[(https://www.bigdatauniversity.com)](https://www.bigdatauniversity.com)

# Classification with Python

In this notebook we try to practice all the classification algorithms that we learned in this course.

We load a dataset using Pandas library, and apply the following algorithms, and find the best one for this specific dataset by accuracy evaluation methods.

Lets first load required libraries:

```
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
```

## 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 |

Lets download the dataset

```
In [2]: !wget -O loan_train.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-course
        s-data/CognitiveClass/ML0101ENv3/labs/loan_train.csv
```

```
--2020-05-10 12:27:25--  https://s3-api.us-geo.objectstorage.softlayer.net/cf-co
urses-data/CognitiveClass/ML0101ENv3/labs/loan_train.csv
Resolving s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorag
e.softlayer.net)... 67.228.254.196
Connecting to s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectsto
rage.softlayer.net)|67.228.254.196|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 23101 (23K) [text/csv]
Saving to: 'loan_train.csv'

100%[=====================================>] 23,101      --.-K/s   in 0.09s

2020-05-10 12:27:25 (263 KB/s) - 'loan_train.csv' saved [23101/23101]
```

## Load Data From CSV File

```
In [3]: 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 | education | Gender |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | PAIDOFF | 1000 | 30 | 9/8/2016 | 10/7/2016 | 45 | High School or Below | male |
| 1 | 2 | 2 | PAIDOFF | 1000 | 30 | 9/8/2016 | 10/7/2016 | 33 | Bechalor | female |
| 2 | 3 | 3 | PAIDOFF | 1000 | 15 | 9/8/2016 | 9/22/2016 | 27 | college | male |
| 3 | 4 | 4 | PAIDOFF | 1000 | 30 | 9/9/2016 | 10/8/2016 | 28 | college | female |
| 4 | 6 | 6 | PAIDOFF | 1000 | 30 | 9/9/2016 | 10/8/2016 | 29 | college | male |

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

Out[4]: (346, 10)

## Convert to date time object

```
In [5]: 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 | education | Gender |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 45 | High School or Below | male |
| 1 | 2 | 2 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 33 | Bechalor | female |
| 2 | 3 | 3 | PAIDOFF | 1000 | 15 | 2016-09-08 | 2016-09-22 | 27 | college | male |
| 3 | 4 | 4 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 28 | college | female |
| 4 | 6 | 6 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 29 | college | male |

# Data visualization and pre-processing

Let's see how many of each class is in our data set

```
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 [7]: # notice: installing seaborn might takes a few minutes
        !conda install -c anaconda seaborn -y
```

```
Solving environment: done

## Package Plan ##

  environment location: /opt/conda/envs/Python36

  added / updated specs:
    - seaborn


The following packages will be downloaded:

    package                    |            build
    ---------------------------|-----------------
    openssl-1.1.1g             |       h7b6447c_0         3.8 MB  anaconda
    seaborn-0.10.1             |             py_0         160 KB  anaconda
    certifi-2020.4.5.1         |           py36_0         159 KB  anaconda
    ca-certificates-2020.1.1   |                0         132 KB  anaconda
    ------------------------------------------------------------
                                           Total:         4.2 MB

The following packages will be UPDATED:

    ca-certificates: 2020.1.1-0          --> 2020.1.1-0          anaconda
    certifi:         2020.4.5.1-py36_0   --> 2020.4.5.1-py36_0   anaconda
    openssl:         1.1.1g-h7b6447c_0   --> 1.1.1g-h7b6447c_0   anaconda
    seaborn:         0.9.0-pyh91ea838_1  --> 0.10.1-py_0         anaconda


Downloading and Extracting Packages
openssl-1.1.1g       | 3.8 MB    | #################################### | 100%
seaborn-0.10.1       | 160 KB    | #################################### | 100%
certifi-2020.4.5.1   | 159 KB    | #################################### | 100%
ca-certificates-2020 | 132 KB    | #################################### | 100%
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
```
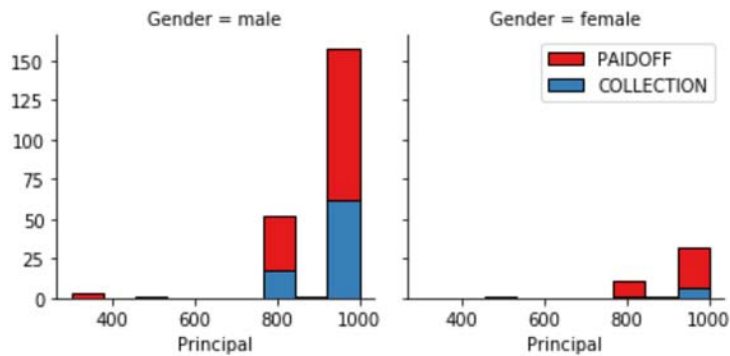
In [8]:
```python
import seaborn as sns

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()
```
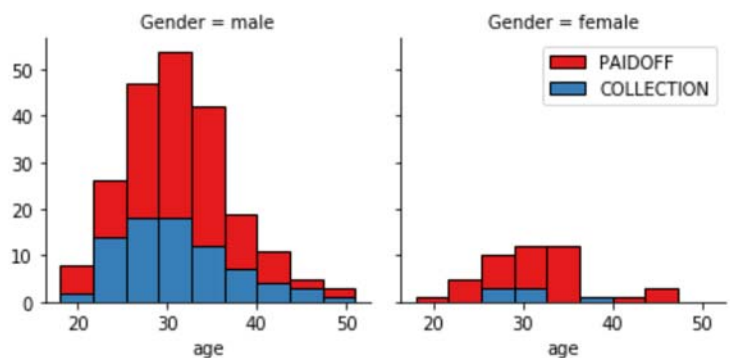


