



(<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 [3]: !wget -O loan_train.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_train.csv

--2020-03-30 06:31:49-- https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_train.csv
Resolving s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorage.softlayer.net)... 67.228.254.196
Connecting to s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorage.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 0.002s

2020-03-30 06:31:49 (12.3 MB/s) - 'loan_train.csv' saved [23101/23101]
```

### Load Data From CSV File

```
In [4]: df = pd.read_csv('loan_train.csv')
df.head()
```

Out[4]:

	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	Bechalar	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 [5]: df.shape
```

Out[5]: (346, 10)

Convert to date time object

```
In [6]: df['due_date'] = pd.to_datetime(df['due_date'])
df['effective_date'] = pd.to_datetime(df['effective_date'])
df.head()
```

Out[6]:

	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	Bechalar	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 [7]: df['loan_status'].value_counts()
```

Out[7]: 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]: # 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		
ca-certificates-2020.1.1	0	132 KB	anaconda
seaborn-0.10.0	py_0	161 KB	anaconda
certifi-2019.11.28	py36_1	157 KB	anaconda
openssl-1.1.1	h7b6447c_0	5.0 MB	anaconda
Total:		5.5 MB	

The following packages will be UPDATED:

```
ca-certificates: 2020.1.1-0 --> 2020.1.1-0 anaconda
certifi:         2019.11.28-py36_0 --> 2019.11.28-py36_1 anaconda
openssl:         1.1.1e-h7b6447c_0 --> 1.1.1-h7b6447c_0 anaconda
seaborn:         0.9.0-pyh91ea838_1 --> 0.10.0-py_0 anaconda
```

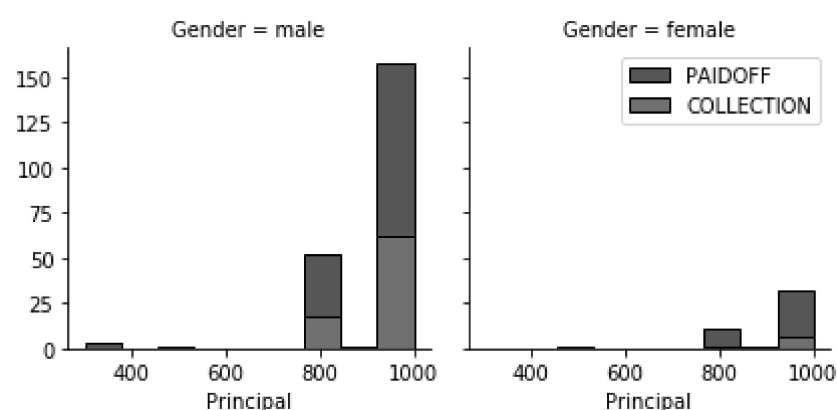
Downloading and Extracting Packages

```
ca-certificates-2020 | 132 KB | ##### | 100%
seaborn-0.10.0       | 161 KB | ##### | 100%
certifi-2019.11.28  | 157 KB | ##### | 100%
openssl-1.1.1       | 5.0 MB | ##### | 100%
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
```

```
In [9]: import seaborn as sns
```

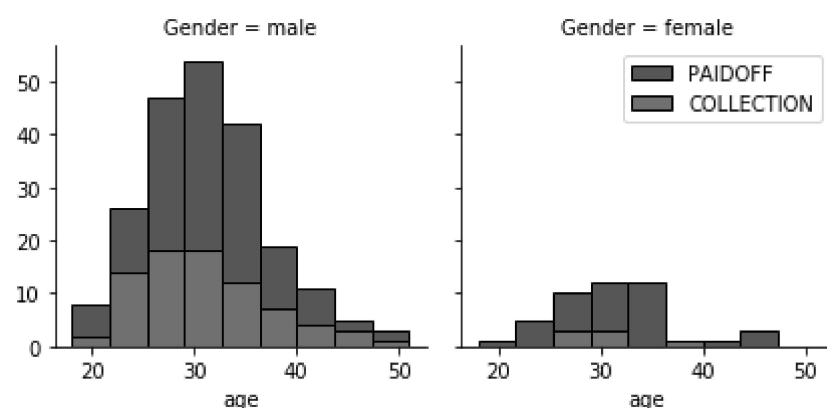
```
bins = np.linspace(df.Principal.min(), df.Principal.max(), 10)
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)
g.map(plt.hist, 'Principal', bins=bins, ec="k")

