

# Introduction to machine learning

Masters M2MO & MIDS

Stéphane Gaïffas



## Today

- Kernels
- Kernel SVM
- Kernel regression

## Supervised learning setting

- We observe a training dataset  $D$  of pairs  $(x_i, y_i)$  for  $i = 1, \dots, n$
- Features  $x_i \in \mathbb{R}^d$  and labels  $y_i \in \mathbb{R}$  (regression) or  $y_i \in \{-1, 1\}$  (binary classification)
- Given a features vector  $x \in \mathbb{R}^d$ , we want to predict the label  $y$

## Features engineering

- Given raw features  $x_1, \dots, x_n \in \mathbb{R}^d$ , we can construct **new** features
- For instance, we can add second order polynomials of the features

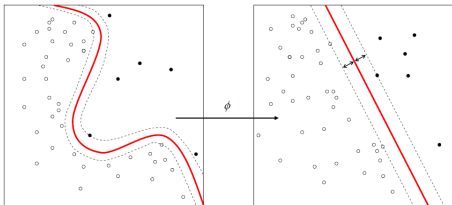
$$x_{i,j}^2, \quad x_{i,j} \times x_{i,k} \quad \text{for any } 1 \leq j, k \leq d$$

- It increases the number of features, hence the dimension of the model weights  $w$  learned from it

## A feature map

- Consider a feature map  $\varphi : \mathbb{R}^d \rightarrow \mathbb{H}$  that "adds" all these new features
- $\mathbb{H}$  is an Hilbert space (eventually infinite dimensional), endowed with an inner product  $\langle \cdot, \cdot \rangle_{\mathbb{H}}$
- The decision boundary  $x \rightarrow \langle w, \varphi(x) \rangle + b = 0$  is **not an hyperplane anymore** (but  $\varphi(x) \rightarrow \langle w, \varphi(x) \rangle + b = 0$  is)

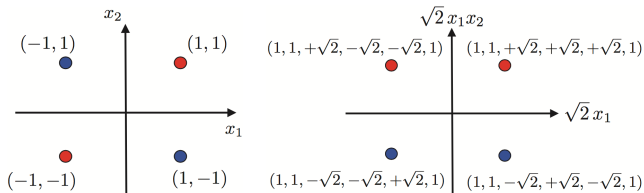
A common belief: **increasing dimension** of features space makes data **almost linearly separable**



The **polynomial** mapping  $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$  for  $x = (x_1, x_2) \in \mathbb{R}^2$

$$\varphi(x) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$$

solves the XOR (Exclusive OR) classification problem



XOR : label  $y_i$  is blue iff one of the coordinates of  $x_i$  equals 1.

- Blue and red points **cannot be linearly separated** in  $\mathbb{R}^2$
- But **they can using the mapping  $\varphi$** , using the hyperplane  $x_1x_2 = 0$

This mapping  $\varphi$  is call **polynomial mapping of order 2**.

Note that for  $x, x' \in \mathbb{R}^2$  we have

$$\begin{aligned}\langle \varphi(x), \varphi(x') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ 1 \end{bmatrix}, \begin{bmatrix} x_1^2 \\ x_1'^2 \\ x_2^2 \\ \sqrt{2}x_1'x_2' \\ \sqrt{2}x_1' \\ \sqrt{2}x_2' \\ 1 \end{bmatrix} \right\rangle \\ &= (x_1x_1' + x_2x_2' + 1)^2 \\ &= (\langle x, x' \rangle + 1)^2\end{aligned}$$

This motivates the definition of

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle = (\langle x, x' \rangle + c)^q$$

where  $q \in \mathbb{N} - \{0\}$  and  $c > 0$ . In this case  $K$  is called the polynomial **kernel** of degree  $q$ .

Given a “raw feature” space  $\mathcal{X}$  (often  $\mathcal{X} = \mathbb{R}^d$ ), a function

$$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

is called a **kernel** over  $\mathcal{X}$ .

**Definition.** We say that a kernel  $K$  is **symmetric** iff

$$K(x, x') = K(x', x)$$

for any  $x, x' \in \mathcal{X}$

**Definition.** We say that a kernel is PDS (positive definite symmetric) iff

- it is symmetric
- for any  $N \in \mathbb{N}$  and any  $\{x_1, \dots, x_N\} \subset \mathcal{X}$  we have

$$\mathbf{K} = [K(x_i, x_j)]_{1 \leq i, j \leq N} \succeq 0$$

meaning that  $\mathbf{K}$  is positive semi-definite (symmetric), or equivalently that

$$u^\top \mathbf{K} u = \sum_{1 \leq i, j \leq N} u_i u_j K(x_i, x_j) \geq 0$$

for any  $u \in \mathbb{R}^N$ , or equivalently that all eigenvalues of  $\mathbf{K}$  are non-negative.

