

Lecture 8-Part2

Logistic Regression

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 %matplotlib inline
```

The Data

Import the Dataset.

```
In [2]: 1 data = pd.read_csv('Downloads/titanic-dataset.csv')
        2 data.head(5)
```

Out[2]:

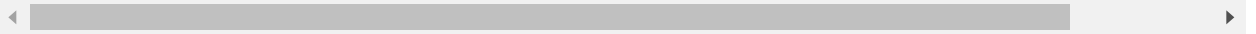
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	Na
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833	C8
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	Na
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C12
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	Na

In [3]:

```
1 data.head()
```

Out[3]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	Na
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C8
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	Na
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C12
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	Na



Missing Data

We can use seaborn to create a simple heatmap to see where we are missing data!

In [4]: 1 data.isnull()

Out[4]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Emb
0	False	False	False	False	False	False	False	False	False	False	True	
1	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	True
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	True
...
886	False	False	False	False	False	False	False	False	False	False	False	True
887	False	False	False	False	False	False	False	False	False	False	False	False
888	False	False	False	False	False	True	False	False	False	False	False	True
889	False	False	False	False	False	False	False	False	False	False	False	False
890	False	False	False	False	False	False	False	False	False	False	False	True

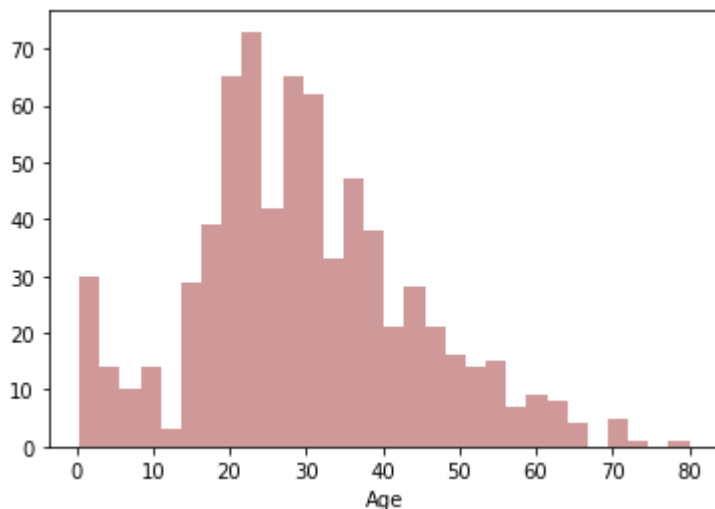
891 rows × 12 columns

In [5]: 1 sns.distplot(data['Age'].dropna(),kde=False,color='darkred',bins=30)

C:\Users\prpou\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

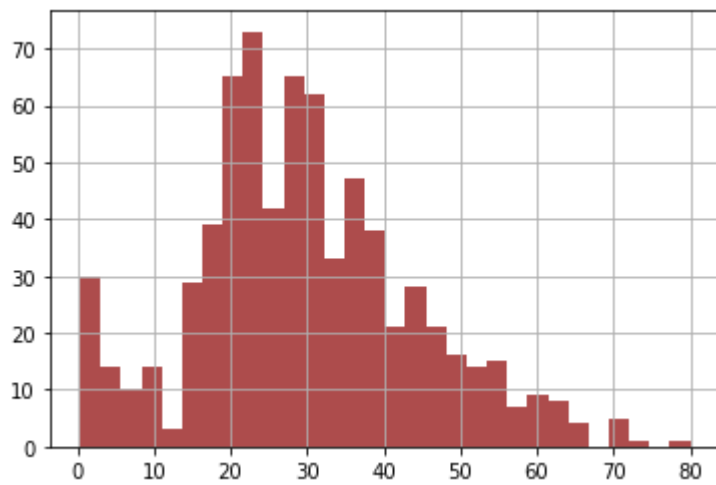
warnings.warn(msg, FutureWarning)

Out[5]: <AxesSubplot:xlabel='Age'>



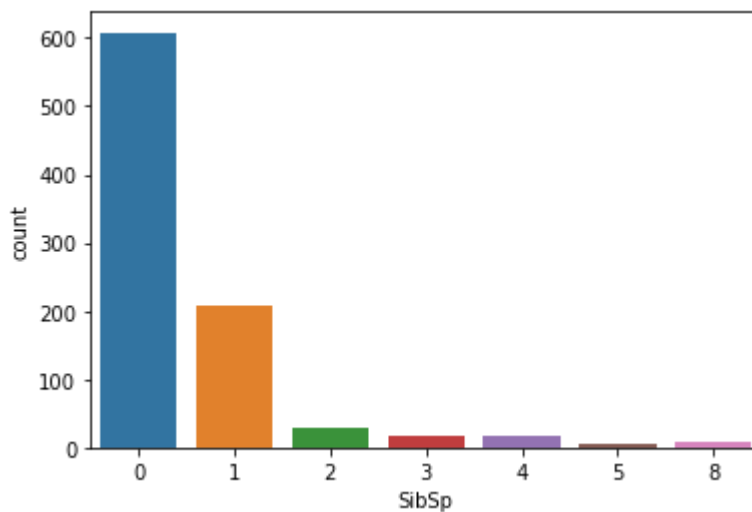
```
In [6]: 1 data['Age'].hist(bins=30,color='darkred',alpha=0.7)
```

Out[6]: <AxesSubplot:>



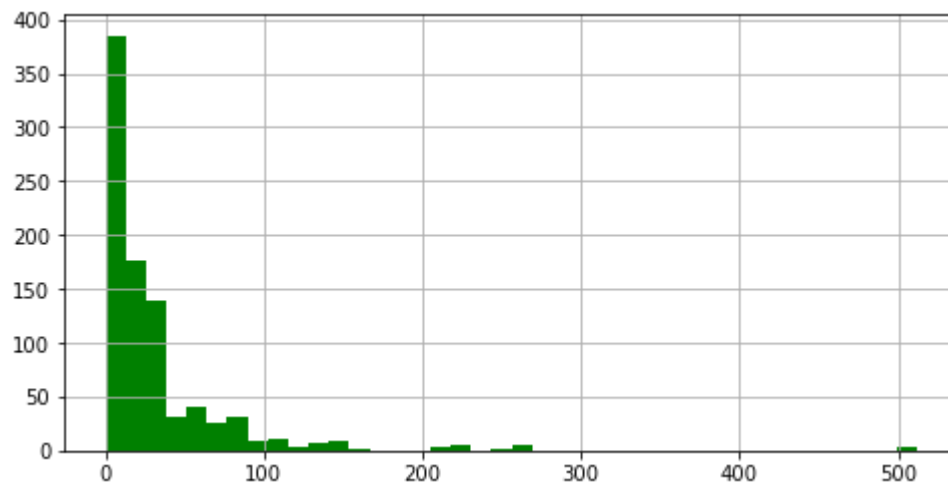
```
In [7]: 1 sns.countplot(x='SibSp',data=data)
```

Out[7]: <AxesSubplot:xlabel='SibSp', ylabel='count'>



