

Intelligent Systems: Reasoning and Recognition

James L. Crowley

ENSIMAG 2 / MoSIG M1

Second Semester 2015/2016

Lesson 5

17 February 2016

Kernel Functions and Support Vector Machines

Contents

Kernel Functions	2
Polynomial Kernel Functions	4
Radial Basis Function (RBF)	5
Kernel Functions for Symbolic Data	6
Support Vector Machines	7
Hard-Margin SVMs - a simple linear classifier	7
Finding for the support vectors.	8

Sources:

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

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

Kernel Functions

A Kernel function transforms the training data so that a non-linear decision surface is transformed to a linear equation in a higher number of dimensions.

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

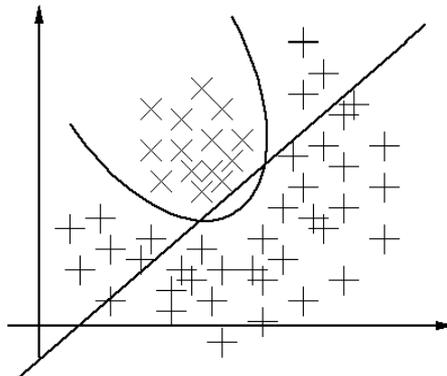
For many domains, it is easier to separate the classes with a linear function if you can transform your feature data into a space with a higher number of dimensions.

One way to do this is to transform the features with a “kernel” function.

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

We use a decision surface $g(\vec{X}) = \vec{W} \cdot \vec{f}(\vec{X}) + b$

The function $\vec{f}(\vec{X})$ provides an implicit non-linear decision surfaces for the original data.



Formally, a "kernel function" is any function that satisfies the Mercer condition.

A function, $k(x, y)$, satisfies Mercer's condition if for all square, integrable functions $f(x)$,

$$\iint k(x, y) f(x) f(y) dx dy \geq 0$$

This condition is satisfied by inner products (dot products) $\vec{W}^T \vec{X} = \sum_{d=1}^D w_d x_d$

Thus $k(\vec{W}, \vec{X}) = \vec{W}^T \vec{X}$ is a valid kernel function. (Known as the linear kernel).

as is $k(\vec{Z}, \vec{X}) = f(\vec{Z})^T f(\vec{X})$

We can learn the discriminant in an inner product space $k(\vec{Z}, \vec{X}) = f(\vec{Z})^T f(\vec{X})$ where \vec{W} will be learned from the training data.

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

The Mercer function can be satisfied by many other functions.

Popular kernel functions include:

- Polynomial Kernels
- Radial Basis Functions
- Fisher Kernels
- Text intersection
- Bayesian Kernels

Kernel functions provide an implicit feature space. We will see that we can learn in the kernel space, and then recognize without explicitly computing the position in this implicit space!

This will allow us to use Kernels for infinite dimensional spaces as well as non-numerical and symbolic data!

Polynomial Kernel Functions

The Polynomial kernel is defined as

$$k(\vec{x}, \vec{z}) = (\vec{z}^T \vec{x} + c)^n$$

where n is the “order” of the kernel, and c is a constant that allows to trade off the influence of the higher order and lower order terms.

Second order or quadratic kernels are a popular form of Polynomial kernel, widely used in Speech Recognition. Higher order kernels tend to “overfit” the training data and thus do not generalize well.

For example for a feature vector with two components $d=2$: $\vec{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

without the c term, the quadratic kernel can be expressed as:

$$k(\vec{x}, \vec{z}) = (\vec{z}^T \vec{x})^2 = f(\vec{z})^T f(\vec{x})$$

$$k(\vec{x}, \vec{z}) = (x_1 z_1 + x_2 z_2)^2 = (x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2)$$

$$k(\vec{x}, \vec{z}) = (z_1^2 + \sqrt{2} z_1 z_2 + z_2^2)(x_1^2 + \sqrt{2} x_1 x_2 + x_2^2)^T = f(\vec{z})^T f(\vec{x})$$

thus

$$f(\vec{X}) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix}$$

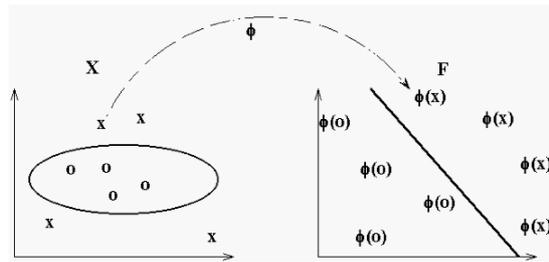
When the c is included the quadratic kernel is expressed as:

$$k(\vec{x}, \vec{z}) = (\vec{z}^T \vec{x} + c)^2 = f(\vec{z})^T f(\vec{x})$$

and $D=2$: and $\vec{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ the implicit feature space has 5 dimensions.

