


```
275: Sex_female Sex_male
65 0 1
445 0 1
659 1 0
591 1 0
690 0 1
... ...
275 1 0
839 0 1
577 1 0
400 0 1
383 1 0
746 rows x 2 columns

In [276]: pclass_dummies = pd.get_dummies(train_rebal['Pclass'], prefix = 'Pclass', drop_first=False)
pclass_dummies

Out[276]:
Pclass_1 Pclass_2 Pclass_3
65 0 0 1
445 1 0 0
659 1 0 0
591 1 0 0
690 1 0 0
... ...
275 1 0 0
839 1 0 0
577 1 0 0
400 0 0 1
383 1 0 0
746 rows x 3 columns

In [277]: embarked_dummies = pd.get_dummies(train_rebal['Embarked'], prefix = 'Embarked', drop_first=False)
embarked_dummies

Out[277]:
Embarked_C Embarked_Q Embarked_S
65 1 0 0
445 0 0 1
659 1 0 0
591 1 0 0
690 0 0 1
... ...
275 0 0 1
839 1 0 0
577 0 0 1
400 0 0 1
383 0 0 1
746 rows x 3 columns

In [278]: age_dummies = pd.get_dummies(train_rebal['Age_c'], prefix = 'Age', drop_first=False)
age_dummies

Out[278]:
Age_Baby/Toddler Age_Child Age_Adult Age_Elderly
65 0 0 1 0
445 0 1 0 0
659 0 0 1 0
591 0 0 1 0
690 0 0 1 0
... ...
275 0 0 1 0
839 0 0 1 0
577 0 0 1 0
400 0 0 1 0
383 0 0 1 0
746 rows x 4 columns

In [279]: ticket_dummies = pd.get_dummies(train_rebal['Ticket2'], prefix = 'Ticket2', drop_first=False)
ticket_dummies

Out[279]:
Ticket2_1 Ticket2_2 Ticket2_3 Ticket2_4 Ticket2_5 Ticket2_6 Ticket2_7 Ticket2_9 Ticket2_A Ticket2_C Ticket2_F Ticket
65 0 1 0 0 0 0 0 0 0 0 0 0
445 0 0 1 0 0 0 0 0 0 0 0 0
659 0 0 1 0 0 0 0 0 0 0 0 0
591 0 0 1 0 0 0 0 0 0 0 0 0
690 1 0 0 0 0 0 0 0 0 0 0 0
... ...
275 1 0 0 0 0 0 0 0 0 0 0 0
839 1 0 0 0 0 0 0 0 0 0 0 0
577 1 0 0 0 0 0 0 0 0 0 0 0
400 0 0 0 0 0 0 0 0 0 0 0 0
383 1 0 0 0 0 0 0 0 0 0 0 0
746 rows x 15 columns

Test Set

In [280]: sex_dummies_test = pd.get_dummies(test['Sex'], prefix = 'Sex', drop_first=False)
pclass_dummies_test = pd.get_dummies(test['Pclass'], prefix = 'Pclass', drop_first=False)
embarked_dummies_test = pd.get_dummies(test['Embarked'], prefix = 'Embarked', drop_first=False)
age_dummies_test = pd.get_dummies(test['Age_c'], prefix = 'Age', drop_first=False)
ticket2_dummies_test = pd.get_dummies(test['Ticket2'], prefix = 'Ticket2', drop_first=False)
```

