



IBM Developer SKILLS NETWORK

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 [6]:  
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

Field	Description
Education	Education of applicant
Gender	The gender of applicant

Let's download the dataset

```
In [7]: !wget -O loan_train.csv https://cf-courses-data.s3.us.cloud-object-storage.appdomain
```

'wget' is not recognized as an internal or external command,
operable program or batch file.

Load Data From CSV File

```
In [8]: df = pd.read_csv('loan_train.csv')
df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	G
0	0	0	PAIDOFF	1000	30	9/8/2016	10/7/2016	45	High School or Below	1
1	2	2	PAIDOFF	1000	30	9/8/2016	10/7/2016	33	Bechelor	1
2	3	3	PAIDOFF	1000	15	9/8/2016	9/22/2016	27	college	1
3	4	4	PAIDOFF	1000	30	9/9/2016	10/8/2016	28	college	1
4	6	6	PAIDOFF	1000	30	9/9/2016	10/8/2016	29	college	1

```
In [9]: df.shape
```

Out[9]: (346, 10)

Convert to date time object

```
In [10]: df['due_date'] = pd.to_datetime(df['due_date'])
df['effective_date'] = pd.to_datetime(df['effective_date'])
df.head()
```

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

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	G
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	college	

Data visualization and pre-processing

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

```
In [11]: df['loan_status'].value_counts()
```

```
Out[11]: 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 [12]: # 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
```

```
## Package Plan ##

environment location: G:\Anaconda3

added / updated specs:
- seaborn
```

The following packages will be downloaded:

package	build
conda-4.12.0	py38haa95532_0
	Total: 14.5 MB

The following packages will be UPDATED:

```
conda          conda-forge::conda-4.11.0-py38haa244f~ --> pkgs/main::conda-4.1
2.0-py38haa95532_0
```

Downloading and Extracting Packages

conda-4.12.0	14.5 MB		0%
conda-4.12.0	14.5 MB		0%
conda-4.12.0	14.5 MB		1%
conda-4.12.0	14.5 MB	2	2%
conda-4.12.0	14.5 MB	4	5%
conda-4.12.0	14.5 MB	7	7%
conda-4.12.0	14.5 MB	9	10%
conda-4.12.0	14.5 MB	#2	12%
conda-4.12.0	14.5 MB	#4	14%
conda-4.12.0	14.5 MB	#6	17%
conda-4.12.0	14.5 MB	#8	19%

conda-4.12.0	14.5 MB	##	21%
conda-4.12.0	14.5 MB	##3	23%
conda-4.12.0	14.5 MB	##6	26%
conda-4.12.0	14.5 MB	##8	29%
conda-4.12.0	14.5 MB	###	31%
conda-4.12.0	14.5 MB	###3	33%
conda-4.12.0	14.5 MB	###4	35%
conda-4.12.0	14.5 MB	###6	36%
conda-4.12.0	14.5 MB	###7	38%
conda-4.12.0	14.5 MB	###9	39%
conda-4.12.0	14.5 MB	####	41%
conda-4.12.0	14.5 MB	####2	42%
conda-4.12.0	14.5 MB	####3	44%
conda-4.12.0	14.5 MB	####4	45%
conda-4.12.0	14.5 MB	####6	46%
conda-4.12.0	14.5 MB	####7	47%
conda-4.12.0	14.5 MB	####8	49%
conda-4.12.0	14.5 MB	####9	50%
conda-4.12.0	14.5 MB	#####	51%
conda-4.12.0	14.5 MB	#####1	52%
conda-4.12.0	14.5 MB	#####2	53%
conda-4.12.0	14.5 MB	#####4	54%
conda-4.12.0	14.5 MB	#####5	55%
conda-4.12.0	14.5 MB	#####6	56%
conda-4.12.0	14.5 MB	#####7	58%
conda-4.12.0	14.5 MB	#####8	58%
conda-4.12.0	14.5 MB	#####9	59%
conda-4.12.0	14.5 MB	#####	61%
conda-4.12.0	14.5 MB	#####1	62%
conda-4.12.0	14.5 MB	#####2	63%
conda-4.12.0	14.5 MB	#####4	64%
conda-4.12.0	14.5 MB	#####5	66%
conda-4.12.0	14.5 MB	#####6	67%
conda-4.12.0	14.5 MB	#####8	68%
conda-4.12.0	14.5 MB	#####9	69%
conda-4.12.0	14.5 MB	#####	71%
conda-4.12.0	14.5 MB	#####2	72%
conda-4.12.0	14.5 MB	#####4	74%
conda-4.12.0	14.5 MB	#####6	76%
conda-4.12.0	14.5 MB	#####8	78%
conda-4.12.0	14.5 MB	#####	81%
conda-4.12.0	14.5 MB	#####3	83%
conda-4.12.0	14.5 MB	#####6	86%
conda-4.12.0	14.5 MB	#####8	89%
conda-4.12.0	14.5 MB	#####1	91%
conda-4.12.0	14.5 MB	#####3	94%
conda-4.12.0	14.5 MB	#####6	96%
conda-4.12.0	14.5 MB	#####8	98%
conda-4.12.0	14.5 MB	#####	100%
conda-4.12.0	14.5 MB	#####	100%

Preparing transaction: ...working... done

Verifying transaction: ...working... done

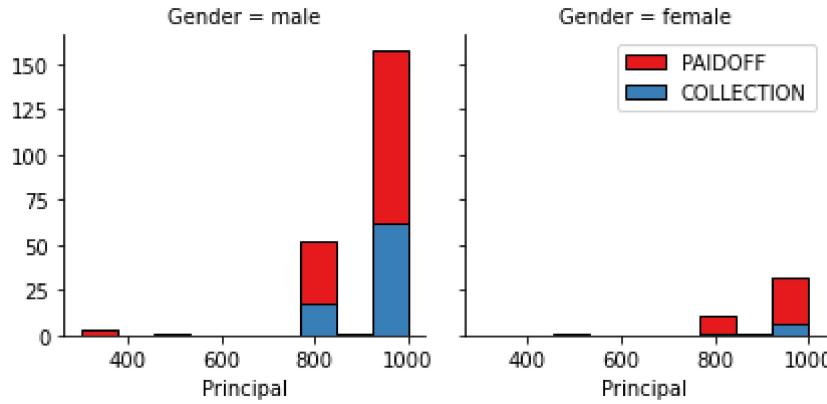
Executing transaction: ...working... done

In [13]:

