

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

In [2]: import seaborn as sns

In [3]: import matplotlib.pyplot as plt
%matplotlib inline

In [4]: import math

In [5]: titanic_data = pd.read_csv("train.csv")

In [6]: titanic_data.head(10)
```

```
Out[6]:
```

	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	Cumings, 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
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	NaN	Q
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	E46	S
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	NaN	S
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	NaN	S
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	NaN	C

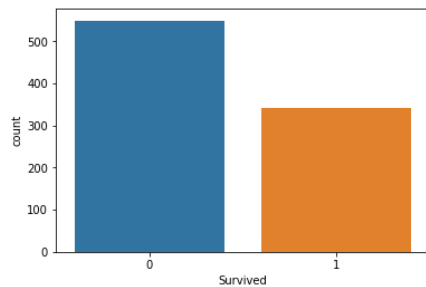
```
In [7]: print("# of passengers in the original data: "+str(len(titanic_data.index)))

# of passengers in the original data: 891
```

Analyzing Data

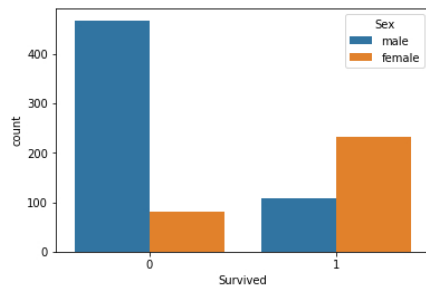
```
In [8]: sns.countplot(x="Survived", data=titanic_data)

Out[8]: <matplotlib.axes._subplots.AxesSubplot at 0x2222dc22400>
```



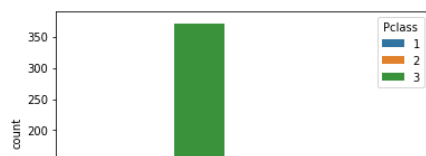
```
In [9]: sns.countplot(x="Survived", hue="Sex", data=titanic_data)

Out[9]: <matplotlib.axes._subplots.AxesSubplot at 0x2222df0c4a8>
```



```
In [10]: sns.countplot(x="Survived", hue="Pclass", data=titanic_data)

Out[10]: <matplotlib.axes._subplots.AxesSubplot at 0x2222df449b0>
```

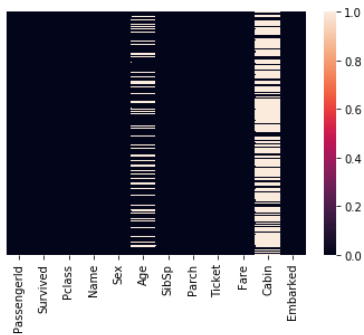


[illegible]

891 rows x 12 columns

```
Out[16]: PassengerId
Survived
Pclass
Name
Sex
Age
SibSp
Parch
Ticket
Fare
Cabin
Embarked
dtype: object
```

```
In [17]: sns.heatmap(titanic_data.isnull(), yticklabels=False)
Out[17]: <matplotlib.axes._subplots.AxesSubplot at 0x2222e1a3320>
```



```
In [18]: titanic_data.head(10)
```

Out[18]:

	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	Cumings, 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
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	NaN	Q
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	E46	S
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	NaN	S
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	NaN	S
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	NaN	C

```
In [19]: titanic_data.drop("Cabin", axis=1, inplace=True)
```

```
In [20]: titanic_data.head(10)
```

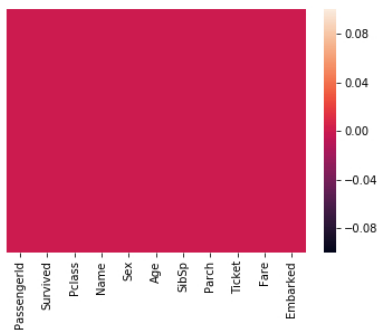
Out[20]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	S
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	Q
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	S
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	S
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	S
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	C

