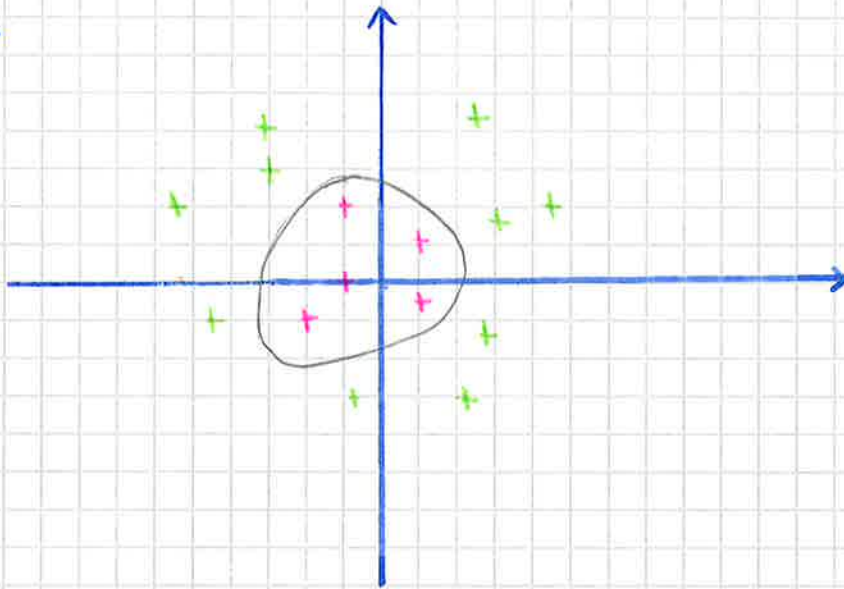


Chapter 5
Kernelized SVM's

Example:



The data cannot be separated by a straight line. However, if, instead of representing each point by $x_n = [x_{n,1}, x_{n,2}]^T$, we represented them by their distance from the origin $r_n = x_{n,1}^2 + x_{n,2}^2$, we could separate them with a straight line. When we have a test point x_{new} , we compute r_{new} and classify it according to it being \leq or \geq to a certain value. In general, we use a transformation $\phi(x_n)$ of the n -th training object.

The most important characteristic of the SVM framework is that we never have actually to perform the transformation. In our objective and decision functions, the data x_n, x_m, x_{new} appear within inner products ($x_n^T x_m, x_n^T x_{new}, \dots$). After applying the transformation, we need to calculate the inner product in the new space ($\phi(x_n)^T \phi(x_m)$). So we do not need to think in terms of transformation. Instead, if we can show that some function $k(x_n, x_m) = \phi(x_n)^T \phi(x_m)$, for some transformation ϕ , we are free to use $k(x_n, x_m)$ in place of any inner product. Function that corresponds to inner products in some space are known as kernel functions.

Our optimisation and decision function (soft margin version) rewritten to include kernel functions are:

$$\begin{cases} \text{argmax}_{\alpha} \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_n \alpha_m t_n t_m k(x_n, x_m) \\ \text{subject to } \sum_{n=1}^N \alpha_n t_n = 0 \text{ and } 0 \leq \alpha_n \leq C, \forall n \\ t_{new} = \text{sign} \left(\sum_{n=1}^N \alpha_n t_n k(x_n, x_{new}) + b \right) \end{cases}$$

Most popular kernels:

$$\begin{cases} \text{Linear: } k(x_n, x_m) = x_n^T x_m \\ \text{Gaussian: } k(x_n, x_m) = \exp \left\{ -\gamma (x_n - x_m)^T (x_n - x_m) \right\} \\ \text{Polynomial: } k(x_n, x_m) = (1 + x_n^T x_m)^\gamma \end{cases}$$

Linear is the one we used in Chapter 3. The Gaussian and polynomial kernels have a parameter that must be set

by the user (usually via cross-validation).

45

We draw a boundary obtained using a Gaussian kernel as the example. For the original SVM, we could compute the decision boundary exactly as it consisted of values of x that satisfied $w^T x + b = 0$. We can no longer compute w as it would be given by $\sum_i \alpha_i t_i \phi(x_i)$ and we do not necessarily know $\phi(x_i)$ (we only know $k(x_i, x_{new}) = \phi(x_i)^T \phi(x_{new})$). Therefore, to draw this decision boundary, we have had to evaluate $\sum_i \alpha_i t_i k(x_i, x_{new})$ over a grid of x_{new} values and then draw "the contour corresponding to $\sum_i \alpha_i t_i k(x_i, x_{new}) = 0$ ".

Tuning γ

Modifying γ changes the (implicit) transformation $\phi(x)$, which will in turn change the kind of decision boundary. For the Gaussian kernel, increasing γ has the effect of increasing the complexity of the boundaries in the original space (many figures du livre de Rogu & Grammi, p. 136).

▷ For γ "too big", the number of support vectors has increased dramatically so the solution can no longer be considered sparse.

Kernelization

The SVM is not the only algorithm that can be kernelized.

Ex: KNN classifier

KNN requires the computation of the distances between x_{new} and each x_i . This distance can be

expressed as : $(x_{\text{new}} - x_u)^T (x_{\text{new}} - x_u)$

$$\|x_{\text{new}} - x_u\|^2 = x_{\text{new}}^T x_{\text{new}} + x_u^T x_u - 2x_{\text{new}}^T x_u$$

We can replace this by a kernelized version:

$$k(x_{\text{new}}, x_{\text{new}}) - 2k(x_{\text{new}}, x_u) + k(x_u, x_u)$$

This gives us a kernelized KNN.

Exercises: p. 72-104 of the python book