

Kernel-Based Learning & Multivariate Modeling

MIRI Master - DMKM Master

Lluís A. Belanche
belanche@cs.upc.edu

Soft Computing Research Group

Universitat Politècnica de Catalunya

2016-2017

Kernel-Based Learning & Multivariate Modeling

Contents by lecture

Sept 14 Introduction to Kernel-Based Learning

Sept 21 The SVM for classification, regression and novelty detection (I)

Sept 28 The SVM for classification, regression and novelty detection (II)

Oct 05 Kernel design (I): theoretical issues

Oct 19 Kernel design (II): practical issues

Oct 26 Kernelizing ML & stats algorithms

Nov 02 Advanced topics

Kernel-Based Learning

Introduction

Desiderata for learning methods (note this is my particular view):

Robust Insensitive to outliers, errors and/or wrong model assumptions

Stable Against variations of the training data samples

Efficient In the computational sense (necessary to handle large datasets)

Effective Flexibility (complexity surplus) & explicit complexity control

Versatile Accept different data types and/or similarity measures

⇒ **Generalize** well to unseen data (as well as possible)

Kernel-Based Learning

Introduction

Kernel-based methods consist of two ingredients:

1. The (right) **kernel function** to convert input data into a similarity matrix
2. The **algorithm** that uses the kernels and data to produce a model

Kernel-Based Learning

Introduction

Examples of kernel functions:

- RBF
- Polynomial
- String
- Wavelet
- Fisher

Examples of learning algorithms:

- Support Vector Machine (SVM) and Relevance Vector Machine (RVM)
- kernel Fisher Discriminant Analysis (FDA)
- kernel PCA, kernel CCA
- kernel regression (logistic, ridge)
- kernel k -means

Kernel-Based Learning

Introduction

Relationship between the two ingredients:

- The real art is to pick an appropriate kernel. Consider the RBF kernel:

$$k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{1}{2} \frac{\|\mathbf{u} - \mathbf{v}\|^2}{\sigma^2}\right), \quad \sigma^2 \in \mathbb{R}$$

- All information of a problem (besides the target values) is tunneled through the kernel matrix $\mathbf{K} = (k_{ij})$, where $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
- if σ^2 is very small, then $\mathbf{K} \approx \mathbf{I}$ (all data are dissimilar): over-fitting
- if σ^2 is very large, then $\mathbf{K} \approx \mathbf{1}$ (all data are similar): under-fitting

Kernel-Based Learning

Linear regression

Problem: We wish to find a function $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ which best models a data set $S = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\} \subset \mathbb{R}^d \times \mathbb{R}$

- If $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$, just add one dimension $\mathbf{x} \leftarrow (1, \mathbf{x})$; $\mathbf{w} \leftarrow (w_0, \mathbf{w})$
- Call $X_{N \times d}$ the matrix of the \mathbf{x}_n and $\mathbf{t} = (t_1, \dots, t_N)^T$
- If S is generated as $(\mathbf{x}, f(\mathbf{x}))$, the \mathbf{x}_n vectors are linearly independent and $N = d$, then there is a unique solution for $X\mathbf{w} = \mathbf{t}$ given by $\mathbf{w}^* = X^{-1}\mathbf{t}$. In any other case, the problem is ill-posed.

Kernel-Based Learning

Ill-posed & well-posedness

A problem is **ill-posed** if the solution may not always exist, is not uniquely determined or is unstable (small variations in the initial conditions of the problem cause the solution to change quite a lot).



Jacques
Hadamard
(1865-1963)

- Learning problems are in general ill-posed.
- The solution is to use an **inductive principle**; one of the most popular is the **regularization** principle.

