

You should have finished watching a video on support vector machines. This algorithm offers very high accuracy compared to other classifiers. It is used in diverse applications such as face detection, intrusion detection, classification of emails, and handwriting recognition. This classifier works by separating data points using a hyperplane with the largest amount of margin.

Let's take a look at using SVM to analyze if we can use specific data measurements to improve the diagnosis of breast cancer. We have a dataset that includes tumors, malignant (cancerous) or benign (non-cancerous), and features obtained from several cell images. Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. The dataset is hosted [here](http://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+%28diagnostic%29) (<http://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+%28diagnostic%29>) but I have put it also in our Github repository.

Briefly:

Ten real-valued features are computed for each cell nucleus:

a) radius (mean of distances from center to points on the perimeter) b) texture (standard deviation of gray-scale values) c) perimeter d) area e) smoothness (local variation in radius lengths) f) compactness (perimeter<sup>2</sup> / area - 1.0) g) concavity (severity of concave portions of the contour) h) concave points (number of concave portions of the contour) i) symmetry j) fractal dimension ("coastline approximation" - 1)

Additionally, 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.

```
In [1]: import pandas as pd
dataset = pd.read_csv('../data/cancer.csv')
dataset.head()
```

Out[1]:

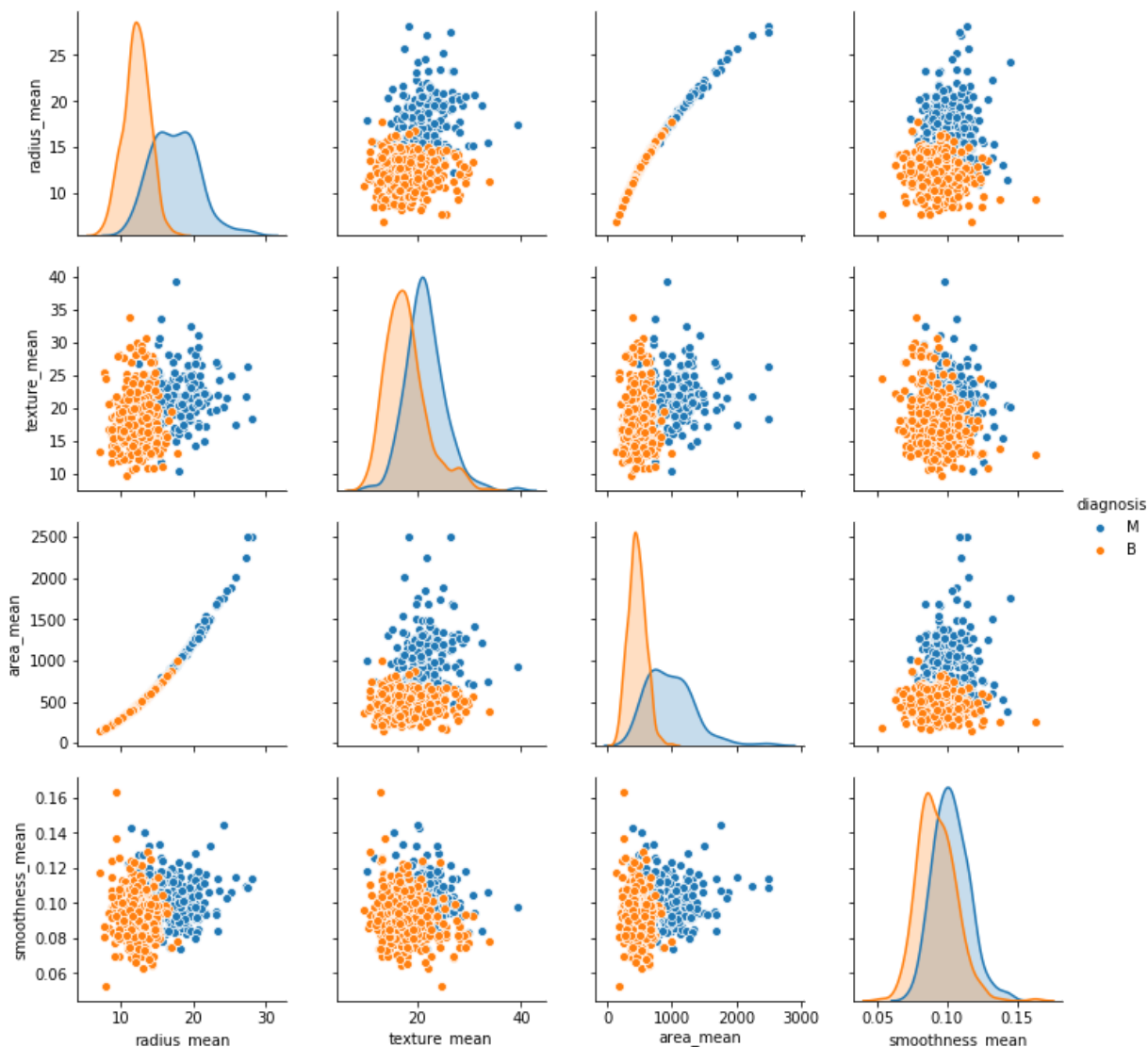
	ID	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactn
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	