For a sample  $x_1, \dots, x_n$  we call  $\mathbf{K} = [K(x_i, x_j)]_{1 \leq i, j \leq n}$  the **Gram matrix** of this sample.



**Definition.** Hadamard product  $\mathbf{A} \odot \mathbf{B}$  between two matrices  $\mathbf{A}$  and  $\mathbf{B}$  (or vectors) with the same dimensions is given by

$$(\mathbf{A} \odot \mathbf{B})_{i,j} = \mathbf{A}_{i,j} \odot \mathbf{B}_{i,j}$$

**Theorem.** The sum, product, pointwise limit and composition with a power series  $\sum_{n \geq 0} a_n x^n$  with  $a_n \geq 0$  for all  $n \geq 0$  preserves the PDS property.

**Proof.** Consider two  $N \times N$  Gram matrices  $\mathbf{K}, \mathbf{K}'$  of PDS kernels  $K, K'$  and take  $u \in \mathbb{R}^N$ . Observe that

$$u^\top (\mathbf{K} + \mathbf{K}') u = u^\top \mathbf{K} u + u^\top \mathbf{K}' u \geq 0$$

So PDS is preserved by the sum and finite sums by recurrence.

Now, to prove that the product  $\mathbf{K} \odot \mathbf{K}'$  is PDS, write  $\mathbf{K} = \mathbf{M}\mathbf{M}^\top$ , where  $\mathbf{M}$  is the square-root of  $\mathbf{K}$  (which is SDP) and note that

$$\begin{aligned} u^\top (\mathbf{K} \odot \mathbf{K}') u &= \sum_{1 \leq i, j \leq N} u_i u_j \mathbf{K}_{i,j} \mathbf{K}'_{i,j} = \sum_{1 \leq i, j \leq N} \sum_{k=1}^N u_i u_j \mathbf{M}_{i,k} \mathbf{M}_{k,j} \mathbf{K}'_{i,j} \\ &= \sum_{k=1}^N z_k^\top \mathbf{K}' z_k \geq 0 \end{aligned}$$

with  $z_k = u \odot \mathbf{M}_{\bullet, k}$ .

This proves that finite products of PDS kernels is PDS.

Assume that  $K_n \rightarrow K$  as  $n \rightarrow +\infty$  pointwise, where  $K_n$  is a sequence of PDS kernels.

It means that any associated sequence of Gram matrices  $\mathbf{K}_n$  and the its limit  $\mathbf{K}$  satisfies  $\mathbf{K}_n \rightarrow \mathbf{K}$  entrywise, so that for any  $u \in \mathbb{R}^N$  we have

$$u^\top \mathbf{K}_n u \rightarrow u^\top \mathbf{K} u$$

so  $u^\top \mathbf{K} u \geq 0$  since  $u^\top \mathbf{K}_n u \geq 0$  for all  $n$ .

This proves stability of PDS property under pointwise limit.

Now, let  $K$  be a kernel such that  $|K(x, x')| < r$  for all  $x, x' \in \mathcal{X}$  and  $\sum_{n \geq 0} a_n x^n$  a power series with radius of convergence  $r$ .

By stability under sum and product, we have that

$$\sum_{k=0}^N a_n K^n$$

is PDS, and

$$\lim_{N \rightarrow +\infty} \sum_{n=0}^N a_n K^n = \sum_{n \geq 0} a_n K^n$$

remains PDS since PDS is kept under pointwise limit.

This concludes the proof of the theorem.

**Theorem.** The following inequality holds for  $K, K'$  two PSD kernels

$$K(x, x')^2 \leq K(x, x)K(x', x')$$

for any  $x, x' \in \mathcal{X}$ . It is called the **Cauchy-Schwartz inequality** for PSD kernels.

