



Classification with Python

In this notebook we try to practice all the classification algorithms that we have 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.

Let's 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 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

Let's download the dataset

```
In [2]: !wget -O loan_train.csv https://cf-courses-data.s3.us.cloud-object-storage.appdomain.c
```

Load Data From CSV File

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

```
Out[3]:
```

	Unnamed: 0.1	Unnamed: 0	loan_status	Principal	terms	effective_date	due_date	age	education	Ge
0	0	0	PAIDOFF	1000	30	9/8/2016	10/7/2016	45	High School or Below	Female
1	2	2	PAIDOFF	1000	30	9/8/2016	10/7/2016	33	Bechelor	Female
2	3	3	PAIDOFF	1000	15	9/8/2016	9/22/2016	27	college	Female
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	Female

```
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: 0.1	Unnamed: 0	loan_status	Principal	terms	effective_date	due_date	age	education	Ge
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below	Female
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechelor	Female
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	Female
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	Female

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

Let's plot some columns to understand data better:

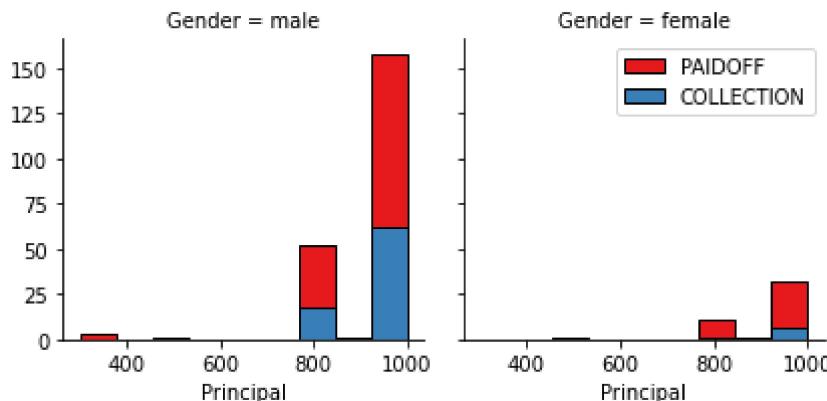
```
In [7]: # notice: installing seaborn might takes a few minutes  
!conda install -c anaconda seaborn -y
```

```
Collecting package metadata (current_repodata.json): ...working... done  
Solving environment: ...working... done
```

```
# All requested packages already installed.
```

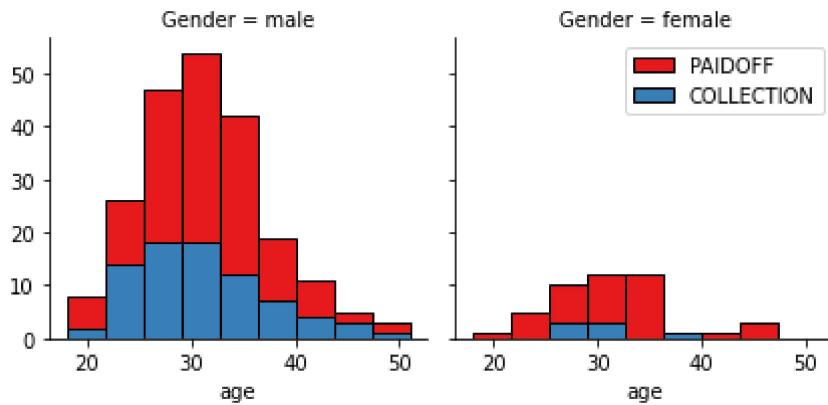
```
In [8]: import seaborn as sns
```

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