5 rows × 32 columns



```
In [3]: import seaborn as sns
sns.pairplot(dataset, hue = 'diagnosis', vars = ['radius_mean', 'texture_mean', 'area_mean', 'smoothness_mean'])
```

Out[3]: <seaborn.axisgrid.PairGrid at 0x13745158048>



To prepare the dataset for SVM algorithms, I'd suggest doing three things:

1. Check for NA values. Previously, I've pretty much ignored this step but for your own datasets you may want to do so to decide how to deal with them.
2. Decide what data goes into your X and y variables.
3. Ensure that the label categories are numeric.
4. Create training and test datasets.

```
In [4]: dataset.isnull().sum()  
dataset.isna().sum()
```

```
Out[4]: ID                0  
diagnosis                0  
radius_mean              0  
texture_mean             0  
perimeter_mean           0  
area_mean                0  
smoothness_mean          0  
compactness_mean         0  
conc_mean                0  
conc_points_mean         0  
symmetry_mean            0  
fractal_mean             0  
radius_se                0  
texture_se               0  
perimeter_se             0  
area_se                  0  
smoothness_se            0  
compactness_se           0  
conc_se                  0  
conc_points_se           0  
symmetry_se              0  
fractal_se               0  
radius_worst             0  
texture_worst            0  
perimeter_worst          0  
area_worst               0  
smoothness_worst         0  
compactness_worst        0  
conc_worst               0  
conc_points_worst        0  
symmetry_worst           0  
fractal_worst            0  
dtype: int64
```

```
In [5]: X = dataset.iloc[:, 2:32].values  
X
```

```
Out[5]: array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01,  
                1.189e-01],  
               [2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01,  
                8.902e-02],  
               [1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01,  
                8.758e-02],  
               ...,  
               [1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01,  
                7.820e-02],  
               [2.060e+01, 2.933e+01, 1.401e+02, ..., 2.650e-01, 4.087e-01,  
                1.240e-01],  
               [7.760e+00, 2.454e+01, 4.792e+01, ..., 0.000e+00, 2.871e-01,  
                7.039e-02]])
```

```
In [6]: Y = dataset.iloc[:,1].values
Y
```

```
Out[6]: array(['M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M',
'M', 'M', 'M', 'M', 'M', 'M', 'B', 'B', 'B', 'M', 'M', 'M', 'M',
'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'B', 'M',
'M', 'M', 'M', 'M', 'M', 'M', 'M', 'B', 'M', 'B', 'B', 'B', 'B',
'B', 'M', 'M', 'B', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'M',
'M', 'B', 'B', 'B', 'B', 'B', 'M', 'B', 'M', 'M', 'B', 'M', 'B',
'M', 'B', 'B', 'M', 'M', 'M', 'B', 'B', 'B', 'M', 'M', 'B', 'B',
'B', 'M', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'M', 'M', 'M', 'B', 'M', 'M', 'B', 'B', 'B', 'M', 'M', 'B', 'M',
'B', 'M', 'M', 'B', 'M', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'B',
'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'M', 'B', 'B', 'B', 'B', 'B', 'M', 'M', 'B', 'M', 'B', 'B', 'M',
'B', 'B', 'M', 'M', 'B', 'M', 'M', 'B', 'B', 'M', 'B', 'B', 'B',
'B', 'B', 'M', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B',
'B', 'B', 'B', 'M', 'B', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M',
'M', 'M', 'M', 'M', 'M', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'M',
'B', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'M', 'B', 'B',
'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B',
'B', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'B', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'B', 'B', 'B', 'M', 'B', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'M',
'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B',
'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B', 'M', 'B', 'M', 'B', 'B',
'B', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'M', 'M', 'B', 'B',
'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'M', 'B', 'M', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B', 'M', 'B', 'B', 'B',
'B', 'B', 'B', 'B', 'M', 'M', 'B', 'M', 'B', 'M', 'B', 'B', 'B',
'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'M', 'B', 'M', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
'B', 'B', 'B', 'M', 'M', 'M', 'M', 'M', 'M', 'M', 'B'], dtype=object)
```

```
In [7]: # You've Learned one way of doing this previously in our Naive Bayes tutorial.
# Here is another nifty way of converting labels to numerical values.
```

```
#Encoding categorical data values
from sklearn.preprocessing import LabelEncoder
labelencoder_Y = LabelEncoder()
Y = labelencoder_Y.fit_transform(Y)
```

```
In [8]: # Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.25, random_state =
0)
```