```
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()
```



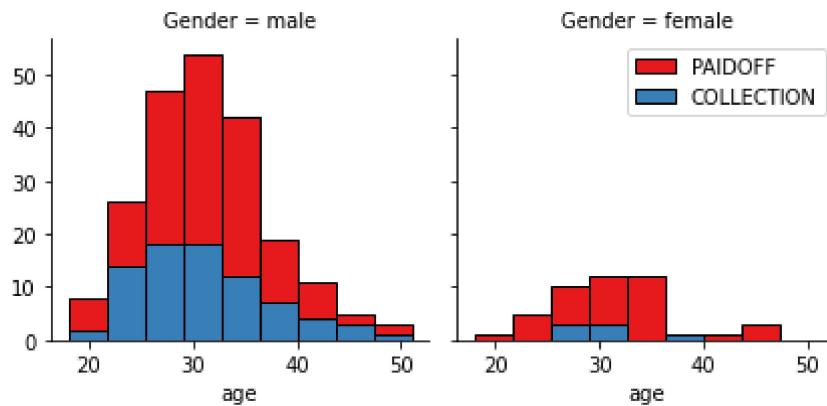
In [14]:

```

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

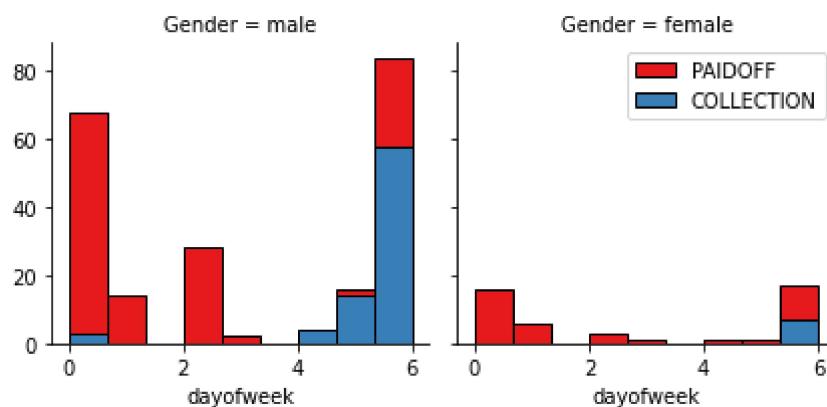
Let's look at the day of the week people get the loan

