

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

```
df= pd.read_csv('/content/drive/MyDrive/Titanic.csv')
df.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	892	0	3	Kelly, Mr. James	male	34.5	0	0	330911	7.8292
1	893	1	3	Wilkes, Mrs. James (Ellen Needs)	female	47.0	1	0	363272	7.0000

```
## statistical info
df.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	418.000000	418.000000	418.000000	332.000000	418.000000	418.000000	417.000000
mean	1100.500000	0.363636	2.265550	30.272590	0.447368	0.392344	35.627188
std	120.810458	0.481622	0.841838	14.181209	0.896760	0.981429	55.907576
min	892.000000	0.000000	1.000000	0.170000	0.000000	0.000000	0.000000
25%	996.250000	0.000000	1.000000	21.000000	0.000000	0.000000	7.895800
50%	1100.500000	0.000000	3.000000	27.000000	0.000000	0.000000	14.454200
75%	1204.750000	1.000000	3.000000	39.000000	1.000000	0.000000	31.500000
max	1308.000000	1.000000	3.000000	76.000000	8.000000	9.000000	512.329200

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  -
0   PassengerId  418 non-null    int64
1   Survived     418 non-null    int64
2   Pclass       418 non-null    int64
3   Name         418 non-null    object
4   Sex          418 non-null    object
5   Age         332 non-null    float64
6   SibSp        418 non-null    int64
7   Parch        418 non-null    int64
8   Ticket       418 non-null    object
9   Fare         417 non-null    float64
10  Cabin        91 non-null     object
11  Embarked     418 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 39.3+ KB
```

```
## categorical attributes
# sns.countplot(df['Survived'])

