

# Intelligent Systems: Reasoning and Recognition

James L. Crowley

ENSIMAG 2 / MoSIG M1

Second Semester 2015/2016

Lesson 6

19 February 2016

## Support Vector Machines using Kernels

### Contents

Kernel Functions .....	2
Support Vector Machines with Kernels .....	3
Soft Margin SVM's - Non-separable training data.....	7

#### Sources:

"Neural Networks for Pattern Recognition", C. M. Bishop, Oxford Univ. Press, 1995.

"A Computational Biology Example using Support Vector Machines", Suzy Fei, 2009 (available on line).

## Kernel Functions

Linear discriminant functions can provide very efficient 2-class classifiers, provided that the class features can be separated by a linear decision surface.

For many domains, it is possible to find a “kernel” function, that transforms the data into a space where the two classes are separate.

Instead of a decision surface:  $g(\vec{X}) = \vec{W}^T \vec{X} + b$

We will use a decision surface of the form:

$$g(\vec{X}) = \vec{W}^T f(\vec{X}) + b$$

where  $\vec{W} = f(\vec{Z}) = \sum_{m=1}^M a_m y_m f(\vec{X}_m)$  is learned from the transformed training data.

## Support Vector Machines with Kernels

Let us assume that a training data composed of  $M$  training samples  $\{\vec{X}_m\}$  and their indicator variable,  $\{y_m\}$ , where  $y_m$  is -1 or +1.

We will seek a linear decision surface  $g(\vec{X}) = \vec{W}^T f(\vec{X}) + b$  such that the training data fall into two separable classes. That is

$$\forall m: y_m(\vec{W}^T f(\vec{X}) + b) > 0$$

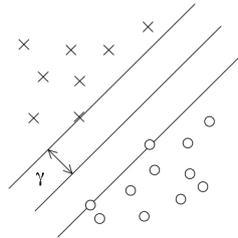
If we assume that the data is separable, then for all training samples:

$$y_m g(\vec{X}_m) > 0$$

For any training sample  $\vec{X}_m$  the perpendicular distance to the decision surface is:

$$d_m = \frac{y_m g(\vec{X}_m)}{\|\vec{W}\|} = \frac{y_m(\vec{W}^T f(\vec{X}_m) + b)}{\|\vec{W}\|}$$

The margin is the smallest distance from the decision surface:



$$\gamma = \min\{y_m(\vec{W}^T f(\vec{X}_m) + b)\}$$

For a decision surface,  $(\vec{W}, b)$ , the support vectors are the subset  $\{\vec{X}_s\}$  of the training sample,  $\{\vec{X}_s\} \subset \{\vec{X}_m\}$  that minimize the margin,  $\gamma_m$ ,

$$\gamma = \min_n \{\gamma_n\} = \min_m \left\{ \frac{1}{\|\vec{W}\|} y_m (\vec{W}^T \vec{X}_m + b) \right\}$$

We will seek to maximize the margin by finding the  $\{\vec{X}_s\}$  training samples that maximize:

$$\arg \max_{\vec{w}, b} \left\{ \frac{1}{\|\vec{W}\|} \min_m \{y_m (\vec{W}^T f(\vec{X}_m) + b)\} \right\}$$

The factor  $\frac{1}{\|\vec{W}\|}$  can be removed from the optimization because  $\|\vec{W}\|$  does not depend on  $m$ .

Direct solution would be very difficult because we do not know how many support vectors will be required.

Fortunately the problem can be converted to an equivalent problem.

Note that rescaling the problem changes nothing. Thus we will scale the equation such for the sample that is closest to the decision surface (smallest margin):

$$y_m (\vec{W}^T f(\vec{X}_m) + b) = 1 \quad \text{that is:} \quad y_m g(\vec{X}_m) = 1$$

For all other sample points:

$$y_m (\vec{W}^T f(\vec{X}_m) + b) > 1$$

This is known as the Canonical Representation for the decision hyperplane.

The training sample where  $y_m (\vec{w}^T f(\vec{X}_m) + b) = 1$  are said to be the "active" constraint. All other training samples are "inactive".

By definition there is always at least one active constraint.

When the margin is maximized, there will be  $D+1$  active constraints.

Thus the optimization problem is to maximize  $\arg \min_{\vec{w}, b} \left\{ \frac{1}{2} \|\vec{W}\|^2 \right\}$  subject to the active constraints.

