

By: Dishan Fapot

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 [1]: 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 [2]: from sklearn.datasets import load_breast_cancer
```

```
In [3]: cancer = load_breast_cancer()
```

The data set is presented in a dictionary form:

```
In [4]: cancer.keys()
```

```
Out[4]: dict_keys(['data', 'target', 'target_names', 'DESCR', 'feature_names', 'filename'])
```

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'])`

```

The mean, standard error, and "worst" or largest (mean of the three
largest values) of these features were computed for each image,
resulting in 30 features. For instance, field 3 is Mean Radius, field
d
13 is Radius SE, field 23 is Worst Radius.

- class:
  - WDBC-Malignant
  - WDBC-Benign

:Summary Statistics:

=====
                                Min      Max
=====
radius (mean):                  6.981    28.11
texture (mean):                 9.71     39.28
perimeter (mean):              43.79    188.5
area (mean):                   143.5    2501.0
smoothness (mean):             0.053    0.163

```

In [5]: `cancer['feature_names']`

```

Out[5]: 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 [6]: 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):
#   Column                                Non-Null Count  Dtype
---  -
0   mean radius                           569 non-null    float64
1   mean texture                           569 non-null    float64
2   mean perimeter                         569 non-null    float64
3   mean area                             569 non-null    float64
4   mean smoothness                       569 non-null    float64
5   mean compactness                      569 non-null    float64
6   mean concavity                         569 non-null    float64
7   mean concave points                   569 non-null    float64
8   mean symmetry                         569 non-null    float64
9   mean fractal dimension                 569 non-null    float64
10  radius error                           569 non-null    float64
11  texture error                          569 non-null    float64
12  perimeter error                       569 non-null    float64
13  area error                            569 non-null    float64
14  smoothness error                      569 non-null    float64
15  compactness error                     569 non-null    float64
16  concavity error                       569 non-null    float64
17  concave points error                  569 non-null    float64
18  symmetry error                        569 non-null    float64
19  fractal dimension error                569 non-null    float64
20  worst radius                          569 non-null    float64
21  worst texture                         569 non-null    float64
22  worst perimeter                       569 non-null    float64
23  worst area                            569 non-null    float64
24  worst smoothness                      569 non-null    float64
25  worst compactness                     569 non-null    float64
26  worst concavity                       569 non-null    float64
27  worst concave points                  569 non-null    float64
28  worst symmetry                        569 non-null    float64
29  worst fractal dimension                569 non-null    float64
dtypes: float64(30)
memory usage: 133.5 KB
```

```
In [7]: cancer['target']
```

```
Out[7]: array([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, 1, 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, 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, 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, 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, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
               1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
               0, 0, 1, 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, 1, 0, 1, 1, 0, 1, 0, 1, 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, 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, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 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, 0, 1, 1, 0, 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, 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, 0, 0, 0, 0, 0, 0, 1])
```

```
In [8]: df_target = pd.DataFrame(cancer['target'], columns=['Cancer'])
```

Now let's actually check out the dataframe!

```
In [10]: df_target.head()
```

```
Out[10]:
```

	Cancer
0	0
1	0
2	0
3	0
4	0

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 [11]: from sklearn.model_selection import train_test_split
```

```
In [12]: X_train, X_test, y_train, y_test = train_test_split(df_feat, np.ravel(df_target),
```

Train the Support Vector Classifier

```
In [13]: from sklearn.svm import SVC
```

```
In [14]: model = SVC()
```

```
In [15]: model.fit(X_train,y_train)
```

```
Out[15]: SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
            decision_function_shape='ovr', degree=3, gamma='scale', 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 [16]: predictions = model.predict(X_test)
```

```
In [17]: from sklearn.metrics import classification_report,confusion_matrix
```

```
In [18]: print(confusion_matrix(y_test,predictions))
```

```
[[ 56  10]
 [  3 102]]
```

```
In [19]: print(classification_report(y_test,predictions))
```

	precision	recall	f1-score	support
0	0.95	0.85	0.90	66
1	0.91	0.97	0.94	105
accuracy			0.92	171
macro avg	0.93	0.91	0.92	171
weighted avg	0.93	0.92	0.92	171

We can search for and use the best parameters (C and gamma) 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 [31]: param_grid = {'C': [0.1, 1, 10, 100, 1000], 'gamma': [1, 0.1, 0.01, 0.001, 0.0001], 'kernel': ['rbf']}
```

```
In [32]: 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 [33]: 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 [34]: # May take awhile!
grid.fit(X_train, y_train)

[CV] .... C=1000, gamma=0.0001, kernel=rbf, score=0.924, total= 0.0s
[CV] C=1000, gamma=0.0001, kernel=rbf .....
[CV] .... C=1000, gamma=0.0001, kernel=rbf, score=0.962, total= 0.0s

[Parallel(n_jobs=1)]: Done 125 out of 125 | elapsed: 6.9s finished
```

```
Out[34]: GridSearchCV(cv=None, error_score=nan,
                    estimator=SVC(C=1.0, break_ties=False, cache_size=200,
                                class_weight=None, coef0=0.0,
                                decision_function_shape='ovr', degree=3,
                                gamma='scale', kernel='rbf', max_iter=-1,
                                probability=False, random_state=None, shrinking=True,
                                tol=0.001, verbose=False),
                    iid='deprecated', n_jobs=None,
                    param_grid={'C': [0.1, 1, 10, 100, 1000],
                                'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
                                'kernel': ['rbf']},
                    pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                    scoring=None, verbose=3)
```

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

```
In [35]: grid.best_params_
```

```
Out[35]: {'C': 1, 'gamma': 0.0001, 'kernel': 'rbf'}
```

```
In [36]: grid.best_estimator_
```

```
Out[36]: SVC(C=1, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma=0.0001, kernel='rbf',
max_iter=-1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False)
```

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

```
In [37]: grid_predictions = grid.predict(X_test)
```

```
In [38]: print(confusion_matrix(y_test,grid_predictions))
```

```
[[ 59   7]
 [  4 101]]
```

```
In [39]: print(classification_report(y_test,grid_predictions))
```

	precision	recall	f1-score	support
0	0.94	0.89	0.91	66
1	0.94	0.96	0.95	105
accuracy			0.94	171
macro avg	0.94	0.93	0.93	171
weighted avg	0.94	0.94	0.94	171

Great job!