$$f(\vec{X}) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 c \\ \sqrt{2} x_1 x_2 \\ \sqrt{2} x_2 c \\ x_2^2 \end{pmatrix}$$

This Kernel can be used to map a hyperbolic surface to a plane



In general when for a d dimensional feature vector

$$k(\vec{x}, \vec{z}) = (\vec{z}^T \vec{x} + c)^n \text{ yields } f(\vec{X}) = \begin{pmatrix} x_1^2 \\ \vdots \\ x_d^2 \\ \sqrt{2}x_1x_2 \\ \vdots \\ \sqrt{2}x_dx_{d-1} \\ \sqrt{2}x_1c \\ \vdots \\ \sqrt{2}x_dc \\ c \end{pmatrix}$$

Radial Basis Function (RBF)

Also known as a Gaussian Kernel, Radial Basis Function (RBF) kernels are often used in Computer Vision. The RBF Kernel function has the form:

$$g(\vec{X}) = \sum_{n=1}^N \vec{W}_n^T f(\|\vec{X} - \vec{X}_n\|)$$

Where the points N samples \vec{X}_n can be derived from the training data.

The term $\|\vec{X} - \vec{X}_n\|$ is the Euclidean distance from the set of points $\{\vec{X}_n\}$.

The distance can be normalized by dividing by σ

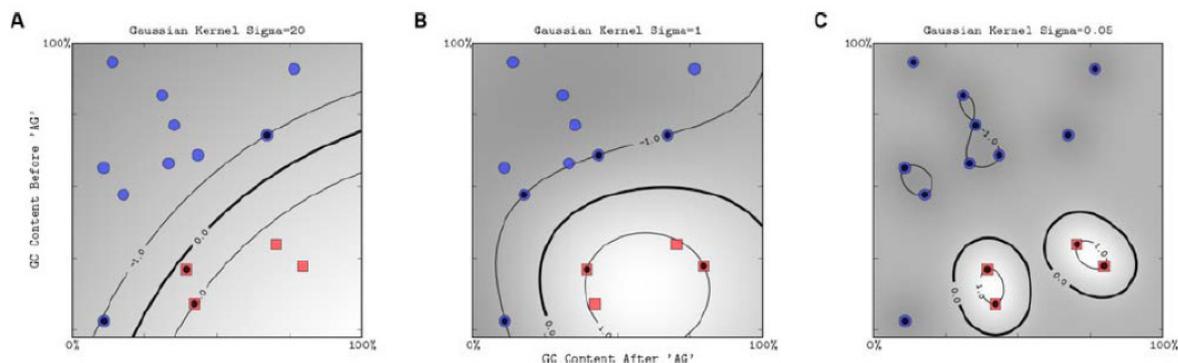
$$g(\vec{X}) = \sum_{n=1}^N \vec{W}_n^T f\left(\frac{\|\vec{X} - \vec{X}_n\|}{\sigma}\right)$$

The sigma parameter acts as a smoothing parameter that determines the influence of each of the points, \vec{X}_n .

The zero-crossings in the distances define the decision surface.

Depending on σ , this can provide a good fit or an over fit to the data. If σ is large compared to the distance between the classes, this can give an overly flat discriminant surface. If σ is small compared to the distance between classes, this will over-fit the samples.

A good choice for σ will be comparable to the distance between the closest members of the two classes.



(images from "A Computational Biology Example using Support Vector Machines", Suzy Fei, 2009)

Each Radial Basis Function is a dimension in a high dimensional basis space.

Kernel Functions for Symbolic Data

Kernel functions can be defined over graphs, sets, strings and text!

Consider for example, a non-vector space composed of a set of words $\{W\}$. We can select a subset of discriminant words $\{S\} \subset \{W\}$

Now given a set of words (a probe), $\{A\} \subset \{W\}$

We can define a kernel function of A and S using the intersection operation.

$$k(A, S) = 2^{|A \cap S|}$$

where $| \cdot |$ denotes the cardinality (the number of elements) of a set.

Support Vector Machines

Support Vector Machines (SVM), also known as maximum margin classifiers are popular for problems of classification, regression and novelty detection. The solution of the model parameters corresponds to a convex optimization problem. SVM's use a minimal subset of the training data (the “support vectors”) to define the “best” decision surface between two classes. We will use the two class problem, $K=2$, to illustrate the principle. Multi-class solutions are possible.

The simplest case, the hard margin SVM, require that the training data be completely separated by at least one hyper-plane. This is generally achieved by using a Kernel to map the features into a high dimensional space where the two classes are separable.

To illustrate the principle, we will first examine a simple linear SVM where the data are separable. We will then generalize with Kernels and with soft margins.