Standardize Numeric Variables

Numeric variables are normalized using Min-Max scaling

Numeric variables include: Age, Fare, Fam

```
In [281]: #normalized variables are added back into original df as new columns 'Var_mm'
train_rebal['Age_mm'] = MinMaxScaler().fit_transform(train_rebal[['Age']])
train_rebal['Fare_mm'] = MinMaxScaler().fit_transform(train_rebal[['Fare']])
train_rebal['Fam_mm'] = MinMaxScaler().fit_transform(train_rebal[['Fam']])
train_rebal.head(10)

Out[281]:
PassengerId Survived Pclass Name Sex Age SibSp Parch Fare Embarked Age_c Fam Fare2 Ticket
65 66 1 3 Moubarek, male 28.25 1 1 2661 15.2458 C Adult 2 Mid
445 446 1 1 Dodge, male 4.00 0 2 33638 81.8583 S Child 2 Max
659 660 0 1 Newell, Mr. male 58.00 0 2 35273 113.2750 C Adult 2 Max
591 592 1 1 Stephenson, female 52.00 1 0 36947 78.2667 C Adult 1 Max
690 691 1 1 Dick, Mr. male 31.00 1 0 17474 57.0000 S Adult 1 Max
397 398 0 2 McKane, Mr. male 46.00 0 0 28403 26.0000 S Adult 0 High
810 811 0 3 Alexander, male 26.00 0 0 3474 7.8875 S Adult 0 Low
273 274 0 1 Natsch, Mr. male 37.00 0 1 PC 29.7000 C Adult 1 High
587 588 1 1 Wilhelm, male 60.00 1 1 13567 79.2000 C Adult 2 Max
673 674 1 2 Withlms, male 31.00 0 0 244270 13.0000 S Adult 0 Mid
```

Models

Logistic Regression

```
In [282]: train_rebal.head()

Out[282]:
PassengerId Survived Pclass Name Sex Age SibSp Parch Ticket Fare Embarked Age_c Fam Fare2 Ticket
65 66 1 3 Moubarek, male 28.25 1 1 2661 15.2458 C Adult 2 Mid
445 446 1 1 Dodge, male 4.00 0 2 33638 81.8583 S Child 2 Max
659 660 0 1 Newell, Mr. male 58.00 0 2 35273 113.2750 C Adult 2 Max
591 592 1 1 Stephenson, female 52.00 1 0 36947 78.2667 C Adult 1 Max
690 691 1 1 Dick, Mr. male 31.00 1 0 17474 57.0000 S Adult 1 Max
```

```
In [283]: train_rebal.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 746 entries, 65 to 383
Data columns (total 17 columns):
# Column Non-Null Count Dtype
---
0 PassengerId 746 non-null int64
1 Survived 746 non-null int64
2 Pclass 746 non-null int64
3 Name 746 non-null object
4 Sex 746 non-null object
5 Age 746 non-null float64
6 SibSp 746 non-null int64
7 Parch 746 non-null int64
8 Ticket 746 non-null object
9 Fare 746 non-null float64
10 Embarked 746 non-null object
11 Age_c 746 non-null category
12 Fam 746 non-null int64
13 Fare2 733 non-null category
14 Ticket2 746 non-null object
15 Title 746 non-null object
16 Age_mm 746 non-null float64
dtypes: category(2), float64(3), int64(6), object(6)
memory usage: 111.3+ KB

*Note that corr() applies to numeric variables only.
```

```
In [284]: train_rebal.corr()

Out[284]:
PassengerId Survived Pclass Age SibSp Parch Fare Fam Age_mm
PassengerId 1.000000 -0.035589 0.010040 0.001134 -0.013229 0.0095082 0.006857 0.040526 0.001134
Survived -0.035589 1.000000 -0.347120 -0.074136 -0.032727 0.108903 0.233654 0.034530 -0.074136
Pclass 0.010040 -0.347120 1.000000 -0.348372 0.062213 0.012752 -0.252360 0.048712 -0.348372
Age 0.001134 -0.074136 -0.348372 1.000000 -0.211754 -0.205865 0.129484 -0.250335 1.000000
SibSp -0.013229 -0.032727 0.062213 -0.211754 1.000000 0.386197 0.131973 0.877294 -0.211754
Parch 0.0095082 0.108903 0.012752 -0.205865 0.386197 1.000000 0.194591 0.781525 -0.205865
Fare 0.006857 0.233654 -0.052360 0.129484 0.131973 0.194591 1.000000 0.190509 0.129484
Fam 0.040526 0.034530 0.048712 -0.250335 0.877294 0.781525 0.190509 1.000000 -0.250335
Age_mm 0.001134 -0.074136 -0.348372 1.000000 -0.211754 -0.205865 0.129484 -0.250335 1.000000
```

