

# Machine Learning

## Learning with kernels

Lluís A. Belanche

`belanche@cs.upc.edu`



Soft Computing Research Group  
Dept. de Ciències de la Computació (Computer Science)  
Universitat Politècnica de Catalunya

# Learning with kernels (I): The SVM

## Linear regression revisited

- **Problem:** We wish to find a function  $y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$  which best models a data set  $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\} \subset \mathbb{R}^d \times \mathbb{R}$
- Then we minimize the regularized (aka penalized) empirical error:

$$E_{\text{emp}}^\lambda(y) = \sum_{n=1}^N \left(t_n - y(\mathbf{x}_n)\right)^2 + \lambda \sum_{i=1}^d w_i^2 = \|\mathbf{t} - X\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

The parameter  $\lambda > 0$  defines a trade-off between the fit to the data and the complexity of the vector  $\mathbf{w}$

# Learning with kernels (I): The SVM

## Linear regression revisited

Setting  $\frac{\partial E_{\text{emp}}^\lambda(y)}{\partial \mathbf{w}} = 0$ , we obtain the (regularized) normal equations:

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

with solution

$$\hat{\mathbf{w}} = (X^\top X + \lambda I_d)^{-1} X^\top t$$

and therefore

$$y(\mathbf{x}) = \hat{\mathbf{w}}^\top \mathbf{x}$$

# Learning with kernels (I): The SVM

## Linear regression revisited

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

$$\hat{\mathbf{w}} = \sum_{n=1}^N \alpha_n \mathbf{x}_n, \quad \hat{\mathbf{w}} = \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \\ \vdots \\ \hat{w}_d \end{pmatrix}$$

$$y(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\mathbf{x}_n^\top \mathbf{x}), \quad \boldsymbol{\alpha} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{pmatrix}$$

The new vector of parameters is given by  $\boldsymbol{\alpha} = (X X^\top + \lambda I_N)^{-1} \mathbf{t}$

# Learning with kernels (I): The SVM

## Linear regression revisited

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

$$y(\mathbf{x}) = \hat{\mathbf{w}}^\top \mathbf{x} \quad \text{and} \quad y(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\mathbf{x}_n^\top \mathbf{x})$$

The dual form is usually more convenient when  $d \gg N$ :

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

# Learning with kernels (I): The SVM

## Key aspects of kernel methods

How can we achieve non-linear regression?

A **feature map** is a function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^M$ :

$$\phi(\mathbf{x}) = \left( \phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_M(\mathbf{x}) \right)^\top$$

- $\phi(\mathbf{x})$  is called the **feature vector**
- $\{\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d\}$  is the **feature space** (FS), and typically  $M \gg d$ .

# Learning with kernels (I): The SVM

## Key aspects of kernel methods

- Define  $\Phi_{N \times M}$  the matrix of the  $\phi(\mathbf{x}_n)$  as

$$\phi_{nm} = \phi_m(\mathbf{x}_n), n = 1, \dots, N, m = 1, \dots, M.$$

- Suppose we perform ridge regression on the  $\Phi$  matrix
- The new regression function has the **primal** representation:

$$y(\mathbf{x}) = \hat{\mathbf{w}}^\top \phi(\mathbf{x})$$

Note the primal now (explicitly) operates in feature space

# Learning with kernels (I): The SVM

## Key aspects of kernel methods

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

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

- The feature space where  $k$  *implicitly* operates is  $\mathbb{R}^M$
- For some feature maps, computing  $k(\mathbf{u}, \mathbf{v})$  is independent of  $M$



# Learning with kernels (I): The SVM

## Key aspects of kernel methods

Since  $\hat{\mathbf{w}} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$ , the new regression function has the **dual** representation:

$$y(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\phi(\mathbf{x}_n)^\top \phi(\mathbf{x})) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

The new vector of parameters is given by

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda I_N)^{-1} \mathbf{t}, \quad \text{where } \mathbf{K} = (k(\mathbf{x}_n, \mathbf{x}_m))$$

# Learning with kernels (I): The SVM

## Key aspects of kernel methods

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

- The Support Vector Machine (SVM) and the Relevance Vector Machine (RVM)
- Fisher Discriminant Analysis (KFDA), Principal Components Analysis (KPCA), Canonical Correlation Analysis (KCCA), ...
- Kernel (regularized) linear regression
- Kernel k-means, kernel kNN
- (less known or very recent): PLS, Parzen Windows, logistic regression, statistical tests, ...

