In [3]: df.head(10)
```

Out[3]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0700



In [4]:

```
df.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
```

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

Out[5]:

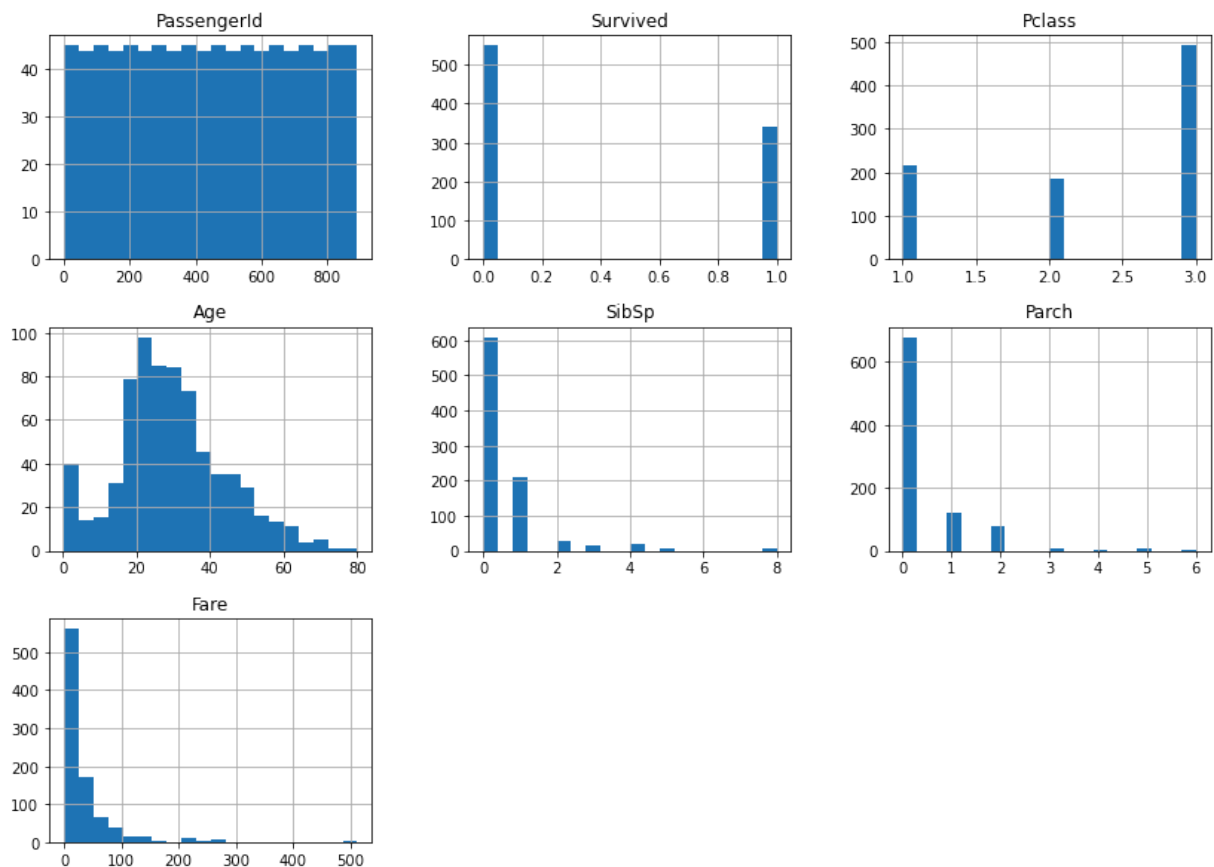
	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

```
In [6]: df.isnull().sum()
```

Out[6]:

PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	177
SibSp	0
Parch	0
Ticket	0
Fare	0
Cabin	687
Embarked	2
dtype:	int64