Kernel-Based Learning

Ridge regression

Under the standard assumptions $t = f(\mathbf{x}) + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$, ML leads to the minimization of the regularized (= penalized) empirical error:

$$E_\lambda(\mathbf{w}) = \sum_{n=1}^N (t_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \sum_{i=0}^d w_i^2 = (\mathbf{t} - X\mathbf{w})^T (\mathbf{t} - X\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

1. Note $E_\lambda(\mathbf{w}) = \|\mathbf{t} - X\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$
2. The parameter $\lambda > 0$ defines a trade-off between the fit to the data and the complexity of the model (length of \mathbf{w} in this case)

Kernel-Based Learning

Ridge regression

Setting $\nabla E_\lambda(\mathbf{w}) = 0$, we obtain the (regularized) normal equations

$$-2X^T(t - X\mathbf{w}) + 2\lambda\mathbf{w} = 0$$

with solution $\mathbf{w}^* = (X^T X + \lambda I_d)^{-1} X^T t$

and therefore $f(\mathbf{x}) = (\mathbf{w}^*)^T \mathbf{x} = \langle \mathbf{w}^*, \mathbf{x} \rangle = t^T X (X^T X + \lambda I_d)^{-1} \mathbf{x}$.

Since X is $N \times d$, the matrix $X^T X$ is $d \times d$

1. The “model size” does not grow with data size (a **parametric** model)
2. $X^T X + \lambda I_d$ always has an inverse, for all $\lambda > 0$

Kernel-Based Learning

Dual representation

It turns out that the regularized solution can be written as:

$$\mathbf{w}^* = \sum_{n=1}^N \alpha_n \mathbf{x}_n$$

In consequence,

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \mathbf{x}_n, \mathbf{x} \rangle$$

1. The vector of parameters $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)^T$ is $\boldsymbol{\alpha} = (XX^T + \lambda I_N)^{-1} \mathbf{t}$
2. The **Gram matrix** XX^T (“matrix of inner products”) is $N \times N$

Kernel-Based Learning

Primal and dual

So we have the **primal** and the **dual** forms for $f(\mathbf{x})$:

$$f(\mathbf{x}) = \langle \mathbf{w}^*, \mathbf{x} \rangle = \sum_{i=0}^d w_i^* x_i \quad \text{and} \quad f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \mathbf{x}_n, \mathbf{x} \rangle$$

The dual form is (apparently) more convenient when $d \gg N$:

- The primal requires the computation and inversion of $X^T X + \lambda I_d$, requiring $O(Nd^2 + d^3)$ operations
- The dual requires the computation and inversion of $XX^T + \lambda I_N$, requiring $O(dN^2 + N^3)$ operations

Kernel-Based Learning

Key aspects of kernel methods

1. Input data is embedded (mapped) into a vector space
2. Linear relations are sought among the elements of this vector space
3. The coordinates of the images (mapped data) are *not* needed: only their pairwise inner products
4. Often these inner products can be computed (efficiently and implicitly) in the input space (via a kernel function) → the PROMISE

Kernel-Based Learning

Feature maps (1)

How can we perform **non-linear** function regression?

First create a *feature map*, a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$:

$$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_D(\mathbf{x}))^T$$

$\phi(\mathbf{x})$ is called the *feature vector* and $\{\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d\}$ is the *feature space* (possibly part of a larger vector space \mathcal{H}), and typically $D \gg d$.

(as a technicality, we could easily define $\phi_0(\mathbf{x}) = 1$ as before, if need be)

Kernel-Based Learning

Feature maps (2)

The non-linear regression function has now the primal representation:

$$f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{j=1}^D w_j \phi_j(\mathbf{x})$$

- This feature space has the structure of \mathbb{R}^D (a vector space)
- In consequence, there is also the dual representation:

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \phi(\mathbf{x}_n), \phi(\mathbf{x}) \rangle$$

(how general is this result???)

Kernel-Based Learning

Feature maps (3)

We will see that it suffices for \mathcal{H} to have the structure of a **Hilbert space**:

- A vector space endowed with an **inner product** whose associated norm defines a complete metric
- Distances, lengths and angles are well-defined for the elements of the space
- Completeness means that all Cauchy sequences defined in \mathcal{H} converge to an element of \mathcal{H}



David
Hilbert
(1862-1943)

Kernel-Based Learning

Feature maps and kernels (1)

Given a feature map $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$, being \mathcal{H} a Hilbert space, we define its associated **kernel function** $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as:

$$k(\mathbf{u}, \mathbf{v}) = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle, \quad \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$$

One key point is that, for some feature maps, computing $k(\mathbf{u}, \mathbf{v})$ is independent of the dimension of \mathcal{H} (it only depends on d).