Prepare separate predictors from response:

```
In [285]: #Predictors
X_logi = pd.DataFrame(train_rebal[['Pclass', 'Fam', 'Age', 'Fare', 'Sex']])
X_logi

Out[285]:
Pclass Fam Age Fare Sex
65 3 2 28.25 15.2458 male
445 1 2 4.00 81.8583 male
659 1 2 58.00 113.2750 male
591 1 1 52.00 78.2667 female
690 1 1 31.00 57.0000 female
... ...
275 1 1 63.00 77.9583 female
839 1 0 28.25 29.7000 male
577 1 1 39.00 55.9000 female
400 3 0 39.00 7.9250 male
383 1 1 35.00 52.0000 female
746 rows x 5 columns
```

```
In [286]: #Response
y_logi = pd.DataFrame(train_rebal[['Survived']])
y_logi

Out[286]:
Survived
65 1
445 1
659 0
591 1
690 1
... ...
275 1
839 1
577 1
400 1
383 1
746 rows x 1 columns
```

Convert categorical variable to numeric:

```
In [287]: X_logi['Sex'].replace(['male', 'female'],
[0, 1], inplace=True)

Convert predictors to dummy variables:
```

```
In [288]: X_logi_d = pd.get_dummies(X_logi)
X_logi_d.head()

Out[288]:
Pclass Fam Age Fare Sex
65 3 2 28.25 15.2458 0
445 1 2 4.00 81.8583 0
659 1 2 58.00 113.2750 0
591 1 1 52.00 78.2667 1
690 1 1 31.00 57.0000 0
746 rows x 6 columns
```

Logistic Regression Model:

```
In [289]: #Add constant to X_logi
X_logi_d = sm.add_constant(X_logi_d)

logreg01 = sm.Logit(y_logi, X_logi_d).fit()
logreg01.summary()

Optimization terminated successfully.
Current function value: 0.428228
Iterations 6
```

Model: Logit Pseudo R-squared: 0.382

Dependent Variable: Survived AIC: 660.9157

Date: 2022-12-09 22:55 BIC: 678.6040

No. Observations: 746 Log-Likelihood: -319.46

Df Model: 5 Df Residuals: 740 LL-Null: -617.09

DF Residuals: 740 LLR p-value: 0.13776

Converged: 1.0000 Scale: 1.0000

No. Iterations: 6.0000

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
const	3.2571	0.5057	6.4405	0.0000	2.2659	4.2483
Pclass	-1.2666	0.1502	-8.4350	0.0000	-1.5609	-0.9723
Fam	-0.2664	0.03726	-3.6709	0.0002	-0.4087	-0.1242
Age	-0.0503	0.0084	-5.9807	0.0000	-0.0668	-0.0338
Fare	0.0009	0.0021	0.4048	0.6856	-0.0033	0.0050
Sex	3.1643	0.2331	13.5721	0.0000	2.7073	3.6212

Logistic Regression Shaving:

Removed Ticket2 because it's showing NaNs for the metrics.

Embarked is showing multicollinearity with very high Std Error, 0 values for z score, and p-value of 1.

Logistic Regression Interpretation:

Model has a pseudo R² of 36.3% which suggest that the model is a good fit. LLR p-value is less than 0.5 significance level. BIC is slightly higher than AIC but they are around the same ball park. All the variable metrics look reasonable.

```
In [290]: lg_pred_tr = logreg01.predict(X_logi_d)
lg_pred_tr_logis = np.where (lg_pred_tr > 0.5, 1, 0)
```