```
In [8]: 1 data['Fare'].hist(color='green',bins=40,figsize=(8,4))
```

Out[8]: <AxesSubplot:>

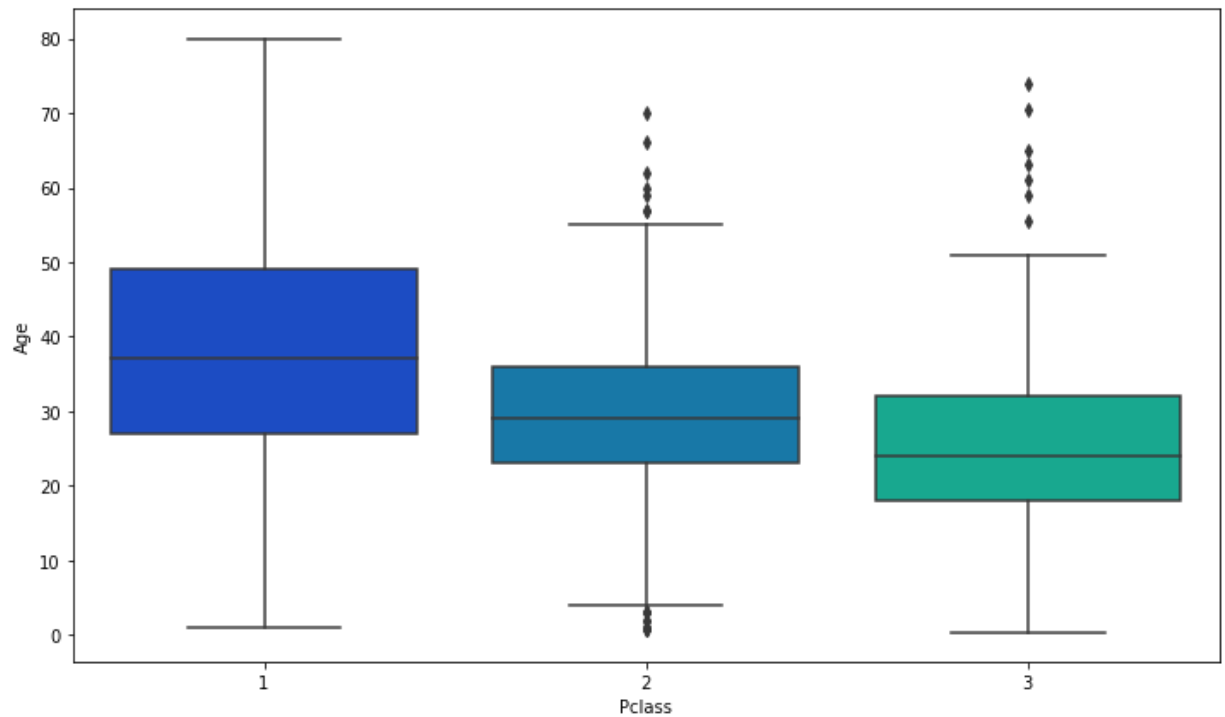


Data Cleaning

We want to fill in missing age data instead of just dropping the missing age data rows. One way to do this is by filling in the mean age of all the passengers (imputation). However we can be smarter about this and check the average age by passenger class. For example:

```
In [9]: 1 plt.figure(figsize=(12, 7))  
2 sns.boxplot(x='Pclass',y='Age',data=data,palette='winter')
```

```
Out[9]: <AxesSubplot:xlabel='Pclass', ylabel='Age'>
```



We can see the wealthier passengers in the higher classes tend to be older, which makes sense. We'll use these average age values to impute based on Pclass for Age.

```

In [10]: 1 def impute_age(cols):
          2     Age = cols[0]
          3     Pclass = cols[1]
          4
          5     if pd.isnull(Age):
          6
          7         if Pclass == 1:
          8             return 37
          9
         10         elif Pclass == 2:
         11             return 29
         12
         13         else:
         14             return 24
         15
         16     else:
         17         return Age

```

Now apply that function!

```

In [11]: 1 data['Age'] = data[['Age', 'Pclass']].apply(impute_age,axis=1)

```

Let's go ahead and drop the Cabin column and the row in Embarked that is NaN.

```

In [14]: 1 data.drop('Cabin',axis = 1,inplace=True)

```

```

In [15]: 1 data.head()

```

Out[15]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Emb
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	

In [16]: 1 data.dropna(inplace=True)

Converting Categorical Features

In [17]: 1 data.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 889 entries, 0 to 890
Data columns (total 11 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   PassengerId     889 non-null    int64
 1   Survived        889 non-null    int64
 2   Pclass          889 non-null    int64
 3   Name            889 non-null    object
 4   Sex             889 non-null    object
 5   Age            889 non-null    float64
 6   SibSp           889 non-null    int64
 7   Parch           889 non-null    int64
 8   Ticket          889 non-null    object
 9   Fare            889 non-null    float64
10   Embarked        889 non-null    object
dtypes: float64(2), int64(5), object(4)
memory usage: 83.3+ KB
```

In [18]: 1 sex = pd.get_dummies(data['Sex'],drop_first=True)
2 embark = pd.get_dummies(data['Embarked'],drop_first=True)

In [19]: 1 data.drop(['Sex', 'Embarked', 'Name', 'Ticket'],axis=1,inplace=True)

In [20]: 1 data = pd.concat([data,sex,embark],axis=1)

In [21]: 1 data.head()

Out[21]:

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	male	Q	S
0	1	0	3	22.0	1	0	7.2500	1	0	1
1	2	1	1	38.0	1	0	71.2833	0	0	0
2	3	1	3	26.0	0	0	7.9250	0	0	1
3	4	1	1	35.0	1	0	53.1000	0	0	1
4	5	0	3	35.0	0	0	8.0500	1	0	1

Logistic Regression model

Train Test Split


```
In [22]: 1 from sklearn.model_selection import train_test_split
```

```
In [23]: 1 X_train, X_test, y_train, y_test = train_test_split(data.drop('Survived',axi
2 data['Survived'], test_s
3 random_state=101)
```

Training and Predicting

```
In [24]: 1 from sklearn.linear_model import LogisticRegression
```

```
In [25]: 1 logmodel = LogisticRegression()
2 logmodel.fit(X_train,y_train)
```

C:\Users\prpou\anaconda3\lib\site-packages\sklearn\linear_model_logistic.py:81
8: 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> (<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 (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,

Out[25]: LogisticRegression()

```
In [26]: 1 predictions = logmodel.predict(X_test)
```

Let's move on to evaluate our model!

Evaluation

We can check precision,recall,f1-score using classification report.

```
In [27]: 1 from sklearn.metrics import classification_report
```

```
In [28]: 1 print(classification_report(y_test,predictions))
```

	precision	recall	f1-score	support
0	0.79	0.91	0.85	163
1	0.81	0.62	0.71	104
accuracy			0.80	267
macro avg	0.80	0.77	0.78	267
weighted avg	0.80	0.80	0.79	267

```
In [ ]: 1
```

```
In [ ]: 1
```