```
In [9]: bins = np.linspace(df.age.min(), df.age.max(), 10)
```

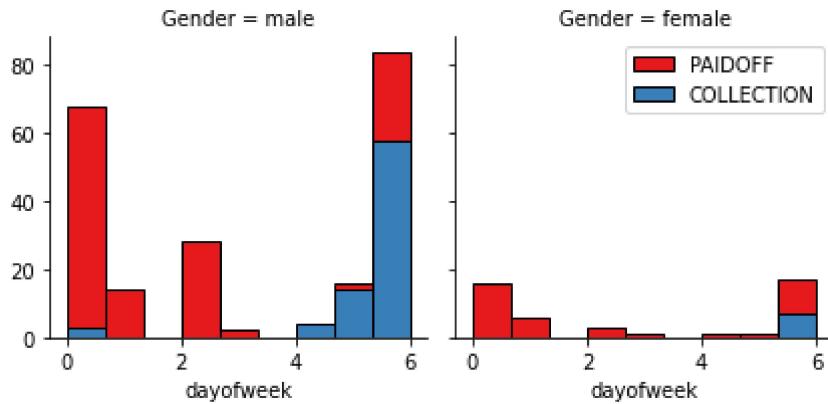
```
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)  
g.map(plt.hist, 'age', bins=bins, ec="k")  
  
g.axes[-1].legend()  
plt.show()
```



Pre-processing: Feature selection/extraction

Let's 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 don't pay it off, so let's use Feature binarization to set a threshold value less than day 4

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

Out[11]:

	Unnamed: 0.1	Unnamed: 0	loan_status	Principal	terms	effective_date	due_date	age	education	Ge
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	Bechelor	fe
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	college	fe
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	college	

Convert Categorical features to numerical values

Let's 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

Let's 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: 0.1	Unnamed: 0	loan_status	Principal	terms	effective_date	due_date	age	education	Ge
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	Bechelor	
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	college	
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	college	

One Hot Encoding

How about education?

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

Out[14]:

education	loan_status	
Bechelor	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

Features before 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	Bechelor
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 [44]: `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[44]:

	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

Let's define feature sets, X:

In [45]:

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

Out[45]:

	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 lables?

In [46]:

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

Out[46]:

```
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 [47]:

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

```
Out[47]: 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 [48]: #Splitting the Data to Train and test
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=42)
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 [49]: #Choosing the best K
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics
Ks = 15
mean_acc = np.zeros((Ks-1))
std_acc = np.zeros((Ks-1))
```

```

ConfustionMx = [];
for n in range(1,Ks):

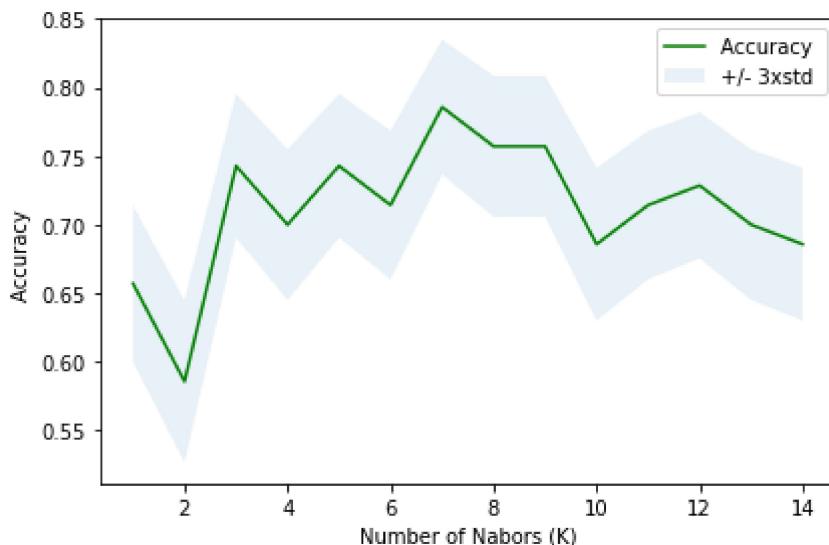
    #Train Model and Predict
    neigh = KNeighborsClassifier(n_neighbors = n).fit(X_train,y_train)
    yhat=neigh.predict(X_test)
    mean_acc[n-1] = metrics.accuracy_score(y_test, yhat)
    std_acc[n-1]=np.std(yhat==y_test)/np.sqrt(yhat.shape[0])

mean_acc