# sns.countplot(df['Pclass'])

# sns.countplot(train['Sex'])

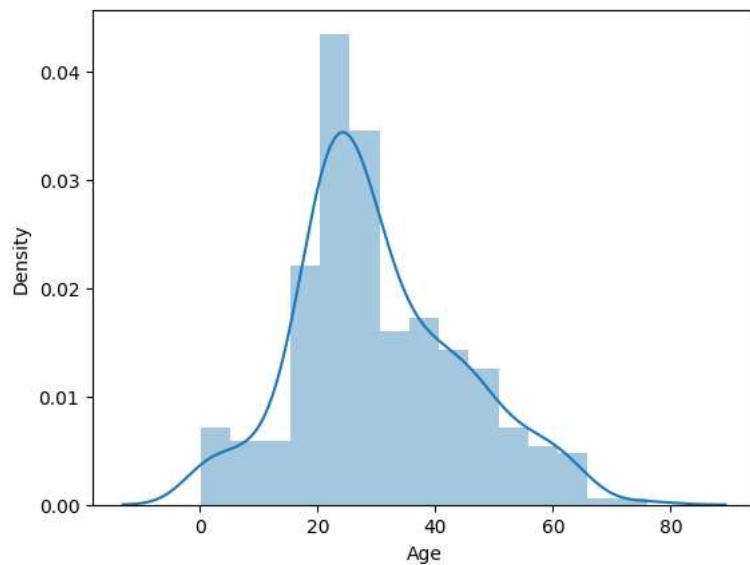
# sns.countplot(train['SibSp'])
```

```
# sns.countplot(train['Parch'])
```

```
# sns.countplot(train['Embarked'])
```

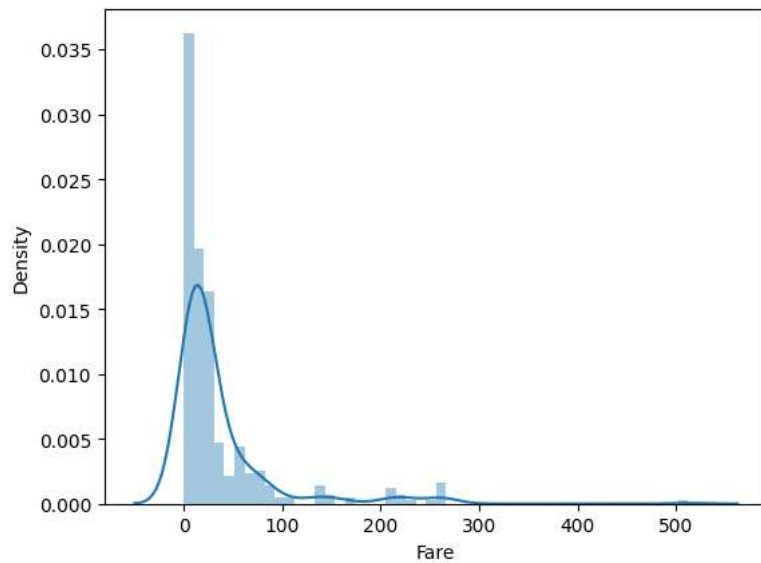
```
## numerical attributes  
sns.distplot(df['Age'])
```

<Axes: xlabel='Age', ylabel='Density'>

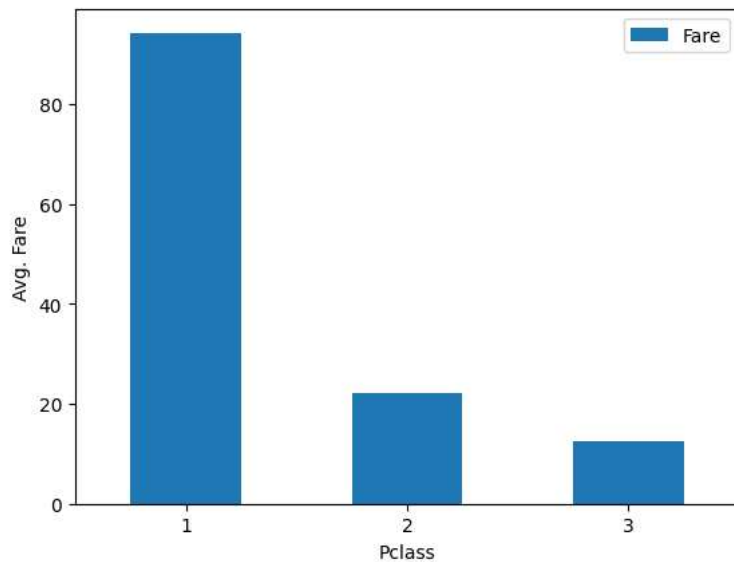


```
sns.distplot(df['Fare'])
```

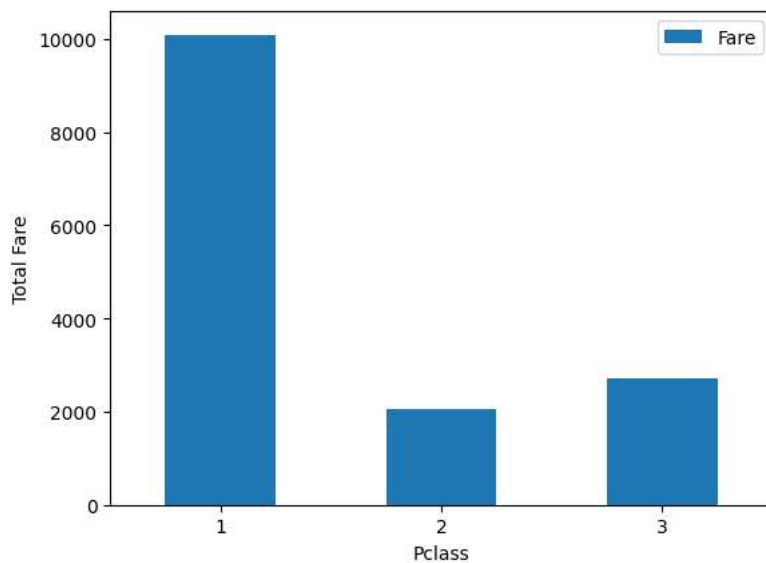
<Axes: xlabel='Fare', ylabel='Density'>



```
class_fare = df.pivot_table(index='Pclass', values='Fare')  
class_fare.plot(kind='bar')  
plt.xlabel('Pclass')  
plt.ylabel('Avg. Fare')  
plt.xticks(rotation=0)  
plt.show()
```

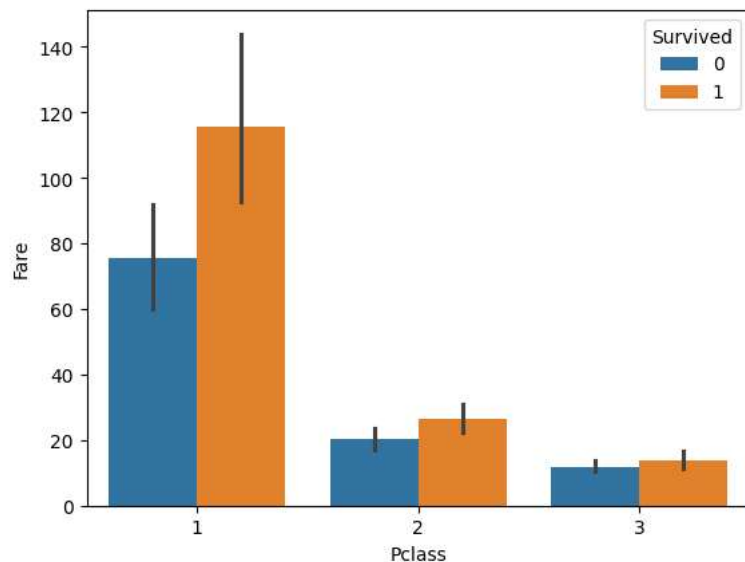


