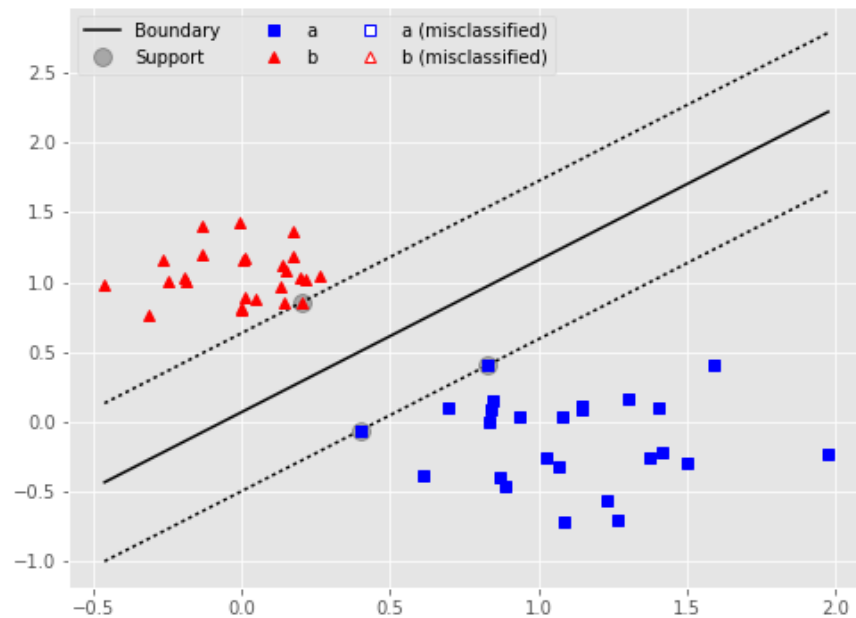


# Introduction to Data Science

## 16 - Classification via Support Vector Machines



### This Lecture

- Classification via Support Vector Machine

### Setup

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import sklearn as sk
5
6 from scipy.stats import norm
7 from pandas import Series, DataFrame
8 from matplotlib.colors import ListedColormap
9
10 plt.style.use('ggplot')
11
12 %matplotlib inline
13 %config InlineBackend.figure_format = 'retina'
```

## Example Problem

This synthetic problem creates two clusters of points which are normally distributed in two dimensions. We're going to start with data which is linearly separable in order to explain SVM better.

```
In [2]: 1 # ensure repeatability of this notebook
2 # (comment out for new results each run)
3 np.random.seed(12345)
4
5 # Get some normally distributed samples
6 def sample_cluster(n, x, y, sigma):
7     x = np.random.randn(n) * sigma + x;
8     y = np.random.randn(n) * sigma + y;
9     return np.array([x, y]).T
10
11 c1 = sample_cluster(25, 1, 0, 0.3)
12 c2 = sample_cluster(25, 0, 1, 0.2)
13 d1 = DataFrame(c1, columns=['x', 'y'])
14 d2 = DataFrame(c2, columns=['x', 'y'])
15 d1['class'] = 'a'
16 d2['class'] = 'b'
17 data = d1.append(d2)
18 data.index = pd.RangeIndex(50)
```

## Novice programmer tip:

In the cell of code above, what might you do if you don't understand why the last line was used, or what it does?

How might you find out the answer?

1. Look at the documentation for [pandas.RangeIndex](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.RangeIndex.html) (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.RangeIndex.html>) (Google for it is usually fastest way to find the webpage).
2. Write some additional code to investigate. See cell below.

```

In [3]: 1  ## Novice programmer tip
        2
        3  # Let's figure out what pd.RangeIndex(50) did, and why we needed it...
        4
        5  # First, let's see what the variable data was and held, prior to the RangeIndex ca.
        6  data = d1.append(d2)
        7  data

```

Out[3]:

	x	y	class
0	0.938588	0.037236	a
1	1.143683	0.090784	a
2	0.844168	0.157132	a
3	0.833281	0.000282	a
4	1.589734	0.403143	a
5	1.418022	-0.214063	a
6	1.027872	-0.249346	a
7	1.084524	-0.711069	a
8	1.230707	-0.558228	a
9	1.373930	-0.258227	a
10	1.302157	0.168044	a
11	0.611134	-0.379780	a
12	1.082497	0.035948	a
13	1.068674	-0.319054	a
14	1.405875	0.099865	a
15	1.265929	-0.707826	a
16	0.399509	-0.059863	a
17	0.888447	-0.462599	a
18	1.500708	-0.291221	a
19	0.868429	-0.392109	a
20	0.838078	0.085905	a
21	1.143096	0.113395	a
22	1.974683	-0.226166	a
23	0.693632	0.099386	a
24	0.826874	0.404923	a
0	0.013975	0.899383	b
1	0.049335	0.875545	b
2	-0.002372	0.815766	b
3	0.200962	0.854757	b
4	0.265439	1.044579	b
5	-0.183852	1.010263	b
6	-0.309821	0.768456	b
7	0.004437	1.163341	b

	x	y	class
8	0.151673	1.086722	b
9	-0.132105	1.202147	b
10	0.172516	1.364975	b
11	-0.002006	0.800496	b
12	0.010002	1.170118	b
13	0.134043	0.973684	b
14	0.170593	1.182483	b
15	-0.191174	1.037642	b
16	-0.004699	1.433892	b
17	-0.460847	0.977014	b
18	-0.130494	1.400739	b
19	-0.243660	1.005922	b
20	-0.266522	1.159051	b
21	0.214925	1.023622	b
22	0.144728	0.850294	b
23	0.138000	1.116994	b
24	0.200309	1.030535	b