In [9]:
```python
bins = np.linspace(df.age.min(), df.age.max(), 10)
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)
g.map(plt.hist, 'age', bins=bins, ec="k")

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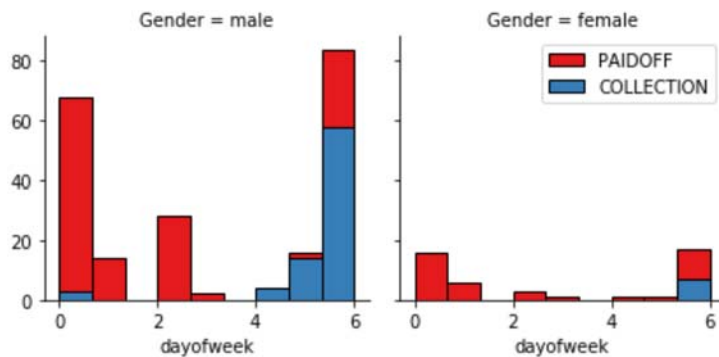


# Pre-processing: Feature selection/extraction

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

```
In [10]:  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()
```



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 | education | Gender | da |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 0 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 45 | High School or Below | male | |
| **1** | 2 | 2 | PAIDOFF | 1000 | 30 | 2016-09-08 | 2016-10-07 | 33 | Bechalor | female | |
| **2** | 3 | 3 | PAIDOFF | 1000 | 15 | 2016-09-08 | 2016-09-22 | 27 | college | male | |
| **3** | 4 | 4 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 28 | college | female | |
| **4** | 6 | 6 | PAIDOFF | 1000 | 30 | 2016-09-09 | 2016-10-08 | 29 | college | male | |

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

Out[13]:

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

# One Hot Encoding

### How about education?

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

```
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]: 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]:
```
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 [88]:
```
X = Feature
X[0:5]
```

Out[88]:

|   | 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'],
      dtype=object)
```

In [125]:
```
df['loan_status_cat']=df['loan_status'].replace(to_replace=['COLLECTION','PAIDOFF
'], value=[0,1])
df.head()
```

Out[125]:

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

In [128]:
```
y2=df['loan_status_cat']
```

## 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]
```

```
/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/preprocessing/data.
py:645: DataConversionWarning: Data with input dtype uint8, int64 were all conve
rted to float64 by StandardScaler.
  return self.partial_fit(X, y)
/opt/conda/envs/Python36/lib/python3.6/site-packages/ipykernel/__main__.py:1: Da
taConversionWarning: Data with input dtype uint8, int64 were all converted to fl
oat64 by StandardScaler.
  if __name__ == '__main__':
```

Out[19]:
```
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

Now, it is your turn, use the training set to build an accurate model. Then use the test set to report the accuracy of the model
You should use the following algorithm:

- K Nearest Neighbor(KNN)
- Decision Tree
- Support Vector Machine
- Logistic Regression

**Notice:**

- You can go above and change the pre-processing, feature selection, feature-extraction, and so on, to make a better model.
- You should use either scikit-learn, Scipy or Numpy libraries for developing the classification algorithms.
- You should include the code of the algorithm in the following cells.