```
In [7]: df.hist(bins=20,figsize=(14,10))
import matplotlib.pyplot as plt
plt.show()
```

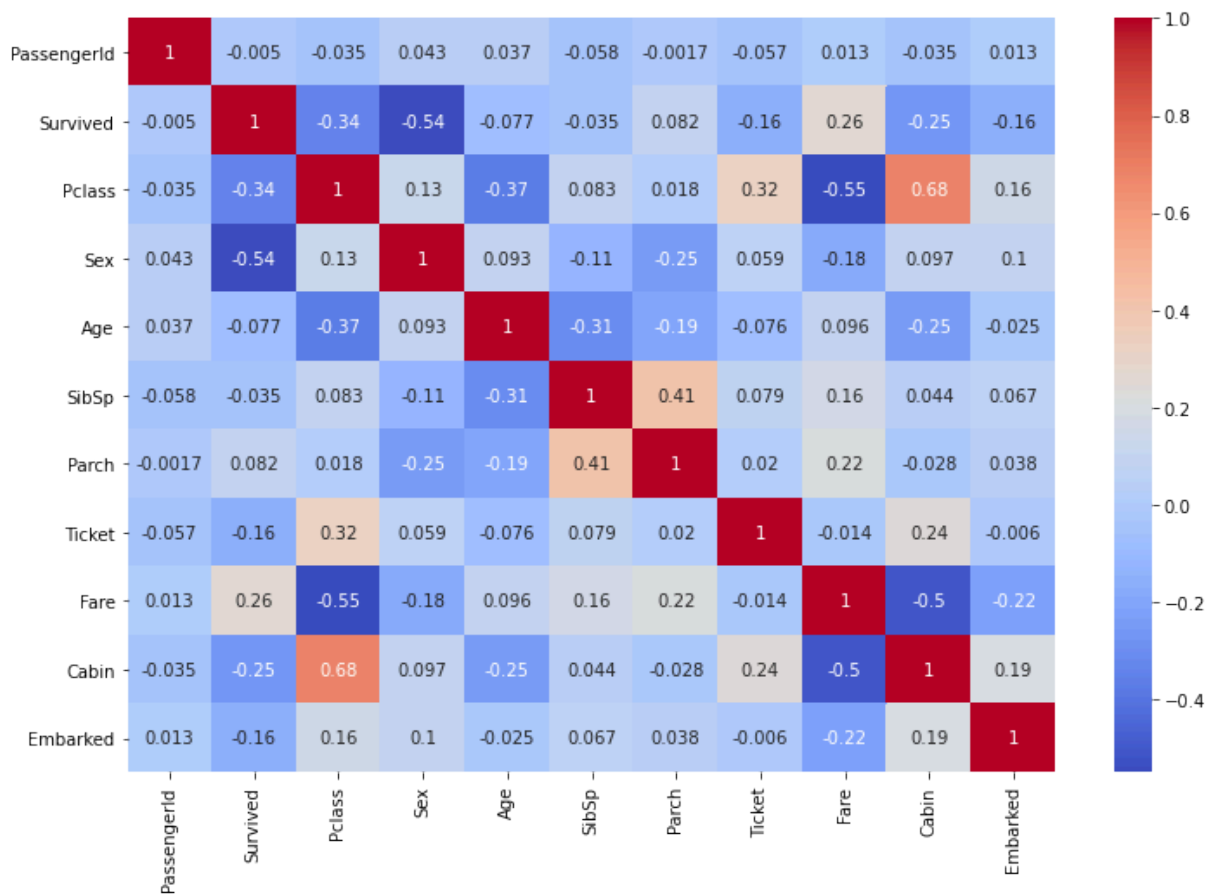


```
In [8]: df.drop(['Name'],axis=1,inplace=True)
```

```
In [9]: from sklearn.preprocessing import LabelEncoder
lb=LabelEncoder()
df['Sex']=lb.fit_transform(df['Sex'])
```

```
In [10]: df['Ticket']=lb.fit_transform(df['Ticket'])
df['Cabin']=lb.fit_transform(df['Cabin'])
df['Embarked']=lb.fit_transform(df['Embarked'])
```

