In [1]	<pre>import pandas as pd import numpy as np import matplotlib.pyplot as plt; import seaborn as sns %matplotlib inline</pre>
	886 887 0 2 Montvila, Rev. Juozas male 27.0 0 0 211536 13.00 NaN S 887 888 1 1 1 Graham, Miss. Margaret Edith female 19.0 0 0 112053 30.00 B42 S
In [4]	888 889 0 3 Johnston, Miss. Catherine Helen "Carrie" 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  sns.heatmap(train.isnull(), yticklabels=False, cbar=False, cmap='viridis') <matplotlib.axessubplots.axessubplot 0x7fd795f337b8="" at=""></matplotlib.axessubplots.axessubplot>
	Survived - Survived - Pclass - Name - Sex - Age - Ticket - Ticket - Fare - Cabin - Embarked -
In [5] Out[5]	Cabin 687 Age 177 Embarked 2 Fare 0 Ticket 0
	Parch 0 SibSp 0 Sex 0 Name 0 Pclass 0 Survived 0 PassengerId 0 dtype: int64  Roughly 20 percent of the Age data is missing. The proportion of Age missing is likely small enough for reasonable replacement with some form of imputation. Looking at the Cabin column, it looks like we are just missing too much of tha
In [6]	data to do something useful with at a basic level. We'll probably drop this later, or change it to another feature like "Cabin Known: 1 or 0"  Let's continue on by visualizing some more of the data! Check out the video for full explanations over these plots, this code is just to serve as reference.  sns.set_style('whitegrid') sns.countplot(x='Survived', data=train, palette='RdBu_r')
	500 400 Tig 300
In [7]	200 100 0 Survived  sns.set_style('whitegrid')
Out[7]	sns.countplot(x='Survived', hue='Sex', data=train, palette='RdBu_r')  : <matplotlib.axessubplots.axessubplot 0x7fd793dc4b38="" at="">  400  Sex  male female</matplotlib.axessubplots.axessubplot>
In [8]	<pre>sns.countplot(x='Survived', hue='Pclass', data=train, palette='rainbow')</pre>
	300 250 150 100
In [9] Out[9]	train['Age'].hist(bins=30,color='darkred',alpha=0.7) <matplotlib.axessubplots.axessubplot 0x7fd793c8b128="" at=""></matplotlib.axessubplots.axessubplot>
	70 60 50 40 30
In [10]	
out[10]	<pre>control in the second sec</pre>
	200 100 0 1 2 3 4 5 8 SibSp
	<pre>sns.countplot(x='Parch', data=train)  countplot(x='Parch', data=train</pre>
	400 200 100 0 1 2 3 4 5 6
In [12] Out[12]	rain['Fare'].hist(color='green',bins=40,figsize=(8,4))  **matplotlib.axessubplots.AxesSubplot at 0x7fd793b74b38>  **matplotlib.axessubplots.AxesSubplot at 0x7fd793b74b38>
	300 250 200 150 100 50
In [13] Out[13]	plt.figure(figsize=(12, 7)) sns.boxplot(x='Pclass', y='Age', data=train, palette='winter')  * <matplotlib.axessubplots.axessubplot 0x7fd793b67fd0="" at=""></matplotlib.axessubplots.axessubplot>
	80 70 60 50
	20 10
	1 2 3 Pclass  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 [14]	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!  train['Age'] = train[['Age', 'Pclass']].apply(impute_age, axis=1)  train['Embarked'] = train['Embarked'].fillna('S')  Now let's check that heat map again!
	sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis') <pre> <matplotlib.axessubplots.axessubplot 0x7fd793a13f28="" at=""> </matplotlib.axessubplots.axessubplot></pre>
	Great! Let's go ahead and drop the Cabin column and the row in Embarked that is NaN.  We will sum of family member
In [18] In [19] Out[19]	<pre>train.drop('Cabin',axis=1,inplace=True) train.head()</pre>
	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
	train.dropna(inplace=True)  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</class>
	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
In [22] In [23]	<pre>embark = pd.get_dummies(train['Embarked'],drop_first=True)</pre>
	train.head()  PassengerId Survived Pclass Age SibSp Parch Fare male Q S  0 1 0 3 22.0 1 0 71.2833 0 0 0  1 2 1 1 38.0 1 0 71.2833 0 0
In [26]	2
In [28] In [29]	train['Survived'], test_size=0.10, random_state=101)  from sklearn.linear_model import LogisticRegression
Out[29] In [30]	<pre>/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.     FutureWarning)  LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,</pre>
Out[30]	331 332 1 45.5 0 0 28.500 1 0 1 700 701 1 18.0 1 0 227.525 0 0 0 748 749 1 19.0 1 0 53.100 1 0 1
	751 752 3 6.0 0 1 12.475 1 0 1  481 482 2 29.0 0 0 0.000 1 0 1  : predictions  : array([0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
	0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
	<pre>from sklearn.metrics import classification_report,confusion_matrix  print(confusion_matrix(y_test,predictions))  [[46    5]     [13    26]]  print(classification_report(y_test,predictions))      precision recall f1-score support</pre>
	0 0.78 0.90 0.84 51 1 0.84 0.67 0.74 39 accuracy 0.80 90 macro avg 0.81 0.78 0.79 90 weighted avg 0.81 0.80 0.80 90
In [36]	dt_model.fit(X_train,y_train)  DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
In [37] In [38]	<pre>min_samples_leaf=1, min_samples_split=2,</pre>
In [39]	[14 25]]
	accuracy
Out[41]	RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
In [43]	<pre>: rf_pre=rf.predict(X_test)  : print(confusion_matrix(y_test,rf_pre))  [[48     3]      [15     24]]  : print(classification_report(y_test,rf_pre))</pre>
	precision recall f1-score support  0 0.76 0.94 0.84 51 1 0.89 0.62 0.73 39  accuracy macro avg 0.83 0.78 0.78 90 weighted avg 0.82 0.80 0.79 90
	<pre>from xgboost import XGBClassifier xgboost = XGBClassifier(n_estimators=1000) xgboost.fit(X_train,y_train)  XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,</pre>
	<pre>nthread=None, objective='binary:logistic', random_state=0,     reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,     silent=None, subsample=1, verbosity=1)  : xg_pred = xgboost.predict(X_test)  : print(confusion_matrix(y_test,xg_pred))  [[44 7]</pre>
In [48]	[144 7] [16 23]]  print(classification_report(y_test, xg_pred))  precision recall f1-score support  0 0.73 0.86 0.79 51 1 0.77 0.59 0.67 39  accuracy 0.74 90
In [49]	macro avg 0.75 0.73 0.73 90 weighted avg 0.75 0.74 0.74 90  import keras from keras.layers import Dense from keras.models import Sequential
In [50]	<pre>ann.add(Dense(units= 32,init= 'uniform', activation = 'relu', input_dim=9)) ann.add(Dense(units= 32,init= 'uniform', activation = 'relu')) ann.add(Dense(units= 1,init= 'uniform', activation = 'sigmoid')) ann.compile(optimizer='adam',</pre>
Jr	/opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:2: UserWarning: Update your `Dense` call to the Keras 2 API: `Dense(units=32, activation="relu", input_dim=9, kernel_initialize="uniform")`  /opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:3: UserWarning: Update your `Dense` call to the Keras 2 API: `Dense(units=32, activation="relu", kernel_initializer="uniform")  This is separate from the ipykernel package so we can avoid doing imports until  /opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:4: UserWarning: Update your `Dense` call to the Keras 2 API: `Dense(units=1, activation="sigmoid", kernel_initializer="uniform after removing the cwd from sys.path.