```
In [21]: titanic_data.dropna(inplace=True)
```

```
In [22]: sns.heatmap(titanic_data.isnull(), yticklabels=False)
```

Out[22]: <matplotlib.axes._subplots.AxesSubplot at 0x2222e256f28>



```
In [23]: titanic_data.isnull().sum()
```

Out[23]:

PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	0
SibSp	0
Parch	0
Ticket	0
Fare	0

```
Embarked      0  
dtype: int64
```

```
In [24]: sex = pd.get_dummies(titanic_data["Sex"], drop_first=True)  
sex.head(5)
```

```
Out[24]:
```

	male
0	1
1	0
2	0
3	0
4	1

```
In [25]: embark = pd.get_dummies(titanic_data["Embarked"], drop_first=True)  
embark.head(5)
```

```
Out[25]:
```

	Q	S
0	0	1
1	0	0
2	0	1
3	0	1
4	0	1

```
In [26]: Pcl = pd.get_dummies(titanic_data["Pclass"], drop_first=True)  
Pcl.head(5)
```

```
Out[26]:
```

	2	3
0	0	1
1	0	0
2	0	1
3	0	0
4	0	1

```
In [27]: titanic_data = pd.concat([titanic_data, sex, embark, Pcl], axis=1)  
titanic_data.head(5)
```

```
Out[27]:
```

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

```
In [28]: titanic_data.drop(['Sex', 'Embarked', 'PassengerId', 'Ticket', 'Name'], axis=1, inplace=True)
```

```
In [29]: titanic_data.head()
```

```
Out[29]:
```

	Survived	Pclass	Age	SibSp	Parch	Fare	male	Q	S	2	3
0	0	3	22.0	1	0	7.2500	1	0	1	0	1
1	1	1	38.0	1	0	71.2833	0	0	0	0	0
2	1	3	26.0	0	0	7.9250	0	0	1	0	1
3	1	1	35.0	1	0	53.1000	0	0	1	0	0
4	0	3	35.0	0	0	8.0500	1	0	1	0	1

```
In [30]: titanic_data.drop('Pclass', axis=1, inplace=True)
```

```
In [31]: titanic_data.head()
```

```
Out[31]:
```

	Survived	Age	SibSp	Parch	Fare	male	Q	S	2	3
0	0	22.0	1	0	7.2500	1	0	1	0	1
1	1	38.0	1	0	71.2833	0	0	0	0	0
2	1	26.0	0	0	7.9250	0	0	1	0	1
3	1	35.0	1	0	53.1000	0	0	1	0	0
4	0	35.0	0	0	8.0500	1	0	1	0	1

#Train

```
In [33]: X = titanic_data.drop("Survived", axis=1)  
y = titanic_data["Survived"]
```

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

```
In [36]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=1)
```

```
In [37]: from sklearn.linear_model import LogisticRegression
```

```
In [38]: logmodel = LogisticRegression()
```

```
In [39]: logmodel.fit(X_train, y_train)
```

```
Out[39]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,  
    intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,  
    penalty='l2', random_state=None, solver='liblinear', tol=0.0001,  
    verbose=0, warm_start=False)
```

```
In [40]: prediction = logmodel.predict(X_test)
```

```
In [41]: from sklearn.metrics import classification_report
```

```
In [44]: classification_report(y_test, prediction)
```

```
Out[44]: '          precision    recall  f1-score   support\n\n 0.75      0.72      0.73      0.73      88\n0.78      0.79      0.78      0.78     126\n1.00      1.00      1.00      1.00     214'
```

```
In [45]: from sklearn.metrics import confusion_matrix
```

```
In [46]: confusion_matrix(y_test, prediction)
```

```
Out[46]: array([[105,  21],  
       [ 25,  63]], dtype=int64)
```

```
In [47]: from sklearn.metrics import accuracy_score
```

```
In [48]: accuracy_score(y_test, prediction)
```

```
Out[48]: 0.7850467289719626
```