Test Dataset Validation

```
In [291]: test.head()

Out[291]:
PassengerId Survived Pclass Name Sex Age SibSp Parch Ticket Fare Embarked Age_c Fam Ticket2
725 726 0 3 Oreskovik, male 20.0 0 0 315084 8.6625 S Adult 0 3
861 862 0 2 Gies, Mr. male 21.0 1 0 28134 11.5000 S Adult 1 2
528 529 0 3 Salonen, Mr. male 39.0 0 0 3101296 7.9250 S Adult 0 3
46 47 0 3 Lemon, Mr. male 28.0 1 0 370371 15.5000 Q Adult 1 3
627 628 1 1 Longley, female 21.0 0 0 13502 77.9583 S Adult 0 1
```

```
In [292]: X_logi_test = pd.DataFrame(test[['Pclass', 'Fam', 'Age', 'Fare', 'Sex']])
y_logi_test = pd.DataFrame(test[['Survived']])

In [293]: X_logi_test['Sex'].replace(['male', 'female'],
[0, 1], inplace=True)
X_logi_d_test = pd.get_dummies(X_logi_test)

In [294]: X_logi_d_test

Out[294]:
Pclass Fam Age Fare Sex
725 3 0 20.0 8.6625 0
861 2 1 21.0 11.5000 0
528 3 1 39.0 7.9250 0
46 3 1 28.0 15.5000 0
627 1 0 21.0 77.9583 1
... ...
360 3 5 40.0 27.0000 0
856 1 2 45.0 164.8667 1
199 2 0 24.0 13.0000 1
451 3 1 28.0 19.9667 0
417 2 2 18.0 13.0000 1
295 rows x 5 columns
```

```
In [295]: X_logi_d_test = sm.add_constant(X_logi_d_test)
logreg01_test = sm.Logit(y_logi_test, X_logi_d_test).fit()
logreg01_test.summary()

Optimization terminated successfully.
Current function value: 0.485070
Iterations 7
```

Model: Logit Pseudo R-squared: 0.281

Dependent Variable: Survived AIC: 298.1912

Date: 2022-12-09 22:55 BIC: 320.3130

No. Observations: 295 Log-Likelihood: -139.46

Df Model: 5 Df Residuals: 289 LL-Null: -198.94

DF Residuals: 289 LLR p-value: 1.8066e-22

Converged: 1.0000 Scale: 1.0000

No. Iterations: 7.0000

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
const	0.4190	0.8941	0.4686	0.6393	-1.3334	2.1714
Pclass	-0.6449	0.2470	-2.6113	0.0090	-1.1289	-0.1609
Fam	-0.2508	0.1225	-2.0473	0.0406	-0.4909	-0.0107
Age	-0.0147	0.0129	-1.1432	0.2530	-0.0400	0.0105
Fare	0.0175	0.0081	2.1679	0.0302	0.0017	0.0333
Sex	2.2155	0.3216	6.8896	0.0000	1.5852	2.8458

Contingency Table

```
Predictions:

In [296]: predictions_prob = logreg01_test.predict(X_logi_d_test)
predictions_prob.head()

Out[296]:
725 0.159946
861 0.226244
528 0.124429
46 0.192446
627 0.354487
dtype: float64

In [297]: cutoff = 0.5

In [298]: ypred_logis = np.where (predictions_prob > cutoff, 1, 0)
ypred_c['Survived'] = pd.DataFrame(ypred_logis)
ypred_c

In [299]: conf_matrix = pd.crosstab(test['Survived'], ypred_logis,
rownames = ['Actual'],
columns = ['Predicted'],
margins = True)
conf_matrix

Out[299]:
Predicted 0 1 All
Actual
0 149 27 176
1 36 83 119
All 185 110 295
```

Evaluation Metrics

```
In [300]: #Baseline
#Accuracy (all negative model) = TNR/QT
logis_baseline = round((176/296)*100, 2)
logis_baseline

Out[300]:
59.46

In [301]: #Accuracy = (TP+FP)/QT
logis_a = round(((150+82)/295)*100, 2)
logis_a

Out[301]:
78.64

In [302]: #Error rate = 1-Accuracy
logis_e = round(100-78.38, 2)
logis_e

Out[302]:
21.62

In [303]: #Sensitivity: Recall = TP/TAP
logis_r = round((82/119)*100, 2)
logis_r

Out[303]:
68.91

In [304]: #Specificity: Specificity = TN/TAN
logis_s = round((150/176)*100, 2)
logis_s

Out[304]:
85.23

In [305]: logis_t = [ ['Metrics', 'Score', '%'],
['Accuracy base', logis_baseline],
['Accuracy', logis_a],
['Error rate', logis_e],
['Sensitivity', logis_r],
['Specificity', logis_s],
]
print(tables(logis_t, headers='firstrow'))

Metrics Score %
-----
Accuracy base 59.46
Accuracy 78.64
Error rate 21.62
Sensitivity 68.91
Specificity 85.23
```

