

In [1]:

```
import itertools
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.ticker import NullFormatter
import pandas as pd
import numpy as np
import matplotlib.ticker as ticker
from sklearn import preprocessing
%matplotlib inline
```

In [2]:

```
df = pd.read_csv('loan_train.csv')
df.head()
```

Out[2]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	educa
0	0	0	PAIDOFF	1000	30	9/8/2016	10/7/2016	45	Scho Be
1	2	2	PAIDOFF	1000	30	9/8/2016	10/7/2016	33	Bech
2	3	3	PAIDOFF	1000	15	9/8/2016	9/22/2016	27	coll
3	4	4	PAIDOFF	1000	30	9/9/2016	10/8/2016	28	coll
4	6	6	PAIDOFF	1000	30	9/9/2016	10/8/2016	29	coll

In [3]:

```
df.shape
```

Out[3]:

(346, 10)

In [4]:

```
df['due_date'] = pd.to_datetime(df['due_date'])
df['effective_date'] = pd.to_datetime(df['effective_date'])
df.head()
```

Out[4]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	educat
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	Schoc Be
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bech
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	coll
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	coll
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	coll



In [5]:

```
df['loan_status'].value_counts()
```

Out[5]:

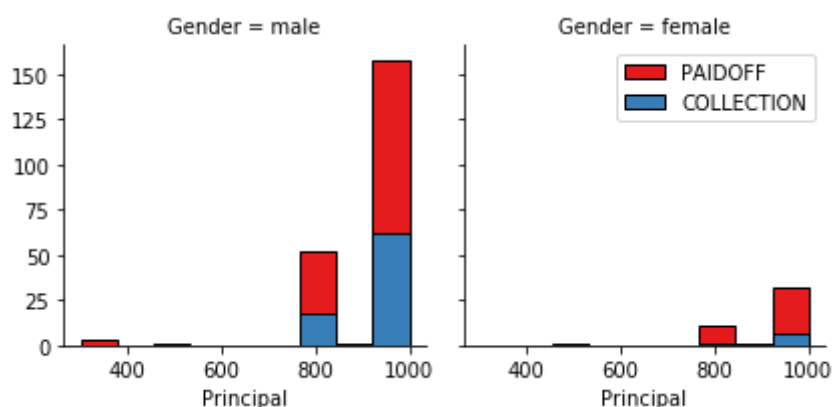
PAIDOFF 260  
COLLECTION 86  
Name: loan\_status, dtype: int64

In [6]:

```
import seaborn as sns
```

```
bins = np.linspace(df.Principal.min(), df.Principal.max(), 10)  
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)  
g.map(plt.hist, 'Principal', bins=bins, ec="k")
```

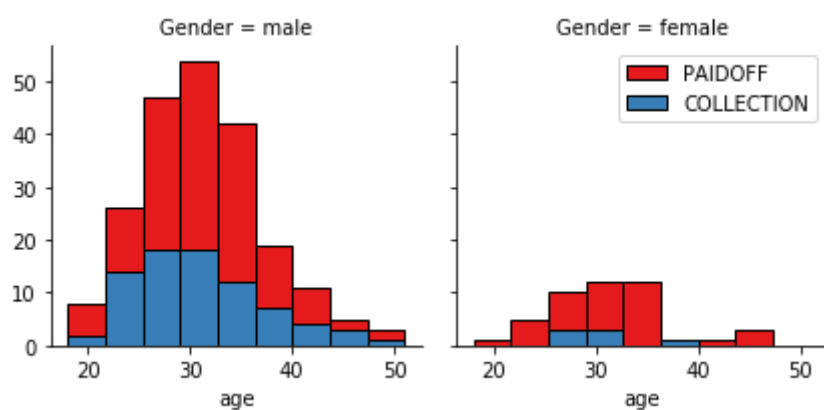
```
g.axes[-1].legend()  
plt.show()
```



In [7]:

