# Classification ML project ¶

```
In [2]:  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
```

## About dataset

This dataset is about past loans. The **Loan_train.csv** data set includes details of 346 customers whose loan are already paid off or defaulted. It includes following fields:

| Field | Description |
| --- | --- |
| Loan_status | Whether a loan is paid off on in collection |
| Principal | Basic principal loan amount at the |
| Terms | Origination terms which can be weekly (7 days), biweekly, and monthly payoff schedule |
| Effective_date | When the loan got originated and took effects |
| Due_date | Since it's one-time payoff schedule, each loan has one single due date |
| Age | Age of applicant |
| Education | Education of applicant |
| Gender | The gender of applicant |

## Load Data From CSV File

In [3]:
```python
df = pd.read_csv('loan_train.csv')
df.head()
```

Out[3]:

| | Unnamed: 0 | Unnamed: 0.1 | loan_status | Principal | terms | effective_date | due_date | age | edu |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | PAIDOFF | 1000 | 30 | 9/8/2016 | 10/7/2016 | 45 | Sch |
| 1 | 2 | 2 | PAIDOFF | 1000 | 30 | 9/8/2016 | 10/7/2016 | 33 | Be |
| 2 | 3 | 3 | PAIDOFF | 1000 | 15 | 9/8/2016 | 9/22/2016 | 27 | ( |
| 3 | 4 | 4 | PAIDOFF | 1000 | 30 | 9/9/2016 | 10/8/2016 | 28 | ( |
| 4 | 6 | 6 | PAIDOFF | 1000 | 30 | 9/9/2016 | 10/8/2016 | 29 | ( |

In [4]:
```python
df.shape
```

Out[4]: (346, 10)

## Convert to date time object

In [5]:
```python
df['due_date'] = pd.to_datetime(df['due_date'])
df['effective_date'] = pd.to_datetime(df['effective_date'])
df.head()
```

Out[5]:

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

# Data visualization and pre-processing
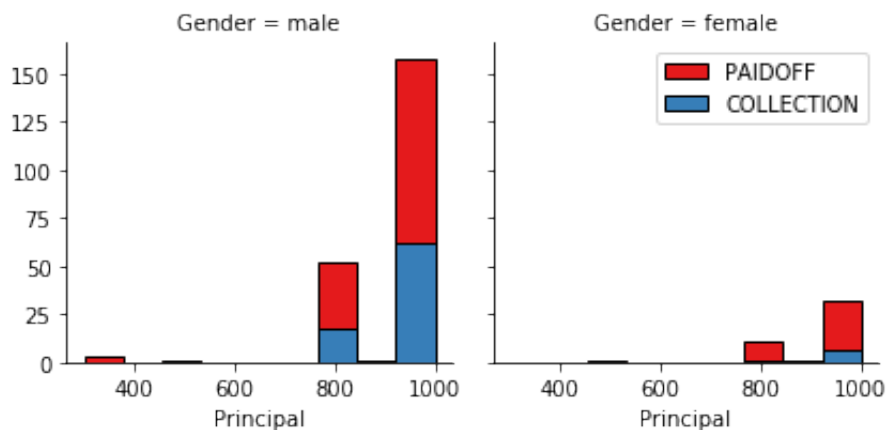
```
In [6]: df['loan_status'].value_counts()
```

```
Out[6]: PAIDOFF         260
        COLLECTION       86
        Name: loan_status, dtype: int64
```

260 people have paid off the loan on time while 86 have gone into collection
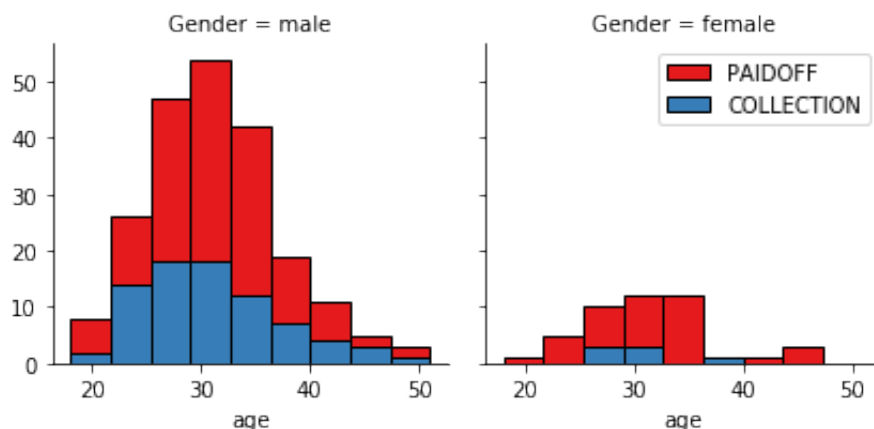
Lets plot some columns to underestand data better:

```
In [8]: import seaborn as sns

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

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

In [9]:
```python
bins=np.linspace(df.age.min(), df.age.max(), 10)
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set
1", col_wrap=2)
g.map(plt.hist, 'age', bins=bins, ec="k")