# K Nearest Neighbor(KNN)

Notice: You should find the best k to build the model with the best accuracy.
**warning:** You should not use the **loan_test.csv** for finding the best k, however, you can split your train_loan.csv into train and test to find the best **k**.

```
In [20]:  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_st
          ate=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 [28]:  from sklearn.neighbors import KNeighborsClassifier
          from sklearn import metrics

          Ks = 10
          mean_acc = np.zeros((Ks-1))
          std_acc = np.zeros((Ks-1))
          ConfustionMx = [];
          for n in range(1,Ks):
              neigh = KNeighborsClassifier(n_neighbors = n).fit(X_train,y_train)
              yhat=neigh.predict(X_test)
              mean_acc[n-1] = metrics.accuracy_score(y_test, yhat)

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

          mean_acc
```

```
Out[28]:  array([0.67142857, 0.65714286, 0.71428571, 0.68571429, 0.75714286,
                 0.71428571, 0.78571429, 0.75714286, 0.75714286])
```

```
In [29]:  plt.plot(range(1,Ks),mean_acc,'g')
          plt.fill_between(range(1,Ks),mean_acc - 1 * std_acc,mean_acc + 1 * std_acc, alpha=
          0.10)
          plt.legend(('Accuracy ', '+/- 3xstd'))
          plt.ylabel('Accuracy ')
          plt.xlabel('Number of Nabors (K)')
          plt.tight_layout()
          plt.show()
```



```
In [26]:  print( "The best accuracy is with", mean_acc.max(), "with k=", mean_acc.argmax()+1)

          The best accuracy is with 0.7857142857142857 with k= 7
```

```
In [155]:  neigh = KNeighborsClassifier(n_neighbors = 7).fit(X_train,y_train)
           neigh
```

```
Out[155]:  KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                        metric_params=None, n_jobs=None, n_neighbors=7, p=2,
                        weights='uniform')
```

# Decision Tree

```
In [156]:  from sklearn.tree import DecisionTreeClassifier
```

```
In [157]:  Tree = DecisionTreeClassifier(criterion="entropy", max_depth = 4)
           Tree.fit(X_train, y_train)
```

```
Out[157]:  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 [158]:
```python
yhat_tree=Tree.predict(X_test)
yhat_tree
```

Out[158]:
```
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)
```

In [71]:
```python
y_test[0:5]
```

Out[71]:
```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
      dtype=object)
```

In [91]:
```python
print("DecisionTrees's Accuracy: ", metrics.accuracy_score(y_testtree, predTree))
```

```
DecisionTrees's Accuracy:  0.7403846153846154
```

In [179]:
```python
!pip install graphviz
!pip install pydotplus
import graphviz
import pydotplus

dot_data = StringIO()
filename = "tree.png"
featureNames = Feature.columns
out=tree.export_graphviz(Tree,feature_names=featureNames, out_file=dot_data, class
_names= np.unique(y_train), filled=True,  special_characters=True,rotate=False)
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png(filename)
img = mpimg.imread(filename)
plt.figure(figsize=(100, 200))
plt.imshow(img,interpolation='nearest')
```