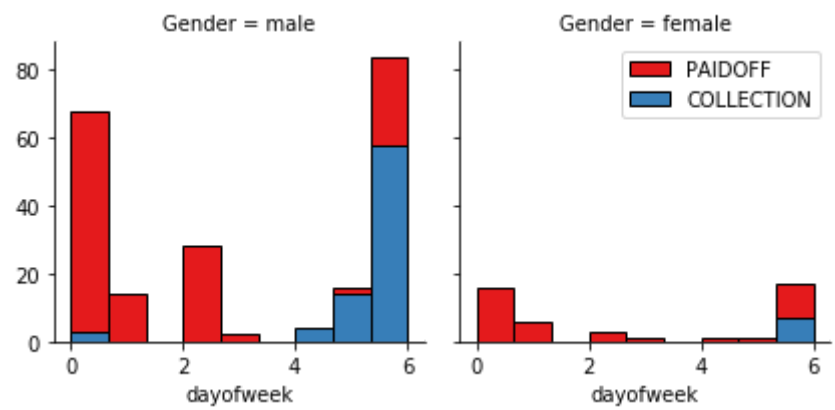
```
bins = np.linspace(df.age.min(), df.age.max(), 10)  
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)  
g.map(plt.hist, 'age', bins=bins, ec="k")
```

```
g.axes[-1].legend()  
plt.show()
```



In [8]:

```
df['dayofweek'] = df['effective_date'].dt.dayofweek
bins = np.linspace(df.dayofweek.min(), df.dayofweek.max(), 10)
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)
g.map(plt.hist, 'dayofweek', bins=bins, ec="k")
g.axes[-1].legend()
plt.show()
```



In [9]:

```
df['weekend'] = df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
df.head()
```

Out[9]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	educat
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	Schoc Be
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bech
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	coll
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	coll
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	coll

In [10]:

```
df.groupby(['Gender'])['loan_status'].value_counts(normalize=True)
```

Out[10]:

```
Gender  loan_status
female  PAIDOFF      0.865385
        COLLECTION   0.134615
male    PAIDOFF      0.731293
        COLLECTION   0.268707
Name: loan_status, dtype: float64
```

In [11]:

```
df['Gender'].replace(to_replace=['male','female'], value=[0,1],inplace=True)
df.head()
```

Out[11]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechalor
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	coll
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	coll
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	coll

In [12]:

```
df.groupby(['education'])['loan_status'].value_counts(normalize=True)
```

Out[12]:

```
education  loan_status
Bechalor   PAIDOFF      0.750000
           COLLECTION   0.250000
High School or Below  PAIDOFF      0.741722
                    COLLECTION   0.258278
Master or Above      COLLECTION   0.500000
                    PAIDOFF      0.500000
college          PAIDOFF      0.765101
                COLLECTION   0.234899
Name: loan_status, dtype: float64
```

In [13]:

```
df[['Principal', 'terms', 'age', 'Gender', 'education']].head()
```

Out[13]:

	Principal	terms	age	Gender	education
0	1000	30	45	0	High School or Below
1	1000	30	33	1	Bechalar
2	1000	15	27	0	college
3	1000	30	28	1	college
4	1000	30	29	0	college

In [14]:

```
Feature = df[['Principal', 'terms', 'age', 'Gender', 'weekend']]
Feature = pd.concat([Feature, pd.get_dummies(df['education'])], axis=1)
Feature.drop(['Master or Above'], axis = 1, inplace=True)
Feature.head()
```

Out[14]:

	Principal	terms	age	Gender	weekend	Bechalar	High School or Below	college
0	1000	30	45	0	0	0	1	0
1	1000	30	33	1	0	1	0	0
2	1000	15	27	0	0	0	0	1
3	1000	30	28	1	1	0	0	1
4	1000	30	29	0	1	0	0	1

In [15]:

```
X = Feature
X[0:5]
```

Out[15]:

	Principal	terms	age	Gender	weekend	Bechalar	High School or Below	college
0	1000	30	45	0	0	0	1	0
1	1000	30	33	1	0	1	0	0
2	1000	15	27	0	0	0	0	1
3	1000	30	28	1	1	0	0	1
4	1000	30	29	0	1	0	0	1

In [16]:

```
y = df['loan_status'].values  
y[0:5]
```

Out[16]:

```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],  
      dtype=object)
```

In [17]:

```
import warnings  
warnings.filterwarnings('ignore')
```

In [18]:

```
X= preprocessing.StandardScaler().fit(X).transform(X)  
X[0:5]
```

Out[18]:

```
array([[ 0.51578458,  0.92071769,  2.33152555, -0.42056004, -1.20577805,  
        -0.38170062,  1.13639374, -0.86968108],  
       [ 0.51578458,  0.92071769,  0.34170148,  2.37778177, -1.20577805,  
        2.61985426, -0.87997669, -0.86968108],  
       [ 0.51578458, -0.95911111, -0.65321055, -0.42056004, -1.20577805,  
        -0.38170062, -0.87997669,  1.14984679],  
       [ 0.51578458,  0.92071769, -0.48739188,  2.37778177,  0.82934003,  
        -0.38170062, -0.87997669,  1.14984679],  
       [ 0.51578458,  0.92071769, -0.3215732 , -0.42056004,  0.82934003,  
        -0.38170062, -0.87997669,  1.14984679]])
```