The factor of  $\frac{1}{2}$  is a convenience for later analysis.

To solve this problem, we will use Lagrange Multipliers,  $a_n \geq 0$ , with one multiplier for each constraint. This gives a Lagrangian function:

$$L(\vec{W}, b, \vec{a}) = \frac{1}{2} \|\vec{W}\|^2 - \sum_{m=1}^M a_m \{y_m (\vec{W}^T f(\vec{X}_m) + b) - 1\}$$

Setting the derivatives to zero, we obtain:

$$\frac{\partial L}{\partial \vec{W}} = 0 \Rightarrow \vec{W} = \sum_{m=1}^M a_m y_m f(\vec{X}_m)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{m=1}^M a_m y_m = 0$$

Eliminating  $\vec{w}, b$  from  $L(\vec{w}, b, \vec{a})$  we obtain:

$$L(\mathbf{a}) = \sum_{m=1}^M a_m - \frac{1}{2} \sum_{m=1}^M \sum_{n=1}^M a_n a_m y_n y_m k(\vec{X}_n, \vec{X}_m)$$

with constraints:

$$a_m \geq 0 \text{ for } m=1, \dots, M$$

$$\sum_{m=1}^M a_m y_m = 0$$

where the kernel function is :  $k(\vec{X}_1, \vec{X}_2) = f(\vec{X}_1)^T f(\vec{X}_2)$

The solution takes the form of a quadratic programming problem in  $D_k$  variables (the dimension of the Kernel space). This would normally take  $O(D_k^3)$  computations.

In going to the dual formulation, we have converted this to a dual problem over  $M$  data points, requiring  $O(M^3)$  computations.

This can appear to be a problem, but the solution only depends on a small number of points  $M_s \ll M$ .

To classify a new observed point, we evaluate:

$$g(\vec{X}) = \sum_{m=1}^M a_m y_m k(\vec{X}_m, \vec{X}) + b$$

The solution to optimization problems of this form satisfy the "Karush-Kuhn-Tucker" condition, requiring:

$$\begin{aligned} a_m &\geq 0 \\ y_m g(\vec{X}_m) - 1 &\geq 0 \\ a_m \{y_m g(\vec{X}_m) - 1\} &\geq 0 \end{aligned}$$

For every observation in the training set,  $\{\vec{X}_m\}$ , either

$$a_m = 0 \quad \text{or} \quad y_m g(\vec{X}_m) = 1$$

Any point for which  $a_m = 0$  does not contribute to  $g(\vec{X}) = \sum_{m=1}^M a_m y_m k(\vec{X}_m, \vec{X}) + b$

and thus is not used! (is not active) .

The remaining  $M_s$  samples for which  $a_m \neq 0$  are the Support vectors.

These points lie on the margin at  $y_m g(\vec{X}_m) = 1$  of the maximum margin hyperplane.

Once the model is trained, all other points can be discarded!

Let us define the support vectors as the set  $\{\vec{X}_s\}$ .

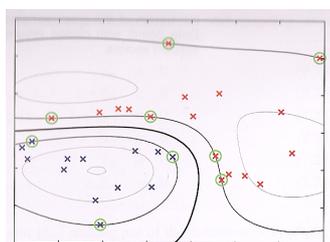
Now that we have solved for  $\{\vec{X}_s\}$  and  $\mathbf{a}$ , we can solve for  $b$ :

we note that for any active training sample  $m$  in  $\{\vec{X}_s\}$

$$y_m \left( \sum_{n \in S} a_n y_n k(\vec{X}_n, \vec{X}_m) + b \right) = 1$$

averaging over all support vectors in  $\{\vec{X}_s\}$  gives:

$$b = \frac{1}{M_s} \sum_{m \in S} \left( y_m - \sum_{n \in S} a_n y_n k(\vec{X}_n, \vec{X}_m) \right)$$



From Bishop p 331.

## Soft Margin SVM's - Non-separable training data.

So far we have assumed that the data are linearly separable in  $f(\vec{X})$ .

For many problems some training data may overlap.

The problem is that the error function goes to  $\infty$  for any point on the wrong side of the decision surface. This is called a "hard margin" SVM.