```

Out[49]: array([0.65714286, 0.58571429, 0.74285714, 0.7 , 0.74285714,
 0.71428571, 0.78571429, 0.75714286, 0.75714286, 0.68571429,
 0.71428571, 0.72857143, 0.7 , 0.68571429])

In [22]: `#Plotting the best K`
`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=0.1)`
`plt.legend(['Accuracy ',' +/- 3xstd'])`
`plt.ylabel('Accuracy ')`
`plt.xlabel('Number of Nabors (K)')`
`plt.tight_layout()`
`plt.show()`



In [23]: `#The best K is 7 as shown in the figure`
`print("The best accuracy was with", mean_acc.max(), "with k=", mean_acc.argmax()+1)`
`neigh = KNeighborsClassifier(n_neighbors = 7).fit(X_train,y_train)`
`yhatkm=neigh.predict(X_test)`

The best accuracy was with 0.7857142857142857 with k= 7

Decision Tree

In [24]: `#impporting decision tree from sklearn`
`from sklearn.tree import DecisionTreeClassifier`

In [25]: `Tree = DecisionTreeClassifier(criterion="entropy", max_depth = 15)`
`Tree.fit(X_train,y_train)`

```
Out[25]: DecisionTreeClassifier(criterion='entropy', max_depth=15)
```

```
In [26]: predTree = Tree.predict(X_test)
print("DecisionTrees's Accuracy: ", metrics.accuracy_score(y_test, predTree))

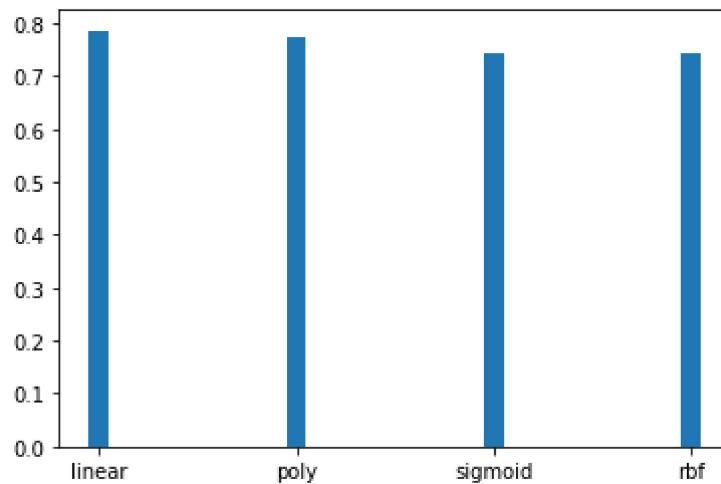
DecisionTrees's Accuracy: 0.7428571428571429
```

Support Vector Machine

```
In [27]: #Choosing best kernel
from sklearn import svm
kernels=['linear','poly','sigmoid','rbf']
mean_acc2 = np.zeros((len(kernels)))
for i,kernel in enumerate(kernels):
    clf = svm.SVC(kernel=kernel)
    clf.fit(X_train, y_train)
    yhat2 = clf.predict(X_test)
    mean_acc2[i] = metrics.accuracy_score(y_test, yhat2)
```

```
In [28]: #Plotting a graph between mean_acc vs kernels
mean_acc2=list(mean_acc2)
plt.bar(kernels,mean_acc2,width=0.1)
```

```
Out[28]: <BarContainer object of 4 artists>
```



From graph linear kernel has max accuracy so-

```
In [29]: clf = svm.SVC(kernel='linear')
clf.fit(X_train, y_train)
yhatsvm = clf.predict(X_test)
```

```
In [30]: print("SVM's Accuracy: ", metrics.accuracy_score(y_test, yhatsvm))
```

```
SVM's Accuracy: 0.7857142857142857
```