```
In [4]: 1  ## Novice programmer tip (contintued)
        2
        3  # Ah. We see that the data.append() statement appended d1 and d2, as requested,
        4  # and also kept the index names for d1 and d2, such that the 0 to 24 labels of
        5  # both are now used twice. What happens if we try to extract rows with a specific
        6  # index name (number, in this case)? Let's use the .loc() method to do so.
        7
        8  data.loc[5]
```

```
Out[4]:
```

	x	y	class
5	1.418022	-0.214063	a
5	-0.183852	1.010263	b

```
In [5]: 1  ## Novice programmer tip (contintued)
        2
        3  # Yep, it gave us the two rows that use 5 as their index.
        4  # If we want to be able to uniquely index each row, we need to give
        5  # the data DataFrame new indices. Seems likely that that's what
        6  # RangeIndex() does. Let's confirm this...
        7  data.index = pd.RangeIndex(50)
        8  data
```

```
Out[5]:
```

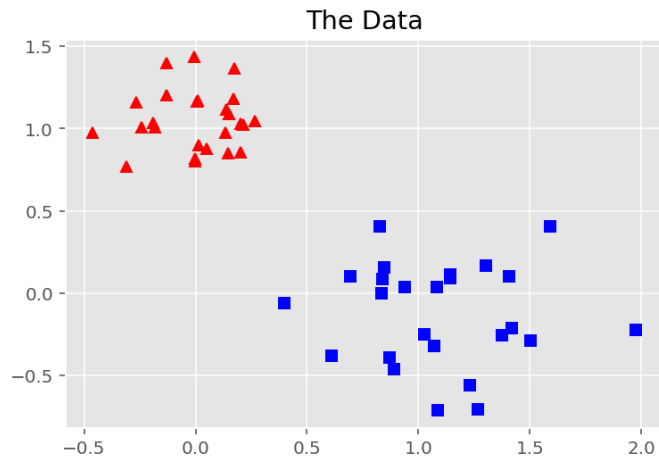
	x	y	class
0	0.938588	0.037236	a
1	1.143683	0.090784	a
2	0.844168	0.157132	a
3	0.833281	0.000282	a
4	1.589734	0.403143	a
5	1.418022	-0.214063	a
6	1.027872	-0.249346	a
7	1.084524	-0.711069	a
8	1.230707	-0.558228	a
9	1.373930	-0.258227	a
10	1.302157	0.168044	a
11	0.611134	-0.379780	a
12	1.082497	0.035948	a
13	1.068674	-0.319054	a
14	1.405875	0.099865	a
15	1.265929	-0.707826	a
16	0.399509	-0.059863	a
17	0.888447	-0.462599	a
18	1.500708	-0.291221	a
19	0.868429	-0.392109	a
20	0.838078	0.085905	a
21	1.143096	0.113395	a
22	1.974683	-0.226166	a
23	0.693632	0.099386	a
24	0.826874	0.404923	a
25	0.013975	0.899383	b
26	0.049335	0.875545	b
27	-0.002372	0.815766	b
28	0.200962	0.854757	b
29	0.265439	1.044579	b
30	-0.183852	1.010263	b
31	-0.309821	0.768456	b
32	0.004437	1.163341	b

	x	y	class
33	0.151673	1.086722	b
34	-0.132105	1.202147	b
35	0.172516	1.364975	b
36	-0.002006	0.800496	b
37	0.010002	1.170118	b
38	0.134043	0.973684	b
39	0.170593	1.182483	b
40	-0.191174	1.037642	b
41	-0.004699	1.433892	b
42	-0.460847	0.977014	b
43	-0.130494	1.400739	b
44	-0.243660	1.005922	b
45	-0.266522	1.159051	b
46	0.214925	1.023622	b
47	0.144728	0.850294	b
48	0.138000	1.116994	b
49	0.200309	1.030535	b

```
In [6]: 1  ## Novice programmer tip (contintued)
        2
        3  # Yep. The RangeIndex(50) command created new indices for the DataFrame,
        4  # from 0 to 49 (50-1).
        5
        6  # So now we understand!
```

**Back to the lecture...**

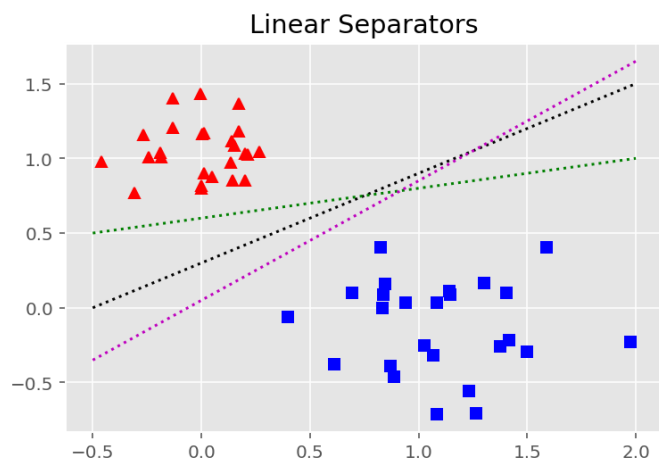
```
In [7]: 1 plt.plot(c1[:,0], c1[:,1], 'bs', label='a')
2 plt.plot(c2[:,0], c2[:,1], 'r^', label='b')
3 plt.title('The Data')
4 plt.show()
```



Note that this data is well separated, and there are many possible linear separators.

Which one is best?