K-means Clustering

```
In [306]: #Predictors
X_kmeans = train_rebal[['Pclass', 'Fam', 'Age', 'Fare']]
X_kmeans

Out[306]:
Pclass Fam Age Fare
65 3 2 28.25 15.2458
445 1 2 4.00 81.8583
659 1 2 58.00 113.2750
591 1 1 52.00 78.2667
690 1 1 31.00 57.0000
... ...
275 1 1 63.00 77.9583
839 1 0 28.25 29.7000
577 1 1 39.00 55.9000
400 3 0 39.00 7.9250
383 1 1 35.00 52.0000
746 rows x 4 columns
```

Standardize predictors using Z-score transformation:

```
*Note: only numeric data

In [307]: X_kms = pd.DataFrame(stats.zscore(X_kmeans),
columns=['Pclass', 'Fam', 'Age', 'Fare'])

K-means clustering model:
```

```
In [308]: kmsans01 = KMeans(n_clusters = 4).fit(X_kms)
kmsans01

Out[308]:
KMeans(n_clusters=4)

Investigate:
```

```
In [309]: cluster = kmsans01.labels_

In [310]: Cluster1 = X_kmeans.loc[cluster == 0]
Cluster2 = X_kmeans.loc[cluster == 1]
Cluster3 = X_kmeans.loc[cluster == 2]
Cluster4 = X_kmeans.loc[cluster == 3]

In [311]: Cluster1.describe()

Out[311]:
Pclass Fam Age Fare
count 103.000000 103.000000 103.000000 103.000000
mean 2.611650 3.533981 11.902913 32.178321
std 0.564176 2.066614 11.782431 19.785354
min 1.000000 1.000000 0.420000 7.854200
25% 2.000000 2.000000 2.000000 19.258300
50% 3.000000 3.000000 8.000000 27.750000
75% 3.000000 5.000000 19.500000 36.877100
max 3.000000 10.000000 48.000000 120.000000
```

```
In [312]: Cluster2.describe()

Out[312]:
Pclass Fam Age Fare
count 232.000000 232.000000 232.000000 232.000000
mean 1.168103 0.741379 30.114224 58.230280
std 0.374767 0.854027 12.762562 38.327296
min 1.000000 0.000000 16.000000 0.000000
25% 1.000000 0.000000 28.250000 26.550000
50% 1.000000 1.000000 36.000000 52.000000
75% 1.000000 1.000000 48.250000 83.475000
max 2.000000 4.000000 71.000000 153.462500
```

```
In [313]: Cluster3.describe()

Out[313]:
Pclass Fam Age Fare
count 390.000000 390.000000 390.000000 390.000000
mean 2.782051 0.276923 28.014744 11.188694
std 0.413353 0.591453 7.940711 6.942283
min 1.000000 0.000000 11.000000 0.000000
25% 3.000000 0.000000 23.500000 7.750000
50% 3.000000 0.000000 28.250000 13.000000
75% 3.000000 0.000000 36.750000 30.500000
max 3.000000 2.000000 70.500000 73.500000
```

```
In [314]: Cluster4.describe()

Out[314]:
Pclass Fam Age Fare
count 210 21.000000 21.000000 21.000000
mean 1.195281 30.107143 281.309333
std 0.0 1.935877 11.629319 118.925800
min 1.0 0.000000 15.000000 15.550000
25% 1.0 0.000000 23.000000 21.137500
50% 1.0 1.000000 28.250000 24.752080
75% 1.0 4.000000 35.000000 263.000000
max 1.0 5.000000 64.000000 512.329200
```