Out[51]	<pre>: ann.fit(X_train,y_train, batch_size=32, nb_epoch=300,verbose= 0)  /opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:1: UserWarning: The `nb_epoch` argument in `fit` has been renamed `epochs`.    """Entry point for launching an IPython kernel.  : <keras.callbacks.callbacks.history 0x7fd782815e10="" at="">  : ann_pred = ann.predict(X_test)    ann_pred = [ 1 if y&gt;=0.5 else 0 for y in ann_pred]    print(ann_pred)</keras.callbacks.callbacks.history></pre>
<sub>[</sub> 54]	print(classification_report(y_test, ann_pred))  precision recall f1-score support  0 0.78 0.90 0.84 51 1 0.84 0.67 0.74 39  accuracy 0.80 90 macro avg 0.81 0.78 0.79 90 weighted avg 0.81 0.80 0.80 90
In [55] In [56] Out[56]	Now we will use test dataset  test = pd.read_csv('/input/test.csv')  sns.heatmap(test.isnull())
ac[56]	<pre></pre>
	220 240 250 250 250 250 250 250 250 250 250 25
	<pre>test.drop('Cabin', axis=1,inplace=True)  test['Fare'].fillna(test['Fare'].median(), inplace=True)  test.info() <class 'pandas.core.frame.dataframe'=""></class></pre>
	RangeIndex: 418 entries, 0 to 417 Data columns (total 10 columns): PassengerId
	Parch 418 non-null int64  Ticket 418 non-null object  Fare 418 non-null float64  Embarked 418 non-null object  dtypes: float64(2), int64(4), object(4)  memory usage: 32.8+ KB  test.head()
Out[60]	PassengerId         Pclass         Pclass         Name         Sex         Age         SibSp         Parch         Ticket         Fare         Embarked           0         892         3         Kelly, Mr. James         male         34.5         0         33091         7.8292         Q           1         893         3         Wilkes, Mrs. James (Ellen Needs)         female         47.0         1         0         363272         7.0000         S           2         894         2         Myles, Mr. Thomas Francis         male         62.0         0         240276         9.6875         Q           3         895         3         Wirz, Mr. Albert         male         27.0         0         315154         8.6625         S           4         896         3         Hirvonen, Mrs. Alexander (Helga E Lindqvist)         female         22.0         1         3101298         12.2875         S
	4 896 3 Hirvonen, Mrs. Alexander (Helga E Lindqvist) female 22.0 1 1 3101298 12.2875 S  : test['Age'] = test[['Age', 'Pclass']].apply(impute_age, axis=1)  : sex_test = pd.get_dummies(test['Sex'], drop_first=True) embark_test= pd.get_dummies(test['Embarked'], drop_first=True)
	<pre>test.drop(['Sex', 'Embarked', 'Name', 'Ticket'], axis=1, inplace=True) test = pd.concat([test, sex_test, embark_test], axis=1)</pre>
In [64]	
In [64] In [65] Out[65]	PassengerId         Pclass         Age         SibSp         Parch         Fare         male         Q         S           0         892         3         34.5         0         0         7.8292         1         1         0           1         893         3         47.0         1         0         7.0000         0         0         1           2         894         2         62.0         0         9.6875         1         1         0           3         895         3         27.0         0         8.6625         1         0         1           4         896         3         22.0         1         1 2.2875         0         0         1
In [64] In [65] Out[65]	Passengerid   Pclass   Age   SibSp   Parch   Fare   male   Q   S     1
<pre>In [64] In [65] Out[65] In [66] Out[66]</pre>	Palae   Pala
In [64] In [65] Out[65] In [66] Out[66] In [67] In [68] In [69]	Passengerid   Polas   Reg   Sin   Sin   Parch   Fare   male   Q   S
In [64] In [65] Out [65] In [66] Out [67] In [68] In [69] In [70]	Passemgerid   Polars   Age   SibSp   Part   Fare   male   Q   S
In [64] In [65] Out [65] In [66] Out [66] In [67] In [68] In [69] In [70] In [71]	Passengered Pollas Age 50:50 Parch   Fare rule Q S