

# Lecture 19: FOSML

①

- We will relate kernel methods to SVM.
- Since SVM's are only about hyperplane projection, which are inner products, there is intuition that kernels should adapt naturally to this setting. This is because kernels can be realized as inner products, admittedly in some (perhaps infinite) high-dim. Hilbert space.

• Consider the following optimization problem:

$$\textcircled{*} \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j) \text{ subject to}$$

- $0 \leq \alpha_i \leq C, \forall i=1, \dots, m$
- $\sum_{i=1}^m \alpha_i y_i = 0,$

• Recalling some SVM results, we may write a hypothesis solution  $h$  as

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

where  $b = y_i - \sum_{j=1}^m \alpha_j y_j K(x_j, x_i)$  for any  $x_i$  with  $\alpha_i \in (0, C)$ .

• Rewriting  $\textcircled{*}$  in vector form using the kernel matrix  $K \in \mathbb{R}^{m \times m}$  associated to

$K$  and the sample  $\{x_i\}_{i=1}^m$  yields

$$\max_{\alpha} 2 \mathbb{1}^T \alpha - (\alpha \circ y)^T K (\alpha \circ y) \text{ subject to}$$

- $0 \leq \alpha \leq C$
- $\alpha^T y = 0$

Here,  $\alpha \circ y = (\alpha_1 y_1, \dots, \alpha_m y_m)$  is the Hadamard product of  $\alpha, y$ .

In vector notation, the solution is the same but with  $b = y_i - (\alpha_i y)^T \underbrace{K e_i}_{(0, \dots, 0, 1, \dots, 0)}$   
 for any  $x_i$  with  $0 < \alpha_i < C$  ↑  
i<sup>th</sup> coordinate

So, the solution can be written as a linear combination of kernel evaluations, at least if we ignore the offsets  $b$ .

This is a general property of a range of optimization problems:

Theorem (Representer Theorem): Let  $K: X \times X \rightarrow \mathbb{R}$  be a PDS kernel, with  $\mathcal{H}$  its associated RKHS. For any  $G: \mathbb{R}^d \rightarrow \mathbb{R}$  and any loss function  $L: \mathbb{R}^m \rightarrow \mathbb{R}$ ,  
<sup>non-decreasing</sup>

$$\arg \min_{h \in \mathcal{H}} F(h) = \arg \min_{h \in \mathcal{H}} G(\|h\|_{\mathcal{H}}) + L(h(x_1), \dots, h(x_m))$$

admits a solution of the form  $h^* = \sum_{i=1}^m \alpha_i K(x_i, \cdot)$ . If  $G$  is in fact strictly increasing, then all solutions have this form.

Proof: Let  $\mathcal{H}_1 = \text{span}(\{K(x_i, \cdot)\}_{i=1}^m)$ . This is a subspace of  $\mathcal{H}$ , and thus  $\mathcal{H} = \mathcal{H}_1 \oplus \mathcal{H}_1^\perp$ , where  $\mathcal{H}_1^\perp$  is the orthogonal complement of  $\mathcal{H}_1$ , i.e.  
 $\mathcal{H}_1^\perp = \{h_0 \in \mathcal{H} \mid \langle h_0, h \rangle_{\mathcal{H}} = 0 \ \forall h \in \mathcal{H}_1\}$ .

So, any  $h \in \mathcal{H}$  may be written as  $h = h_1 + h_1^\perp$ , where  $h_1 \in \mathcal{H}_1, h_1^\perp \in \mathcal{H}_1^\perp$ .

Since  $G$  is non-decreasing,

$$G(\|h_1\|_{\mathcal{H}}) \leq G(\sqrt{\|h_1\|_{\mathcal{H}}^2 + \|h_1^\perp\|_{\mathcal{H}}^2}) = G(\|h\|_{\mathcal{H}}).$$

By the reproducing property, and the fact that  $h^\perp$  is orthogonal to all elements of  $\mathcal{H}_1^\perp$ ,

$$\begin{aligned}
h(x_i) &= \langle h, K(x_i, \cdot) \rangle_{\mathcal{H}} \\
&= \langle h_1 + h^\perp, K(x_i, \cdot) \rangle_{\mathcal{H}} \\
&= \langle h_1, K(x_i, \cdot) \rangle_{\mathcal{H}} \\
&= h_1(x_i)
\end{aligned}$$

Hence,  $L(h(x_1), \dots, h(x_n)) = L(h_1(x_1), \dots, h_1(x_n))$ . Thus,  $F(h_1) \leq F(h)$ .