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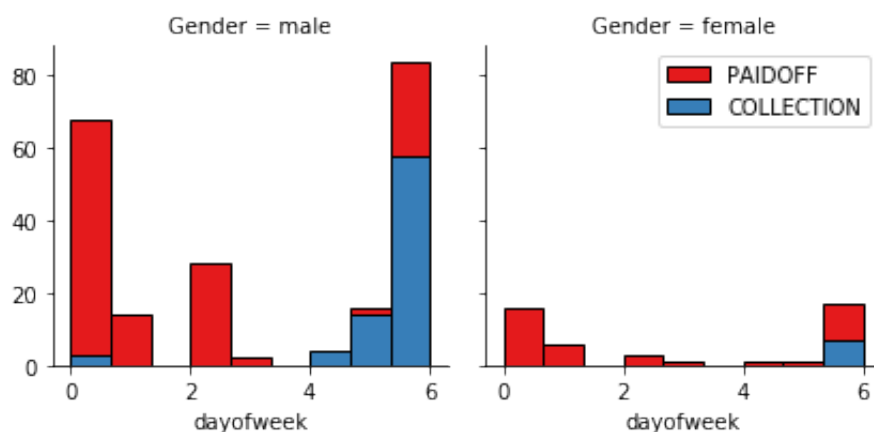


# Pre-processing: Feature selection/extraction

## Lets look at the day of the week people get the loan

In [10]:
```python
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="Set
1", col_wrap=2)
g.map(plt.hist, 'dayofweek', bins=bins, ec="k")
g.axes[-1].legend()
plt.show()
```

We see that people who get the loan at the end of the week dont pay it off, so lets use Feature binarization to set a threshold values less then day 4

```
In [11]:  df['weekend']= df['dayofweek'].apply(lambda x: 1 if (x>3)  else 0)
          df.head()
```

Out[11]:

| | Unnamed: 0 | Unnamed: 0.1 | loan_status | Principal | terms | effective_date | due_date | age | edu |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 0 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 45 | Sch |
| **1** | 2 | 2 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 33 | Be |
| **2** | 3 | 3 | PAIDOFF | 1000 | 15 | 2016-09-08 | 2016-09-22 | 27 | c |
| **3** | 4 | 4 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 28 | c |
| **4** | 6 | 6 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 29 | c |

# Convert Categorical features to numerical values

Lets look at gender:

```
In [12]:  df.groupby(['Gender'])['loan_status'].value_counts(normalize=True)
```

```
Out[12]:  Gender   loan_status
          female   PAIDOFF        0.865385
                   COLLECTION     0.134615
          male     PAIDOFF        0.731293
                   COLLECTION     0.268707
          Name: loan_status, dtype: float64
```

86 % of female pay there loans while only 73 % of males pay there loan

Lets convert male to 0 and female to 1:

In [13]:
```python
df['Gender'].replace(to_replace=['male','female'], value=[0,1],inpl
ace=True)
df.head()
```

Out[13]:

| | Unnamed: 0 | Unnamed: 0.1 | loan_status | Principal | terms | effective_date | due_date | age | educ |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 0 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 45 | Sch |
| **1** | 2 | 2 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 33 | Be |
| **2** | 3 | 3 | PAIDOFF | 1000 | 15 | 2016-09-08 | 2016-09-22 | 27 | c |
| **3** | 4 | 4 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 28 | c |
| **4** | 6 | 6 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 29 | c |

# One Hot Encoding

**How about education?**

In [14]:
```python
df.groupby(['education'])['loan_status'].value_counts(normalize=Tru
e)
```

Out[14]:
```
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
```

**Feature befor One Hot Encoding**

In [15]:
```python
df[['Principal','terms','age','Gender','education']].head()
```

Out[15]:

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

**Use one hot encoding technique to conver categorical varables to binary variables and append them to the feature Data Frame**

In [16]:
```python
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[16]:

|   | Principal | terms | age | Gender | weekend | Bechalor | 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 |

# Feature selection

Lets defind feature sets, X:

```
In [17]:  X = Feature
          X[0:5]
```

Out[17]:

| | Principal | terms | age | Gender | weekend | Bechalor | 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 |

What are our lables?

```
In [18]:  y = df['loan_status'].values
          y[0:5]
```

Out[18]:  array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'], dty
          pe=object)

# Normalize Data

Data Standardization give data zero mean and unit variance (technically should be done after train test split )

```
In [19]:  X = preprocessing.StandardScaler().fit(X).transform(X)
          X[0:5]
```

Out[19]:  array([[ 0.51578458,  0.92071769,  2.33152555, -0.42056004, -1.205
          77805,
                  -0.38170062,  1.13639374, -0.86968108],
                 [ 0.51578458,  0.92071769,  0.34170148,  2.37778177, -1.205
          77805,
                   2.61985426, -0.87997669, -0.86968108],
                 [ 0.51578458, -0.95911111, -0.65321055, -0.42056004, -1.205
          77805,
                  -0.38170062, -0.87997669,  1.14984679],
                 [ 0.51578458,  0.92071769, -0.48739188,  2.37778177,  0.829
          34003,
                  -0.38170062, -0.87997669,  1.14984679],
                 [ 0.51578458,  0.92071769, -0.3215732 , -0.42056004,  0.829
          34003,
                  -0.38170062, -0.87997669,  1.14984679]])
```