# Learning with kernels (I): The SVM

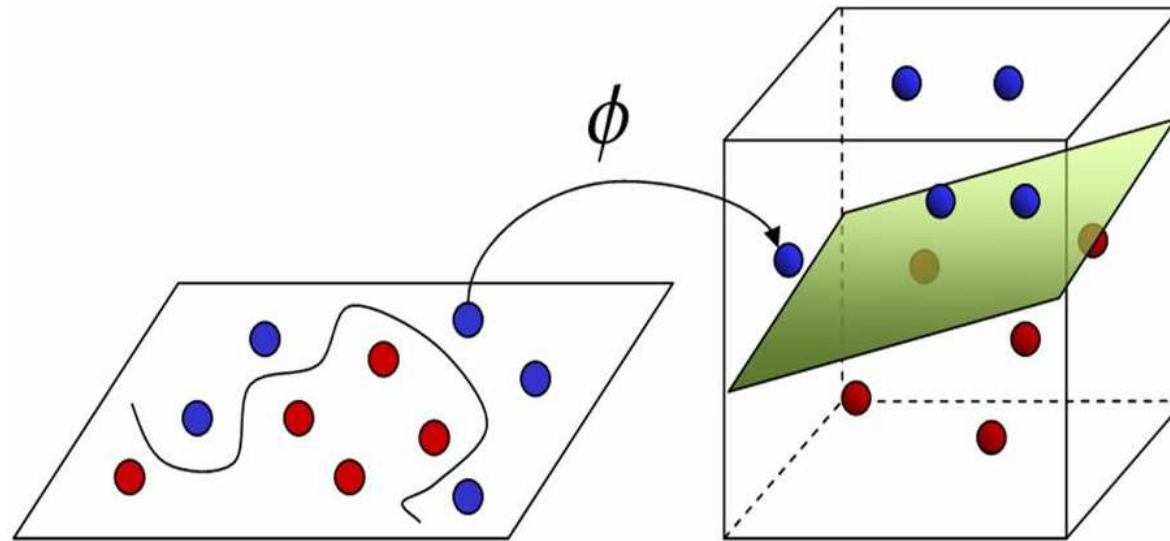
## Key aspects of kernel methods

Kernel-based methods consist of two ingredients:

1. The kernel function (this is non-trivial)
2. The algorithm taking kernels as input
  - Data items are embedded into a vector space (feature space FS)
  - Linear relations are sought among the elements of the FS
  - The coordinates of these images are not needed: only their pairwise inner products
  - These inner products can sometimes be computed efficiently and implicitly in the input space (kernel function)
  - The solution vector is expressed as a linear combination of the kernel centered at the data

# Learning with kernels (II): Kernel functions

A kernel function implicitly defines a map  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  from an input space of objects  $\mathcal{X}$  into some Hilbert space  $\mathcal{H}$  (called the *feature space*). The “kernel trick” consists in performing the mapping and the inner product simultaneously by defining its associated kernel function:



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

# Learning with kernels (II): Kernel functions

## The Kernel Trick

- Suppose we take  $k(\mathbf{u}, \mathbf{v}) = \langle \mathbf{u}, \mathbf{v} \rangle^d$  (a simple choice).
- What is the underlying mapping  $\phi$  here?  
 $\implies$  Answer: this choice of kernel corresponds to a map  $\phi$  leading into the space spanned by all products of exactly  $d$  dimensions of  $\mathbb{R}^n$ .
- Let us take, for instance,  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^2$ , and take  $d = 2$ :

$$\begin{aligned} k(\mathbf{u}, \mathbf{v}) &= \langle \mathbf{u}, \mathbf{v} \rangle^2 = \left[ \left\langle \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}, \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} \right\rangle \right]^2 \\ &= (u_1 v_1 + u_2 v_2)^2 = (u_1 v_1)^2 + 2u_1 v_1 u_2 v_2 + (u_2 v_2)^2 \\ &= u_1^2 v_1^2 + (\sqrt{2} u_1 u_2)(\sqrt{2} v_1 v_2) + u_2^2 v_2^2 \\ &= \left\langle \begin{pmatrix} u_1^2 \\ \sqrt{2} u_1 u_2 \\ u_2^2 \end{pmatrix}, \begin{pmatrix} v_1^2 \\ \sqrt{2} v_1 v_2 \\ v_2^2 \end{pmatrix} \right\rangle = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle \end{aligned}$$