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

```
Out[9]: (426, 30)
```

```
In [10]: #Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
X_test
```

```
Out[10]: array([[ -0.21395901,  0.3125461 , -0.14355187, ...,  1.37043754,
         1.08911166,  1.53928319],
        [ -0.26750714,  1.461224 , -0.32955207, ..., -0.84266106,
        -0.71577388, -0.88105993],
        [ -0.03922298, -0.86770223, -0.10463112, ..., -0.505318 ,
        -1.20298225, -0.92494342],
        ...,
        [ -0.51270124, -1.69096186, -0.54095317, ..., -0.12632201,
         0.33773512, -0.42872244],
        [ -0.17732081, -2.01395163, -0.17345939, ..., -0.62875108,
        -0.29500302, -0.65432858],
        [  1.5305829 , -0.26300709,  1.57961296, ...,  1.6694843 ,
         1.18085869,  0.48889253]])
```

```
In [11]: #Using SVC method of svm class to use Support Vector Machine Algorithm
from sklearn.svm import SVC
classifier = SVC(kernel = 'linear', random_state = 0)
classifier.fit(X_train, Y_train)

#Using SVC method of svm class to use Kernel SVM Algorithm
from sklearn.svm import SVC
classifier = SVC(kernel = 'rbf', random_state = 0)
classifier.fit(X_train, Y_train)
```

```
Out[11]: 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=0,
         shrinking=True, tol=0.001, verbose=False)
```

```
In [12]: Y_pred = classifier.predict(X_test)
```

```
In [13]: from sklearn.metrics import confusion_matrix
cm = confusion_matrix(Y_test, Y_pred)
```

```
In [14]: cm
```

```
Out[14]: array([[88,  2],
        [ 3, 50]], dtype=int64)
```

Discuss the success of this model for predicting malignant and benign tumors.

Conceptual Review Questions:

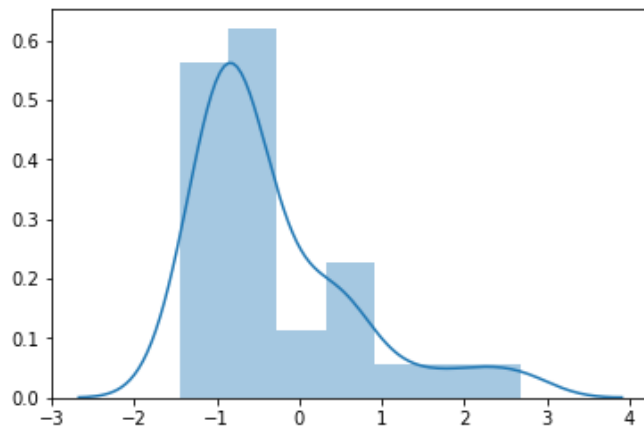
1. What is a hyperplane in SVM and how is it used?
2. What are two tuning parameters that can be used for SVM and how do they impact the model in terms of overfitting?
3. What are the pros and cons of SVMs?

Answers:

- 1) the line (2D) or plane that separates the dataset into two or more categories.
- 2) C and gamma values. High C and gamma values tend to overfit the model: high C values penalize the slack variables more which results in low margin; high gamma values results in lower radius of model influence, which would have low bias but high variance outputs.
- 3) Pros: good at dealing with high dimensional data, which can be hard to visualize; does not need a large dataset  
Cons: picking the right kernel and parameters (C and gamma) can be computational intensive (this can be overcome by developing several models with various kernels/parameters using loop function).

```
In [15]: sns.distplot(X_train[1,:])
```

```
Out[15]: <matplotlib.axes._subplots.AxesSubplot at 0x137484592e8>
```



```
In [ ]:
```