

# FoSM: Lecture 20

①

- Given  $\{x_i\}_{i=1}^m \subset \mathbb{R}^D$ , one needs to evaluate  $K(x_i, x_j)$  for all  $(x_i, x_j)$  pairs, in order to use, e.g., kernel SVM.
- Letting  $C_K = \text{cost of evaluating } K(x_i, x_j) \text{ for one } (x_i, x_j) \text{ pair}$ , this is  $O(C_K \cdot m^2)$ , which is quadratic in  $m$ . This is not acceptable when  $m$  is large.
- Moreover, given the kernel matrix / fast construction of  $K \in \mathbb{R}^{m \times m}$ , an evaluation of a single test point is also slow. Indeed, we classify using
$$h(x) = \sum_{i=1}^m \alpha_i K(x_i, x) + b,$$

so to classify a single  $x$  is  $O(C_K \cdot m)$ .

The interesting case here is for data s.t.  $\tilde{D} < m$ , where  $\dim(H_K) = \tilde{D}$ , where  $K(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle_{H_K}$ ,  $\Phi: \mathbb{R}^D \rightarrow H_K$ . Indeed, otherwise we can handle the primal problem and incur complexity  $O(\tilde{D}m)$ . But, ~~usually~~ typically  $\tilde{D} \gg m$  and it is easy to find cases where  $\tilde{D} = +\infty$  (i.e. Gaussian kernels).

Our goal is to discuss some ways around this. One interesting way is to approximate  $\Phi$  with  $\Psi$ , s.t. (1)  $\Psi: \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$ ,  $D' < \tilde{D}$

$$(2) \langle \Psi(x_i), \Psi(x_j) \rangle \approx \langle \Phi(x_i), \Phi(x_j) \rangle.$$

In this sense, we drop the dimension of the embedding space ((1)), while preserving in an approximate sense the desirable aspects of  $\Phi$  related to our kernel  $K$ , namely  $\langle \Psi(x_i), \Psi(x_j) \rangle \approx K(x_i, x_j)$ .

To do so, we will need a little harmonic analysis: (2)

Theorem (Bochner): Let  $K(x, x') = G(x - x')$  be a continuous kernel defined over a locally compact set  $X$ . Then the following are equivalent:

(a)  $K$  is positive definite

(b)  $\exists \mu$ , a non-negative measure, s.t.  $G(x) = \int_X \exp(-i \langle \omega, x \rangle) d\mu(\omega)$ .

Condition (b) says  $G$  is the Fourier transform of a non-negative measure  $\mu$ .

Note that the condition  $K(x, x') = G(x - x')$  gives such kernels the ~~monitored~~ shift-invariant, since  $K(x+z, x'+z) = K(x, x') \quad \forall x, x', z \in X$ .

Proposition: Let  $K(x, x') = G(x, x')$  be a shift-invariant, continuous kernel with associated measure  $\mu \geq 0$ . Suppose  $K(x, x) = 1 \quad \forall x$ , i.e.  $\mu$  is a probability measure. Then

$$K(x, x') = \mathbb{E}_{\omega \sim \mu} \left( \begin{bmatrix} \cos(\omega \cdot x) & \sin(\omega \cdot x) \end{bmatrix}^T \begin{bmatrix} \cos(\omega \cdot x') & \sin(\omega \cdot x') \end{bmatrix} \right).$$

Proof: Note that since  $K, \mu$  are  $\mathbb{R}$ -valued, it suffices to consider the real part of  $\exp(-i \langle \omega, x \rangle)$ , since the imaginary part necessarily averages out to 0. So, noting  $\operatorname{Re}(e^{i \langle \omega, x \rangle}) = \cos(\langle \omega, x \rangle)$ , we get

$$K(x, x') = \int_X \cos(\langle \omega, x - x' \rangle) d\mu(\omega)$$

Noting that  $\cos(a-b) = \cos(a)\cos(b) + \sin(a)\sin(b)$ , we get

(5)

$$\int_X \cos(\langle w, x - x' \rangle) d\mu(w)$$

$$= \int_X \cos(\langle w, x \rangle - \langle w, x' \rangle) d\mu(w)$$

$$= \int_X \left( \cos(\langle w, x \rangle) \cos(\langle w, x' \rangle) + \sin(\langle w, x \rangle) \sin(\langle w, x' \rangle) \right) d\mu(w)$$

$$= \mathbb{E}_{w \sim \mu} \left( \begin{bmatrix} \cos(w \cdot x) & \sin(w \cdot x) \end{bmatrix}^T \cdot \begin{bmatrix} \cos(w \cdot x') & \sin(w \cdot x') \end{bmatrix} \right)$$

This suggests a scheme for estimating an approximate kernel mapping.