Test Dataset Validation

```
In [315]: X_kmeans_test = test[['Pclass', 'Fam', 'Age', 'Fare']]
X_kmeans_test

Out[315]:
Pclass Fam Age Fare
725 3 0 20.0 8.6625
861 2 1 21.0 11.5000
528 3 0 39.0 7.9250
46 3 1 28.0 15.5000
627 1 0 21.0 77.9583
... ...
360 3 5 40.0 27.0000
856 1 2 45.0 164.8667
199 2 0 24.0 13.0000
451 3 1 28.0 19.9667
417 2 2 18.0 13.0000
295 rows x 4 columns
```

```
In [316]: X_kms_test = pd.DataFrame(stats.zscore(X_kmeans_test),
columns=['Pclass', 'Fam', 'Age', 'Fare'])
kmeans_test = KMeans(n_clusters = 4).fit(X_kms_test)
cluster_test = kmeans_test.labels_

In [317]: Cluster1_test = X_kmeans_test.loc[cluster_test == 0]
Cluster2_test = X_kmeans_test.loc[cluster_test == 1]
Cluster3_test = X_kmeans_test.loc[cluster_test == 2]
Cluster4_test = X_kmeans_test.loc[cluster_test == 3]

In [318]: Cluster1_test.describe()

Out[318]:
Pclass Fam Age Fare
count 179.000000 179.000000 179.000000 179.000000
mean 2.793296 0.402235 26.642458 10.611668
std 0.406077 0.681193 9.255007 32.000000
min 2.000000 0.000000 1.000000 0.000000
25% 3.000000 0.000000 21.500000 7.750000
50% 3.000000 0.000000 28.000000 8.662500
75% 3.000000 1.000000 29.000000 14.477100
max 3.000000 2.000000 61.000000 73.500000
```

```
In [319]: Cluster2_test.describe()

Out[319]:
Pclass Fam Age Fare
count 76.000000 76.000000 76.000000 76.000000
mean 1.289474 0.447368 39.677632 41.242270
std 0.484858 0.640723 13.868599 24.971153
min 1.000000 0.000000 16.000000 0.000000
25% 1.000000 0.000000 28.000000 26.000000
50% 1.000000 0.000000 36.750000 30.500000
75% 2.000000 1.000000 49.250000 60.044800
max 3.000000 3.000000 80.000000 91.079200
```

```
In [320]: Cluster3_test.describe()

Out[320]:
Pclass Fam Age Fare
count 26.000000 26.000000 26.000000 26.000000
mean 2.846154 0.576892 20.653846 35.033019
std 0.367946 2.035077 12.699425 14.440364
min 2.000000 3.000000 1.000000 18.750000
25% 3.000000 4.000000 8.250000 25.850025
50% 3.000000 6.000000 24.500000 31.250000
75% 3.000000 8.000000 28.000000 39.687500
max 3.000000 10.000000 40.000000 69.550000
```

```
In [321]: Cluster4_test.describe()

Out[321]:
Pclass Fam Age Fare
count 14.0 14.000000 14.000000 14.000000
mean 1.0 1.500000 29.351429 165.804164
std 0.0 1.224745 16.176330 46.532749
min 1.0 0.000000 0.920000 110.883300
25% 1.0 0.250000 19.000000 133.650000
50% 1.0 1.500000 32.000000 151.550000
75% 1.0 2.750000 39.750000 211.960425
max 1.0 3.000000 50.000000 247.520800
```

CART

Predictor Variables: Embarked, Age_c, Sex, Fam, Fare, Fare2, Pclass, Ticket

Build Model


```
In [322]: #save target variable as y CART
y_dt = train_rebal[['Survived']]

#specify levels of target variable
y_dt_names = ["Did not survive", "Survived"]

In [323]: #combine all predictor variables, with dummies for categorical attributes
X_dt = pd.concat([sex_dummies,
                  embarked_dummies,
                  age_dummies,
                  pclass_dummies,
                  train_rebal['Fare'],
                  train_rebal['Fam']], axis=1)
X_dt
```

	Sex_female	Sex_male	Embarked_C	Embarked_Q	Embarked_S	Age_Baby/Toddler	Age_Child	Age_Adult	Age_Elderly	Pclass_1
65	0	1	1	0	0	0	0	1	0	0
445	0	1	0	0	0	1	0	1	0	0
659	0	1	1	0	0	0	0	0	1	0
591	1	0	1	0	0	0	0	0	1	0
690	0	1	0	0	1	0	0	0	1	0
...
275	1	0	0	0	1	0	0	0	1	0
839	0	1	1	0	0	0	0	0	1	0
577	1	0	0	0	1	0	0	0	1	0
400	0	1	0	0	1	0	0	0	1	0
383	1	0	0	0	1	0	0	0	1	0

746 rows x 14 columns

```
In [323]: #Build model
#DecisionTreeClassifier(criterion = "gini").fit(X_dt, y_dt)
export_graphviz(cart01, out_file = "cart01.dot", class_names=y_dt_names)

In [325]: #save predictions
pred_CART = cart01.predict(X_dt)
```