```
Requirement already satisfied: graphviz in /opt/conda/envs/Python36/lib/python3.
6/site-packages (0.14)
Requirement already satisfied: pydotplus in /opt/conda/envs/Python36/lib/python
3.6/site-packages (2.0.2)
Requirement already satisfied: pyparsing>=2.0.1 in /opt/conda/envs/Python36/lib/
python3.6/site-packages (from pydotplus) (2.3.1)
```

Out[179]:   <matplotlib.image.AxesImage at 0x7f7193754470>



# Support Vector Machine

In [159]:
```python
from sklearn import svm
```

In [161]:
```python
model_svm = svm.SVC(kernel='rbf')
model_svm.fit(X_train, y_train)
yhatsvm = model_svm.predict(X_test)
yhatsvm
```

/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/svm/base.py:196: Fu
tureWarning: The default value of gamma will change from 'auto' to 'scale' in ve
rsion 0.22 to account better for unscaled features. Set gamma explicitly to 'aut
o' or 'scale' to avoid this warning.
  "avoid this warning.", FutureWarning)

Out[161]:
```
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)
```

# Logistic Regression

In [162]:
```python
from sklearn.linear_model import LogisticRegression
LR = LogisticRegression(C=0.01, solver='sag').fit(X_train,y_train)
LR
```

Out[162]:
```
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='sag',
          tol=0.0001, verbose=0, warm_start=False)
```

In [164]:
```python
yhat_lr = LR.predict(X_test)
yhat_lr
```

Out[164]:
```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF',
       'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'], dtype=object)
```

# Model Evaluation using Test set

In [165]: 
```python
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:

In [166]: 
```python
!wget -O loan_test.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-course
s-data/CognitiveClass/ML0101ENv3/labs/loan_test.csv
```

```
--2020-05-10 14:58:21--  https://s3-api.us-geo.objectstorage.softlayer.net/cf-co
urses-data/CognitiveClass/ML0101ENv3/labs/loan_test.csv
Resolving s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorag
e.softlayer.net)... 67.228.254.196
Connecting to s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectsto
rage.softlayer.net)|67.228.254.196|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 3642 (3.6K) [text/csv]
Saving to: 'loan_test.csv'

100%[======================================>] 3,642       --.-K/s   in 0s

2020-05-10 14:58:21 (275 MB/s) - 'loan_test.csv' saved [3642/3642]
```

## Load Test set for evaluation

In [167]: 
```python
test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

Out[167]:

|   | Unnamed: 0 | Unnamed: 0.1 | loan_status | Principal | terms | effective_date | due_date | age | education | Gender |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | PAIDOFF | 1000 | 30 | 9/8/2016 | 10/7/2016 | 50 | Bechalor | female |
| 1 | 5 | 5 | PAIDOFF | 300 | 7 | 9/9/2016 | 9/15/2016 | 35 | Master or Above | male |
| 2 | 21 | 21 | PAIDOFF | 1000 | 30 | 9/10/2016 | 10/9/2016 | 43 | High School or Below | female |
| 3 | 24 | 24 | PAIDOFF | 1000 | 30 | 9/10/2016 | 10/9/2016 | 26 | college | male |
| 4 | 35 | 35 | PAIDOFF | 800 | 15 | 9/11/2016 | 9/25/2016 | 29 | Bechalor | male |

In [168]:
```python
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]
```

/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/preprocessing/data.
py:645: DataConversionWarning: Data with input dtype uint8, int64 were all conve
rted to float64 by StandardScaler.
  return self.partial_fit(X, y)
/opt/conda/envs/Python36/lib/python3.6/site-packages/ipykernel/__main__.py:9: Da
taConversionWarning: Data with input dtype uint8, int64 were all converted to fl
oat64 by StandardScaler.

Out[168]:
```
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 [169]:
```python
test_y = test_df['loan_status'].values
test_y[0:5]
```

Out[169]:
```
array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
      dtype=object)
```

In [170]:
```python
knn_yhat=neigh.predict(test_X)
print("Avg F1-score of KNN: %.4f" % f1_score(test_y, knn_yhat, average='weighted
'))
print("Jaccard score of KNN: %.4f" % jaccard_similarity_score(test_y, knn_yhat))
```

Avg F1-score of KNN: 0.6328
Jaccard score of KNN: 0.6667

In [172]:
```python
tree_yhat=Tree.predict(test_X)
print("Avg F1-score of Decision Tree: %.4f" % f1_score(test_y, tree_yhat, average=
'weighted'))
print("Jaccard score of Decision Tree: %.4f" % jaccard_similarity_score(test_y, tr
ee_yhat))
```

