

ML0101EN-Clas-SVM-cancer-py-v1

February 28, 2019

SVM (Support Vector Machines)

In this notebook, you will use SVM (Support Vector Machines) to build and train a model using human cell records, and classify cells to whether the samples are benign or malignant.

SVM works by mapping data to a high-dimensional feature space so that data points can be categorized, even when the data are not otherwise linearly separable. A separator between the categories is found, then the data is transformed in such a way that the separator could be drawn as a hyperplane. Following this, characteristics of new data can be used to predict the group to which a new record should belong.

Table of contents

```
<ol>
  <li><a href="#load_dataset">Load the Cancer data</a></li>
  <li><a href="#modeling">Modeling</a></li>
  <li><a href="#evaluation">Evaluation</a></li>
  <li><a href="#practice">Practice</a></li>
</ol>
```

```
In [1]: import pandas as pd
import pylab as pl
import numpy as np
import scipy.optimize as opt
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
%matplotlib inline
import matplotlib.pyplot as plt
```

Load the Cancer data

The example is based on a dataset that is publicly available from the UCI Machine Learning Repository (Asuncion and Newman, 2007)[<http://mllearn.ics.uci.edu/MLRepository.html>]. The dataset consists of several hundred human cell sample records, each of which contains the values of a set of cell characteristics. The fields in each record are:

Field name	Description
ID	Clump thickness
Clump	Clump thickness
UnifSize	Uniformity of cell size
UnifShape	Uniformity of cell shape

Field name	Description
MargAdh	Marginal adhesion
SingEpiSize	Single epithelial cell size
BareNuc	Bare nuclei
BlandChrom	Bland chromatin
NormNucl	Normal nucleoli
Mit	Mitoses
Class	Benign or malignant

For the purposes of this example, we're using a dataset that has a relatively small number of predictors in each record. To download the data, we will use `!wget` to download it from IBM Object Storage.

Did you know? When it comes to Machine Learning, you will likely be working with large datasets. As a business, where can you host your data? IBM is offering a unique opportunity for businesses, with 10 Tb of IBM Cloud Object Storage: [Sign up now for free](#)

In [2]: *#Click here and press Shift+Enter*

```
!wget -O cell_samples.csv https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-d
--2019-02-28 23:11:03-- https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/Cogni
Resolving s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorage.softlayer.net).
Connecting to s3-api.us-geo.objectstorage.softlayer.net (s3-api.us-geo.objectstorage.softlayer.net).
HTTP request sent, awaiting response... 200 OK
Length: 20675 (20K) [text/csv]
Saving to: cell_samples.csv

cell_samples.csv  100%[=====>]  20.19K  --.-KB/s   in 0.02s

2019-02-28 23:11:03 (976 KB/s) - cell_samples.csv saved [20675/20675]
```

0.0.1 Load Data From CSV File

```
In [3]: cell_df = pd.read_csv("cell_samples.csv")
        cell_df.head()
```

```
Out[3]:
```

	ID	Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	\
0	1000025	5	1	1	1	2	1	
1	1002945	5	4	4	5	7	10	
2	1015425	3	1	1	1	2	2	
3	1016277	6	8	8	1	3	4	
4	1017023	4	1	1	3	2	1	

	BlandChrom	NormNucl	Mit	Class
0	3	1	1	2
1	3	2	1	2
2	3	1	1	2

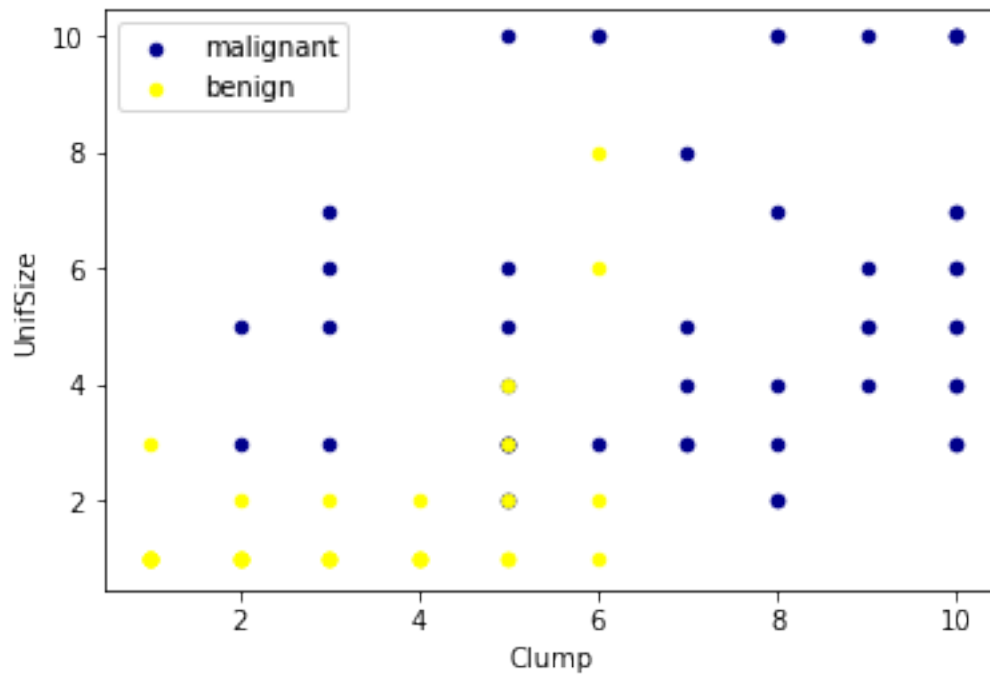
3	3	7	1	2
4	3	1	1	2

The ID field contains the patient identifiers. The characteristics of the cell samples from each patient are contained in fields Clump to Mit. The values are graded from 1 to 10, with 1 being the closest to benign.