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



```
In [10]: 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")

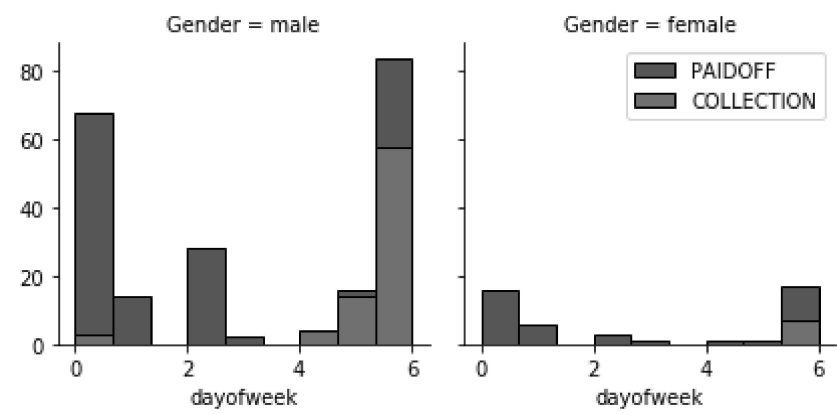
g.axes[-1].legend()
plt.show()
```



## Pre-processing: Feature selection/extraction

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

```
In [11]: 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 [12]: df['weekend'] = df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
df.head()
```

Out[12]:

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

Convert Categorical features to numerical values

Lets look at gender:

```
In [13]: df.groupby(['Gender'])['loan_status'].value_counts(normalize=True)
```

Out[13]:

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

Out[14]:

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

One Hot Encoding

How about education?

```
In [15]: df.groupby(['education'])['loan_status'].value_counts(normalize=True)
```

```
Out[15]: 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 [16]: df[['Principal','terms','age','Gender','education']].head()
```

```
Out[16]:
```

	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 variables to binary variables and append them to the feature Data Frame

```
In [17]: 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[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

Feature selection

Lets definnd feature sets, X:

```
In [18]: X = Feature
X[0:5]
```

```
Out[18]:
```

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

```
Out[19]: array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
              dtype=object)
```

Normalize Data

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

```
In [20]: 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 converted to float64 by StandardScaler.
    return self.partial_fit(X, y)
/opt/conda/envs/Python36/lib/python3.6/site-packages/ipykernel/__main__.py:1: DataConversionWarning: Data with input dtype uint8, int64 were all converted to float64 by StandardScaler.
    if __name__ == '__main__':

Out[20]: 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 [21]: # split train_loan
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 [22]: # import library
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics

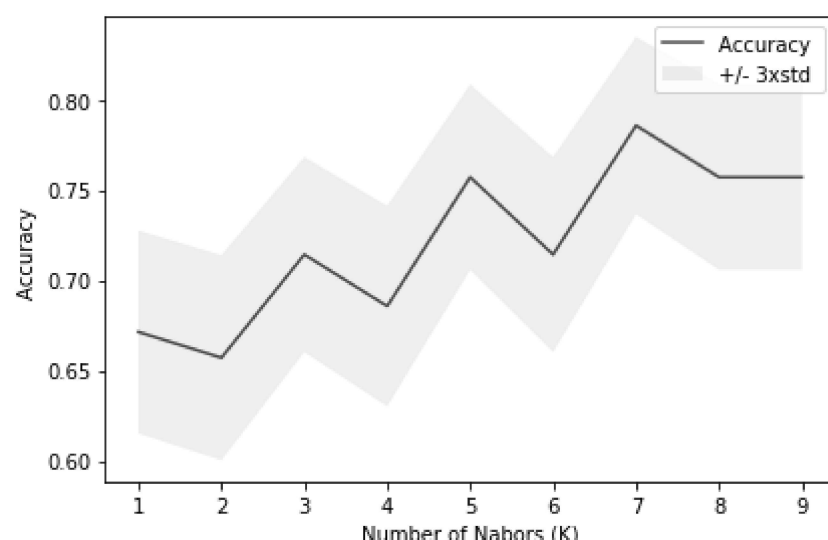
# try with 10 different values of k to find the best one
Ks = 10
mean_acc = np.zeros((Ks-1))
std_acc = np.zeros((Ks-1))
ConfustionMx = [];
for n in range(1,Ks):
    #Train Model and Predict
    neigh = KNeighborsClassifier(n_neighbors = n).fit(X_train,y_train)
    yhat=neigh.predict(X_test)
    mean_acc[n-1] = metrics.accuracy_score(y_test, yhat)
    std_acc[n-1]=np.std(yhat==y_test)/np.sqrt(yhat.shape[0])

