

Lecture 10-11- Feature Maps, Representer Theorem and Kernels.

Lecturer: Lorenzo Rosasco

In this class we introduce the concepts of feature map and kernel that allow to generalize Regularization Networks, and not only, well beyond linear models. Our starting point will be again Tikhonov regularization,

$$\min_{w \in \mathbb{R}^D} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_w(x_i)) + \lambda \|w\|^2. \quad (10.1)$$

10.1 Feature Maps

A feature map is a map

$$\Phi : X \rightarrow F$$

from the input space into a new space called feature space where there is a scalar product $\Phi(x)^T \Phi(x')$. The feature space can be infinite dimensional and the following notation is used for the scalar product $\langle \Phi(x), \Phi(x') \rangle_F$.

10.1.1 Beyond Linear Models.

The simplest case is when $F = \mathbb{R}^p$, and we can view the entries $\Phi(x)^j$, $j = 1, \dots, p$ as novel measurements on the input points. For illustrative purposes consider $X = \mathbb{R}^2$. An example of feature map could be $x = (x_1, x_2) \mapsto \Phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$. With this choice if we now consider

$$f_w(x) = w^T \Phi(x) = \sum_{j=1}^p w^j \Phi(x)^j$$

we effectively have that the function is no longer linear but it is a polynomial of degree 2. Clearly the same reasoning holds for much more general choice of measurements (features), in fact *any* finite set of measurements. Although seemingly simple, the above observation allows to consider very general models. Figure 10.1 gives a geometric interpretation of the potential effect of considering a feature map. Points which are not easily classified by a linear model, can be easily classified by a *linear model in the feature space*. Indeed, the model is no longer linear in the original input space.

10.1.2 Computations.

While feature maps allow to consider nonlinear models, the computations are essentially the same as in the linear case. Indeed, it is easy to see that the computations considered for linear models, under different loss functions, remain unchanged, as long as we change $x \in \mathbb{R}^D$ into $\Phi(x) \in \mathbb{R}^p$. For example, for least squares we simply need to replace the n by D matrix X_n with a new n by p matrix Φ_n , where each row is the image of an input point in the feature space as defined by the feature map.

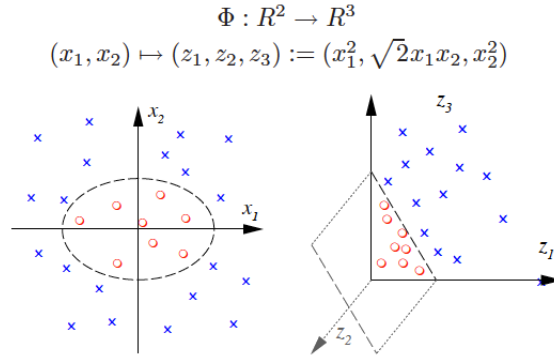


Figure 10.1. A pictorial representation of the potential effect of considering a feature map in a simple two dimensional example.

10.2 Representer Theorem

In this section we discuss how the above reasoning can be further generalized. The key result is that the solution of regularization problems of the form (10.1), can always be written as

$$\hat{w}^T = \sum_{i=1}^n x_i^T c_i, \quad (10.2)$$

where x_1, \dots, x_n are the inputs in the training set and $c = (c_1, \dots, c_n)$ a set of coefficients. The above result is an instance of the so called representer theorem. We first discuss this result in the context of RLS.

10.2.1 Representer Theorem for RLS

The result follows noting that the following equality holds,

$$(X_n^T X_n + \lambda n I)^{-1} X_n^T = X_n^T (X_n X_n^T + \lambda n I)^{-1}, \quad (10.3)$$

so that we have,

$$w = X_n^T \underbrace{(X_n X_n^T + \lambda n I)^{-1} Y_n}_c = \sum_{i=1}^n x_i^T c_i.$$

Equation (10.3) follows from considering the SVD of X_n , that is $X_n = U \Sigma V^T$. Indeed we have $X_n^T = V \Sigma U^T$ so that

$$(X_n^T X_n + \lambda n I)^{-1} X_n^T = V (\Sigma^2 + \lambda)^{-1} \Sigma U^T$$

and

$$X_n^T (X_n X_n^T + \lambda n I)^{-1} = V \Sigma (\Sigma^2 + \lambda)^{-1} U^T.$$

10.2.2 Representer Theorem Implications

Using Equation (10.2) it is possible to show how the vector c of coefficients can be computed considering different loss functions. In particular, for the square loss the vector of coefficients satisfies the following linear system

$$(K_n + \lambda nI)c = Y_n.$$

where K_n is the n by n matrix with entries $(K_n)_{i,j} = x_i^T x_j$. The matrix K_n is called the *kernel matrix* and is symmetric and positive semi-definite.

10.3 Kernels

One of the main advantages of using the representer theorem is that the solution of the problem depends on the input points only through inner products $x^T x'$. Kernel methods can be seen as replacing the inner product with a more general function $K(x, x')$. In this case, the representer theorem (10.2), that is $f_w(x) = w^T x = \sum_{i=1}^n x_i^T x c_i$, becomes

$$\hat{f}(x) = \sum_{i=1}^n K(x_i, x) c_i. \quad (10.4)$$

and we can promptly derive kernel versions of the Regularization Networks induced by different loss functions.

The function K is often called a kernel and to be admissible it should *behave like* an inner product. More precisely it should be: 1) symmetric, and 2) positive definite, that is the kernel matrix K_n should be positive semi-definite for any set of n input points. While the symmetry property is typically easy to check, positive semi-definiteness is trickier. Popular examples of positive definite kernels include:

- the linear kernel $K(x, x') = x^T x'$,
- the polynomial kernel $K(x, x') = (x^T x' + 1)^d$,
- the Gaussian kernel $K(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$,

where the last two kernels have a tuning parameter, the degree and Gaussian width, respectively.

A positive definite kernel is often called a *reproducing kernel* and is a key concept in the theory of reproducing kernel Hilbert spaces.

We end noting that there are some basic operations that can be used to build new kernels. In particular it is easy to see that, if K_1, K_2 are reproducing kernels then $K_1 + K_2$ is also a kernel.