



## 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 metrics.

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
import 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 or 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-goe-objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENR3/1a/loan_train.csv
--2020-08-27 23:03:56-- https://s3-api-us-goe-objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENR3/1a/loan_train.csv
Resolving s3-api-us-goe-objectstorage.softlayer.net (s3-api-us-goe-objectstorage.softlayer.net)... 67.228.254.196
Connecting to s3-api-us-goe-objectstorage.softlayer.net (s3-api-us-goe-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

2020-08-27 23:03:57 (10.7 MB/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	9/6/2016	10/7/2016	45	High School or Below	male
1	2	2	PAIDOFF	1000	30	9/6/2016	10/7/2016	33	Bachelor	female
2	3	3	PAIDOFF	1000	15	9/6/2016	9/22/2016	27	college	male
3	4	4	PAIDOFF	1000	30	9/6/2016	10/6/2016	28	college	female
4	6	6	PAIDOFF	1000	30	9/9/2016	10/6/2016	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	Bachelor	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

Lets see how many of each class is in our data set

```
In [6]: df['loan_status'].value_counts()
```

```
Out[6]:
```

```
PAIDOFF    269
COLLECTION    86
Name: loan_status, dtype: int64
```

269 people have paid off the loan on time while 86 have gone into collection

Lets plot some columns to understand 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
ca-certificates-2020.7.22	py36_0
certifi-2020.6.20	py36_0
seaborn-0.10.1	py36_0
Total:	
4.2 MB	

The following packages will be UPDATED:

ca-certificates	2020.6.24-0 --> 2020.7.22-0	anaconda
certifi	2020.6.20-py36_0 --> 2020.6.20-py36_0	anaconda
openssl	1.1.1g-h7b6447c_0 --> 1.1.1g-h7b6447c_0	anaconda
seaborn	0.9.0-py36ha831_1 --> 0.10.1-py36_0	anaconda

Downloading and Extracting Packages

openssl-1.1.1.g	3.4 MB	#####	100%
ca-certificates-2020.7.22	132 KB	#####	100%
certifi-2020.6.20	140 KB	#####	100%
seaborn-0.10.1	160 KB	#####	100%

Preparing transaction: done  
Verifying transaction: done  
Executing transaction: done

```
In [14]: 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 [15]: 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 [16]: 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 don't pay it off, so lets use Feature binarization to set a threshold values less than day 4

```
In [17]: df['weekend'] = df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
df.head()
```

```
Out[17]:
```

	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	Bachelor	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 [18]: df.groupby(['Gender'])['loan_status'].value_counts(normalize=True)
```

```
Out[18]:
```

```
Gender  loan_status
female  PAIDOFF      0.863385
        COLLECTION  0.136615
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 [19]: df['Gender'].replace(to_replace=['male','female'], value=[0,1],inplace=True)
df.head()
```

```
Out[19]:
```

	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	Bachelor	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 [20]: df.groupby(['education'])['loan_status'].value_counts(normalize=True)
```

```
Out[20]:
```

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

```
In [21]: df[['Principal','terms','age','Gender','education']].head()
```

```
Out[21]:
```

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

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

```
In [22]: 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[22]:
```

	Principal	terms	age	Gender	weekend	Bachelor	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 define feature sets, X:

```
In [23]: X = Feature
X[0:5]
```

```
Out[23]:
```

	Principal	terms	age	Gender	weekend	Bachelor	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 labels?

```
In [24]: y = df['loan_status'].values
y[0:5]
```

```
Out[24]: array([ 0, 1, 0, 1, 0])
dtype: object
```

### Normalize Data

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

```
In [25]: X = preprocessing.StandardScaler().fit(X).transform(X)
X[0:5]
```

```
Out[25]:
```

```
array([[ 0.51578458,  0.92071769,  2.33152555, -0.42054004, -1.20577805,
        -0.38170942,  1.1363974, -0.86846108],
       [ 0.51578458,  0.92071769,  0.34010449, -1.23778177, -1.20577805,
        -0.6185424, -0.87997669, -0.86968108],
       [ 0.51578458,  0.89911111, -0.65221055, -0.42054004, -1.20577805,
        -0.38170942, -0.87997669, -1.16946479],
       [ 0.51578458,  0.92071769, -0.48739188,  2.3778177,  0.82934003,
        -0.38170942, -0.87997669, -1.16946479],
       [ 0.51578458,  0.92071769, -0.3215732, -0.42054004,  0.82934003,
        -0.38170942, -0.87997669, -1.16946479]])
```

### Classification

Now, it's 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 sklearn, 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)

Warning: 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 [26]: from sklearn.neighbors import KNeighborsClassifier

k = 6
#Train Model and Predict
model_knn = KNeighborsClassifier(n_neighbors = k)
model_knn.fit(X,y)
model_knn
```

```
Out[26]:
```