We will relax this by adding a "slack" variable,  $z_n$  for each training sample.

$$z_m \geq 1$$

We will define

$$z_m = 0 \quad \text{for training samples on the correct side of the margin, and}$$

$$z_m = |y_m - g(\vec{X}_m)| \quad \text{for other training samples.}$$

For a sample inside the margin, but on the correct side of the decision surface:

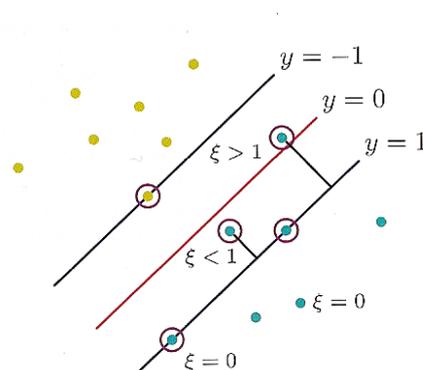
$$0 < z_m \leq 1$$

For a sample on the decision surface:

$$z_m = 1$$

For a sample on the wrong side of the decision surface:

$$z_m > 1$$



Soft margin SVM: Bishop p 332 (note use of  $\xi_n$  in place  $z_n$ )

This is called a soft margin SVM.

To softly penalize points on the wrong side, we minimize :

$$C \sum_{m=1}^M z_m + \frac{1}{2} \|\vec{w}\|^2$$

where  $C > 0$  controls the tradeoff between slack variables and the margin.

because any misclassified point  $z_m > 1$ , the upper bound on the number of misclassified points is  $\sum_{m=1}^M z_m$ .

$C$  is an inverse factor. (note  $C = \infty$  is the SVM with hard margins)

To solve for the SVM we write the Lagrangian:

$$L(\vec{W}, b, z, \vec{a}, \mu) = \frac{1}{2} \|\vec{W}\|^2 + C \sum_{m=1}^M z_m - \sum_{m=1}^M a_m \{y_m g(\vec{X}_m) - 1 + z_m\} - \sum_{m=1}^M \mu_m z_m$$

where  $\{a_m \geq 0\}$  and  $\{\mu_m \geq 0\}$  are the Lagrange multipliers.

The KKT conditions are

$$\begin{aligned} a_m &\geq 0 \\ y_m g(\vec{X}_m) - 1 + z_m &\geq 0 \\ a_m \{y_m g(\vec{X}_m) - 1 + z_m\} &\geq 0 \\ \mu_m &\geq 0 \\ z_m &\geq 1 \\ \mu_m z_m &= 0 \end{aligned}$$

We optimize for  $\vec{w}$ ,  $b$ , and  $\{z_m\}$ , using  $g(\vec{X}) = \vec{W}^T f(\vec{X}) + b$

Solving the derivatives of  $L(\vec{W}, b, \vec{a})$  for zero gives

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \vec{W} = \sum_{m=1}^M a_m y_m f(\vec{X}_m)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{m=1}^M a_m y_m = 0$$

$$\frac{\partial L}{\partial z_n} = 0 \Rightarrow a_m = C - \mu_n$$

using these to eliminate  $w$ ,  $b$  and  $\{S_m\}$  from  $L(w, b, a)$  we obtain

$$L(\vec{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{m=1}^M \sum_{n=1}^M a_m a_n y_m y_n k(\vec{X}_m, \vec{X}_n)$$

This appears to be the same as before, except that the constraints are different.

$$0 \leq a_m \leq C \quad \text{and} \quad \sum_{m=1}^M a_m y_m = 0$$

(referred to as a "box" constraint). The solution is a quadratic programming problem, with complexity  $O(M^3)$ . However, as before, a large subset of training samples have  $a_m = 0$ , and thus do not contribute to the optimization.

For the remaining points  $y_m g(\vec{X}_m) = 1 - S_m$

For samples ON the margin  $a_m < C$  hence  $\mu_m > 0$  requiring that  $S_m = 0$

For samples INSIDE the margin:  $a_m = C$  and  $S_m \leq 1$  if correctly classified and  $S_m > 1$  if misclassified.

as before to solve for  $b$  we note that :

$$y_m \left( \sum_{n \in S} a_n y_n k(\vec{X}_n, \vec{X}_m) + b \right) = 1$$

Averaging over all support vectors in  $S$  gives:

$$b = \frac{1}{M_N} \sum_{m \in N} \left( y_m - \sum_{n \in S} a_n y_n k(\vec{X}_m, \vec{X}_n) \right)$$

where  $\mathcal{N}$  denotes the set of support vectors such that  $0 < a_n < C$ .