- Therefore,  $\phi : \mathbb{R}^2 \longrightarrow \mathbb{R}^3$  with  $\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$

# Kernel-Based Learning

## Key aspects of kernel methods

- 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  $k(\mathbf{u}, \mathbf{v}) = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle_{\mathcal{H}}$ ,  $\mathbf{u}, \mathbf{v} \in \mathcal{X}$
- $\mathcal{X}$  can be any space,  $\mathcal{H}$  is any **Hilbert space**:
  - An abstract complete vector space possessing the structure of an inner product
  - Examples would be  $\mathbb{R}^M$  or the  $l_2$  space of square-summable sequences

---

In our previous discussion,  $\mathcal{X} = \mathbb{R}^d$  and  $\mathcal{H} = \mathbb{R}^M$

# Learning with kernels (II): Kernel functions

## Characterization of Kernels

**Definition:** A symmetric function  $k$  is called **positive semi-definite** in  $\mathcal{X}$  if:

for every  $N \in \mathbb{N}$ , and every choice  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}$ ,

the matrix  $\mathbf{K} = (k_{ij})$ , where  $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  is **positive semi-definite**.

---

**Theorem:**  $k$  admits the existence of a map  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  s.t.

$\mathcal{H}$  is a Hilbert space and  $k(\mathbf{u}, \mathbf{v}) = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle_{\mathcal{H}}$

if and only if  $k$  is a positive semi-definite symmetric function in  $\mathcal{X}$ .

# Learning with kernels (II): Kernel functions

## On positive semi-definiteness

There are many equivalent characterizations of the psd (*positive semi-definite*) property for real symmetric matrices. Here are some:

1.  $A_{d \times d}$  is psd if and only if all of its eigenvalues are non-negative.
2.  $A_{d \times d}$  is psd if and only if the determinants of all of its leading principal minors are non-negative.
3.  $A_{d \times d}$  is psd if and only if there is a psd matrix  $B$  such that  $BB^T = A$  (this matrix  $B$  is unique and called the square root of  $A$ ).
4.  $A_{d \times d}$  is psd if and only if,  $\forall \mathbf{c} \in \mathbb{R}^d$ ,  $\mathbf{c}^T A \mathbf{c} \geq 0$ .



# Learning with kernels (II): Kernel functions

## General linear kernel

If  $A_{d \times d}$  is a psd matrix, then the function  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  given by  $k(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T A \mathbf{v}$  is a kernel.

*Proof.* Since  $A$  is psd we can write it in the form  $A = BB^T$ . For every  $N \in \mathbb{N}$ , and every choice  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ , we form the matrix  $K = (k_{ij})$ , where  $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T A \mathbf{x}_j$ . Then for every  $\mathbf{c} \in \mathbb{R}^N$ :

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j k_{ij} = \sum_{i=1}^N \sum_{j=1}^N c_i c_j \mathbf{x}_i^T A \mathbf{x}_j = \sum_{i=1}^N \sum_{j=1}^N c_i c_j (B^T \mathbf{x}_i)^T (B^T \mathbf{x}_j)$$

$$= \left\| \sum_{i=1}^N c_i (B^T \mathbf{x}_i) \right\|^2 \geq 0. \quad \text{Note that } \phi(\mathbf{x}) = B^T \mathbf{x}$$

# Learning with kernels (II): Kernel functions

## Summary of kernel properties

If  $k, k'$  are kernels on  $\mathbb{R}^d$ ,  $k''$  is a kernel on  $\mathbb{R}^N$ ,  $a, b \geq 0$ , and  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^N$ , then the following functions are kernels on  $\mathbb{R}^d$ :