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                    metric_params=None, n_jobs=1, n_neighbors=6, p=2,
                    weights='uniform')
```

```
In [ ]:
```

```
In [ ]:
```

### Decision Tree

```
In [28]: from sklearn.tree import DecisionTreeClassifier

model_tree = DecisionTreeClassifier(criterion='entropy', max_depth = 4)
model_tree.fit(X,y)
model_tree
```

```
Out[28]:
```

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

```
In [ ]:
```

```
In [ ]:
```

### Support Vector Machine

```
In [29]: from sklearn import svm

model_svm = svm.SVC(kernel='linear')
model_svm.fit(X, y)
model_svm
```

```
Out[29]:
```

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

```
In [ ]:
```

```
In [ ]:
```

### Logistic Regression

```
In [30]: from sklearn.linear_model import LogisticRegression

model_lr = LogisticRegression(C=0.0001, solver='liblinear')
model_lr.fit(X,y)
model_lr
```

```
Out[30]:
```

```
LogisticRegression(C=0.0001, class_weight=None, dual=False,
                    fit_intercept=True, intercept_scaling=1, max_iter=10,
                    multi_class='warn', n_jobs=None, penalty='l2', random_state=None,
                    solver='liblinear', tol=0.0001, verbose=0, warm_start=False)
```

```
In [ ]:
```

```
In [ ]:
```

### Model Evaluation using Test set

```
In [31]: 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 [32]: wget -O loan_test.csv https://s3-api-us-goe-objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENR3/1ab/loan_test.csv
--2020-08-27 23:13:12-- https://s3-api-us-goe-objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENR3/1ab/loan_test.csv
Resolving s3-api-us-goe-objectstorage.softlayer.net (s3-api-us-goe-objectstorage.softlayer.net)... 67.228.254.196
Connecting to s3-api-us-goe-objectstorage.softlayer.net (s3-api-us-goe-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

2020-08-27 23:13:12 (339 KB/s) - 'loan_test.csv' saved [3642/3642]
```

### Load Test set for evaluation

```
In [33]: test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

```
Out[33]:
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	1	1	PAIDOFF	1000	30	9/6/2016	10/7/2016	50	Bachelor	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	Bachelor	male

```
In [34]: 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_df[['Principal','terms','age','Gender','education']].head()
```

```
Out[34]:
```

	Principal	terms	age	Gender	education
0	1000	30	50	1	Bachelor
1	300	7	35	0	Master or Above
2	1000	30	43	1	High School or Below
3	1000	30	26	0	college
4	800	15	29	0	Bachelor

```
In [35]: 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)
Feature_test.head()
```

```
Out[35]:
```

	Principal	terms	age	Gender	weekend	Bachelor	High School or Below	college
0	1000	30	50	1	0	1	0	0
1	300	7	35	0	1	0	0	0
2	1000	30	43	1	1	0	1	0
3	1000	30	26	0	1	1	0	1
4	800	15	29	0	1	1	0	0

```
In [36]: X_test = Feature_test
X_test = preprocessing.StandardScaler().fit(X_test).transform(X_test)
X_test[0:5]
```

```
Out[36]:
```

```
array([[ 0.51578458,  0.80635838, 1.5928902, -0.42054004, -1.20577805,
        -0.68129585,  1.18518519,  0.82934003],
       [ 0.51578458,  0.80635838,  0.75145986,  0.7052023,
        -0.42054004,  0.82934003,  0.82934003],
       [ 0.51578458,  0.80635838,  0.75145986,  0.7052023,
        -0.42054004,  0.82934003,  0.82934003],
       [ 0.51578458,  0.80635838,  0.75145986,  0.7052023,
        -0.42054004,  0.82934003,  0.82934003],
       [ 0.51578458,  0.80635838,  0.75145986,  0.7052023,
        -0.42054004,  0.82934003,  0.82934003]])
```

```
In [37]: m = "Jaccard",f1_score,("LogLoss")
report = report.set_index(["Algorithm", "Jaccard", "F1-score", "LogLoss"])
report
```

```
Out[37]:
```

	Algorithm	Jaccard	F1-score	LogLoss
KNN	0.6851851851851852	0.681295852331905	nan	nan
Decision Tree	0.7777777777777778	0.7283950617283951	nan	nan
SVM	0.7402740274027403	0.63041		