

Support Vector Machines with Python

Welcome to the Support Vector Machines with Python Lecture Notebook! Remember to refer to the video lecture for the full background information on the code here!

Import Libraries

```
In [51]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

Get the Data

We'll use the built in breast cancer dataset from Scikit Learn. We can get with the load function:

```
In [52]: from sklearn.datasets import load_breast_cancer
```

```
In [54]: cancer = load_breast_cancer()
```

The data set is presented in a dictionary form:

```
In [55]: cancer.keys()
```

```
Out[55]: dict_keys(['DESCR', 'target', 'data', 'target_names', 'feature_names'])
```

We can grab information and arrays out of this dictionary to set up our data frame and understanding of the features:

```
In [4]: print(cancer['DESCR'])
```

Separating plane described above was obtained using Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, pp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes.

The actual linear program used to obtain the separating plane in the 3-dimensional space is that described in: [K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

```
ftp ftp.cs.wisc.edu
cd math-prog/cpo-dataset/machine-learn/WDBC/
```

```
In [56]: cancer['feature_names']
```

```
Out[56]: array(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
               'mean smoothness', 'mean compactness', 'mean concavity',
               'mean concave points', 'mean symmetry', 'mean fractal dimension',
               'radius error', 'texture error', 'perimeter error', 'area error',
               'smoothness error', 'compactness error', 'concavity error',
               'concave points error', 'symmetry error', 'fractal dimension error',
               'worst radius', 'worst texture', 'worst perimeter', 'worst area',
               'worst smoothness', 'worst compactness', 'worst concavity',
               'worst concave points', 'worst symmetry', 'worst fractal dimension'],
              dtype='<U23')
```

Set up DataFrame

```
In [12]: df_feat = pd.DataFrame(cancer['data'], columns=cancer['feature_names'])
df_feat.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 30 columns):
mean radius           569 non-null float64
mean texture          569 non-null float64
mean perimeter        569 non-null float64
mean area             569 non-null float64
mean smoothness       569 non-null float64
mean compactness      569 non-null float64
mean concavity        569 non-null float64
mean concave points   569 non-null float64
mean symmetry         569 non-null float64
mean fractal dimension 569 non-null float64
radius error          569 non-null float64
texture error         569 non-null float64
perimeter error       569 non-null float64
area error            569 non-null float64
smoothness error      569 non-null float64
compactness error     569 non-null float64
concavity error       569 non-null float64
concave points error  569 non-null float64
symmetry error        569 non-null float64
fractal dimension error 569 non-null float64
worst radius          569 non-null float64
worst texture         569 non-null float64
worst perimeter       569 non-null float64
worst area            569 non-null float64
worst smoothness      569 non-null float64
worst compactness     569 non-null float64
worst concavity       569 non-null float64
worst concave points  569 non-null float64
worst symmetry        569 non-null float64
worst fractal dimension 569 non-null float64
dtypes: float64(30)
memory usage: 133.4 KB
```

```
In [14]: cancer['target']
```

```
Out[14]: array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1,
1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0,
1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1,
1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1,
0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 1,
0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1,
0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 0,
0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1,
1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1,
1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0,
1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1,
0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1,
1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1,
1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1,
0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1,
1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1])
```

```
In [16]: df_target = pd.DataFrame(cancer['target'], columns=['Cancer'])
```

Now let's actually check out the dataframe!

```
In [8]: df.head()
```

```
Out[8]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius	worst texture	worst perimeter	worst area	
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	25.38	17.33	184.60	2019.0	
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	24.99	23.41	158.80	1956.0	
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	23.57	25.53	152.50	1709.0	
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	14.91	26.50	98.87	567.7	
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	22.54	16.67	152.20	1575.0	

5 rows × 30 columns

Exploratory Data Analysis

We'll skip the Data Viz part for this lecture since there are so many features that are hard to interpret if you don't have domain knowledge of cancer or tumor cells. In your project you will have more to visualize for the data.

Train Test Split

```
In [57]: from sklearn.model_selection import train_test_split
```

```
In [58]: X_train, X_test, y_train, y_test = train_test_split(df_feat, np.ravel(df_target), test_size=0.30, random_state=1)
```

Train the Support Vector Classifier

```
In [59]: from sklearn.svm import SVC
```

```
In [60]: model = SVC()
```

```
In [61]: model.fit(X_train,y_train)
```

```
Out[61]: SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,  
            decision_function_shape=None, degree=3, gamma='auto', kernel='rbf',  
            max_iter=-1, probability=False, random_state=None, shrinking=True,  
            tol=0.001, verbose=False)
```

