

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
matplotlib inline

In [1]: train = pd.read_csv("../input/train.csv")

In [2]: train.tail()

Out[3]:
   PassengerId  Survived  Pclass    Name     Sex  Age  SibSp  Parch    Ticket   Fare Cabin Embarked
886           887         0       2  Monville, Rev. Jozias      male  27.0   0   0   211536  11.00  NaN      S
887           888         1       1    Graham, Miss. Margaret Edith female  19.0   0   0  112053  30.00  B42      S
888           889         0       3  Johnson, Mrs. Catherine Helen "Catie" female  NaN   1   2   W/C. 6607  23.45  NaN      S
889           890         1       1    Behr, Mr. Karl Howell      male  26.0   0   0  111369  30.00  C148      C
890           891         0       3        Dooley, Mr. Patrick      male  32.0   0   0  370376  7.75  NaN      Q

In [4]: sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis')

Out[4]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b3737d>

In [5]: train.isnull().sum().sort_values(ascending=False)

Out[5]:
Cabin          887
Age            372
Embarked       2
Fare           0
Parch          0
SibSp          0
Sex            0
Name           0
Pclass         0
Survived       0
PassengerId    0
dtype: int64

In [6]: sns.set_style('whitegrid')
sns.countplot(x='Survived',data=train,palette='RdBu_r')

Out[6]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b0886d>

In [7]:
sns.set_style('whitegrid')
sns.countplot(x='Survived',hue='Sex',data=train,pallette='RdBu_r')

Out[7]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b4b3bd>

In [8]:
sns.set_style('whitegrid')
sns.countplot(x='Survived',hue='Pclass',data=train,pallette='rainbow')

Out[8]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b3c4ad>

In [9]: train['Age'].hist(bins=39,color='darkred',alpha=0.7)

Out[9]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b8b12d>

In [10]:
sns.countplot(x='SibSp',data=train)

Out[10]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b0e3cb>

In [11]:
sns.countplot(x='Parch',data=train)

Out[11]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b955db>

In [12]: train['Fare'].hist(color='green',bins=40,figure=plt.figure(figsize=(8,4)))

Out[12]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b74b3d>

In [13]:
plt.figure(figsize=(12,7))
sns.boxplot(x='Pclass',y='Age',data=train,pallette='winter')

Out[13]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b07f0d>

In [14]:
def impute_age(cols):
    Age = cols[0]
    Pclass = cols[1]

    if pd.isnull(Age):

        if Pclass == 1:
            return 37
        elif Pclass == 2:
            return 29
        else:
            return 24
    else:
        return Age

Now apply that function!

In [15]: train['Age'] = train[['Age','Pclass']].apply(impute_age,axis=1)

In [16]: train['Embarked'] = train['Embarked'].fillna('S')

Now let's check that heat map again!

In [17]: sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis')

Out[17]:
<matplotlib.axes._subplots.AxesSubplot at 8x7f4793b31f2d>

In [18]: train.drop("Cabin",axis=1,inplace=True)

Out[18]:
   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 21371  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  Heikinen, Miss. Laina      female  26.0   0   0  STON/O2 3101282  7.9250      S
3             4         1       1  Futerle, 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  8.6500      S

In [19]: train.droptna(inplace=True)

In [20]: train.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 891 entries, 0 to 890
Data columns (total 11 columns):
PassengerId    891 non-null int64
Survived       891 non-null int64
Pclass         891 non-null int64
Name           891 non-null object
Sex            891 non-null object
Age            891 non-null float64
SibSp          891 non-null int64
Parch          891 non-null int64
Ticket         891 non-null object
Fare           891 non-null float64
Embarked       891 non-null object
dtypes: float64(2), int64(4), object(4)
memory usage: 83.5+ KB

In [21]: sex = pd.get_dummies(train['Sex'],drop_first=True)
embark = pd.get_dummies(train['Embarked'],drop_first=True)

In [22]: train.drop(['Sex','Embarked','Name','Ticket'],axis=1,inplace=True)

In [23]: train = pd.concat([train,sex,embark],axis=1)

In [24]: train.head()

Out[25]:
   PassengerId  Survived  Pclass  Age  SibSp  Parch  Fare  male  Q  S
0             1         0       3  22.0   1.0   0.0  7.2500   1.0  0.1
1             2         1       1  38.0   1.0   0.0  71.2833   0.0  0.0
2             3         1       3  26.0   0.0   0.0  7.9250   0.0  0.1
3             4         1       1  35.0   1.0   0.0  53.1000   0.0  0.1
4             5         0       3  35.0   0.0   0.0  8.6500   1.0  0.1