```
In [11]: corr_matrix=df.corr()
plt.figure(figsize=(12,8))
import seaborn as sns
sns.heatmap(corr_matrix,annot=True,cmap='coolwarm')
plt.show()
```



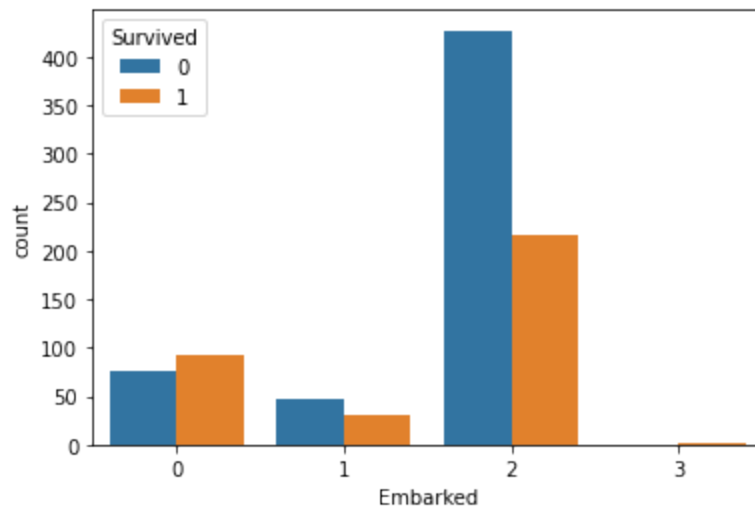
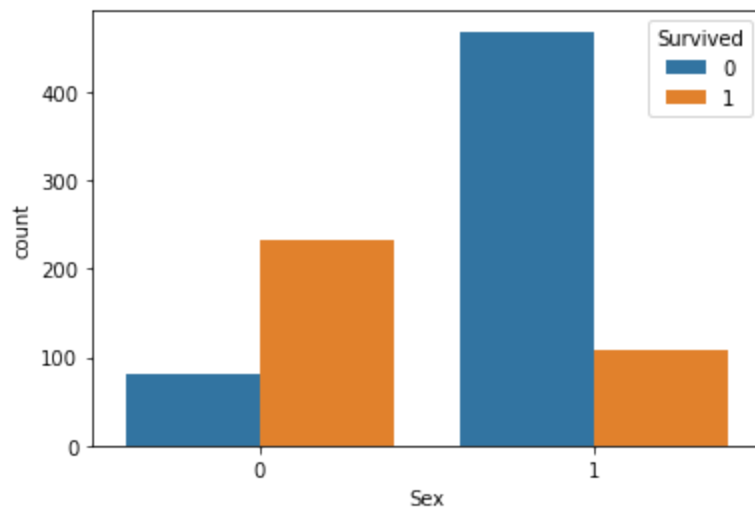
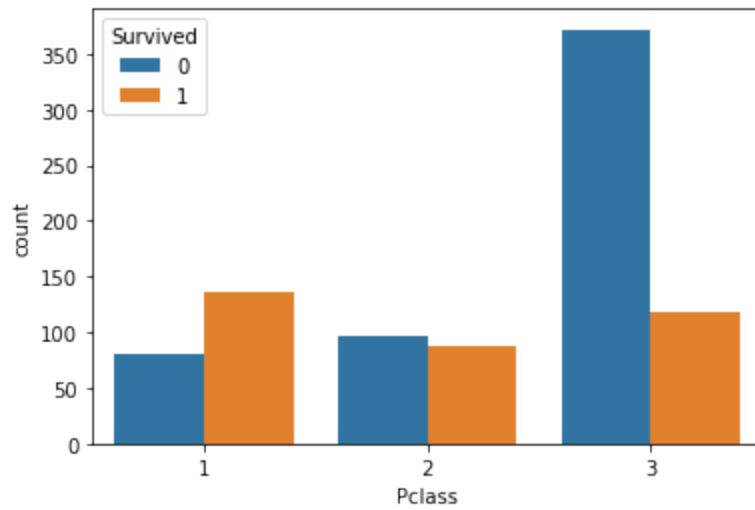
```
In [12]: sns.pairplot(df,hue='Survived',diag_kind='hist')
```

```
Out[12]: <seaborn.axisgrid.PairGrid at 0x1fcc97dc8b0>
```



```
In [13]: feature_e=['Pclass', 'Sex', 'Embarked']
```

```
In [14]: for feature in feature_e:
sns.countplot(x=feature, hue='Survived', data=df)
plt.show()
```