Predictions and Evaluations

Now let's predict using the trained model.

```
In [27]: predictions = model.predict(X_test)
```

```
In [45]: from sklearn.metrics import classification_report, confusion_matrix
```

```
In [46]: print(confusion_matrix(y_test, predictions))
```

```
[[ 0 66]  
 [ 0 105]]
```

```
In [62]: print(classification_report(y_test, predictions))
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	66
1	0.61	1.00	0.76	105
avg / total	0.38	0.61	0.47	171

```
/Users/marci/anaconda/lib/python3.5/site-packages/sklearn/metrics/classification.py:1074: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples.
'precision', 'predicted', average, warn_for)
```

Woah! Notice that we are classifying everything into a single class! This means our model needs to have its parameters adjusted (it may also help to normalize the data).

We can search for parameters using a GridSearch!

Gridsearch

Finding the right parameters (like what C or gamma values to use) is a tricky task! But luckily, we can be a little lazy and just try a bunch of combinations and see what works best! This idea of creating a 'grid' of parameters and just trying out all the possible combinations is called a Gridsearch, this method is common enough that Scikit-learn has this functionality built in with GridSearchCV! The CV stands for cross-validation which is the

GridSearchCV takes a dictionary that describes the parameters that should be tried and a model to train. The grid of parameters is defined as a dictionary, where the keys are the parameters and the values are the settings to be tested.

```
In [63]: param_grid = {'C': [0.1, 1, 10, 100, 1000], 'gamma': [1, 0.1, 0.01, 0.001, 0.0001], 'kernel': ['rbf']}
```

```
In [64]: from sklearn.model_selection import GridSearchCV
```

One of the great things about GridSearchCV is that it is a meta-estimator. It takes an estimator like SVC, and creates a new estimator, that behaves exactly the same - in this case, like a classifier. You should add refit=True and choose verbose to whatever number you want, higher the number, the more verbose (verbose just means the text output describing the process).

```
In [65]: grid = GridSearchCV(SVC(),param_grid,refit=True,verbose=3)
```

What fit does is a bit more involved than usual. First, it runs the same loop with cross-validation, to find the best parameter combination. Once it has the best combination, it runs fit again on all data passed to fit (without cross-validation), to build a single new model using the best parameter setting.

```
In [40]: # May take awhile!
grid.fit(X_train,y_train)

[CV] ..... gamma=0.01, C=100, kernel=rbf, score=0.636364 - 0.0s
[CV] gamma=0.001, C=100, kernel=rbf .....
[CV] ..... gamma=0.001, C=100, kernel=rbf, score=0.894737 - 0.0s
[CV] gamma=0.001, C=100, kernel=rbf .....
[CV] ..... gamma=0.001, C=100, kernel=rbf, score=0.932331 - 0.0s
[CV] gamma=0.001, C=100, kernel=rbf .....
[CV] ..... gamma=0.001, C=100, kernel=rbf, score=0.916667 - 0.0s
[CV] gamma=0.0001, C=100, kernel=rbf .....
[CV] ..... gamma=0.0001, C=100, kernel=rbf, score=0.917293 - 0.0s
[CV] gamma=0.0001, C=100, kernel=rbf .....
[CV] ..... gamma=0.0001, C=100, kernel=rbf, score=0.977444 - 0.0s
[CV] gamma=0.0001, C=100, kernel=rbf .....
[CV] ..... gamma=0.0001, C=100, kernel=rbf, score=0.939394 - 0.0s
[CV] gamma=1, C=1000, kernel=rbf .....
[CV] ..... gamma=1, C=1000, kernel=rbf, score=0.631579 - 0.0s
[CV] gamma=1, C=1000, kernel=rbf .....
[CV] ..... gamma=1, C=1000, kernel=rbf, score=0.631579 - 0.0s
[CV] gamma=1, C=1000, kernel=rbf .....
[CV] ..... gamma=1, C=1000, kernel=rbf, score=0.636364 - 0.0s
-----
```

You can inspect the best parameters found by GridSearchCV in the `best_params_` attribute, and the best estimator in the `best_estimator_` attribute:

```
In [41]: grid.best_params_
```

```
Out[41]: {'C': 10, 'gamma': 0.0001, 'kernel': 'rbf'}
```

```
In [ ]: grid.best_estimator_
```

Then you can re-run predictions on this grid object just like you would with a normal model.


```
In [48]: grid_predictions = grid.predict(X_test)
```

```
In [49]: print(confusion_matrix(y_test,grid_predictions))
```

```
[[ 60   6]
 [  3 102]]
```

```
In [50]: print(classification_report(y_test,grid_predictions))
```

	precision	recall	f1-score	support
0	0.95	0.91	0.93	66
1	0.94	0.97	0.96	105
avg / total	0.95	0.95	0.95	171

Great job!