So,  $h_1$  a minimizer of  $F(h)$  necessarily lies in  $\mathcal{H}_1$ . Moreover, if  $F$  is strictly increasing, we have that  $F(h_1) < F(h)$  if  $\|h^\perp\|_{\mathcal{H}} > 0 \Rightarrow$  any minimizer must have  $\|h^\perp\|_{\mathcal{H}} = 0$ , i.e. must lie in  $\mathcal{H}_1$ .

We are now in a position to prove learning bounds for KSVM:

Theorem (Rademacher complexity for kernel hypothesis classes): Let  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a PDS kernel and let  $\Phi: \mathcal{X} \rightarrow \mathcal{H}$  be a feature mapping associated to  $K$ . Let  $S \subseteq \{x \mid K(x, x) \leq r^2\}$ ,  $r > 0$ , be a sample of  $\mathcal{S}$  size  $m$ . Let  $\mathcal{H} = \{x \mapsto \langle w, \Phi(x) \rangle_{\mathcal{H}} \mid \|w\|_{\mathcal{H}} \leq \Delta\}$ ,  $\Delta \geq 0$ . Then for  $K$  the associated Gram matrix,

$$\hat{R}_S(\mathcal{H}) \leq \frac{\Delta \sqrt{\text{Tr}(K)}}{m} \leq \sqrt{\frac{r^2 \Delta^2}{m}}$$

Proof: 
$$\hat{R}_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_S \left( \sup_{\|w\|_{\mathcal{H}} \leq \Delta} \left\langle w, \sum_{i=1}^m z_i \Phi(x_i) \right\rangle_{\mathcal{H}} \right)$$

$$\leq \frac{1}{m} \cdot \Delta \cdot \mathbb{E}_S \left( \left\| \sum_{i=1}^m z_i \Phi(x_i) \right\|_{\mathcal{H}} \right)$$

$$\begin{aligned}
&\leq \frac{\Delta}{m} \left( \mathbb{E} \left( \left\| \sum_{i=1}^m \mathcal{C}_i \Phi(x_i) \right\|_{\mathcal{H}}^2 \right) \right)^{1/2} \\
&= \frac{\Delta}{m} \left( \mathbb{E} \left( \left\langle \sum_{i=1}^m \mathcal{C}_i \Phi(x_i), \sum_{j=1}^m \mathcal{C}_j \Phi(x_j) \right\rangle \right) \right)^{1/2} \\
&= \frac{\Delta}{m} \left( \mathbb{E} \left( \sum_{i,j=1}^m \mathcal{C}_i \mathcal{C}_j \langle \Phi(x_i), \Phi(x_j) \rangle \right) \right)^{1/2} \\
&= \frac{\Delta}{m} \left( \mathbb{E} \left( \sum_{i,j=1}^m \mathcal{C}_i^2 \langle \Phi(x_i), \Phi(x_j) \rangle \right) \right)^{1/2} \\
&= \frac{\Delta}{m} \left( \sum_{i=1}^m \|\Phi(x_i)\|_{\mathcal{H}}^2 \right)^{1/2} \\
&= \frac{\Delta}{m} \left( \sum_{i=1}^m K(x_i, x_{\mathcal{Y}_i}) \right)^{1/2} \\
&= \frac{\Delta}{m} \cdot \sqrt{\text{Tr}(K)} \\
&\leq \frac{\Delta}{m} \cdot \sqrt{m \cdot r^2} \\
&= \sqrt{\frac{\Delta^2 r^2}{m}} \quad \blacksquare
\end{aligned}$$

Remark: By the Khintchine-Kahane inequality, the empirical Rademacher complexity can be lower bounded by  $\frac{1}{\sqrt{2}} \sqrt{\text{Tr}(K)} \cdot \frac{\Delta}{m}$ , which suggests the argument made in the above theorem is nearly tight. (5)