**Proof.** Take  $x, x' \in \mathcal{X}$  and consider the Gram matrix

$$\mathbf{K} = \begin{bmatrix} K(x, x) & K(x, x') \\ K(x', x) & K(x', x') \end{bmatrix}.$$

Since  $K$  is PSD, then  $\mathbf{K} \succeq 0$ , which entails that

$$0 \leq \det \mathbf{K} = K(x, x)K(x', x') - K(x, x')^2$$

**Theorem** [Reproducing kernel Hilbert space]. Let  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a PDS kernel. Then, there is a Hilbert space  $\mathbb{H}$  endowed with an inner product  $\langle \cdot, \cdot \rangle$  and a mapping  $\varphi : \mathcal{X} \rightarrow \mathbb{H}$  such that

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle$$

and such that the **reproducing property** holds:

$$h(x) = \langle h, K(x, \cdot) \rangle$$

for any  $h \in \mathbb{H}$  and  $x \in \mathcal{X}$ .

**Proof.** Available on the moodle. Think of  $\mathbb{H}$  as containing limits of functions

$$\sum_{i=1}^N a_i K(x_i, \cdot)$$

for any  $a_1, \dots, a_N \in \mathbb{R}$  and  $x_1, \dots, x_N \in \mathcal{X}$ .

**Remark.** Stresses the fact that a PDS kernel is some kind of similarity measure, since it is actually an inner product

- We say that  $\mathbb{H}$  is a **reproducing kernel Hilbert space** associated to the kernel  $K$ .
- The Hilbert space  $\mathbb{H}$  is called the **features space** associated to  $K$
- The corresponding mapping  $\varphi : \mathcal{X} \rightarrow \mathbb{H}$  is called the **features mapping**
- $\mathbb{H}$  is endowed with an inner product  $\langle h, h' \rangle$  for  $h, h' \in \mathbb{H}$  and a norm  $\|h\| = \sqrt{\langle h, h \rangle}$
- The feature space might is not unique in general

### In summary

- Choose a kernel  $K$  you think relevant, if it's PDS, then there is a mapping  $\varphi$  and a RKHS  $\mathbb{H}$  for it
- Feature engineering becomes kernel engineering with kernel methods

**Definition.** The **normalized kernel**  $K'$  associated to a kernel  $K$  is given by

$$K'(x, x') = \frac{K(x, x')}{\sqrt{K(x, x)K(x', x')}}}$$

if  $K(x, x)K(x', x') > 0$  and  $K(x, x') = 0$  otherwise.

**Theorem.** If  $K$  is a PDS kernel, its normalized kernel  $K'$  is PDS.

**Remark.** We have that  $K(x, x')$  is the cosine of the angle between  $\varphi(x)$  and  $\varphi(x')$  if  $K$  is a normalized kernel (if none is zero).  
Once again,  $K(x, x')$  is a similarity measure between  $x$  and  $x'$



**Proof.** Let  $x_1, \dots, x_N \in \mathcal{X}$  and  $c \in \mathbb{R}^N$ . If  $K(x_i, x_i) = 0$  or  $K(x_j, x_j) = 0$  then  $K(x_i, x_j) = 0$  using Cauchy-Schwartz, so  $K'(x_i, x_j) = 0$ .

So, we can assume  $K(x_i, x_i) > 0$  for all  $i = 1, \dots, N$  and write the following:

$$\begin{aligned} \sum_{1 \leq i, j \leq N} \frac{c_i c_j K(x_i, x_j)}{\sqrt{K(x_i, x_i) K(x_j, x_j)}} &= \sum_{1 \leq i, j \leq N} \frac{c_i c_j \langle \varphi(x_i), \varphi(x_j) \rangle}{\|\varphi(x_i)\| \|\varphi(x_j)\|} \\ &= \left\| \sum_{i=1}^N \frac{c_i \varphi(x_i)}{\|\varphi(x_i)\|} \right\|^2 \geq 0 \end{aligned}$$

which proves the theorem.

**Remark.** If  $K$  is a normalized kernel, then

$$\|\varphi(x)\| = \langle \varphi(x), \varphi(x) \rangle = K(x, x) = 1$$

for any  $x \in \mathcal{X}$

**The polynomial kernel.** For  $c > 0$  and  $q \in \mathbb{N} - \{0\}$  we define the polynomial kernel

$$K(x, x') = (\langle x, x' \rangle + c)^q.$$

It is a PDS kernel

**Proof.** It is the power of the PDS kernel  $(x, x') \mapsto \langle x, x' \rangle + b$ .

We already computed its mapping  $\varphi(x)$ : it contains all the monomials of degree less than  $q$  of the coordinates of  $x$