The Class field contains the diagnosis, as confirmed by separate medical procedures, as to whether the samples are benign (value = 2) or malignant (value = 4).

Lets look at the distribution of the classes based on Clump thickness and Uniformity of cell size:

```
In [4]: ax = cell_df[cell_df['Class'] == 4][0:50].plot(kind='scatter', x='Clump', y='UnifSize',
cell_df[cell_df['Class'] == 2][0:50].plot(kind='scatter', x='Clump', y='UnifSize', color
plt.show()
```



0.1 Data pre-processing and selection

Lets first look at columns data types:

```
In [5]: cell_df.dtypes
```

```
Out[5]: ID          int64
        Clump       int64
        UnifSize    int64
        UnifShape    int64
```

```

MargAdh      int64
SingEpiSize  int64
BareNuc       object
BlandChrom    int64
NormNucl     int64
Mit          int64
Class        int64
dtype: object

```

It looks like the **BareNuc** column includes some values that are not numerical. We can drop those rows:

```

In [6]: cell_df = cell_df[pd.to_numeric(cell_df['BareNuc'], errors='coerce').notnull()]
        cell_df['BareNuc'] = cell_df['BareNuc'].astype('int')
        cell_df.dtypes

```

```

Out[6]: ID      int64
        Clump    int64
        UnifSize int64
        UnifShape int64
        MargAdh  int64
        SingEpiSize int64
        BareNuc  int64
        BlandChrom int64
        NormNucl int64
        Mit      int64
        Class    int64
        dtype: object

```

```

In [7]: feature_df = cell_df[['Clump', 'UnifSize', 'UnifShape', 'MargAdh', 'SingEpiSize', 'BareNuc']]
        X = np.asarray(feature_df)
        X[0:5]

```

```

Out[7]: array([[ 5,  1,  1,  1,  2,  1,  3,  1,  1],
               [ 5,  4,  4,  5,  7, 10,  3,  2,  1],
               [ 3,  1,  1,  1,  2,  2,  3,  1,  1],
               [ 6,  8,  8,  1,  3,  4,  3,  7,  1],
               [ 4,  1,  1,  3,  2,  1,  3,  1,  1]])

```

We want the model to predict the value of Class (that is, benign (=2) or malignant (=4)). As this field can have one of only two possible values, we need to change its measurement level to reflect this.

```

In [8]: cell_df['Class'] = cell_df['Class'].astype('int')
        y = np.asarray(cell_df['Class'])
        y [0:5]

```

```

Out[8]: array([2, 2, 2, 2, 2])

```

0.2 Train/Test dataset

Okay, we split our dataset into train and test set:

```
In [9]: 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: (546, 9) (546,)
```

```
Test set: (137, 9) (137,)
```

Modeling (SVM with Scikit-learn)

The SVM algorithm offers a choice of kernel functions for performing its processing. Basically, mapping data into a higher dimensional space is called kernelling. The mathematical function used for the transformation is known as the kernel function, and can be of different types, such as:

- 1.Linear
- 2.Polynomial
- 3.Radial basis function (RBF)
- 4.Sigmoid

Each of these functions has its characteristics, its pros and cons, and its equation, but as there's no easy way of knowing which function performs best with any given dataset, we usually choose different functions in turn and compare the results. Let's just use the default, RBF (Radial Basis Function) for this lab.

```
In [10]: from sklearn import svm
        clf = svm.SVC(kernel='rbf')
        clf.fit(X_train, y_train)
```

```
/home/jupyterlab/conda/lib/python3.6/site-packages/sklearn/svm/base.py:196: FutureWarning: The d
"avoid this warning.", FutureWarning)
```

```
Out[10]: 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)
```

After being fitted, the model can then be used to predict new values:

```
In [11]: yhat = clf.predict(X_test)
        yhat [0:5]
```

```
Out[11]: array([2, 4, 2, 4, 2])
```