Given  $D \geq 1$ , let  $\Psi: X \rightarrow \mathbb{R}^{2D}$   
 $x \mapsto \frac{1}{D} \left( \cos(w_1 \cdot x), \sin(w_1 \cdot x), \dots, \cos(w_D \cdot x), \sin(w_D \cdot x) \right)^T$

where  $\{w_i\}_{i=1}^D \stackrel{\text{iid}}{\sim} \mu$ , where the FT of  $\mu$  is  $\kappa$ .

Why do this? Well,  $\langle \Psi(x), \Psi(x') \rangle$   
 $= \frac{1}{D} \sum_{i=1}^D \left[ \cos(w_i \cdot x), \sin(w_i \cdot x) \right]^T \left[ \cos(w_i \cdot x'), \sin(w_i \cdot x') \right]$

which is the empirical version of

$$\mathbb{E}_{w \sim \mu} \left( \begin{bmatrix} \cos(w \cdot x) & \sin(w \cdot x) \end{bmatrix}^T \cdot \begin{bmatrix} \cos(w \cdot x') & \sin(w \cdot x') \end{bmatrix} \right)$$

In fact, as we might expect, this empirical expectation improves as  $D \rightarrow \infty$ . (4)

Lemma: Let  $K$  be a continuously differentiable kernel function that satisfies ~~the~~ the above conditions, i.e.  $K$  is shift-invariant and has associated non-negative measure  $\mu$ . Assume moreover that  $X$  is compact and let  $D = \dim(X)$ . Let  $R$  be s.t.  $X \subset B_R(0)$  and  $\sigma_{\text{wgn}}^2 = \mathbb{E}(\|w\|^2) < \infty$ . Then for

$\Psi: X \rightarrow \mathbb{R}^D$  defined as  $X \mapsto \frac{1}{\sqrt{D}}(\cos(w_1 \cdot x), \sin(w_1 \cdot x), \dots, \cos(w_D \cdot x), \sin(w_D \cdot x))$ ,

let  $r \in (0, 2R]$  and let  $\varepsilon > 0$ . Then

$$\mathbb{P}_{\text{wgn}} \left( \sup_{x, x' \in X} \left| \langle \Psi(x), \Psi(x') \rangle - K(x, x') \right| \geq \varepsilon \right)$$

$$\leq 2 \cdot \mathcal{N}(2R, r) \exp(-D\varepsilon^2/8) + \frac{4r\sigma_p}{\varepsilon},$$

where  $\mathcal{N}(R, r)$  is the minimum number of balls of radius  $r$  needed to cover a ball of radius  $R$ .

Proof: Let  $Z = \{x - x' \mid x, x' \in X\}$ . (Clearly  $X \subset B_R(0) \Rightarrow Z \subset B_{2R}(0)$ .)

Since  $X$  is closed,  $Z$  is also closed and hence compact. Let  $B := \mathcal{N}(2R, r)$ , and let  $\{z_j\}_{j=1}^B$  be the centers of the minimal cover. In particular,  $\forall z \in Z$ ,

$$\exists j \text{ s.t. } \exists \delta \text{ s.t. } z = z_j + \delta, \|\delta\|_2 \leq r.$$

Define  $S: \mathcal{Z} \rightarrow \mathbb{R}$

$$z \mapsto \langle \Psi(x), \Psi(x') \rangle - \mathcal{K}(x, x'),$$

where  $z = x - x'$ ; this map is well-defined because  $\mathcal{K}$  is shift invariant,

same with  $\langle \Psi(x), \Psi(x') \rangle$  by construction. Since  $S$  is  $\mathcal{C}^1$  over  $\mathcal{Z}$ ,

$S$  is  $L$ -Lipschitz, with  $L = \sup_{z \in \mathcal{Z}} \|\nabla S(z)\|_2$ . Moreover, if  $L < \frac{\epsilon}{2r}$  for

all  $j \in \{1, \dots, B\}$ , we have  $|S(z_j)| < \frac{\epsilon}{2}$ , then it follows that

$$|S(z)| = |S(z_j + \delta)|$$

$$\leq L \cdot \|z_j - (z_j + \delta)\|_2 + |S(z_j)|$$

$$\leq rL + \frac{\epsilon}{2}$$

$$< \epsilon.$$

So, it suffices to bound the probability that  $L \geq \frac{\epsilon}{2r}$  and  $|S(z_j)| \geq \frac{\epsilon}{2}$ .