# accuracy
print(mean_acc)

# Plot model accuracy for Different number of Neighbors
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()

# result
print( "The best accuracy was with", mean_acc.max(), "with k=", mean_acc.argmax()+1 )
```

```
[0.67142857 0.65714286 0.71428571 0.68571429 0.75714286 0.71428571
 0.78571429 0.75714286 0.75714286]
```



The best accuracy was with 0.7857142857142857 with k= 7

```
In [23]: # train model with k=7
neigh = KNeighborsClassifier(n_neighbors = 7).fit(X_train,y_train)
neigh
```

```
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]: # import library
from sklearn.tree import DecisionTreeClassifier
# create an instance of the DecisionTreeClassifier called loanTree
loanTree = DecisionTreeClassifier(criterion="entropy", max_depth = 4)
# training
loanTree.fit(X_train,y_train)
loanTree
```

```
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=None,
                                splitter='best')
```

In [ ]:

In [ ]:

# Support Vector Machine

```
In [25]: # import library
from sklearn import svm
# training
clf = svm.SVC()
clf.fit(X_train, y_train)

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

Out[25]: 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 [ ]:

In [ ]:
```

# Logistic Regression

```
In [26]: # import library
from sklearn.linear_model import LogisticRegression
# training
LR = LogisticRegression(C=0.01).fit(X_train,y_train)
LR

/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)

Out[26]: 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)
```

# Model Evaluation using Test set

```
In [27]: 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 [28]: !wget -O loan_test.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_test.csv

--2020-03-30 07:05:57-- https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_test.csv
Resolving s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorage.softlayer.net)... 67.228.254.196
Connecting to s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorage.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

2020-03-30 07:05:57 (387 MB/s) - 'loan_test.csv' saved [3642/3642]
```

## Load Test set for evaluation



```
In [29]: test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

Out[29]:

	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	Bechalar	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	Bechalar	male

```
In [30]: # Pre-processing Loan_test
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'] = 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)
Feature_test = test_df[['Principal','terms','age','Gender','weekend']]
Feature_test = pd.concat([Feature_test,pd.get_dummies(test_df['education'])], axis=1)
Feature_test.drop(['Master or Above'], axis = 1,inplace=True)
test_X = Feature_test
test_X = preprocessing.StandardScaler().fit(test_X).transform(test_X)
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 converted to float64 by StandardScaler.  
return self.partial\_fit(X, y)  
/opt/conda/envs/Python36/lib/python3.6/site-packages/ipykernel/\_\_main\_\_.py:11: DataConversionWarning: Data with input dtype uint8, int64 were all converted to float64 by StandardScaler.