Test Dataset Validation

```
In [326]: #first dataset prep
#combine all predictor variables, with dummies for categorical attributes
X_dt_test = pd.concat([sex_dummies_test,
                      embarked_dummies_test,
                      age_dummies_test,
                      pclass_dummies_test,
                      test['Fare'],
                      test['Fam']], axis=1)
X_dt_test
```

	Sex_female	Sex_male	Embarked_C	Embarked_Q	Embarked_S	Age_Baby/Toddler	Age_Child	Age_Adult	Age_Elderly	Pclass_1
725	0	1	0	0	1	0	0	1	0	0
881	0	1	0	0	1	0	0	0	1	0
528	0	1	0	0	1	0	0	0	1	0
46	0	1	0	1	0	0	0	0	1	0
627	1	0	0	0	1	0	0	0	1	0
...
360	0	1	0	0	1	0	0	0	1	0
856	1	0	0	0	1	0	0	0	1	0
199	1	0	0	0	1	0	0	0	1	0
451	0	1	0	0	1	0	0	0	1	0
417	1	0	0	0	1	0	0	0	1	0

295 rows x 14 columns

```
In [327]: #use CART model to make predictions on test dataset
pred_CART_test = cart01.predict(X_dt_test)
```

Model Performance

```
In [328]: #Show contingency table for actual vs. predicted
ypredCART = pd.crosstab(test['Survived'], pred_CART_test,
                        rownames = ['Actual'],
                        colnames = ['Predicted'])
ypredCART['Total'] = ypredCART.sum(axis=1)
ypredCART.loc['Total'] = ypredCART.sum()
ypredCART
```

	Predicted 0	1	Total
Actual			
0	151	25	176
1	42	77	119
Total	193	102	295

Compared to baseline model: (1) An All Negative Model would have an accuracy rate of 60% and (2) An All Positive Model would have an accuracy rate of 40%.

Accuracy: (153+74)/295 = 77%

Error Rate: 100% - 77% = 23%

Sensitivity: TP/Total Positive = 74/119 = 62%

Specificity: TN/Total Negative = 153/176 = 87%

Random Forest

Build Model

```
In [329]: # re-format target variable as a 1D-array
rfy_dt = np.ravel(y_dt)

In [330]: #run Random Forest algorithm
rf01 = RandomForestClassifier(n_estimators = 100, criterion = "gini").fit(X_dt, rfy_dt)
```

Test Dataset Validation

```
In [331]: #use Random Forest model to make predictions on test dataset
pred_rf_test = rf01.predict(X_dt_test)
```

Model Performance

```
In [332]: #Show contingency table for actual vs. predicted
ypredRF = pd.crosstab(test['Survived'], pred_rf_test,
                      rownames = ['Actual'],
                      colnames = ['Predicted'])
ypredRF['Total'] = ypredRF.sum(axis=1)
ypredRF.loc['Total'] = ypredRF.sum()
ypredRF
```

	Predicted 0	1	Total
Actual			
0	147	29	176
1	43	76	119
Total	190	105	295

Accuracy: (148+75)/295 = 76%

Error Rate: 100% - 76% = 24%

Sensitivity: TP/Total Positive = 74/119 = 62%

Specificity: TN/Total Negative = 148/176 = 84%

Naive Bayes

Build Model

```
In [333]: ymb2 = train_rebal[['Survived']]
yyy = test[['Survived']]
ymb = np.ravel(ymb2)
ymb = pd.concat([sex_dummies,
                embarked_dummies,
                age_dummies,
                pclass_dummies,
                train_rebal['Fare'],
                train_rebal['Fam']], axis=1)
xnbtest = pd.concat([sex_dummies_test,
                    embarked_dummies_test,
                    age_dummies_test,
                    pclass_dummies_test,
                    test['Fare'],
                    test['Fam']], axis=1)
nb_01 = MultinomialNB().fit(xnb, ymb)
yprednb2 = nb_01.predict(xnbtest)
yprednb2 = pd.crosstab(test['Survived'], yprednb2,
                      rownames = ['Actual'],
                      colnames = ['Predicted'])
yprednb2['Total'] = yprednb2.sum(axis=1)
yprednb2.loc['Total'] = yprednb2.sum()
yprednb2
```