In [15]:

```

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 [16]:

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

Out[16]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	G
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	1
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	1
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 [17]:

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

Out[17]:

```
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 [18]:

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

Out[18]:

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	G
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	

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	G
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	college	

One Hot Encoding

How about education?

```
In [19]: df.groupby(['education'])['loan_status'].value_counts(normalize=True)
```

```
Out[19]: 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 [20]: df[['Principal', 'terms', 'age', 'Gender', 'education']].head()
```

	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 [21]: 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()
```

	Principal	terms	age	Gender	weekend	Bechelor	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 [22]:

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

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

What are our labels?

In [23]:

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

Out[23]: 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 [24]:

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

Out[24]: 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.

In [27]:

```
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,)

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 [25]:

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
```

In [28]:

```
for k in range(1, 10):
    knn_model = KNeighborsClassifier(n_neighbors = k).fit(x_train, y_train)
    knn_yhat = knn_model.predict(x_test)
    print("For K = {} accuracy = {}".format(k,accuracy_score(y_test,knn_yhat)))
```

For K = 1 accuracy = 0.6714285714285714
 For K = 2 accuracy = 0.6571428571428571
 For K = 3 accuracy = 0.7142857142857143
 For K = 4 accuracy = 0.6857142857142857
 For K = 5 accuracy = 0.7571428571428571
 For K = 6 accuracy = 0.7142857142857143
 For K = 7 accuracy = 0.7857142857142857
 For K = 8 accuracy = 0.7571428571428571
 For K = 9 accuracy = 0.7571428571428571

In [29]:

```
print("We can see that the KNN model is the best for K=7")
```

We can see that the KNN model is the best for K=7

In [30]:

```
best_knn_model = KNeighborsClassifier(n_neighbors = 7).fit(x_train, y_train)
best_knn_model
```

Out[30]:

KNeighborsClassifier(n_neighbors=7)

In [32]:

```
from sklearn.metrics import f1_score
from sklearn.metrics import jaccard_score

print("Train set Accuracy (Jaccard): ", jaccard_score(y_train, best_knn_model.predict(x_train)))
print("Test set Accuracy (Jaccard): ", jaccard_score(y_test, best_knn_model.predict(x_test)))

print("Train set Accuracy (F1): ", f1_score(y_train, best_knn_model.predict(x_train)))
print("Test set Accuracy (F1): ", f1_score(y_test, best_knn_model.predict(x_test)), accuracy_score(y_test, best_knn_model.predict(x_test)))
```

```
Train set Accuracy (Jaccard): 0.7782426778242678
Test set Accuracy (Jaccard): 0.765625
Train set Accuracy (F1): 0.8000194668761034
Test set Accuracy (F1): 0.7766540244416351
```

Decision Tree

```
In [33]: # importing Libraries
from sklearn.tree import DecisionTreeClassifier
```



```
In [34]: for d in range(1,10):
    dt = DecisionTreeClassifier(criterion = 'entropy', max_depth = d).fit(x_train, y
    dt_yhat = dt.predict(x_test)
    print("For depth = {} the accuracy score is {}".format(d, accuracy_score(y_te
```



```
For depth = 1 the accuracy score is 0.7857142857142857
For depth = 2 the accuracy score is 0.7857142857142857
For depth = 3 the accuracy score is 0.6142857142857143
For depth = 4 the accuracy score is 0.6142857142857143
For depth = 5 the accuracy score is 0.6428571428571429
For depth = 6 the accuracy score is 0.7714285714285715
For depth = 7 the accuracy score is 0.7571428571428571
For depth = 8 the accuracy score is 0.7571428571428571
For depth = 9 the accuracy score is 0.6571428571428571
```



```
In [35]: print("The best value of depth is d = 2 ")
```



```
The best value of depth is d = 2
```



```
In [36]: ## Creating the best model for decision tree with best value of depth 2
best_dt_model = DecisionTreeClassifier(criterion = 'entropy', max_depth = 2).fit(x_t
best_dt_model
```