**The RBF kernel** (Radial Basis Function). For  $\gamma > 0$  it is given by

$$K(x, x') = \exp(-\gamma \|x - x'\|_2^2)$$

**Theorem.** The RBF kernel is a PDS and normalized kernel.

**Proof.** First remark that

$$\begin{aligned} \exp(-\gamma \|x - x'\|_2^2) &= \frac{\exp(2\gamma \langle x, x' \rangle)}{\exp(\gamma \|x\|_2^2) \exp(\gamma \|x'\|_2^2)} \\ &= \frac{K'(x, x')}{\sqrt{K'(x, x) K'(x', x')}} \end{aligned}$$

with  $K'(x, x') = \exp(2\gamma \langle x, x' \rangle)$  and that  $K'$  is PDS since

$$K'(x, x') = \sum_{n \geq 0} \frac{(2\gamma \langle x, x' \rangle)^n}{n!}$$

namely a series of the PDS kernel  $(x, x') \mapsto 2\gamma \langle x, x' \rangle$ .

**The tanh kernel.** Also called the sigmoid kernel

$$K'(x, x') = \tanh(a\langle x, x' \rangle + c) = \frac{e^{a\langle x, x' \rangle + c} - e^{-a\langle x, x' \rangle - c}}{e^{a\langle x, x' \rangle + c} + e^{-a\langle x, x' \rangle - c}}$$

for  $a, c > 0$ . It is again a PDS kernel (same argument as for the RBF kernel).

**Remark.** By far, the RBF kernel is the most widely used: uses as a similarity measure the Euclidean norm More mathematical details are covered in the exercises.

**Kernel based algorithms** how to use kernels for classification and regression?

- Let's recall the primal and dual formulation of the SVM

**Linear SVM.** Primal problem is

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, s \in \mathbb{R}^n} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n s_i$$

subject to  $y_i(\langle x_i, w \rangle + b) \geq 1 - s_i$  and  $s_i \geq 0$  for all  $i = 1, \dots, n$

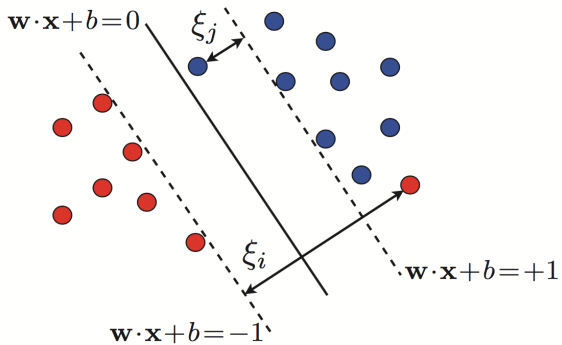
or equivalently

$$\operatorname{argmin}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \ell(y_i, \langle x_i, w \rangle + b)$$

where  $\ell(y, y') = \max(0, 1 - yy') = (1 - yy')_+$  is the hinge loss

Label prediction given by

$$y = \operatorname{sgn}(\langle x, w \rangle + b)$$



**Kernel SVM:** replace  $x_i$  by  $\varphi(x_i)$ . In the primal this leads to

$$\operatorname{argmin}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \ell(y_i, \langle \varphi(x_i), w \rangle + b)$$

Label prediction is given by

$$y = \operatorname{sgn}(\langle \varphi(x), w \rangle + b)$$

In the primal, you need to compute  $\varphi(x)$ !

Dual problem is

$$\max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

subject to  $0 \leq \alpha_i \leq C$  and  $\sum_{i=1}^n \alpha_i y_i = 0$  for all  $i = 1, \dots, n$

and the label prediction using dual variables

$$x \mapsto \operatorname{sgn}(\langle w, x \rangle + b) = \operatorname{sgn}\left(\sum_{i=1}^n \alpha_i y_i \langle x, x_i \rangle + b\right)$$

depends only on the features  $x_i$  via their inner products  $\langle x_i, x_j \rangle$



**Fundamental remark.** The dual problem depends only on the features via their inner products

Given some kernel  $K$ , let's replace the "raw" inner products  $\langle x_i, x_j \rangle$  by the "new" inner products  $K(x_i, x_j) = \langle \varphi(x_i), \varphi(x_j) \rangle$