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

Hard-Margin SVMs - a simple linear classifier.

The simplest case is a linear classifier trained from separable data.

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

Where the decision rule is: IF $\vec{W}^T \vec{X} + b > 0$ THEN C_1 else C_2

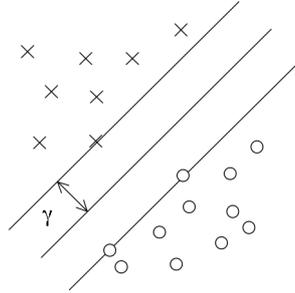
For a hard margin SVM we assume that the two classes are separable for all of the training data:

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

We will use a subset S of the training samples, $\{\vec{X}_s\} \subset \{\vec{X}_m\}$ composed of M_s training samples to define the “best” decision surface $g(\vec{X}) = \vec{W}^T \vec{X} + b$. The minimum number of support vectors depends on the number of features ($M_s = D+1$). The M_s selected training samples $\{\vec{X}_s\}$ are called the support vectors. For example, in a 2D feature space we need only 3 training samples to serve as support vectors.

Thus to define the classifier we will look for the subset of S training samples that maximizes the separation between the two classes.

The separation is defined using the margin: γ .



Finding for the support vectors.

Assume that we have M training samples $\{\vec{X}_m\}$ and their indicator variable $\{y_m\}$, where, y_m is -1 or $+1$.

Assume that we have normalized the coefficients of the hyperplane such at

$$\|\vec{W}\| = 1$$

Then the distance of any sample point \vec{X}_m from the hyper-plane, \vec{W} .

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

The margin is the minimum distance

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

A D dimensional decision surface is defined by at least D points. At least one additional point is required to define the margin. Thus we seek a subset of $M_s = D + 1$ training samples, $\{\vec{X}_s\} \subset \{\vec{X}_m\}$ to define a decision surface and its margin.

Our algorithm must choose $D + 1$ training samples $\{\vec{X}_s\}$ from the M Training samples in $\{\vec{X}_m\}$ such that the margin is as large as possible. This is equivalent to a search for a pair of parallel surfaces a distance γ from the decision surface.

We will use these samples as support vectors.

For the $D+1$ support vectors $\{\vec{X}_s\}$

$$d_s = \gamma$$

For all other training samples:

$$d_m \geq \gamma$$

To find the support vectors, we can arbitrarily define the margin as $\gamma = 1$ and then renormalize $\|\vec{W}\|$ once the support vectors have been discovered. We will look for two separate hyper-planes that “bound” the decision surface, such that for points on these surfaces:

$$\vec{W}^T \vec{X} + b = 1 \quad \text{and} \quad \vec{W}^T \vec{X} + b = -1$$

The distance between these two planes is $\frac{2}{\|\vec{W}\|}$

We will add the constraint that for all training samples that are support vectors

$$\vec{W}^T \vec{X}_m + b \geq 1$$

while for all other samples:

$$\vec{W}^T \vec{X}_m + b \leq -1$$

This can be written as: $y_m (\vec{W}^T \vec{X}_m + b) \geq 1$

This gives we have an optimization algorithm that minimizes $\|\vec{W}\|$ subject to

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

If we note that minimizing $\|\vec{W}\|$ is equivalent minimizing $\frac{1}{2}\|\vec{W}\|^2$, we can set this up as a quadratic optimization problem, and use Lagrange Multipliers.

Our problem is then to find $\arg\min_{w,b} \{\|\vec{W}\|^2\}$ such that $y_m (\vec{W}^T \vec{X}_m + b) \geq 1$

We look for a surface, (\vec{W}, b) such that

$$\arg\min_{\vec{w}, b} \left\{ \frac{1}{2} \|\vec{W}\|^2 \right\} \quad \text{subject to} \quad y_m (\vec{W}^T \vec{X}_m + b) \geq 1$$

by searching for :

$$\arg\min_{\vec{w}, b} \left\{ \max_{\alpha_m} \left\{ \frac{1}{2} \|\vec{W}\|^2 - \sum_{n=1}^N \alpha_n [y_n (\vec{W}^T \vec{X}_n + b) - 1] \right\} \right\}$$

for a subset of $M_s \geq D+1$ samples, $\alpha_m \geq 0$. These are the samples on the margins.

For all other samples where $(\alpha_m < 0)$ we set $\alpha_m = 0$.

The normal of the decision surface is then:

$$\vec{W} = \sum_{m=1}^M \alpha_m y_m \vec{X}_m$$

and the offset can be found by solving for:

$$b = \frac{1}{M_s} \sum_{m \in S} \vec{W}^T \vec{X}_m - y_m$$

The solution can be generalized for use with non-linear decision surfaces using kernels.