```
Out[30]: array([[ 0.49362588,  0.92844966,  3.05981865,  1.97714211, -4.12310563,
                2.39791576, -0.79772404, -0.86135677],
               [-3.56269116, -1.70427745,  0.53336288, -0.50578054, -4.12310563,
               -0.41702883, -0.79772404, -0.86135677],
               [ 0.49362588,  0.92844966,  1.88080596,  1.97714211, -4.12310563,
               -0.41702883,  1.25356634, -0.86135677],
               [ 0.49362588,  0.92844966, -0.98251057, -0.50578054,  0.24253563,
               -0.41702883, -0.79772404,  1.16095912],
               [-0.66532184, -0.78854628, -0.47721942, -0.50578054,  0.24253563,
                2.39791576, -0.79772404, -0.86135677]])
```

```
In [33]: # Pre-processing Loan_test (cont)
test_y = test_df['loan_status'].values
test_y[0:5]
```

```
Out[33]: array(['PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF', 'PAIDOFF'],
               dtype=object)
```

KNN

```
In [34]: # predicted y
yhat_knn = neigh.predict(test_X)

# jaccard
jaccard_knn = jaccard_similarity_score(test_y, yhat_knn)
print("KNN Jaccard index: ", jaccard_knn)

# f1_score
f1_score_knn = f1_score(test_y, yhat_knn, average='weighted')
print("KNN F1-score: ", f1_score_knn)

KNN Jaccard index:  0.6851851851851852
KNN F1-score:  0.6453810131971051
```

Decision Tree

```
In [35]: # predicted y
yhat_dt = loanTree.predict(test_X)

# jaccard
jaccard_dt = jaccard_similarity_score(test_y, yhat_dt)
print("DT Jaccard index: ", jaccard_dt)

# f1_score
f1_score_dt = f1_score(test_y, yhat_dt, average='weighted')
print("DT F1-score: ", f1_score_dt)

DT Jaccard index:  0.5185185185185185
DT F1-score:  0.5385802469135802
```

# SVM

```
In [36]: # predicted y
yhat_svm = clf.predict(test_X)

# jaccard
jaccard_svm = jaccard_similarity_score(test_y, yhat_svm)
print("SVM Jaccard index: ", jaccard_svm)

# f1_score
f1_score_svm = f1_score(test_y, yhat_svm, average='weighted')
print("SVM F1-score: ", f1_score_svm)

SVM Jaccard index:  0.8148148148148148
SVM F1-score:  0.7861952861952862
```

# Logistic regression

```
In [38]: # predicted y
yhat_lg = LR.predict(test_X)
yhat_lg_prob = LR.predict_proba(test_X)

# jaccard
jaccard_lg = jaccard_similarity_score(test_y, yhat_lg)
print("LR Jaccard index: ", jaccard_lg)

# f1_score
f1_score_lg = f1_score(test_y, yhat_lg, average='weighted')
print("LR F1-score: ", f1_score_lg)

# logloss
logloss_lg = log_loss(test_y, yhat_lg_prob)
print("LR log loss: ", logloss_lg)

LR Jaccard index:  0.7407407407407407
LR F1-score:  0.6304176516942475
LR log loss:  0.6037871272191607

/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/metrics/classification.py:1143: UndefinedMetricWarning:
F-score is ill-defined and being set to 0.0 in labels with no predicted samples.
'precision', 'predicted', average, warn_for)
```

# Report

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

Algorithm		Jaccard	F1-score	LogLoss
KNN		0.6851851851851852	0.6453810131971051	NA
Decision Tree		0.5185185185185185	0.5385802469135802	NA
SVM		0.8148148148148148	0.7861952861952862	NA
LogisticRegression		0.7407407407407407	0.6304176516942475	0.6037871272191607

## 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 \(http://cocl.us/ML0101EN-SPSSModeler\)](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 \(https://cocl.us/ML0101EN\\_DSX\)](https://cocl.us/ML0101EN_DSX).

## Thanks for completing this lesson!

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

[Saeed Aghabozorgi \(https://ca.linkedin.com/in/saeedaghabozorgi\)](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.

---

Copyright © 2018 [Cognitive Class \(https://cocl.us/DX0108EN\\_CC\)](https://cocl.us/DX0108EN_CC). This notebook and its source code are released under the terms of the [MIT License \(https://bigdatauniversity.com/mit-license/\)](https://bigdatauniversity.com/mit-license/).