```
Out[36]: DecisionTreeClassifier(criterion='entropy', max_depth=2)
```



```
In [37]: ## Evaluation Metrics
# f1 score
from sklearn.metrics import f1_score
from sklearn.metrics import jaccard_score

print("Train set Accuracy (Jaccard): ", jaccard_score(y_train, best_dt_model.predict(x
print("Test set Accuracy (Jaccard): ", jaccard_score(y_test, best_dt_model.predict(x
```



```
print("Train set Accuracy (F1): ", f1_score(y_train, best_dt_model.predict(x_train)),
print("Test set Accuracy (F1): ", f1_score(y_test, best_dt_model.predict(x_test), av
```



```
Train set Accuracy (Jaccard): 0.7427536231884058
Test set Accuracy (Jaccard): 0.7857142857142857
Train set Accuracy (F1): 0.6331163939859591
Test set Accuracy (F1): 0.6914285714285714
```

Support Vector Machine

```
In [38]: #importing svm
from sklearn import svm
```

```
from sklearn.metrics import f1_score
```

In [39]:

```
for k in ('linear', 'poly', 'rbf', 'sigmoid'):
    svm_model = svm.SVC(kernel=k).fit(x_train, y_train)
    svm_yhat = svm_model.predict(x_test)
    print("For kernel: {}, the f1 score is: {}".format(k, f1_score(y_test, svm_yhat, average='weighted')))
```

For kernel: linear, the f1 score is: 0.6914285714285714
 For kernel: poly, the f1 score is: 0.7064793130366899
 For kernel: rbf, the f1 score is: 0.7275882012724117
 For kernel: sigmoid, the f1 score is: 0.6892857142857144

In [40]:

```
print("We can see the rbf has the best f1 score of 0.7275882012724117")
```

We can see the rbf has the best f1 score of 0.7275882012724117

In [41]:

```
## building best SVM with kernel = rbf
best_svm = svm.SVC(kernel='rbf').fit(x_train, y_train)
best_svm
```

Out[41]:

SVC()

In [43]:

```
from sklearn.metrics import jaccard_score
```

In [46]:

```
from sklearn.metrics import f1_score

print("Train set Accuracy (Jaccard): ", jaccard_score(y_train, best_svm.predict(x_train)))
print("Test set Accuracy (Jaccard): ", jaccard_score(y_test, best_svm.predict(x_test)))

print("Train set Accuracy (F1): ", f1_score(y_train, best_svm.predict(x_train), average='weighted'))
print("Test set Accuracy (F1): ", f1_score(y_test, best_svm.predict(x_test), average='weighted'))
```

Train set Accuracy (Jaccard): 0.7560975609756098
 Test set Accuracy (Jaccard): 0.7272727272727273
 Train set Accuracy (F1): 0.7682165861513688
 Test set Accuracy (F1): 0.7275882012724117

Logistic Regression

In [47]:

```
# importing Libraries
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import log_loss
```

In [48]:

```
for k in ('lbfgs', 'saga', 'liblinear', 'newton-cg', 'sag'):
    lr_model = LogisticRegression(C=0.01, solver=k).fit(x_train, y_train)
    lr_yhat = lr_model.predict(x_test)
    y_prob = lr_model.predict_proba(x_test)
    print('When Solver is {}, logloss is : {}'.format(k, log_loss(y_test, y_prob)))
```

When Solver is lbfgs, logloss is : 0.4920179847937498
 When Solver is saga, logloss is : 0.49201994093949303
 When Solver is liblinear, logloss is : 0.5772287609479654
 When Solver is newton-cg, logloss is : 0.492017801467927
 When Solver is sag, logloss is : 0.49201107252812804

```
In [49]: print("We can see that the best solver is liblinear of 0.5772287609479654")
```

We can see that the best solver is liblinear of 0.5772287609479654

```
In [50]: # Best Logistic regression model with liblinear solver
```

```
best_lr_model = LogisticRegression(C = 0.01, solver = 'liblinear').fit(x_train, y_train)
best_lr_model
```