```
class_fare = df.pivot_table(index='Pclass', values='Fare', aggfunc=np.sum)
class_fare.plot(kind='bar')
plt.xlabel('Pclass')
plt.ylabel('Total Fare')
plt.xticks(rotation=0)
plt.show()
```



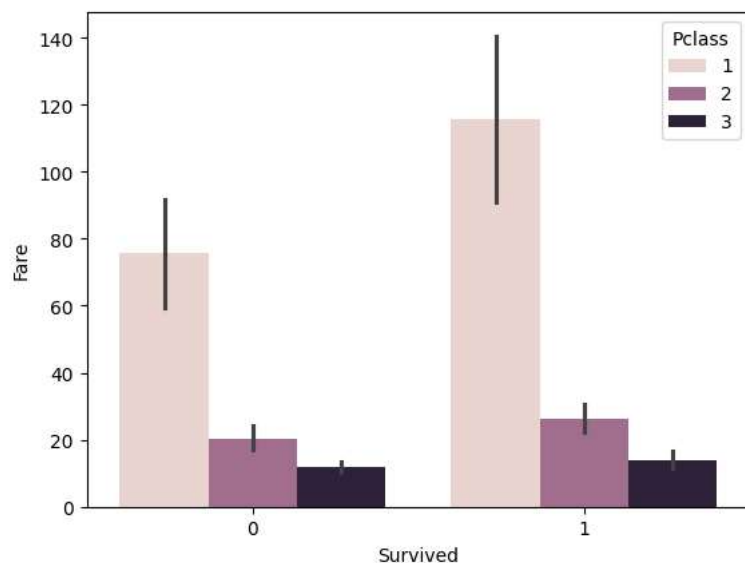
```
sns.barplot(data=df, x='Pclass', y='Fare', hue='Survived')
```

<Axes: xlabel='Pclass', ylabel='Fare'>



```
sns.barplot(data=df, x='Survived', y='Fare', hue='Pclass')
```

<Axes: xlabel='Survived', ylabel='Fare'>



```
## find the null values
df.isnull().sum()
```

```
PassengerId    0
Survived        0
Pclass          0
Name            0
Sex             0
Age            86
SibSp           0
Parch           0
Ticket          0
Fare            1
Cabin          327
Embarked        0
dtype: int64
```

```
# drop or delete the column
df = df.drop(columns=['Cabin'], axis=1)
```

```
df['Age'].mean()
```

```
30.272590361445783
```

```
# fill missing values using mean of the numerical column
df['Age'] = df['Age'].fillna(df['Age'].mean())
df['Fare'] = df['Fare'].fillna(df['Fare'].mean())

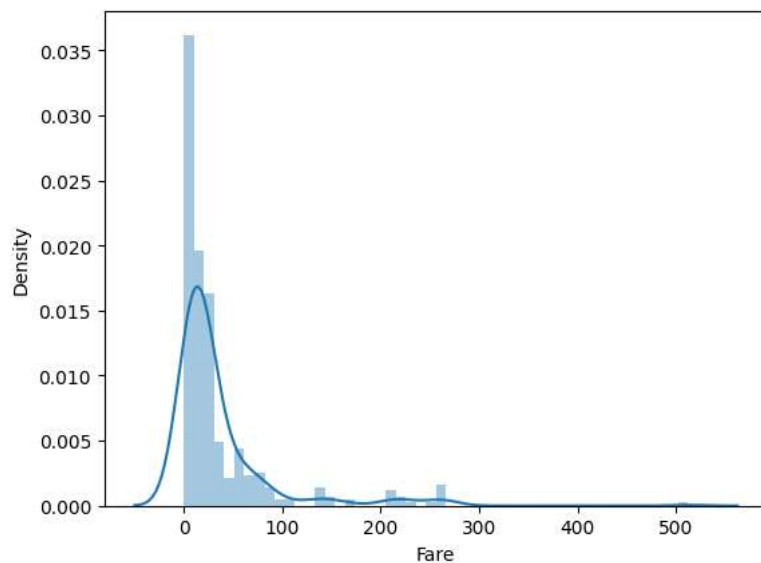
df['Embarked'].mode()[0]