→ the PROMISE!!!

Kernel-Based Learning

Feature maps and kernels (2)

Our regression function has now the dual representation:

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \phi(\mathbf{x}_n), \phi(\mathbf{x}) \rangle = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

This dual representation:

- ... is a **non-linear model** (in the input space)
- ... is a **linear model** (in the feature space)

We now have a non-parametric model (complexity grows with data size)

Kernel-Based Learning

Feature maps and kernels (3)

The new vector of parameters $\alpha = (\alpha_1, \dots, \alpha_N)^T$ is now given by

$$\alpha = (\mathbf{K} + \lambda I_N)^{-1} \mathbf{t},$$

where (as introduced earlier) $\mathbf{K} = (k_{ij})$, with $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$.

So we can do ridge regression based only on \mathbf{K} (and throw away X)

Kernel-Based Learning

Kernel ridge regression

What if we take the (simplest) choice $\phi(\mathbf{x}) = \mathbf{x}$? In this case $d = D$ and $k(\mathbf{u}, \mathbf{v}) = \langle \mathbf{u}, \mathbf{v} \rangle$. The regularized solution reads

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \mathbf{x}_n, \mathbf{x} \rangle$$

where

$$\boldsymbol{\alpha} = (X X^T + \lambda I_N)^{-1} \mathbf{t},$$

(so $\mathbf{K} = X X^T$ in this particularly simple case)

This means we have generalized (the dual of) standard ridge regression via a kernel function (we have **kernelized**)

Kernel-Based Learning

Kernelizing ...

Many (classical and new) learning algorithms can be kernelized:

1. They require solving a problem where the data appear in the form of pairwise inner products (or pairwise Euclidean distances)
2. The solution is expressed as a linear combination of the kernel function centered at the data (ideally we only use *some* of the data: **sparsity**)
3. Examples include SVMs, ridge regression, perceptrons, FDA, PLS [supervised], as well as PCA, k-means, Parzen Windows [unsupervised]

Kernel-Based Learning

General feature maps (4)

A feature map is of the general form $\phi : \mathcal{X} \rightarrow \mathcal{H}$

The associated kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is:

$$k(\mathbf{u}, \mathbf{v}) = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle, \quad \mathbf{u}, \mathbf{v} \in \mathcal{X}$$

\mathcal{X} can be any space, \mathcal{H} is always a Hilbert space

In our starting development, $\mathcal{X} = \mathbb{R}^d$

Kernel-Based Learning

Regularization-based learning algorithms

The key for this is the **regularization framework**; consider the regularized empirical error for mapped data:

$$E_{\lambda}(\mathbf{w}) = \sum_{n=1}^N (t_n - \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle)^2 + \lambda \|\mathbf{w}\|^2, \quad \lambda > 0$$

which we now generalize to arbitrary **loss** functions $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$:

$$E_{\lambda}(\mathbf{w}) = \sum_{n=1}^N L(t_n, \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle) + \lambda \|\mathbf{w}\|^2, \quad \lambda > 0$$

Kernel-Based Learning

Regularization-based learning algorithms

Theorem. If L is differentiable w.r.t. its second argument and \mathbf{w}^* is a minimizer of $E_\lambda(\mathbf{w})$, then we can represent:

$$\mathbf{w}^* = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$$

and therefore

$$f(\mathbf{x}) = \langle \mathbf{w}^*, \phi(\mathbf{x}) \rangle = \sum_{j=1}^D w_j^* \phi_j(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}), \quad \alpha = (\mathbf{K} + \lambda I_N)^{-1} \mathbf{t}$$

This result is usually called the **Representer Theorem**.