Logistic Regression

```
In [31]: #Logistic regression cannot deal with categorial string variables so we must change th
```

```

y_testlog=y_test.copy()
y_trainlog=y_train.copy()
y_testlog[y_testlog == 'PAIDOFF']=1
y_testlog[y_testlog == 'COLLECTION']=0
y_trainlog[y_trainlog == 'PAIDOFF']=1
y_trainlog[y_trainlog == 'COLLECTION']=0
y_trainlog = y_trainlog.astype('float64')
y_testlog = y_testlog.astype('float64')
print("Ytrain:",y_train[0:5],"Ytrain Log:",y_trainlog[0:5])
print("Ytest:",y_test[0:5],"Ytest Log:",y_testlog[0:5])

```

```

Ytrain: ['PAIDOFF' 'COLLECTION' 'PAIDOFF' 'PAIDOFF' 'PAIDOFF'] ,Ytrain Log: [1. 0. 1.
1. 1.]
Ytest: ['PAIDOFF' 'PAIDOFF' 'PAIDOFF' 'PAIDOFF' 'PAIDOFF'] ,Ytest Log: [1. 1. 1. 1.
1.]

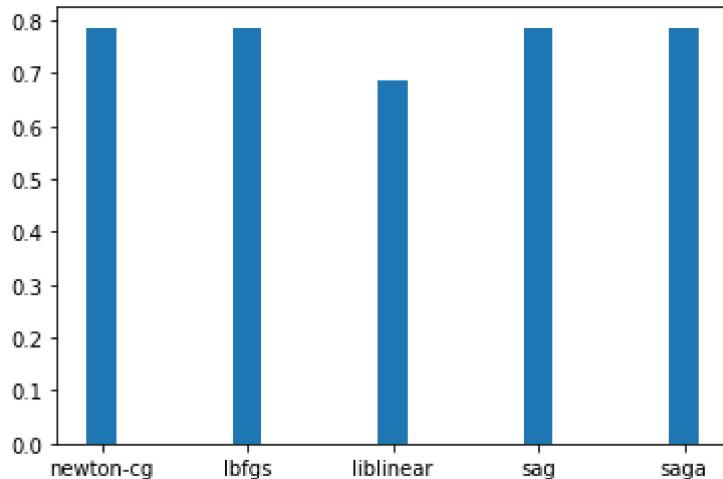
```

In [32]:

```

from sklearn.linear_model import LogisticRegression
solvers=['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']
mean_acc3 = np.zeros((len(solvers)))
for i,solver in enumerate(solvers):
    LR = LogisticRegression(C=0.01, solver=solver).fit(X_train,y_trainlog)
    yhat3=LR.predict(X_test)
    mean_acc3[i] = metrics.accuracy_score(y_testlog, yhat3)
mean_acc3=list(mean_acc3)
ax=plt.bar(solvers,mean_acc3,width=0.2)

```



In [33]:

```

LR = LogisticRegression(C=0.01, solver='sag').fit(X_train,y_trainlog)
yhatLR=LR.predict(X_test)
print("LR's Accuracy: ", metrics.accuracy_score(y_testlog, yhatLR))

```

LR's Accuracy: 0.7857142857142857

Model Evaluation using Test set

In [34]:

```

from sklearn.metrics import jaccard_score
from sklearn.metrics import f1_score
from sklearn.metrics import log_loss

```

First, download and load the test set:

In [35]:

```
!wget -O loan_test.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data
```

Load Test set for evaluation

```
In [36]: test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

```
Out[36]:
```

	Unnamed: 0.1	Unnamed: 0	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	1	1	PAIDOFF	1000	30	9/8/2016	10/7/2016	50	Bechelor	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	School or Below	High School
3	24	24	PAIDOFF	1000	30	9/10/2016	10/9/2016	26	college	Female
4	35	35	PAIDOFF	800	15	9/11/2016	9/25/2016	29	Bechelor	Female