Avg F1-score of Decision Tree: 0.7367
Jaccard score of Decision Tree: 0.7222

In [174]:
```python
svm_yhat=model_svm.predict(test_X)
print("Avg F1-score of SVM: %.4f" % f1_score(test_y, svm_yhat, average='weighted
'))
print("Jaccard score of SVM: %.4f" % jaccard_similarity_score(test_y, svm_yhat))
```

Avg F1-score of SVM: 0.7584
Jaccard score of SVM: 0.7963

In [177]:
```python
LR_yhat=model_svm.predict(test_X)
LR_yhat_prob = LR.predict_proba(test_X)
print("Avg F1-score of LR: %.4f" % f1_score(test_y, LR_yhat, average='weighted'))
print("Jaccard score of LR: %.4f" % jaccard_similarity_score(test_y, LR_yhat))
print("LogLoss score of LR: %.2f" % log_loss(test_y, LR_yhat_prob))
```

```
Avg F1-score of LR: 0.7584
Jaccard score of LR: 0.7963
LogLoss score of LR: 0.52
```

In [180]:
```python
jc1=jaccard_similarity_score(test_y, knn_yhat)
fs1=f1_score(test_y, knn_yhat, average='weighted')

jc2=jaccard_similarity_score(test_y, tree_yhat)
fs2=f1_score(test_y, tree_yhat, average='weighted')

jc3=jaccard_similarity_score(test_y, svm_yhat)
fs3=f1_score(test_y, svm_yhat, average='weighted')

jc4=jaccard_similarity_score(test_y, LR_yhat)
fs4=f1_score(test_y, LR_yhat, average='weighted')
ll4=log_loss(test_y, LR_yhat_prob)

list_jc = [jc1, jc2, jc3, jc4]
list_fs = [fs1, fs2, fs3, fs4]
list_ll = ['NA', 'NA', 'NA', ll4]


import pandas as pd

# fomulate the report format
df = pd.DataFrame(list_jc, index=['KNN','Decision Tree','SVM','Logistic Regression'])
df.columns = ['Jaccard']
df.insert(loc=1, column='F1-score', value=list_fs)
df.insert(loc=2, column='LogLoss', value=list_ll)
df.columns.name = 'Algorithm'
df
```

Out[180]:

| Algorithm | Jaccard | F1-score | LogLoss |
|---|---|---|---|
| KNN | 0.666667 | 0.632840 | NA |
| Decision Tree | 0.722222 | 0.736682 | NA |
| SVM | 0.796296 | 0.758350 | NA |
| Logistic Regression | 0.796296 | 0.758350 | 0.516366 |

# Report

You should be able to report the accuracy of the built model using different evaluation metrics:

| Algorithm | Jaccard | F1-score | LogLoss |
|---|---|---|---|
| KNN | 0.6667 | 0.6328 | NA |
| Decision Tree | 0.7222 | 0.7367 | NA |
| SVM | 0.7963 | 0.7584 | NA |
| LogisticRegression | 0.7963 | 0.7584 | 0.52 |

# Want to learn more?

IBM SPSS Modeler is a comprehensive analytics platform that has many machine learning algorithms. It has been designed to bring predictive intelligence to decisions made by individuals, by groups, by systems – by your enterprise as a whole. A free trial is available through this course, available here: SPSS Modeler (http://cocl.us/ML0101EN-SPSSModeler)

Also, you can use Watson Studio to run these notebooks faster with bigger datasets. Watson Studio is IBM's leading cloud solution for data scientists, built by data scientists. With Jupyter notebooks, RStudio, Apache Spark and popular libraries pre-packaged in the cloud, Watson Studio enables data scientists to collaborate on their projects without having to install anything. Join the fast-growing community of Watson Studio users today with a free account at Watson Studio (https://cocl.us /ML0101EN_DSX)

## Thanks for completing this lesson!

**Author: Saeed Aghabozorgi (https://ca.linkedin.com/in/saeedaghabozorgi)**

Saeed Aghabozorgi (https://ca.linkedin.com/in/saeedaghabozorgi), PhD is a Data Scientist in IBM with a track record of developing enterprise level applications that substantially increases clients' ability to turn data into actionable knowledge. He is a researcher in data mining field and expert in developing advanced analytic methods like machine learning and statistical modelling on large datasets.

---