```
In [15]: df=df.drop(['PassengerId','Ticket','Cabin'],axis=1)
```

```
In [16]: df.isnull().sum()
```

```
Out[16]: Survived      0
         Pclass       0
         Sex          0
         Age         177
         SibSp        0
         Parch        0
         Fare         0
         Embarked     0
         dtype: int64
```

```
In [17]: from sklearn.impute import SimpleImputer
         imputer=SimpleImputer(strategy='median')
```

```
In [18]: df['Age']=imputer.fit_transform(df[['Age']])
```

```
In [19]: df.isnull().sum()
```

```
Out[19]: Survived      0
         Pclass       0
         Sex          0
         Age          0
         SibSp        0
         Parch        0
         Fare         0
         Embarked     0
         dtype: int64
```

```
In [20]: df.columns
```

```
Out[20]: Index(['Survived', 'Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare',
               'Embarked'],
              dtype='object')
```

```
In [21]: x=df.iloc[:,1:]
```

```
In [22]: x.columns
```

```
Out[22]: Index(['Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked'], dtype='object')
```

```
In [23]: y=df.iloc[:,0]
```

```
In [24]: from sklearn.tree import DecisionTreeClassifier
```

```
In [25]: dt=DecisionTreeClassifier()
```

```
In [26]: from sklearn.model_selection import cross_val_score
```

```
In [27]: cv_score=cross_val_score(dt,x,y,cv=5)
```

```
In [28]: cv_score
```

```
Out[28]: array([0.73743017, 0.78089888, 0.80337079, 0.74157303, 0.79213483])
```



```
In [29]: import numpy as np
```

```
In [30]: np.mean(cv_score)
```

```
Out[30]: 0.7710815391375305
```

```
In [31]: from sklearn.linear_model import LogisticRegression
```

```
In [32]: lr=LogisticRegression()
```

```
In [33]: cv_score1=cross_val_score(lr,x,y,cv=5)
```

```
C:\Users\CompLab30\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:76
3: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```

```
C:\Users\CompLab30\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:76
```

```
3: ConvergenceWarning: lbfgs failed to converge (status=1):
```

```
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```

```
In [34]: np.mean(cv_score1)
```

```
Out[34]: 0.7890025735986442
```

```
In [35]: from sklearn.model_selection import train_test_split
```

```
In [36]: xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=0.25,random_state=123)
```

```
In [37]: dt_final=DecisionTreeClassifier()
```

```
In [38]: dt_final.fit(xtrain,ytrain)
```

```
Out[38]: DecisionTreeClassifier()
```

```
In [39]: pred=dt_final.predict(xtest)
```

```
In [40]: lr_final=LogisticRegression()
```

```
In [41]: lr_final.fit(xtrain,ytrain)
```

```
Out[41]: LogisticRegression()
```

```
In [42]: pred_lr=lr_final.predict(xtest)
```

```
In [43]: df['Survived'].value_counts()
```

```
Out[43]: 0    549
         1    342
         Name: Survived, dtype: int64
```

```
In [44]: from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

```
In [45]: accuracy_score(ytest,pred),accuracy_score(ytest,pred_lr)
```

```
Out[45]: (0.8071748878923767, 0.7982062780269058)
```

```
In [46]: confusion_matrix(ytest,pred),confusion_matrix(ytest,pred_lr)
```

```
Out[46]: (array([[118, 21],
                [ 22, 62]], dtype=int64),
         array([[116, 23],
                [ 22, 62]], dtype=int64))
```

```
In [47]: print(classification_report(ytest,pred),classification_report(ytest,pred_lr))
```

	precision	recall	f1-score	support
0	0.84	0.85	0.85	139
1	0.75	0.74	0.74	84
accuracy			0.81	223
macro avg	0.79	0.79	0.79	223
weighted avg	0.81	0.81	0.81	223

	precision	recall	f1-score	support
0	0.84	0.83	0.84	139
1	0.73	0.74	0.73	84
accuracy			0.80	223
macro avg	0.78	0.79	0.79	223
weighted avg	0.80	0.80	0.80	223

```
In [ ]:
```