**The kernel trick.** Once again, to train the SVM with a kernel, you don't need to know or compute the  $\varphi(x_i)$

## The kernel SVM

$$\max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

subject to  $0 \leq \alpha_i \leq C$  and  $\sum_{i=1}^n \alpha_i y_i = 0$  for all  $i = 1, \dots, n$

and the label prediction using dual variables

$$x \mapsto \text{sgn} \left( \sum_{i=1}^n \alpha_i y_i K(x, x_i) + b \right)$$

with the intercept given by

$$b = y_i - \sum_{j=1}^n \alpha_j y_j K(x_j, x_i)$$

for any  $i$  such that  $0 < \alpha_i < C$  (cf previous lecture)

This proves that the hypothesis solution writes

$$h(x) = \text{sgn} \left( \sum_{i=1}^n \alpha_i y_i K(x, x_i) + b \right),$$

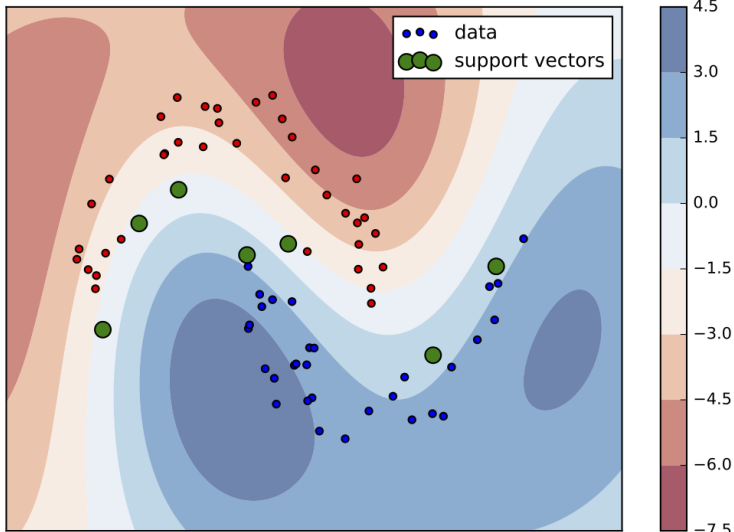
namely a combination of functions  $K(x_i, \cdot)$  where  $x_i$  are the support vectors.

For the RBF kernel, the decision function is

$$x \mapsto \sum_{i: \alpha_i \neq 0} \alpha_i y_i \exp \left( -\gamma \|x - x_i\|_2^2 \right) + b$$

It is a mixture of Gaussian “densities”. Let’s recall that the  $x_i$  with  $\alpha_i \neq 0$  are the support vectors

$$x \mapsto \sum_{i:\alpha_i \neq 0} \alpha_i y_i \exp(-\gamma \|x - x_i\|_2^2) + b$$



The kernel trick is not only for the SVM

**Representer theorem.** If  $K$  is a PDS kernel and  $\mathbb{H}$  its corresponding RKHS, we have that for any increasing function  $g$  and any function  $L : \mathbb{R}^n \rightarrow \mathbb{R}$  that the optimization problem

$$\operatorname{argmin}_{h \in \mathbb{H}} g(\|h\|) + L(h(x_1), \dots, h(x_n))$$

admits only solutions of the form

$$h = \sum_{i=1}^n \alpha_i K(x_i, \cdot).$$

## Kernel ridge regression.

- Consider this time a continuous label  $y_i \in \mathbb{R}$ , features  $x_i \in \mathcal{X}$  for  $i = 1, \dots, n$  and a features mapping  $\varphi : \mathcal{X} \rightarrow \mathbb{H}$  with PDS kernel  $K$
- Kernel ridge regression considers the problem

$$\operatorname{argmin}_w \left\{ \sum_{i=1}^n \ell(y_i, \langle w, \varphi(x_i) \rangle) + \frac{\lambda}{2} \|w\|_2^2 \right\}$$

where  $\lambda$  is a penalization parameter, and  $\ell(y, y') = \frac{1}{2}(y - y')^2$  is the least-squares loss

- Can be written as

$$\operatorname{argmin}_w F(w) \quad \text{with} \quad F(w) = \|y - \mathbf{X}w\|_2^2 + \lambda \|w\|_2^2$$

with  $\mathbf{X}$  the matrix with rows containing the  $\varphi(x_i)$  and  $y = [y_1 \cdots y_n] \in \mathbb{R}^n$

- This problem is strongly convex, and admits a global minimum iff

$$\nabla F(w) = 0 \quad \text{namely} \quad (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})w = \mathbf{X}^T y$$