In [26]: from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(train.drop(['Survived'],axis=1),train['Survived'],test_size=0.18,random_state=191)

In [27]: from sklearn.linear_model import LogisticRegression
logmodel = LogisticRegression()
logmodel.fit(X_train,y_train)

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:422: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)

Out[28]:
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, max_iter=100, multi_class='warn', n_jobs=None, penalty='l2',
                    random_state=None, solver='warn', tol=0.0001, verbose=0, warm_start=False)

In [29]: predictions = logmodel.predict(X_test)
X_test.head()

Out[30]:
   PassengerId  Pclass  Age  SibSp  Parch  Fare  male  Q  S
331           332   1  45.5   0.0   0.0  8.5500   1.0  0.1
700           701   1  18.0   1.0   0.0  227.525   0.0  0.0
748           749   1  19.0   1.0   0.0  53.100   1.0  0.1
751           752   3  6.0   0.0   1.0  12.475   1.0  0.1
481           482   2  29.0   0.0   0.0  0.000   1.0  0.1

In [31]: predictions

Out[31]:
array([0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0,
       1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0,
       0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0,
       1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0,
       0, 1])

In [32]: from sklearn.metrics import classification_report,confusion_matrix

In [33]: print(confusion_matrix(y_test,predictions))

[[44  5]
 [13 26]]

In [34]: print(classification_report(y_test,predictions))

              precision    recall  f1-score   support

0               0.78        0.90        0.84         51
1               0.84        0.67        0.74         39

 accuracy          0.88         0.88         90
 macro avg         0.81        0.78        0.79         90
weighted avg         0.81        0.80        0.88         90

In [35]: from sklearn.tree import DecisionTreeClassifier

In [36]: dt_model=DecisionTreeClassifier()
dt_model.fit(X_train,y_train)

Out[36]:
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=None, splitter='best')

In [37]: dt_pred = dt_model.predict(X_test)

In [38]: print(confusion_matrix(y_test,dt_pred))

[[44  9]
 [14 25]]

In [39]: print(classification_report(y_test,dt_pred))

              precision    recall  f1-score   support

0               0.75        0.82        0.79         51
1               0.74        0.64        0.68         39

 accuracy          0.74        0.73        0.74         90
 macro avg         0.74        0.73        0.73         90
weighted avg         0.74        0.74        0.74         90

In [40]: from sklearn.ensemble import RandomForestClassifier

In [41]: rf = RandomForestClassifier(n_estimators=500)
rf.fit(X_train,y_train)

Out[41]:
RandomForestClassifier(boosting=False, class_weight=None, criterion='gini',
                        max_depth=None, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, n_estimators=500,
                        n_jobs=None, oob_score=False, random_state=None,
                        verbose=0, warm_start=False)

In [42]: rf_pre_rf_predict(X_test)

In [43]: print(confusion_matrix(y_test,rf_pre))

[[44  1]
 [15 24]]

In [44]: print(classification_report(y_test,rf_pre))

              precision    recall  f1-score   support

0               0.76        0.94        0.84         51
1               0.89        0.62        0.73         39

 accuracy          0.83        0.78        0.88         90
 macro avg         0.82        0.78        0.79         90
weighted avg         0.82        0.80        0.79         90

In [45]: from xgboost import XGBClassifier
xgbost = XGBClassifier(n_estimators=1000)
xgbost.fit(X_train,y_train)

Out[45]:
XGBClassifier(base_score=1, booster='gbtree', colsample_bylevel=1,
               colsample_bynode=1, colsample_bynode=1, gamma=0,
               learning_rate=0.1, max_delta_step=0, max_depth=3,
               min_child_weight=1, missing=None, n_estimators=1000,
               nthread=-1, objective='binary:logistic', random_state=0,
               reg_alpha=0.0, reg_lambda=1.0, scale_pos_weights=1.0, seed=None,
               silent=None, subsample=1, verbosity=1)

In [46]: xg_pred = xgboost.predict(X_test)

In [47]: print(confusion_matrix(y_test,xg_pred))

[[44  7]
 [14 23]]

In [48]: print(classification_report(y_test,xg_pred))

              precision    recall  f1-score   support

0
```