1.  $k(u, v) + k'(u, v)$

2.  $ak(u, v) + b$

3.  $k(u, v)k'(u, v)$

4.  $k''(\phi(u), \phi(v))$

# Learning with kernels (II): Kernel functions

## Normalization

If  $k$  is a kernel, then so is:

$$k_n(\mathbf{u}, \mathbf{v}) := \frac{k(\mathbf{u}, \mathbf{v})}{\sqrt{k(\mathbf{u}, \mathbf{u})} \sqrt{k(\mathbf{v}, \mathbf{v})}}$$

Moreover,

$$\begin{aligned} |k_n(\mathbf{u}, \mathbf{v})| &\leq 1 \\ k_n(\mathbf{u}, \mathbf{u}) &= 1 \end{aligned}$$

# Learning with kernels (II): Kernel functions

## Polynomial combinations

1. If  $k$  is a kernel and  $p$  is a polynomial of degree  $m$  with positive coefficients, then the function

$$k_p(u, v) = p(k(u, v))$$

is also a kernel.

2. The special case where  $k$  is linear and  $p(z) = (az + 1)^m$  leads to the so-called **polynomial kernel**

# Learning with kernels (II): Kernel functions

## Polynomial combinations

Consider the kernel family:

$$\left\{ k_i(\mathbf{u}, \mathbf{v}) = \alpha_i \left( \langle \mathbf{u}, \mathbf{v} \rangle + a_i \right)^{\beta_i} \mid \beta_i \in \mathbb{N}, \alpha_i > 0, a_i \geq 0 \right\}$$

For any  $q > 0 \in \mathbb{N}$ ,

$$\sum_{i=0}^q k_i(\mathbf{u}, \mathbf{v})$$

is a kernel.

# Learning with kernels (II): Kernel functions

## Polynomial combinations

Consider the particular case  $a_i = 0$ ,  $\beta_i = i$  and  $\alpha_i = \frac{\alpha^i}{i!}$ , for some real  $\alpha > 0$ , and take the limit  $q \rightarrow \infty$ .

The obtained series is convergent for all  $\alpha$  and the resulting kernel is:

$$\sum_{i=0}^{\infty} \frac{\alpha^i}{i!} (\langle \mathbf{u}, \mathbf{v} \rangle)^i = e^{\alpha \langle \mathbf{u}, \mathbf{v} \rangle}$$

Assume that  $u, v \in \mathbb{R}$ ; then  $\exp(\alpha uv) = \langle \phi(u), \phi(v) \rangle$  with

$\phi(z) = \left( \sqrt{\frac{\alpha^i}{i!}} z^i \right)_{i=0}^{\infty}$ , and therefore we have designed a feature space of infinite dimension!

# Learning with kernels (II): Kernel functions

## Translation invariant and radial kernels

We say that a kernel  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is:

**Translation invariant** if it has the form  $k(\mathbf{u}, \mathbf{v}) = T(\mathbf{u} - \mathbf{v})$ , where  $T : \mathbb{R}^d \rightarrow \mathbb{R}$  is a differentiable function.

**Radial** if it has the form  $k(\mathbf{u}, \mathbf{v}) = t(\|\mathbf{u} - \mathbf{v}\|)$ , where  $t : [0, \infty) \rightarrow [0, \infty)$  is a differentiable function.

Radial kernels fulfill  $k(\mathbf{u}, \mathbf{u}) = t(0)$ .

# Learning with kernels (II): Kernel functions

## The Gaussian kernel

Consider the function  $t(z) = \exp(-\gamma z^2)$ ,  $\gamma > 0$ . The resulting radial kernel is known as the **Gaussian RBF kernel**:

$$k(\mathbf{u}, \mathbf{v}) = \exp(-\gamma \|\mathbf{u} - \mathbf{v}\|^2)$$

Note that some people call it the RBF kernel *par excellence*!