'S'

# fill missing values using mode of the categorical column
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

```
sns.distplot(df['Fare'])
```

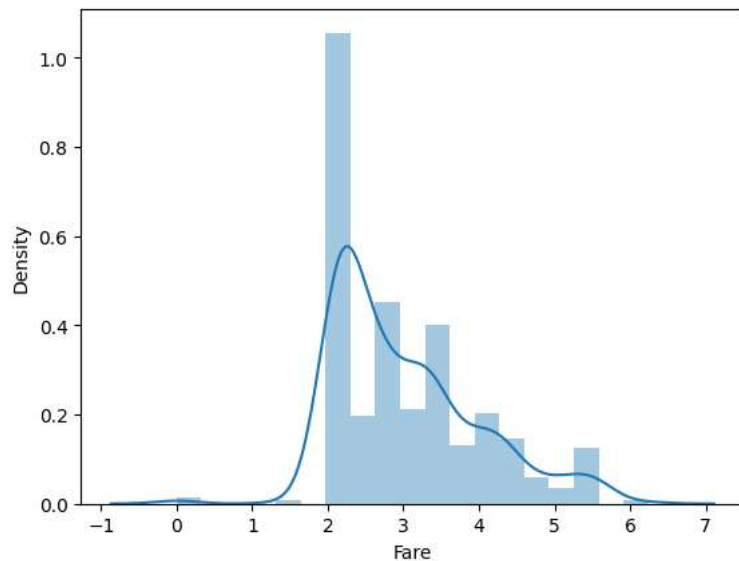
<Axes: xlabel='Fare', ylabel='Density'>



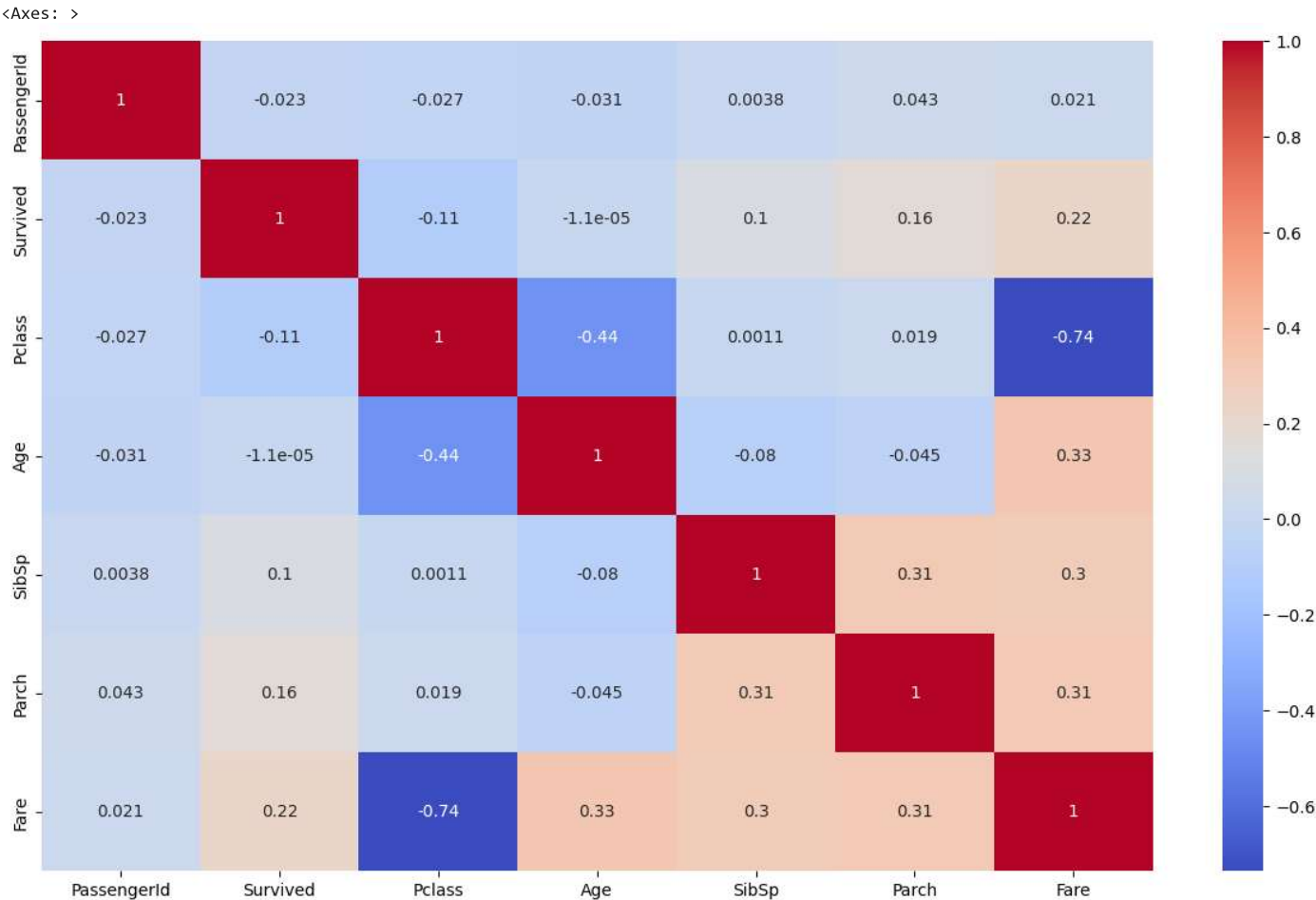
```
df['Fare'] = np.log(df['Fare']+1)
```

```
sns.distplot(df['Fare'])
```

