

# Support Vector Machines with Python

## Import Libraries

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
import warnings; warnings.simplefilter('ignore')
```

## Get the Data

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

```
In [3]: from sklearn.datasets import load_breast_cancer
```

```
In [4]: cancer = load_breast_cancer()
```

The data set is presented in a dictionary form:

```
In [5]: cancer.keys()
```

```
Out[5]: 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 [6]: print(cancer['DESCR'])

.. _breast_cancer_dataset:

Breast cancer wisconsin (diagnostic) dataset
-----

**Data Set Characteristics:**

: Number of Instances: 569

: Number of Attributes: 30 numeric, predictive attributes and the class

: Attribute Information:
  - radius (mean of distances from center to points on the perimeter)
  - texture (standard deviation of gray-scale values)
  - perimeter
  - area
  - smoothness (local variation in radius lengths)
  - compactness (perimeter^2 / area - 1.0)
  - concavity (severity of concave portions of the contour)
```

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

Out[7]: 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 [8]: 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 [9]: cancer['target']
```

```
Out[9]: 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, 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, 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, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0,
              0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0,
              0, 0, 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, 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, 1, 0, 1, 1, 0,
              1, 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, 1, 0, 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, 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 [10]: df_target = pd.DataFrame(cancer['target'], columns=['Cancer'])
```

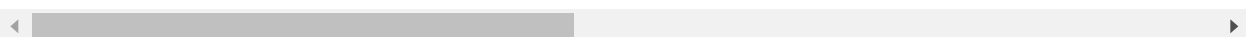
Now let's actually check out the dataframe!

```
In [11]: df_feat.head()
```

```
Out[11]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	d
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	

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

```
In [13]: X_train, X_test, y_train, y_test = train_test_split(df_feat,
np.ravel(df_target), test_size=0.30, random_state=101)
```

## Train the Support Vector Classifier

```
In [14]: from sklearn.svm import SVC
```

```
In [15]: model = SVC()
```

```
In [16]: model.fit(X_train,y_train)
```

```
Out[16]: 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)
```

## Predictions and Evaluations

Now let's predict using the trained model.

```
In [17]: predictions = model.predict(X_test)
```

```
In [18]: from sklearn.metrics import classification_report, confusion_matrix
```

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

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

```
In [20]: 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
micro avg	0.61	0.61	0.61	171
macro avg	0.31	0.50	0.38	171
weighted avg	0.38	0.61	0.47	171

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

```
In [22]: param_grid = [{'C': [0.1,1, 10, 100, 1000], 'gamma': [1,0.1,0.01,0.001,0.0001],
           'kernel': ['rbf']},{'C': [0.1,1, 10, 100, 1000], 'kernel': ['linear']}]
```

```
In [23]: 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 [24]: 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 [25]: # May take awhile!
grid.fit(X_train,y_train)
```

```
Fitting 3 folds for each of 30 candidates, totalling 90 fits
[CV] C=0.1, gamma=1, kernel=rbf .....
[CV] C=0.1, gamma=1, kernel=rbf, score=0.631578947368421, total= 0.0s
[CV] C=0.1, gamma=1, kernel=rbf .....
[CV] C=0.1, gamma=1, kernel=rbf, score=0.631578947368421, total= 0.0s
[CV] C=0.1, gamma=1, kernel=rbf .....
[CV] C=0.1, gamma=1, kernel=rbf, score=0.6363636363636364, total= 0.0s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] C=0.1, gamma=0.1, kernel=rbf, score=0.631578947368421, total= 0.0s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] C=0.1, gamma=0.1, kernel=rbf, score=0.631578947368421, total= 0.0s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] C=0.1, gamma=0.1, kernel=rbf, score=0.6363636363636364, total= 0.0s
[CV] C=0.1, gamma=0.01, kernel=rbf .....
[CV] C=0.1, gamma=0.01, kernel=rbf, score=0.631578947368421, total= 0.0s

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent worke
rs.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.0s remaining: 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 [26]: grid.best_params_
```

```
Out[26]: {'C': 100, 'kernel': 'linear'}
```

```
In [27]: grid.best_estimator_
```

```
Out[27]: SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
kernel='linear', 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 [28]: grid_predictions = grid.predict(X_test)
```

```
In [29]: print(confusion_matrix(y_test,grid_predictions))
```

```
[[ 60  6]
 [ 2 103]]
```

```
In [30]: print(classification_report(y_test,grid_predictions))
```

	precision	recall	f1-score	support
0	0.97	0.91	0.94	66
1	0.94	0.98	0.96	105
micro avg	0.95	0.95	0.95	171
macro avg	0.96	0.95	0.95	171
weighted avg	0.95	0.95	0.95	171

## Great job!

```
In [ ]:
```