# Classification Modeling

# K Nearest Neighbor(KNN)

In [21]:
```python
# 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: (276, 8) (276,)
Test set: (70, 8) (70,)
```

In [45]:
```python
# Modeling
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[45]:
```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minko
wski',
               metric_params=None, n_jobs=1, n_neighbors=3, p=2,
               weights='uniform')
```

In [46]:
```python
# just for sanity chaeck
yhat = kNN_model.predict(X_test)
yhat[0:5]
```

Out[46]:
```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'], dty
pe=object)
```

In [67]:
```python
# Best k
Ks=15
mean_acc=np.zeros((Ks-1))
std_acc=np.zeros((Ks-1))
ConfustionMx=[];
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[67]:
```
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 [68]:
```python
# 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[68]:
```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
           metric_params=None, n_jobs=1, n_neighbors=7, p=2,
           weights='uniform')
```

# Decision Tree

In [84]:
```python
from sklearn.tree import DecisionTreeClassifier
DT_model = DecisionTreeClassifier(criterion="entropy", max_depth = 4)
DT_model.fit(X_train,y_train)
DT_model
```

Out[84]:
```
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=None,
            splitter='best')
```

In [85]:
```python
yhat = DT_model.predict(X_test)
yhat
```

Out[85]:
```
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 [77]:
```python
from sklearn import svm
SVM_model = svm.SVC()
SVM_model.fit(X_train, y_train)
```

Out[77]:
```
SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
  decision_function_shape='ovr', degree=3, gamma='auto', kernel='r
bf',
  max_iter=-1, probability=False, random_state=None, shrinking=Tru
e,
  tol=0.001, verbose=False)
```

In [82]:
```python
yhat = SVM_model.predict(X_test)
yhat
```

Out[82]:
```
array(['COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', '
PAIDOFF',
       'COLLECTION', 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'
,
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COL
LECTION',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COL
LECTION',
       'COLLECTION', 'PAIDOFF', 'COLLECTION', 'COLLECTION', 'PAIDO
FF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
DOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COLLECTION', '
PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COLLECTION', 'PAIDOFF', '
PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
DOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
DOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'COL
LECTION',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
DOFF'], dtype=object)
```

# Logistic Regression

```
In [80]: from sklearn.linear_model import LogisticRegression
         LR_model = LogisticRegression(C=0.01).fit(X_train,y_train)
         LR_model
```

```
Out[80]: LogisticRegression(C=0.01, class_weight=None, dual=False, fit_inte
         rcept=True,
                   intercept_scaling=1, max_iter=100, multi_class='ovr', n_
         jobs=1,
                   penalty='l2', random_state=None, solver='liblinear', tol
         =0.0001,
                   verbose=0, warm_start=False)
```

```
In [81]: yhat = LR_model.predict(X_test)
         yhat
```

```
Out[81]: array(['COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', '
         PAIDOFF',
                 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
         DOFF',
                 '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', 'PAI
         DOFF',
                 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
         DOFF',
                 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAI
         DOFF',
                 'COLLECTION', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', '
         PAIDOFF',
                 'PAIDOFF'], dtype=object)
```

# Model Evaluation using Test set

```
In [93]: from sklearn.metrics import jaccard_similarity_score
         from sklearn.metrics import f1_score
         from sklearn.metrics import log_loss
```

First, download and load the test set:

## Load Test set for evaluation

```
In [100]:   test_df = pd.read_csv('loan_test.csv')
            test_df.head()
```

Out[100]:

| | Unnamed: 0 | Unnamed: 0.1 | loan_status | Principal | terms | effective_date | due_date | age | edu |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 1 | PAIDOFF | 1000 | 30 | 9/8/2016 | 10/7/2016 | 50 | Be |
| **1** | 5 | 5 | PAIDOFF | 300 | 7 | 9/9/2016 | 9/15/2016 | 35 | Ma |
| **2** | 21 | 21 | PAIDOFF | 1000 | 30 | 9/10/2016 | 10/9/2016 | 43 | Sc |
| **3** | 24 | 24 | PAIDOFF | 1000 | 30 | 9/10/2016 | 10/9/2016 | 26 | ( |
| **4** | 35 | 35 | PAIDOFF | 800 | 15 | 9/11/2016 | 9/25/2016 | 29 | Be |

In [101]:
```python
## Preprocessing
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[101]:
```
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 [102]:
```python
test_y = test_df['loan_status'].values
test_y[0:5]
```

Out[102]:
```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'], dtype=object)
```

In [103]:
```python
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 [104]:
```python
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 [105]:
```python
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 [106]:
```python
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
```

# Report

| Algorithm | Jaccard | F1-score | LogLoss |
|---|---|---|---|
| KNN | 0.67 | 0.63 | NA |
| Decision Tree | 0.72 | 0.74 | NA |
| SVM | 0.80 | 0.76 | NA |
| LogisticRegression | 0.74 | 0.66 | 0.57 |