```
In [8]: 1 plt.plot(c1[:,0], c1[:,1], 'bs', label='a')
2 plt.plot(c2[:,0], c2[:,1], 'r^', label='b')
3 plt.plot([-0.5, 2.0], [0, 1.5], 'k:')
4 plt.plot([-0.5, 2.0], [0.5, 1.0], 'g:')
5 plt.plot([-0.5, 2.0], [-0.35, 1.65], 'm:')
6 plt.title('Linear Separators')
7 plt.show()
```



## Maximum Margin Classifier

One answer to the question is, where can we draw the line such that the nearest exemplar(s) in each class are equidistant from the line, and where the distance is as maximal as possible?

It turns out this produces a neat little convex quadratic program, which we can feed to any QP solver.

Downside: it can be expensive with lots of data!

(Math omitted - beyond the scope of this course)

To show what the maximum margin classifier looks like on our data, we're going to create a linear SVM classifier, and fit it using all of the data.

```
In [9]: 1 from sklearn import svm
        2 model = svm.SVC(kernel='linear', C=10)
        3 model.fit(data[['x','y']], data['class'])
```

```
Out[9]: SVC(C=10, kernel='linear')
```

The `plot_predicted` function below is what we've been using to visualize our data points (correctly and incorrectly classified), together with lines that show us the decision boundary and the support vectors.

```
In [10]: 1 def plot_predicted(model, data):
        2     predicted = model.predict(data[['x','y']])
        3     correct = data[data['class'] == predicted]
        4     correct_a = correct[correct['class'] == 'a']
        5     correct_b = correct[correct['class'] == 'b']
        6     incorrect = data[data['class'] != predicted]
        7     incorrect_a = incorrect[incorrect['class'] == 'a']
        8     incorrect_b = incorrect[incorrect['class'] == 'b']
        9
       10     plt.plot(correct_a['x'], correct_a['y'], 'bs', label='a')
       11     plt.plot(correct_b['x'], correct_b['y'], 'r^', label='b')
       12     plt.plot(incorrect_a['x'], incorrect_a['y'], 'bs', markerfacecolor='w', label='a')
       13     plt.plot(incorrect_b['x'], incorrect_b['y'], 'r^', markerfacecolor='w', label='b')
       14     plt.legend(ncol=2)
```

The rather complicated `plot_linear_separator` function below extracts the relevant data from the model to plot the linear decision function and the parallel "maximum margin" that was found.



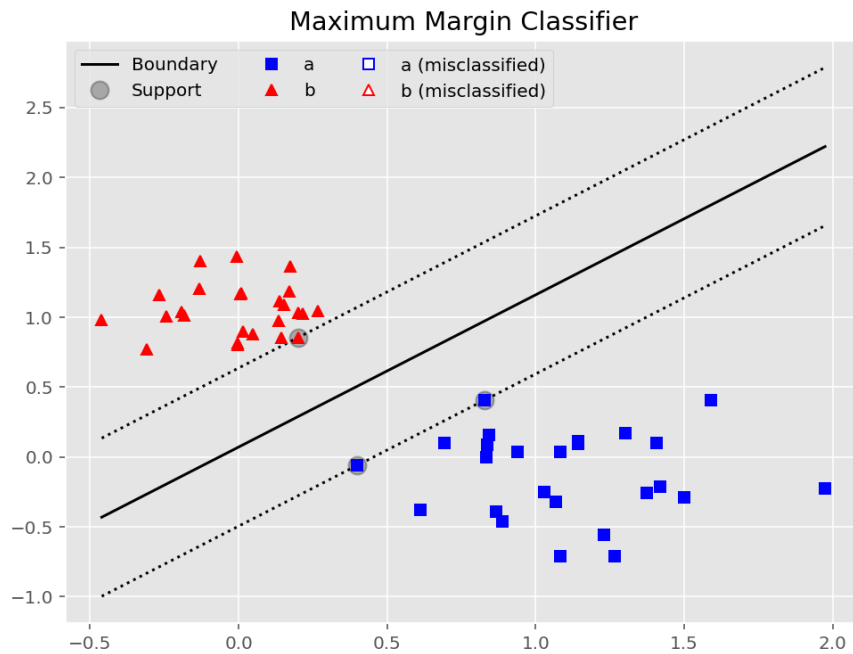
```

In [11]: 1 def plot_linear_separator(model, data, ax=None):
2         # This code modified from Scikit-learn documentation on SVM
3         if ax is None:
4             plt.figure(figsize=(8,6))
5         else:
6             plt.sca(ax)
7
8         # get the separating hyperplane as  $ax + y + c = 0$ 
9         w = model.coef_[0]
10        a = w[0] / w[1]
11        c = (model.intercept_[0]) / w[1]
12        xx = np.linspace(data['x'].min(), data['x'].max())
13        yy = -a * xx - c
14
15        # find the support vectors that define the maximal separation
16        # there ought to be a better way...
17        spos = 0
18        sneg = 0
19        sposdist = 0
20        snegdist = 0
21        for s in model.support_vectors_:
22            # find the orthogonal point
23            ox = (s[0] - a * s[1] - a * c) / (a * a + 1)
24            oy = (a * (a * s[1] - s[0]) - c) / (a * a + 1)
25
26            # find the squared distance
27            d = (s[0] - ox)**2 + (s[1] - oy)**2
28
29            if s[1] > oy and d > sposdist:
30                spos = s
31                sposdist = d
32            if s[1] < oy and d > snegdist:
33                sneg = s
34                snegdist = d
35
36
37        # plot the parallels to the separating hyperplane that pass through the
38        # support vectors
39        yy_pos = -a * xx + (spos[1] + a * spos[0])
40        yy_neg = -a * xx + (sneg[1] + a * sneg[0])
41
42        # plot the separator and the maximum margin lines
43        plt.plot(xx, yy, 'k-', label='Boundary')
44        plt.plot(xx, yy_pos, 'k:')
45        plt.plot(xx, yy_neg, 'k:')
46
47        plt.plot(model.support_vectors_[0], model.support_vectors_[1], 'ko', mar
48
49        # plot the points
50        plot_predicted(model, data)
51        plt.legend(ncol=3)

```

And finally, here's the plot, showing our data:

```
In [13]: 1 plot_linear_separator(model, data)
2 plt.title('Maximum Margin Classifier')
3 plt.show()
```



The solid line in the plot above is the decision boundary.

The dotted lines show the margin area between the classifier and the nearest points in the two clusters.

Note that the dotted lines pass through points in the clusters; these points are called the *support vectors* of the classifier.

## Non-Separable Data

So what happens when the data is not linearly separable?

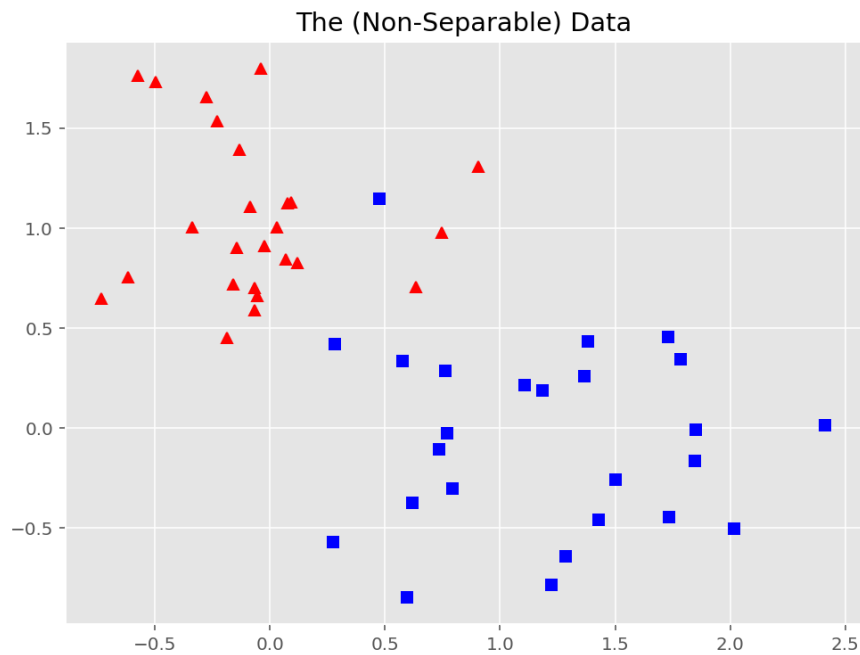
Our QP will break, because there is no feasible solution.

So, the clever fix is to relax the QP to allow points to be *misclassified*; but to get the best classifier possible, a penalty is attached to each misclassified point.

```

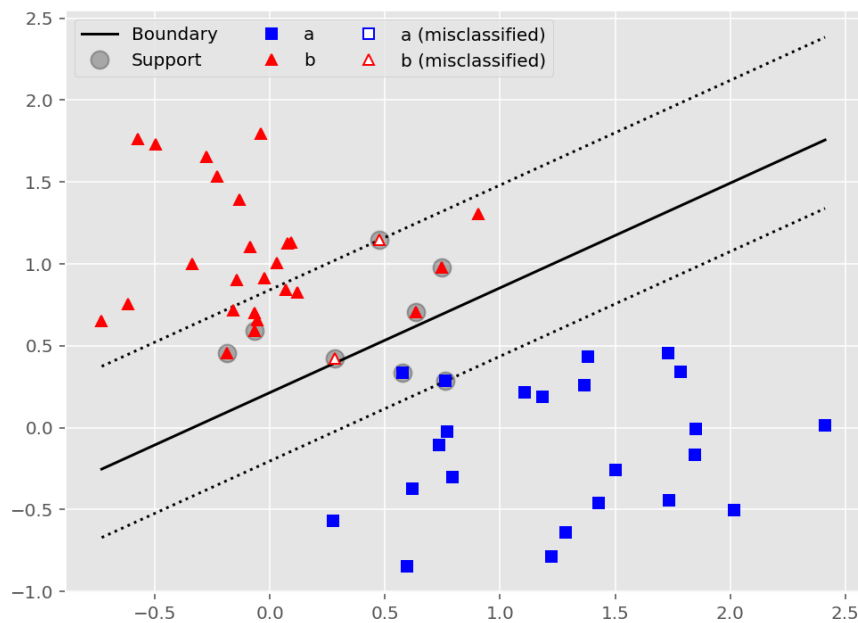
In [14]: 1 c1 = sample_cluster(25, 1, 0, 0.5)
          2 c2 = sample_cluster(25, 0, 1, 0.4)
          3 d1 = DataFrame(c1, columns=['x', 'y'])
          4 d2 = DataFrame(c2, columns=['x', 'y'])
          5 d1['class'] = 'a'
          6 d2['class'] = 'b'
          7 data = d1.append(d2)
          8 data.index = pd.RangeIndex(50)
          9
         10 plt.figure(figsize=(8,6))
         11 plt.plot(c1[:,0], c1[:,1], 'bs', label='a')
         12 plt.plot(c2[:,0], c2[:,1], 'r^', label='b')
         13 plt.title('The (Non-Separable) Data')
         14 plt.show()

```



Let's see what the support vector classifier does with this data:

```
In [16]: 1 model = svm.SVC(kernel='linear', C=10) # call the support vector classifier function
2 model.fit(data[['x','y']], data['class'])
3 plot_linear_separator(model, data)
```



## Effect of C parameter

The strength of the penalty term for the QP is controlled by a new parameter: C.

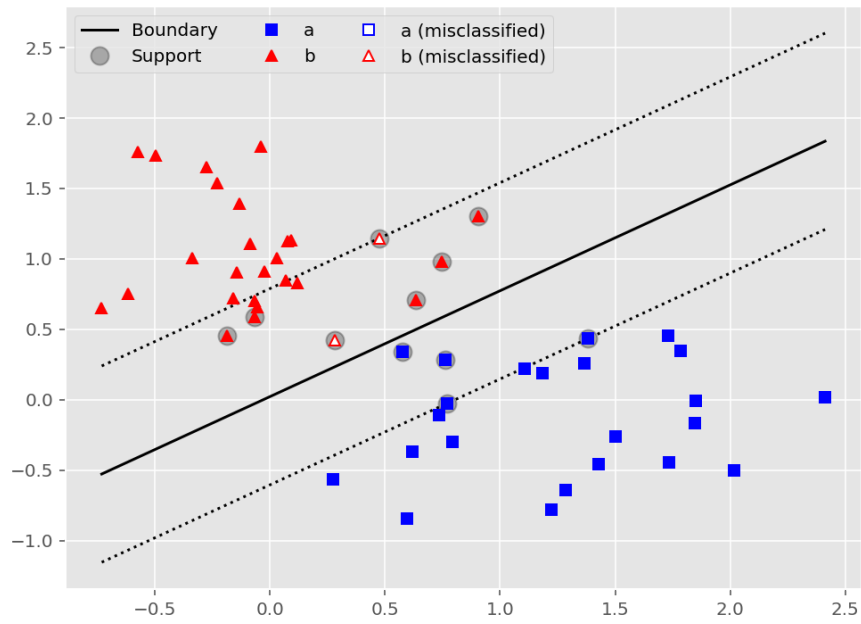
A larger C means a stronger penalty, i.e., gives the QP incentive to reduce misclassifications.

Above we used C = 10.

Let's see the effects of different choices for C:

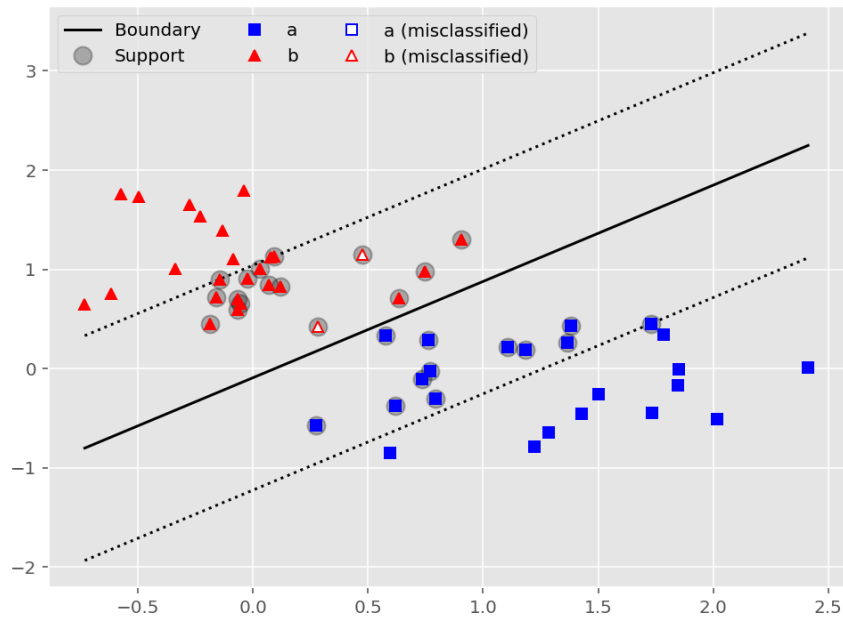
C = 1

```
In [17]: 1 model = svm.SVC(kernel='linear', C=1)
2 model.fit(data[['x','y']], data['class'])
3 plot_linear_separator(model, data)
```



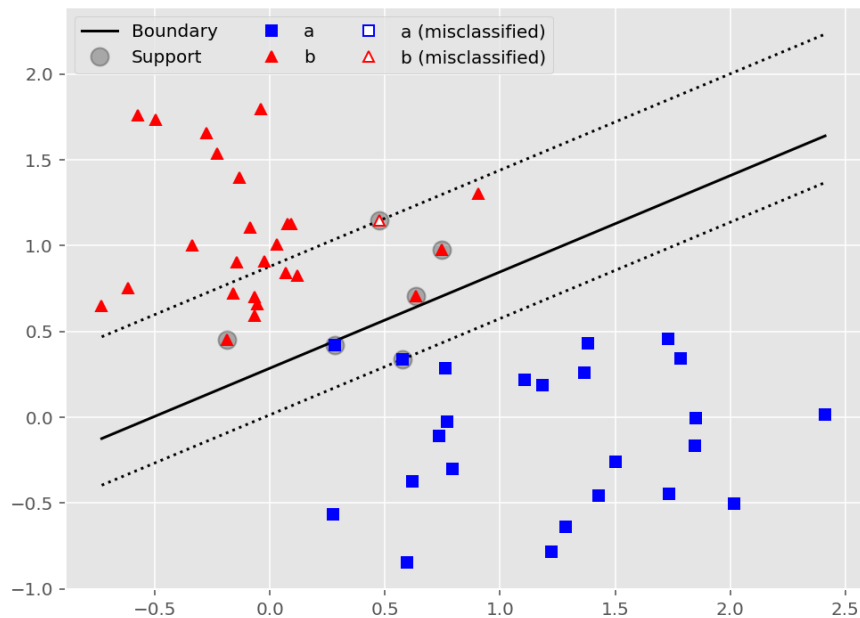
C = 0.1