```
Out[50]: LogisticRegression(C=0.01, solver='liblinear')
```

```
In [53]: ## Evaluation Metrics
```

```
# jaccard score and f1 score
```

```
from sklearn.metrics import f1_score
```

```
print("Train set Accuracy (Jaccard): ", jaccard_score(y_train, best_lr_model.predict(x_train)))
print("Test set Accuracy (Jaccard): ", jaccard_score(y_test, best_lr_model.predict(x_test)))
```

```
print("Train set Accuracy (F1): ", f1_score(y_train, best_lr_model.predict(x_train)))
print("Test set Accuracy (F1): ", f1_score(y_test, best_lr_model.predict(x_test)), average='macro')
```

Train set Accuracy (Jaccard): 0.7351778656126482

Test set Accuracy (Jaccard): 0.6764705882352942

Train set Accuracy (F1): 0.7341146337750953

Test set Accuracy (F1): 0.6670522459996144

Model Evaluation using Test set

```
In [54]: 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 [ ]: !wget -O loan_test.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-
```

Load Test set for evaluation

```
In [55]: test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

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

In [56]:

```
# data processing
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)

Feature1 = test_df[['Principal','terms','age','Gender','weekend']]
Feature1 = pd.concat([Feature1,pd.get_dummies(test_df['education'])], axis=1)
Feature1.drop(['Master or Above'], axis = 1,inplace=True)

x_loan_test = Feature1
x_loan_test = preprocessing.StandardScaler().fit(x_loan_test).transform(x_loan_test)

y_loan_test = test_df['loan_status'].values
```

In [57]:

```
# Jaccard

# KNN
knn_yhat = best_knn_model.predict(x_loan_test)
jacc1 = round(jaccard_score(y_loan_test, knn_yhat, pos_label = "PAIDOFF"), 2)

# Decision Tree
dt_yhat = best_dt_model.predict(x_loan_test)
jacc2 = round(jaccard_score(y_loan_test, dt_yhat, pos_label = "PAIDOFF"), 2)

# Support Vector Machine
svm_yhat = best_svm.predict(x_loan_test)
jacc3 = round(jaccard_score(y_loan_test, svm_yhat, pos_label = "PAIDOFF"), 2)

# Logistic Regression
lr_yhat = best_lr_model.predict(x_loan_test)
jacc4 = round(jaccard_score(y_loan_test, lr_yhat, pos_label = "PAIDOFF"), 2)

jss = [jacc1, jacc2, jacc3, jacc4]
jss
```

Out[57]: [0.65, 0.74, 0.78, 0.74]

In [58]:

```
# F1_score

# KNN
knn_yhat = best_knn_model.predict(x_loan_test)
f1 = round(f1_score(y_loan_test, knn_yhat, average = 'weighted'), 2)

# Decision Tree
dt_yhat = best_dt_model.predict(x_loan_test)
f2 = round(f1_score(y_loan_test, dt_yhat, average = 'weighted'), 2)

# Support Vector Machine
svm_yhat = best_svm.predict(x_loan_test)
f3 = round(f1_score(y_loan_test, svm_yhat, average = 'weighted'), 2)

# Logistic Regression
lr_yhat = best_lr_model.predict(x_loan_test)
f4 = round(f1_score(y_loan_test, lr_yhat, average = 'weighted'), 2)
```

```
f1_list = [f1, f2, f3, f4]
f1_list
```

Out[58]: [0.63, 0.63, 0.76, 0.66]

In [59]:

```
# Log Loss

# Logistic Regression
lr_prob = best_lr_model.predict_proba(x_loan_test)
ll_list = ['NA', 'NA', 'NA', round(log_loss(y_loan_test, lr_prob), 2)]
ll_list
```

Out[59]: ['NA', 'NA', 'NA', 0.57]

In [60]:

```
columns = ['KNN', 'Decision Tree', 'SVM', 'Logistic Regression']
index = ['Jaccard', 'F1-score', 'Logloss']

accuracy_df = pd.DataFrame([jss, f1_list, ll_list], index = index, columns = columns)
accuracy_df1 = accuracy_df.transpose()
accuracy_df1.columns.name = 'Algorithm'
accuracy_df1
```

Out[60]:

Algorithm	Jaccard	F1-score	Logloss
KNN	0.65	0.63	NA
Decision Tree	0.74	0.63	NA
SVM	0.78	0.76	NA
Logistic Regression	0.74	0.66	0.57

Report

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

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.