```
In [37]: 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)
X_test_df = test_df[['Principal','terms','age','Gender','weekend']]
X_test_df = pd.concat([X_test_df,pd.get_dummies(test_df['education'])], axis=1)
X_test_df.drop(['Master or Above'], axis = 1, inplace=True)
y_test_df = test_df['loan_status'].values
finaltest=preprocessing.StandardScaler().fit(X_test_df).transform(X_test_df)
finalresultlog=y_test_df.copy()
finalresultlog[finalresultlog == 'PAIDOFF']=1
finalresultlog[finalresultlog == 'COLLECTION']=0
finalresultlog=finalresultlog.astype('float64')
```

Making a dataframe for report

```
In [38]: report = pd.DataFrame(columns=['Jaccard', 'F1-score', 'LogLoss'])
report['Algorithm']=['KNN', 'Decision Tree', 'SVM', 'LogisticRegression']
report=report.set_index('Algorithm',)
```

```
In [39]: finalkmeans=neigh.predict(finaltest)
finalDT=Tree.predict(finaltest)
finalsvm=clf.predict(finaltest)
finallog=LR.predict(finaltest)
```

Filling data (f_score, jaccard_score and log_loss) in dataframe

```
In [42]: report.loc['LogisticRegression','Jaccard']=jaccard_score(finalresultlog, finallog, average='weighted')
report.loc['LogisticRegression','F1-score']=f1_score(finalresultlog, finallog, average='weighted')
report.loc['LogisticRegression','LogLoss']=log_loss(finalresultlog, LR.predict_proba(finaltest))
report.loc['SVM','Jaccard']=jaccard_score(y_test_df,finalsvm,average='weighted')
report.loc['SVM','F1-score']=f1_score(y_test_df,finalsvm, pos_label='PAIDOFF', average='weighted')
report.loc['Decision Tree','Jaccard']=jaccard_score(y_test_df,finalDT,average='weighted')
```

```
report.loc['Decision Tree', 'F1-score']=f1_score(y_test_df,finalDT, pos_label='PAIDOFF')
report.loc['KNN', 'Jaccard']=jaccard_score(y_test_df,finalkmeans,average='weighted')
report.loc['KNN', 'F1-score']=f1_score(y_test_df,finalkmeans, pos_label='PAIDOFF', average='weighted')
```

```
C:\Users\lenovo\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1370:
UserWarning: Note that pos_label (set to 'PAIDOFF') is ignored when average != 'binary' (got 'weighted'). You may use labels=[pos_label] to specify a single positive class.
    warnings.warn(
C:\Users\lenovo\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1370:
UserWarning: Note that pos_label (set to 'PAIDOFF') is ignored when average != 'binary' (got 'weighted'). You may use labels=[pos_label] to specify a single positive class.
    warnings.warn(
C:\Users\lenovo\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1370:
UserWarning: Note that pos_label (set to 'PAIDOFF') is ignored when average != 'binary' (got 'weighted'). You may use labels=[pos_label] to specify a single positive class.
    warnings.warn(
```

Report

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

In [43]: report

Out[43]:

Jaccard F1-score LogLoss

Algorithm

	KNN	0.532716	0.660225	NaN
	Decision Tree	0.586826	0.725253	NaN
	SVM	0.548697	0.630418	NaN
	LogisticRegression	0.548697	0.630418	0.516366

Algorithm Jaccard F1-score LogLoss

KNN	?	?	NA
Decision Tree	?	?	NA
SVM	?	?	NA
LogisticRegression	?	?	?

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](#)

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](#)

Thanks for completing this lesson!

Author: [Saeed Aghabozorgi](#)

[Saeed Aghabozorgi](#), 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.

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-10-27	2.1	Lakshmi Holla	Made changes in import statement due to updates in version of sklearn library
2020-08-27	2.0	Malika Singla	Added lab to GitLab

© IBM Corporation 2020. All rights reserved.