```
In [18]: 1 model = svm.SVC(kernel='linear', C=0.1) # will have the largest margin hyperplane
2 model.fit(data[['x','y']], data['class'])
3 plot_linear_separator(model, data)
```

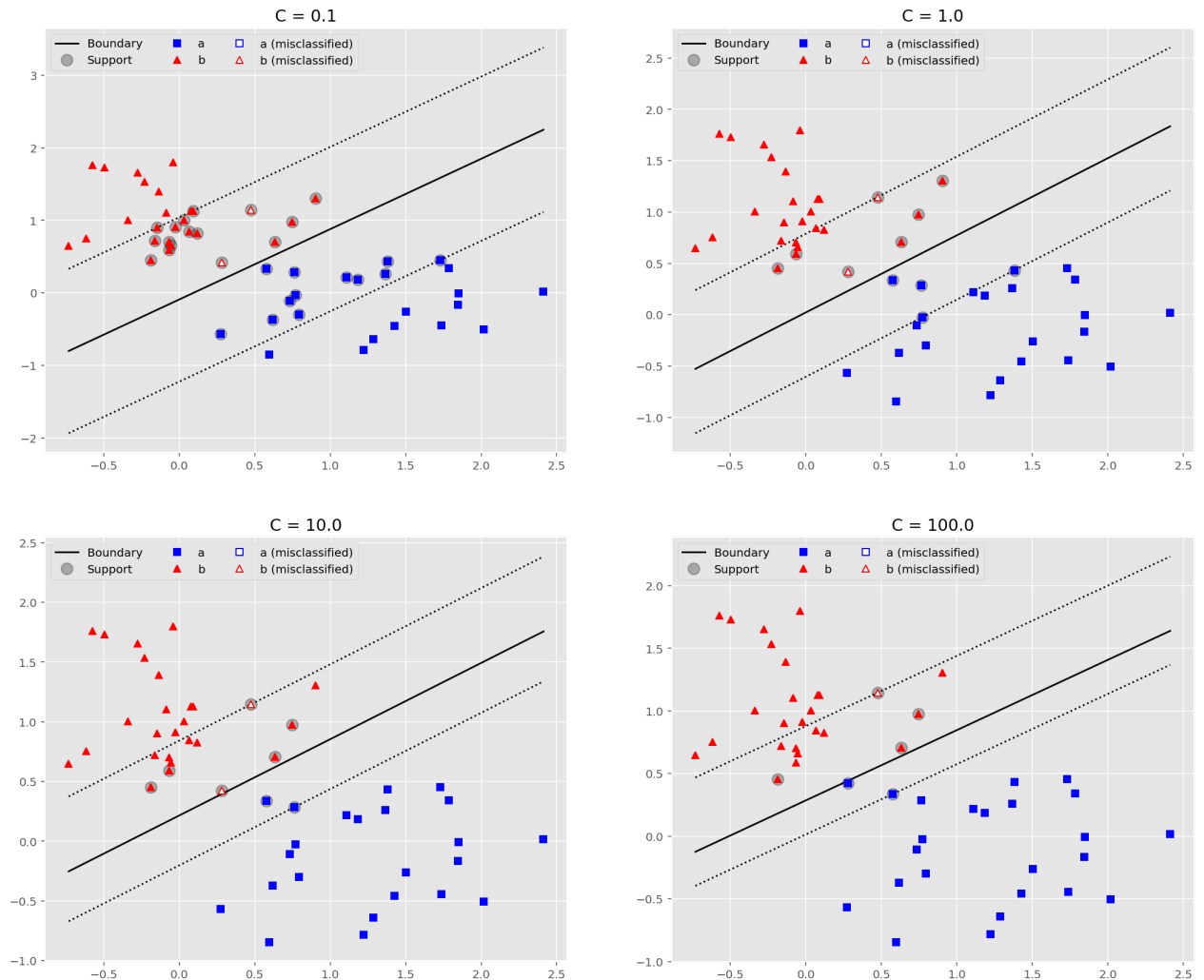


C = 100

```
In [19]: 1 model = svm.SVC(kernel='linear', C=100) # will have the smallest margin hyperplane
2 model.fit(data[['x','y']], data['class'])
3 plot_linear_separator(model, data)
```



```
In [23]: 1 plt.figure(figsize=(18,15))
2
3 C_list = [0.1, 1, 10, 100]
4 for i_subplot, C in enumerate(C_list):
5     model = svm.SVC(kernel='linear', C=C)
6     model.fit(data[['x','y']], data['class'])
7     ax = plt.subplot(2,2,i_subplot+1)
8     plot_linear_separator(model, data, ax)
9     plt.title('C = {:.1f}'.format(C))
10
11 # as you can see in the below plots: a larger C value creates a smaller-margin hyp
```



## Non-Linear Data

It turns out that a quirk in the QP formulation for SVMs allows us to efficiently replace the linear separator model with a non-linear model.

This quirk is known as the "kernel trick".

It lets us use different *kernels* without significant added expense.

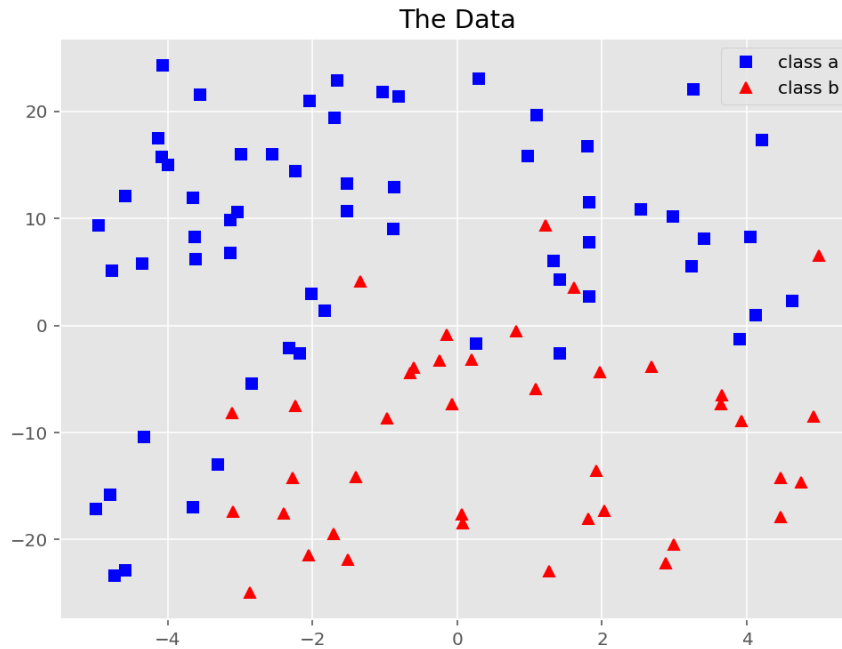
The most popular kernels are linear, polynomial, and radial basis function (RBF). Radial basis functions are basically Gaussian surfaces centered on the data points.

Let's see how this works on our example problem from before:

```
In [22]: 1 def f(X):
2         return 3 + 0.5 * X - X**2 + 0.15 * X**3
3
4         # convenience function for generating samples
5         def sample(n, fn, limits, sigma):
6             width = limits[1] - limits[0]
7             height = limits[3] - limits[2]
8             x = np.random.random(n) * width + limits[0]
9             y = np.random.random(n) * height + limits[2]
10            s = y > fn(x)
11            p = norm.cdf(np.abs(y - fn(x)), scale = sigma) # assigns p with normally distr.
12            r = np.random.random(n) # r is assigned n random variables from [0.0, 1.0).
13
14            def assign(sign, prob, rnum):
15                if sign:
16                    if rnum > prob:
17                        return 'b'
18                    else:
19                        return 'a'
20                else:
21                    if rnum > prob:
22                        return 'a'
23                    else:
24                        return 'b'
25
26            c = [assign(s[i], p[i], r[i]) for i in range(n)]
27
28            return DataFrame({'x' : x, 'y' : y, 'class' : c})
29
```



```
In [24]: 1 data = sample(100, f, [-5, 5, -25, 25], 5)
2 plt.figure(figsize=(8,6))
3 dataa = data[data['class']=='a']
4 datab = data[data['class']=='b']
5 plt.plot(dataa['x'], dataa['y'],'bs', label='class a')
6 plt.plot(datab['x'], datab['y'],'r^', label='class b')
7 plt.legend()
8 plt.title('The Data')
9 plt.show()
```



The "out of the box" default kernel for the Scikit-learn SVC is 'rbf':

```
In [25]: 1 model = svm.SVC()
2 model.fit(data[['x','y']], data['class'])
3
```

Out[25]: SVC()

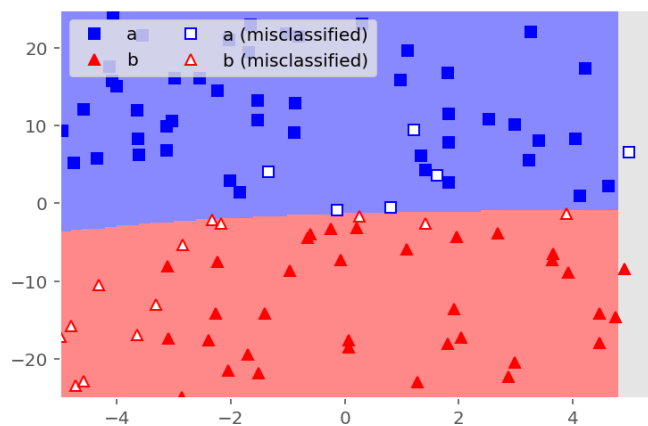
As before, we can visualize the decision boundary by simply plotting all the points in our plane:

```
In [26]: 1 def plot_boundary(model, data):
2     cmap = ListedColormap(['#8888FF', '#FF8888'])
3     xmin, xmax, ymin, ymax = -5, 5, -25, 25
4     grid_size = 0.2
5     xx, yy = np.meshgrid(np.arange(xmin, xmax, grid_size),
6                           np.arange(ymin, ymax, grid_size))
7     pp = model.predict(np.c_[xx.ravel(), yy.ravel()])
8     zz = np.array([{'a':0, 'b':1}[ab] for ab in pp])
9     zz = zz.reshape(xx.shape)
10    plt.figure()
11    plt.pcolormesh(xx, yy, zz, cmap = cmap)
12    plot_predicted(model, data)
13    plt.legend(loc='upper left', ncol=2)
```

```
In [27]: 1 plot_boundary(model, data)
2
```

<ipython-input-26-53d9b4961664>:11: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'go uraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



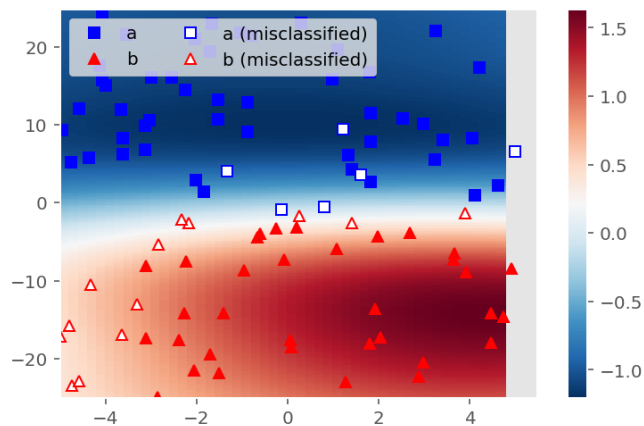
We can also plot the decision function in the plane:

```
In [28]: 1 def plot_decision(model, data):
2         cmap = 'RdBu_r'
3         xmin, xmax, ymin, ymax = -5, 5, -25, 25
4         grid_size = 0.2
5         xx, yy = np.meshgrid(np.arange(xmin, xmax, grid_size),
6                               np.arange(ymin, ymax, grid_size))
7         pp = model.decision_function(np.c_[xx.ravel(), yy.ravel()])
8         zz = pp.reshape(xx.shape)
9         plt.figure()
10        plt.pcolormesh(xx, yy, zz, cmap = cmap)
11        plt.colorbar()
12        plot_predicted(model, data)
13        plt.legend(loc='upper left', ncol=2)
```

```
In [29]: 1 plot_decision(model, data)
2         plt.show()
```

<ipython-input-28-276e3f029014>:10: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'guaranteed', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



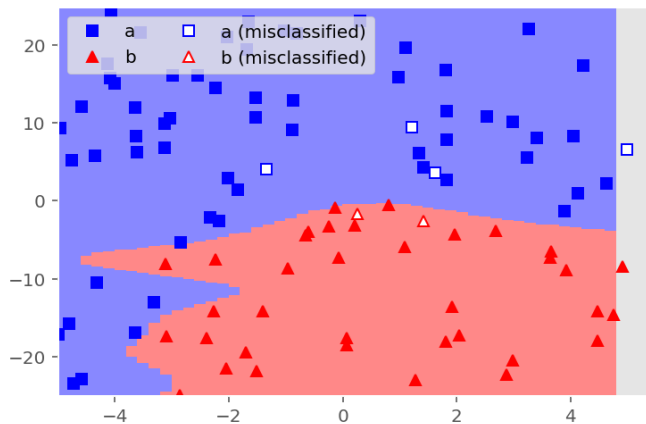
## Effects of *[Math Processing Error]* and C

The RBF kernel has an additional parameter, called `gamma`. A smaller `gamma` results in a shallower, more spread out Gaussian function, and therefore a smoother result:

```
In [30]: 1 model = svm.SVC(gamma=0.1, C = 1)
2 model.fit(data[['x','y']], data['class'])
3 plot_boundary(model, data)
4
```

<ipython-input-26-53d9b4961664>:11: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'go uraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```

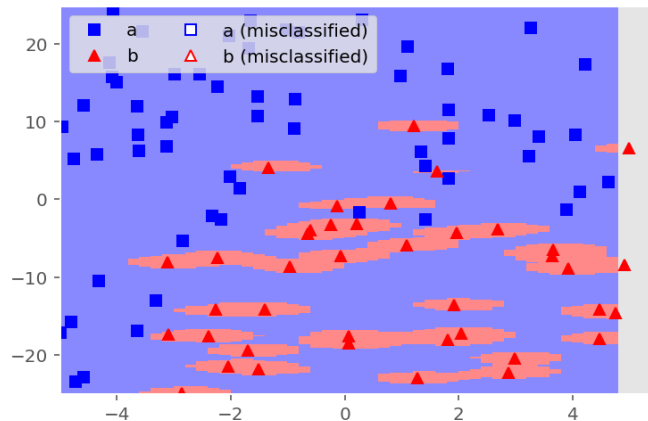


Conversely, a large gamma results in very narrow, spiky Gaussians.

```
In [31]: 1 model = svm.SVC(gamma=2, C = 1)
2 model.fit(data[['x','y']], data['class'])
3 plot_boundary(model, data)
```

<ipython-input-26-53d9b4961664>:11: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'go uraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



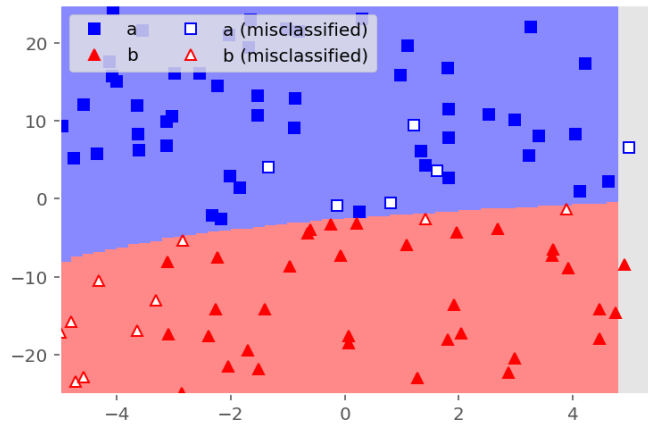
## Polynomial Kernel

There is also a polynomial kernel, which computes separators based on polynomial functions of the inputs. It has an additional parameter, `coef0`, which generally needs to be set to 1. The `degree` parameter determines the degree of the polynomial; typically degree is best kept at 2 or 3.

```
In [32]: 1 model = svm.SVC(kernel='poly', degree=3, coef0 = 1, C = 0.1)
2 model.fit(data[['x', 'y']], data['class'])
3 plot_boundary(model, data)
4 plt.show()
```

<ipython-input-26-53d9b4961664>:11: MatplotlibDeprecationWarning: shading='flat' when X and Y have the same dimensions as C is deprecated since 3.3. Either specify the corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or 'go uraud', or set rcParams['pcolor.shading']. This will become an error two minor releases later.

```
plt.pcolormesh(xx, yy, zz, cmap = cmap)
```



## Next Time

### Model selection

- Cross validation
- Parameter search

```
In [ ]: 1
```