You can also find it as:

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



# Learning with kernels (II): Kernel functions

## Popular choices for the Kernel

**Polynomial kernels** (relation to GLDs)

$$k(\mathbf{u}, \mathbf{v}) = (\langle \mathbf{u}, \mathbf{v} \rangle + 1)^d, \quad d \in \mathbb{N}$$

**Gaussian RBF kernels** (relation to RBFNNs)

$$k(\mathbf{u}, \mathbf{v}) = \exp(-\gamma \|\mathbf{u} - \mathbf{v}\|^2) \quad \gamma > 0 \in \mathbb{R}$$

**Laplacian RBF kernels** (relation to ???)

$$k(\mathbf{u}, \mathbf{v}) = \exp(-\gamma \|\mathbf{u} - \mathbf{v}\|) \quad \gamma > 0 \in \mathbb{R}$$

**Sigmoidal kernels** (relation to MLPs)

$$k(\mathbf{u}, \mathbf{v}) = g(\alpha \langle \mathbf{u}, \mathbf{v} \rangle + \beta)$$

with  $g$  a sigmoidal (e.g., logistic, tanh, ...) and particular choices for  $\alpha, \beta$ .

# Learning with kernels (II): Kernel functions

Euclidean space  $\mathbb{R}^d$ , but not only

- Kernels on sets/bitstrings
- Graph kernels
- Generative kernels (on probability distributions)
- Convolution kernels (on combinatorial structures)
- Tree kernels
- String kernels (text)

... and many others (functional data, categorical data, ...)

# Learning with kernels (II): Kernel functions

## A kernel for set comparison

Given two sets  $A, B$ , consider

$$k(A, B) = \sum_{a \in A} \sum_{b \in B} k_{\text{base}}(a, b)$$

If  $k_{\text{base}}$  is the overlap kernel:

$$k(a, b) = \begin{cases} 1 & \text{if } a = b; \\ 0 & \text{otherwise.} \end{cases}$$

we get  $k(A, B) = |A \cap B|$ . Remarkably,  $k(A, B) = \frac{|A \cap B|}{|A \cup B|}$  is also a kernel.

# Learning with kernels (II): Kernel functions

---

## Kernels for/from graphs (I)

- Consider a graph  $G = (V, E)$ , where the set of vertices (nodes)  $V$  are the data points and  $E$  is the set of edges. Call  $N = |V|$ , the number of nodes
- The idea is to compute a (base) matrix  $B_{N \times N}$  whose entries are the weights of the edges and consider  $B^2 = BB$  ( $B$  need not be symmetric)
- Typical use: **connectivity matrix** of  $G$ : the  $(i, j)$  element of  $B^2$  is the number of paths of length exactly 2 between  $i$  and  $j$

---

Examples:

1. protein-protein interactions
2. people-to-people interactions

---

In 2, the  $(i, j)$  element of  $B^2$  is the number of common friends between data points  $i$  and  $j$  (it can be thought of as a measure of their similarity)

# Learning with kernels (II): Kernel functions

## Kernels for/from graphs (II)

Notes:

1. The entries of  $B$  may be real-valued numbers (e.g., symmetric bounded similarities)
2. Higher powers of  $B$  measure higher-order similarities
3. Only the even powers are guaranteed to be PSD

Consider, for a given  $\lambda \in (0, 1)$ :

$$\sum_{k=0}^{\infty} \frac{1}{k!} \lambda^k B^k = \exp(\lambda B)$$

If  $B$  is symmetric, then  $B = U\Lambda U^T$  is its spectral decomposition, so  $B^2 = (U\Lambda U^T)(U\Lambda U^T) = U\Lambda^2 U^T$ . In general, we have  $B^k = U\Lambda^k U^T$  and therefore:

$$K = \exp(\lambda B) = U \exp(\lambda \Lambda) U^T$$

is an example of a **diffusion** kernel (the name comes from the *heat equation* in physics)

# Learning with kernels (II): Kernel functions

## Afterthoughts

1. Importance of designing kernels that do not constitute explicit inner products between objects, and therefore fully exploit the kernel trick.
2. Possibility of learning the kernel function (or the kernel matrix) from the training data.
3. Need more research for handling special situations –like missing, imprecise or not-applicable (NA) values.
4. Theoretical analyses on the implications of the kernel choice for the success of kernel-based methods.