- Note that  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is always invertible. Thus kernel ridge allows admits a closed-form solution
- Requires to solve a  $D \times D$  linear system, where  $D$  is the dimension of  $\mathbb{H}$
- What if  $D$  is large ?
- Let's us the kernel trick, as we did for SVM

- Representer theorem says that we can find  $\alpha$  such that

$$h(x) = \langle w, \varphi(x) \rangle = \sum_{i=1}^n \alpha_i K(x_i, x) = \sum_{i=1}^n \alpha_i \langle \varphi(x_i), \varphi(x) \rangle$$

for any  $x \in \mathcal{X}$

- This means that

$$w = \mathbf{X}^\top \alpha$$

Now, use the following trick: for any matrix  $\mathbf{X}$ , we have

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1}$$

This entails

$$w = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top y = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} y$$

which gives (note that  $(\mathbf{X} \mathbf{X}^\top)_{i,j} = \langle \varphi(x_i), \varphi(x_j) \rangle = K(x_i, x_j)$ )

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} y$$



**Proof** of the trick. Note that

$$(\mathbf{X}^T \mathbf{X} + \lambda I) \mathbf{X}^T = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda I).$$

Multiplying on the left by  $(\mathbf{X}^T \mathbf{X} + \lambda I)^{-1}$  leads to

$$\mathbf{X}^T = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda I).$$

and then on the right by  $(\mathbf{X} \mathbf{X}^T + \lambda I)^{-1}$  concludes with

$$(\mathbf{X} \mathbf{X}^T + \lambda I)^{-1} \mathbf{X}^T = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T$$

A cute trick. But let's do it like we did for the SVMs  
(just to be sure...)

An alternative formulation of

$$\min_w \sum_{i=1}^n (y_i - \langle w, \varphi(x_i) \rangle)^2 + \lambda \|w\|_2^2$$

is given by

$$\min_w \sum_{i=1}^n (y_i - \langle w, \varphi(x_i) \rangle)^2 \quad \text{subject to} \quad \|w\|_2^2 \leq r^2$$

and also

$$\min_w \sum_{i=1}^n s_i^2 \quad \text{subject to} \quad \|w\|_2^2 \leq r^2 \quad \text{and} \quad s_i = y_i - \langle w, \varphi(x_i) \rangle$$

Which leads to the following Lagrangian

$$L(w, s, \alpha, \lambda) = \min_w \sum_{i=1}^n s_i^2 + \min_w \sum_{i=1}^n \alpha_i (y_i - s_i - \langle w, \varphi(x_i) \rangle) \\ + \lambda (\|w\|_2^2 - r^2)$$

so that the KKT conditions leads to the following properties:

$$\nabla_w L = - \sum_{i=1}^n \alpha_i \varphi(x_i) + 2\lambda w \Rightarrow w = \frac{1}{2\lambda} \sum_{i=1}^n \alpha_i \varphi(x_i)$$

$$\nabla_{s_i} L = 2s_i - \alpha_i \Rightarrow s_i = \alpha_i/2$$

and the slackness complementary conditions:

$$\alpha_i (y_i - s_i - \langle w, \varphi(x_i) \rangle) = 0 \quad \text{and} \quad \lambda (\|w\|_2^2 - r^2) = 0$$

Plugging the expressions of  $w$  and  $s_i$  in functions of  $\alpha$  in  $L$  gives after some algebra the dual objective

$$D(\alpha) = -\lambda \sum_{i=1}^n \alpha_i^2 + 2 \sum_{i=1}^n \alpha_i y_i - \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j \langle \varphi(x_i), \varphi(x_j) \rangle - \lambda r^2$$

(where we replaced  $2\lambda\alpha_i$  by  $\alpha_i$ ) which can be written matricially as

$$\begin{aligned} D(\alpha) &= -\lambda \|\alpha\|_2^2 + 2\langle \alpha, y \rangle - \alpha^\top \mathbf{X}\mathbf{X}^\top \alpha \\ &= 2\langle \alpha, y \rangle - \alpha^\top (\mathbf{K} + \lambda \mathbf{I}) \alpha \end{aligned}$$

with optimum achieved for

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} y$$

(same as before, of course...)

## In summary

- Solving a problem in the dual benefits from the kernel trick
- Allows to construct complex non-linear decision functions
- OK if  $n$  is not too large... (if the  $n \times n$  Gram matrix  $\mathbf{K}$  fits in memory)
- Otherwise, stick to the primal! (and forget about kernels...)
- But don't forget about feature engineering (yes, again !)

**Thank you!**