Evaluation

```
In [16]: from sklearn.metrics import classification_report, confusion_matrix
        import itertools
```

```

In [17]: def plot_confusion_matrix(cm, classes,
                                   normalize=False,
                                   title='Confusion matrix',
                                   cmap=plt.cm.Blues):
    """
    This function prints and plots the confusion matrix.
    Normalization can be applied by setting `normalize=True`.
    """
    if normalize:
        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
        print("Normalized confusion matrix")
    else:
        print('Confusion matrix, without normalization')

    print(cm)

    plt.imshow(cm, interpolation='nearest', cmap=cmap)
    plt.title(title)
    plt.colorbar()
    tick_marks = np.arange(len(classes))
    plt.xticks(tick_marks, classes, rotation=45)
    plt.yticks(tick_marks, classes)

    fmt = '.2f' if normalize else 'd'
    thresh = cm.max() / 2.
    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
        plt.text(j, i, format(cm[i, j], fmt),
                 horizontalalignment="center",
                 color="white" if cm[i, j] > thresh else "black")

    plt.tight_layout()
    plt.ylabel('True label')
    plt.xlabel('Predicted label')

In [18]: # Compute confusion matrix
cnf_matrix = confusion_matrix(y_test, yhat, labels=[2,4])
np.set_printoptions(precision=2)

print (classification_report(y_test, yhat))

# Plot non-normalized confusion matrix
plt.figure()
plot_confusion_matrix(cnf_matrix, classes=['Benign(2)', 'Malignant(4)'], normalize= False

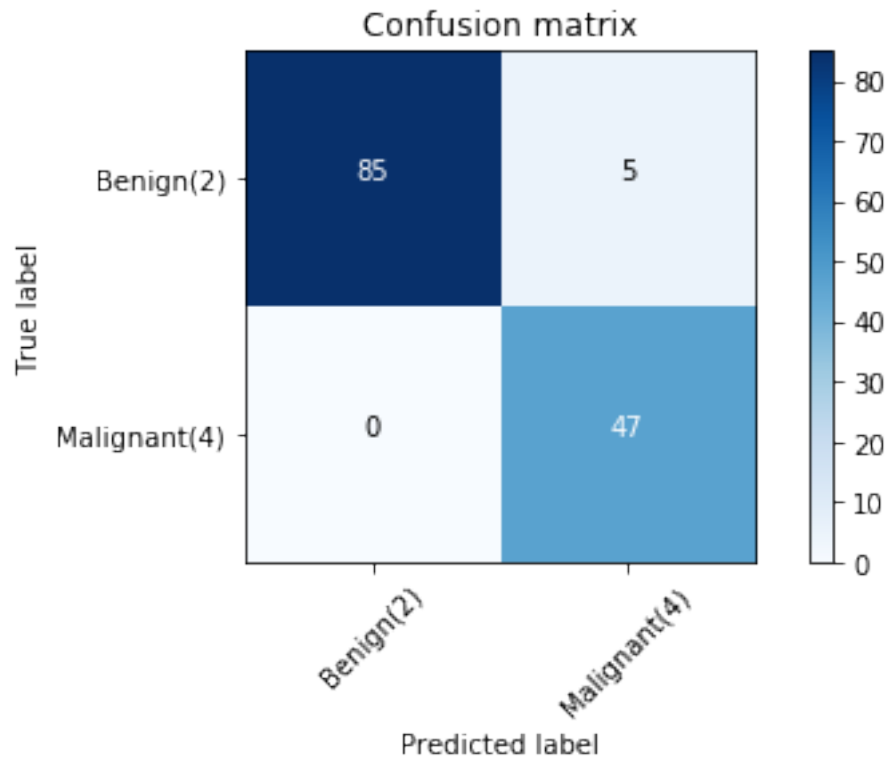
```

	precision	recall	f1-score	support
2	1.00	0.94	0.97	90
4	0.90	1.00	0.95	47

micro avg	0.96	0.96	0.96	137
macro avg	0.95	0.97	0.96	137
weighted avg	0.97	0.96	0.96	137

Confusion matrix, without normalization

```
[[85  5]
 [ 0 47]]
```



You can also easily use the **f1_score** from sklearn library:

```
In [19]: from sklearn.metrics import f1_score
         f1_score(y_test, yhat, average='weighted')
```

```
Out[19]: 0.9639038982104676
```

Lets try jaccard index for accuracy:

```
In [20]: from sklearn.metrics import jaccard_similarity_score
         jaccard_similarity_score(y_test, yhat)
```

```
Out[20]: 0.9635036496350365
```

Practice

Can you rebuild the model, but this time with a **linear** kernel? You can use **kernel='linear'** option, when you define the svm. How the accuracy changes with the new kernel function?

```
In [26]: # write your code here
from sklearn import svm
clf2 = svm.SVC(kernel='linear')
clf2.fit(X_train, y_train)
yhat2 = clf2.predict(X_test)
yhat2 [0:5]

f1_score(y_test, yhat2, average='weighted')
jaccard_similarity_score(y_test, yhat2)

clf2 = svm.SVC(kernel='linear')
clf2.fit(X_train, y_train)
yhat2 = clf2.predict(X_test)
print("Avg F1-score: %.4f" % f1_score(y_test, yhat2, average='weighted'))
print("Jaccard score: %.4f" % jaccard_similarity_score(y_test, yhat2))
```

Avg F1-score: 0.9639

Jaccard score: 0.9635

Double-click **here** for the solution.

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.

Copyright I 2018 Cognitive Class. This notebook and its source code are released under the terms of the MIT License.