# Machine Learning 2 – MAP569

Stéphane Gaïffas



ÉCOLE
**POLYTECHNIQUE**
UNIVERSITÉ PARIS-SACLAY

**Today**

- Kernels
- Kernel SVM
- Kernel regression

**Supervised learning setting**

- We observe a training dataset $D$ of pairs $(x_i, y_i)$ for $i = 1, \ldots, n$
- 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 $x_1, \ldots, x_n \in \mathbb{R}^d$, we can construct **new** features
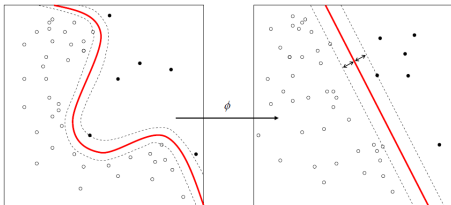- For instance, we can add second order polynomials of the features

$$x_j^2, x_j x_k \quad \text{for any} \quad 1 \leq j, k \leq d$$

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

A **feature map**

- Consider a feature map $\varphi : \mathbb{R}^d \to \mathbb{H}$ 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 $x \to \langle w, \varphi(x) \rangle + b = 0$ is **not an hyperplane anymore** (but $\varphi(x) \to \langle w, \varphi(x) \rangle + b = 0$ is)
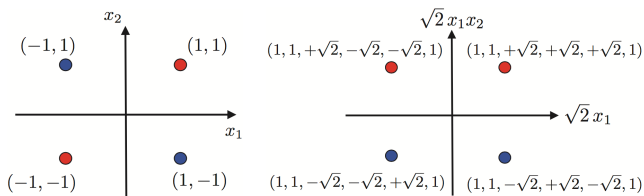
A common belief: **increasing dimension** of features space makes data **almost linearly separable**

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

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

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 $x_1 x_2 = 0$

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

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

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

$$
= (x_1 x_1' + x_2 x_2' + 1)^2
$$

$$
= (\langle x, x' \rangle + 1)^2
$$

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

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

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 symetric) iff

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

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

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

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

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

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

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

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

**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 $\boldsymbol{K}, \boldsymbol{K}'$ of PDS kernels $K, K'$ and take $u \in \mathbb{R}^N$. Observe that

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

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

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

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

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

This proves that finite products of PDS kernels is PDS.

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

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

$$u^\top \boldsymbol{K}_n u \to u^\top \boldsymbol{K} u$$

so $u^\top \boldsymbol{K} u \geq 0$ since $u^\top \boldsymbol{K}_n u \to u$ 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

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

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 PDS 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

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

Since $K$ is PDS, then $\boldsymbol{K} \succeq 0$, which entails that

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

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

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

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

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

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

**Proof**. Available on the moodle

**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 **reproducting kernel Hilbert space** associated to the kernel $K$.
- The Hilbert space $\mathbb{H}$ is called the **features space** associated to $K$
- The corresponding mapping $\varphi : \mathcal{X} \to \mathbb{H}$ 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 engineernig 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 $x_1, \ldots, x_N \in \mathcal{X}$ 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 $i = 1, \ldots, N$ and write the following:

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

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

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

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 Don't worry, you will compute its mapping in PC today :)

**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 $\quad y_i(\langle x_i, w \rangle + b) \geq 1 - s_i \quad$ and $\quad s_i \geq 0 \quad$ for all $\quad i = 1, \ldots, n$
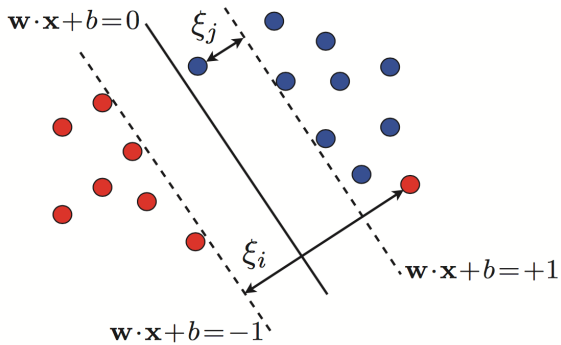
or equivalently

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

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

Label prediction given by

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

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

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

Label prediction is given by

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

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

Dual problem is

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

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

and the label prediction using dual variables

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

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**

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

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

and the label prediction using dual variables

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

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

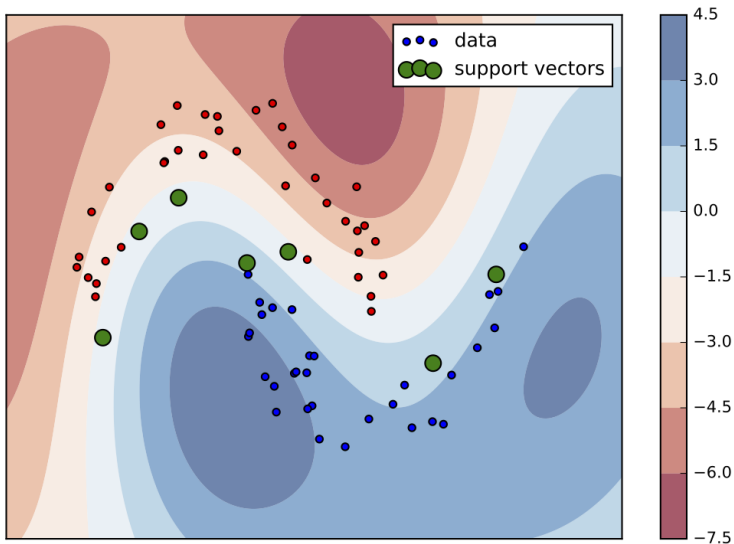$$h(x) = \operatorname{sgn}\Big( \sum_{i=1}^{n} \alpha_i y_i K(x, x_i) + b \Big),$$

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

For the RBF kernel, the decision function is

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

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

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

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 $L : \mathbb{R}^n \to \mathbb{R}$ that the optimization problem

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

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 $i = 1, \ldots, n$ and a features mapping $\varphi : \mathcal{X} \to \mathbb{H}$ 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

$$\operatorname*{argmin}_{w} F(x) \quad \text{with} \quad F(w) = \|y - \boldsymbol{X}w\|_2^2 + \lambda\|w\|_2^2$$

with $\boldsymbol{X}$ 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

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

- Note that $\boldsymbol{X}^\top \boldsymbol{X} + \lambda \boldsymbol{I}$ 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

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

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

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

This entails

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

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

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

**Proof** of the trick. Note that

$$(\boldsymbol{X}^\top \boldsymbol{X} + \lambda \boldsymbol{I})\boldsymbol{X}^\top = \boldsymbol{X}^\top(\boldsymbol{X}\boldsymbol{X}^\top + \lambda \boldsymbol{I}).$$

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

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

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

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

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

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

and also

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

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:

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

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

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

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

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

with optimum achieved for

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

(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 $\boldsymbol{K}$ fits in memory)
- Otherwise, stick to the primal! (and forget about kernels...)
- But don't forget about feature engineering (yes, again !)

**Next week.** We have seen a lot of problem of the form

$$\underset{w}{\arg\min} \, f(w) + g(w)$$

with $f$ a goodness-of-fit function

$$f(w) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \langle w, x_i \rangle)$$

where $\ell$ is some loss and

$$g(w) = \frac{1}{C} \, \text{pen}(w)$$

where pen is some penalization function, examples being
$\text{pen}(w) = \frac{1}{2}\|w\|_2^2$ (ridge) and $\text{pen}(w) = \|w\|_1$ (Lasso)

Next week we'll learn how to solve this kind of problems using **amazing** optimization algorithms

# Thank you!