	Predicted 0	1	Total
Actual			
0	159	17	176
1	43	76	119
Total	202	93	295

Accuracy: 160+77/295 = 80.3%

Error Rate: 100% - 77% = 19.7%

Sensitivity: TP/Total Positive = 77/119 = 65%

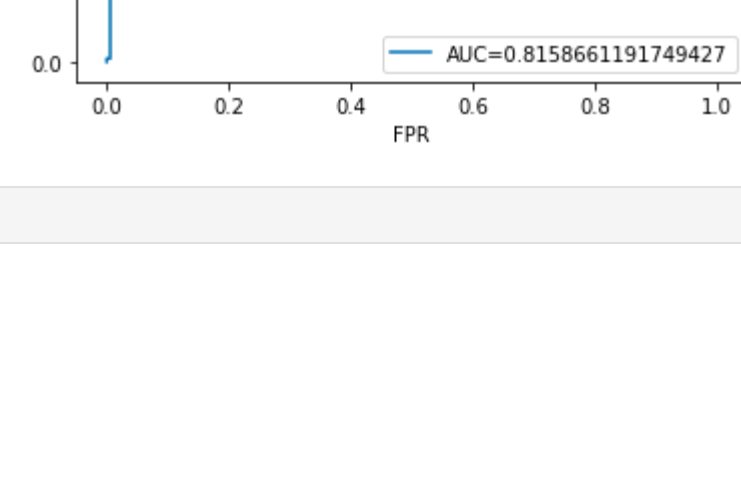
Specificity: TN/Total Negative = 160/176 = 91%

ROC

CART ROC

```
In [337]: print("Accuracy", metrics.accuracy_score(yyy, pred_CART_test))
Accuracy 0.7728813559322034

In [338]: y_pred_proba2 = DecisionTreeClassifier(criterion = "gini").fit(X_dt, y_dt).predict_proba(X_dt_test[:,1])
fpr1, tpr1, _ = metrics.roc_curve(yyy, y_pred_proba2)
auc1 = metrics.roc_auc_score(yyy, y_pred_proba2)
plt.plot(fpr1,tpr1,label="CART auc="+str(auc1))
plt.legend(loc=4)
plt.title('CART')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



RF ROC

```
In [339]: print("Accuracy", metrics.accuracy_score(yyy, pred_rf_test))
Accuracy 0.7599322033898305

In [340]: y_pred_proba2 = RandomForestClassifier(n_estimators = 100, criterion = "gini").fit(X_dt, rfy_dt).predict_proba(X_dt_test[:,1])
fpr2, tpr2, _ = metrics.roc_curve(yyy, y_pred_proba2)
auc2 = metrics.roc_auc_score(yyy, y_pred_proba2)
plt.plot(fpr2,tpr2,label="RF auc="+str(auc2))
plt.legend(loc=4)
plt.title('Random Forest')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



NB ROC

```
In [341]: print("Accuracy", metrics.accuracy_score(yyy, yprednb2))
Accuracy 0.7966101694915254

In [342]: y_pred_proba = MultinomialNB().fit(xnb, ymb).predict_proba(xnbtest[:,1])
fpr, tpr, _ = metrics.roc_curve(yyy, y_pred_proba)
auc = metrics.roc_auc_score(yyy, y_pred_proba)
plt.plot(fpr,tpr,label="NB auc="+str(auc))
plt.legend(loc=4)
plt.title('Naive Bayes')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



Logistic Regression ROC

```
In [344]: print("Accuracy", metrics.accuracy_score(yyy, ypred_logis))
Accuracy 0.7864406779661017

In [343]: #format ROC curve for Logistic Regression model
lg_pred_proba = logreg01.predict(X_logi_d_test)
fpr9, tpr9, _ = metrics.roc_curve(yyy, lg_pred_proba)
auc9 = metrics.roc_auc_score(yyy, lg_pred_proba)
plt.plot(fpr9,tpr9,label="AUC="+str(auc9))
plt.legend(loc=4)
plt.title('Logistic Regression')
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.show()
```



```
In [343]:
```