## Classification

In [19]:

```
# We split the X into train and test to find the best k  
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=4  
)  
print ('Train set:', X_train.shape, y_train.shape)  
print ('Test set:', X_test.shape, y_test.shape)
```

```
('Train set:', (276L, 8L), (276L,))  
('Test set:', (70L, 8L), (70L,))
```

## K Nearest Neighbor(KNN)

In [20]:

```
from sklearn.neighbors import KNeighborsClassifier
k = 3
#Train Model and Predict
kNN_model = KNeighborsClassifier(n_neighbors=k).fit(X_train,y_train)
kNN_model
```

Out[20]:

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                    metric_params=None, n_jobs=None, n_neighbors=3, p=2,
                    weights='uniform')
```

In [21]:

```
yhat = kNN_model.predict(X_test)
yhat[0:5]
```

Out[21]:

```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
      dtype=object)
```

In [22]:

```
Ks=15
mean_acc=np.zeros((Ks-1))
std_acc=np.zeros((Ks-1))
ConfusionMx=[];
for n in range(1,Ks):

    #Train Model and Predict
    kNN_model = KNeighborsClassifier(n_neighbors=n).fit(X_train,y_train)
    yhat = kNN_model.predict(X_test)

    mean_acc[n-1]=np.mean(yhat==y_test);

    std_acc[n-1]=np.std(yhat==y_test)/np.sqrt(yhat.shape[0])
mean_acc
```

Out[22]:

```
array([0.67142857, 0.65714286, 0.71428571, 0.68571429, 0.75714286,
       0.71428571, 0.78571429, 0.75714286, 0.75714286, 0.67142857,
       0.7          , 0.72857143, 0.7          , 0.7          ])
```



In [23]:

```
# Building the model again, using k=7
from sklearn.neighbors import KNeighborsClassifier
k = 7
#Train Model and Predict
kNN_model = KNeighborsClassifier(n_neighbors=k).fit(X_train,y_train)
kNN_model
```

Out[23]:

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                     metric_params=None, n_jobs=None, n_neighbors=7, p=2,
                     weights='uniform')
```

## Decision Tree

In [24]:

```
from sklearn.tree import DecisionTreeClassifier
DT_model = DecisionTreeClassifier(criterion="entropy", max_depth = 4)
DT_model.fit(X_train,y_train)
DT_model
```

Out[24]:

```
DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=
4,
                       max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, presort=False, random_state=Non
e,
                       splitter='best')
```

In [25]:

```
yhat = DT_model.predict(X_test)
yhat
```

Out[25]:

```
array(['COLLECTION', 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COLLECTION',
      'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'COLLECTION', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'COLLECTION', 'COLLECTION', 'PAIDOFF',
      'COLLECTION', 'COLLECTION', 'PAIDOFF', 'COLLECTION', 'PAIDOFF',
      'COLLECTION', 'COLLECTION', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'COLLECTION', 'PAIDOFF',
      'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'COLLECTION', 'PAIDOFF',
      'COLLECTION', 'COLLECTION', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'COLLECTION',
      'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF'], dtype=object)
```

# Support Vector Machine

In [26]:

```
from sklearn import svm
SVM_model = svm.SVC()
SVM_model.fit(X_train, y_train)
```

Out[26]:

```
SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False)
```

In [27]:

```
yhat = SVM_model.predict(X_test)
yhat
```

Out[27]:

```
array(['COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'COLLECTION', 'PAIDOFF', 'COLLECTION',
      'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COLLECTION',
      'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
      dtype=object)
```

In [ ]:

In [ ]:

## Logistic Regression

In [28]:

```
from sklearn.linear_model import LogisticRegression
LR_model = LogisticRegression(C=0.01).fit(X_train,y_train)
LR_model
```

Out[28]:

```
LogisticRegression(C=0.01, 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]:

```
yhat = LR_model.predict(X_test)
yhat
```

Out[29]:

```
array(['COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'COLLECTION', 'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'COLLECTION',
       'COLLECTION', 'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COLLECTION',
       'PAIDOFF', 'PAIDOFF', 'COLLECTION', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF'], dtype=object)
```

In [ ]:

In [ ]:

## Model Evaluation using Test set

In [30]:

```
from sklearn.metrics import jaccard_similarity_score
from sklearn.metrics import f1_score
from sklearn.metrics import log_loss
```

In [31]:

```
test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

Out[31]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education
0	1	1	PAIDOFF	1000	30	9/8/2016	10/7/2016	50	Bech
1	5	5	PAIDOFF	300	7	9/9/2016	9/15/2016	35	Maste Ab
2	21	21	PAIDOFF	1000	30	9/10/2016	10/9/2016	43	Scho Be
3	24	24	PAIDOFF	1000	30	9/10/2016	10/9/2016	26	coll
4	35	35	PAIDOFF	800	15	9/11/2016	9/25/2016	29	Bech

In [32]:

```
test_df['due_date'] = pd.to_datetime(test_df['due_date'])
test_df['effective_date'] = pd.to_datetime(test_df['effective_date'])
test_df['dayofweek'] = test_df['effective_date'].dt.dayofweek
test_df['weekend'] = test_df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
test_df['Gender'].replace(to_replace=['male','female'], value=[0,1],inplace=True)
test_Feature = test_df[['Principal','terms','age','Gender','weekend']]
test_Feature = pd.concat([test_Feature,pd.get_dummies(test_df['education'])], axis=1)
test_Feature.drop(['Master or Above'], axis = 1,inplace=True)
test_X = preprocessing.StandardScaler().fit(test_Feature).transform(test_Feature)
test_X[0:5]
```

Out[32]:

```
array([[ 0.49362588,  0.92844966,  3.05981865,  1.97714211, -1.30384048,
         2.39791576, -0.79772404, -0.86135677],
       [-3.56269116, -1.70427745,  0.53336288, -0.50578054,  0.76696499,
        -0.41702883, -0.79772404, -0.86135677],
       [ 0.49362588,  0.92844966,  1.88080596,  1.97714211,  0.76696499,
        -0.41702883,  1.25356634, -0.86135677],
       [ 0.49362588,  0.92844966, -0.98251057, -0.50578054,  0.76696499,
        -0.41702883, -0.79772404,  1.16095912],
       [-0.66532184, -0.78854628, -0.47721942, -0.50578054,  0.76696499,
         2.39791576, -0.79772404, -0.86135677]])
```

In [33]:

```
test_y = test_df['loan_status'].values
test_y[0:5]
```

Out[33]:

```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
      dtype=object)
```

In [34]:

```
knn_yhat = knn_model.predict(test_X)
print("KNN Jaccard index: %.2f" % jaccard_similarity_score(test_y, knn_yhat))
print("KNN F1-score: %.2f" % f1_score(test_y, knn_yhat, average='weighted') )
```

KNN Jaccard index: 0.67

KNN F1-score: 0.63

In [35]:

```
DT_yhat = DT_model.predict(test_X)
print("DT Jaccard index: %.2f" % jaccard_similarity_score(test_y, DT_yhat))
print("DT F1-score: %.2f" % f1_score(test_y, DT_yhat, average='weighted') )
```

DT Jaccard index: 0.72

DT F1-score: 0.74

In [36]:

```
SVM_yhat = SVM_model.predict(test_X)
print("SVM Jaccard index: %.2f" % jaccard_similarity_score(test_y, SVM_yhat))
print("SVM F1-score: %.2f" % f1_score(test_y, SVM_yhat, average='weighted') )
```

SVM Jaccard index: 0.80

SVM F1-score: 0.76

In [37]:

```
LR_yhat = LR_model.predict(test_X)
LR_yhat_prob = LR_model.predict_proba(test_X)
print("LR Jaccard index: %.2f" % jaccard_similarity_score(test_y, LR_yhat))
print("LR F1-score: %.2f" % f1_score(test_y, LR_yhat, average='weighted') )
print("LR LogLoss: %.2f" % log_loss(test_y, LR_yhat_prob))
```

LR Jaccard index: 0.74

LR F1-score: 0.66

LR LogLoss: 0.57

In [ ]:

Report

	Algorithm	Jaccard	F1-score	LogL
oss				
	KNN	0.67	0.63	N
A				
	Decision Tree	0.72	0.74	N
A				
	SVM	0.80	0.76	N
A				
	LogisticRegression	0.74	0.66	0.
57				