Corollary (Margin Bound for Kernel-Based Hypothesis Classes): Let  $K: X \times X \rightarrow \mathbb{R}$  be a PSD kernel with  $r^2 = \sup_{x \in X} K(x, x)$ . Let  $\Phi: X \rightarrow \mathcal{H}$  be a feature mapping associated to  $K$  and let  $\mathcal{H} = \{x \mapsto \langle w, \Phi(x) \rangle_{\mathcal{H}} \mid \|w\|_{\mathcal{H}} \leq \Delta\}$ . Fix  $\delta, \rho > 0$ . Then  $\forall \epsilon > 0$ , each of the following holds with probability at least  $1 - \delta$ .

(a.)  $R(h) \leq \hat{R}_{S, \rho}(h) + 2 \sqrt{\frac{r^2 \Delta^2}{\rho^2 m}} + \sqrt{\frac{\log(\frac{2}{\delta})}{m}}$

(b.)  $R(h) \leq \hat{R}_{S, \rho}(h) + 2 \sqrt{\frac{\text{Tr}(K) \Delta^2}{m \rho^2}} + 3 \sqrt{\frac{\log(\frac{2}{\delta})}{2m}}$ .

Proof: Apply our SVM margin bound to our Rademacher complexity result for kernel hypothesis classes.

• We have so far considered kernels  $K: X \times X \rightarrow \mathbb{R}$  that induce SPSD Gram matrices  $K \in \mathbb{R}^{m \times m}$ .

• In some cases,  $K$  is constructed as  $K(x, x') = f(-d(x, x'))$ , where  $f: \mathbb{R} \rightarrow \mathbb{R}$  is increasing and  $d: X \times X \rightarrow \mathbb{R}$  is a metric.

• For example,  $f_2(z) = \exp(z^2/2)$ ,  $d(x, x') = \|x - x'\|_2^2$  yields the

### Gaussian Kernel.

- It makes sense to derive conditions on the metric directly.

Defn: A kernel  $K: X \times X \rightarrow \mathbb{R}$  is negative-definite symmetric (NDS) if it is symmetric and if for all  $\{x_1, \dots, x_m\} \subseteq X$ , with  $c \in \mathbb{R}^{m \times 1}$  and  $\mathbb{1}^T c = 0$ ,  $c^T K c \leq 0$

Remark:  $K$  PSS  $\Rightarrow -K$  NDS, but  $K$  NDS  $\not\Rightarrow -K$  PSS.

ex: Let  $D(x, x') = \|x - x'\|_2^2$ . Then  $D: \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  is NDS. Indeed, let  $c \in \mathbb{R}^{m \times 1}$  be s.t.  $\mathbb{1}^T c = 0$ , i.e.  $\sum_{i=1}^m c_i = 0$ . Let  $\{x_i\}_{i=1}^m \subset \mathbb{R}^D$  be arbitrary. Then

$$\begin{aligned}
 & \sum_{i,j=1}^m c_i c_j \|x_i - x_j\|_2^2 && \langle, \rangle = \text{usual } \mathbb{R}^D \text{ inner product} \\
 &= \sum_{i,j=1}^m c_i c_j (\|x_i\|_2^2 + \|x_j\|_2^2 - 2 \langle x_i, x_j \rangle) \\
 &= \sum_{i,j=1}^m c_i c_j (\|x_i\|_2^2 + \|x_j\|_2^2) - 2 \left\langle \sum_{i=1}^m c_i x_i, \sum_{j=1}^m c_j x_j \right\rangle \\
 &= \sum_{i,j=1}^m c_i c_j (\|x_i\|_2^2 + \|x_j\|_2^2) - 2 \left\| \sum_{i=1}^m c_i x_i \right\|_2^2 \\
 &\leq \sum_{i,j=1}^m c_i c_j (\|x_i\|_2^2 + \|x_j\|_2^2) \\
 &= \left( \sum_{j=1}^m c_j \right) \cdot \left( \sum_{i=1}^m c_i \|x_i\|_2^2 \right) + \left( \sum_{i=1}^m c_i \right) \cdot \left( \sum_{j=1}^m c_j \|x_j\|_2^2 \right) = 0.
 \end{aligned}$$