<Axes: xlabel='Fare', ylabel='Density'>



```
corr = df.corr()
plt.figure(figsize=(15, 9))
sns.heatmap(corr, annot=True, cmap='coolwarm')
```



```
## drop unnecessary columns
df = df.drop(columns=['Name', 'Ticket'], axis=1)
df.head()
```

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	892	0	3	male	34.5	0	0	2.178064	Q
1	893	1	3	female	47.0	1	0	2.079442	S
2	894	0	2	male	62.0	0	0	2.369075	Q
3	895	0	3	male	27.0	0	0	2.268252	S
4	896	1	3	female	22.0	1	1	2.586824	S

```
from sklearn.preprocessing import LabelEncoder
cols = ['Sex', 'Embarked']
le = LabelEncoder()

for col in cols:
    df[col] = le.fit_transform(df[col])
df.head()
```

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	892	0	3	1	34.5	0	0	2.178064	1
1	893	1	3	0	47.0	1	0	2.079442	2
2	894	0	2	1	62.0	0	0	2.369075	1
3	895	0	3	1	27.0	0	0	2.268252	2
4	896	1	3	0	22.0	1	1	2.586824	2

```
# input split
X = df.drop(columns=['PassengerId', 'Survived'], axis=1)
y = df['Survived']

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)

from sklearn.linear_model import LogisticRegression
model = LogisticRegression()

model.fit(X_train,y_train)

▼ LogisticRegression
LogisticRegression()

print('Accuracy:', model.score(X_test, y_test))

Accuracy: 1.0

from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()

model.fit(X_train,y_train)

▼ DecisionTreeClassifier
DecisionTreeClassifier()

print('Accuracy:', model.score(X_test, y_test))

Accuracy: 1.0

from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()

model.fit(X_train,y_train)

▼ RandomForestClassifier
RandomForestClassifier()

print('Accuracy:', model.score(X_test, y_test))

Accuracy: 1.0

from sklearn.ensemble import ExtraTreesClassifier
model = ExtraTreesClassifier()

model.fit(X_train,y_train)

▼ ExtraTreesClassifier
ExtraTreesClassifier()

print('Accuracy:', model.score(X_test, y_test))

Accuracy: 1.0

from xgboost import XGBClassifier
model = XGBClassifier()
classify(model)

Accuracy: 1.0
CV Score: 1.0
```

```
from lightgbm import LGBMClassifier
model = LGBMClassifier()
classify(model)
```

[illegible]

```
model = LGBMClassifier()
model.fit(X, y)
```


[illegible]

```
X_test = df.drop(columns=['PassengerId', 'Survived'], axis=1)
```

```
[LightGBM] [Warning] No further splits with positive gain. best gain: -inf
```

```
pred = model.predict(X_test)
```

pred

```
array([0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0,  
       1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1,  
       1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1,  
       1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1,  
       1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0,  
       0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0,  
       1, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1,  
       0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1,  
       1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1,  
       0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1,  
       1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1,  
       0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0,  
       0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0,  
       0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0,  
       0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0,  
       1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0])
```