(a) Bounding probability of  $L \geq \frac{\epsilon}{2r}$ : By linearity of  $\mathbb{E}$ ,

$$\begin{aligned} \mathbb{E}(\nabla \langle \Psi(x), \Psi(x') \rangle) &= \nabla \mathbb{E}(\langle \Psi(x), \Psi(x') \rangle) \\ &= \nabla \mathcal{K}(x, x'). \end{aligned}$$

$$\text{So, } \mathbb{E}(L^2) = \mathbb{E}\left(\sup_{z \in \mathcal{Z}} \|\nabla S(z)\|_2^2\right)$$

$$= \mathbb{E}\left(\sup_{x, x' \in \mathcal{X}} \|\nabla(\langle \Psi(x), \Psi(x') \rangle - \mathcal{K}(x, x'))\|_2^2\right)$$

$$\leq 2 \mathbb{E} \left( \sup_{x, x' \in \mathcal{X}} \|\nabla \langle \Psi(x), \Psi(x') \rangle\|_2^2 \right) + 2 \sup_{x, x' \in \mathcal{X}} \|\nabla \Psi(x, x')\|_2^2 \quad (6)$$

$$= 2 \mathbb{E} \left( \sup_{x, x' \in \mathcal{X}} \|\nabla \langle \Psi(x), \Psi(x') \rangle\|_2^2 \right) + 2 \sup_{x, x' \in \mathcal{X}} \|\mathbb{E}(\nabla \langle \Psi(x), \Psi(x') \rangle)\|_2^2$$

$$\leq 4 \mathbb{E} \left( \sup_{x, x' \in \mathcal{X}} \|\nabla \langle \Psi(x), \Psi(x') \rangle\|_2^2 \right) \quad (\star)$$

We now analyze  $\nabla \langle \Psi(x), \Psi(x') \rangle$  where  $\nabla$  is taken w.r.t.  $z = x - x'$ .

Indeed,  $\nabla \langle \Psi(x), \Psi(x') \rangle$

$$= \nabla \frac{1}{D} \sum_{i=1}^D \cos(w_i \cdot x) \cos(w_i \cdot x') + \sin(w_i \cdot x) \sin(w_i \cdot x')$$

$$= \nabla \frac{1}{D} \sum_{i=1}^D \cos(w_i (x - x'))$$

$$= \frac{1}{D} \sum_{i=1}^D \sin(w_i (x - x')) \cdot w_i \quad (\star\star)$$

Combining  $(\star)$ ,  $(\star\star)$  yields

$$\mathbb{E}(L^2) \leq 4 \mathbb{E} \left( \sup_{x, x' \in \mathcal{X}} \|\nabla \langle \Psi(x), \Psi(x') \rangle\|_2^2 \right)$$

$$\leq 4 \mathbb{E} \left( \sup_{x, x' \in \mathcal{X}} \left\| \frac{1}{D} \sum_{i=1}^D \sin(w_i (x - x')) \cdot w_i \right\|_2^2 \right)$$

$$\leq 4 \mathbb{E} \left( \sup_{x, x' \in \mathcal{X}} \left\| \frac{1}{D} \sum_{i=1}^D \|w_i\|_2^2 \right\|_2^2 \right)$$

$$= 4 \mathbb{E}(\|w_i\|_2^2)$$

$$= 4 \sigma_w^2.$$

So, by Markov's inequality,

$$P\left(L \geq \frac{\epsilon}{2r}\right) \leq \left(\frac{4rZ_\mu}{\epsilon}\right)^2.$$

(b) Bounding probability that  $|S(z_j)| \geq \frac{\epsilon}{2}$ : Note that  $S$  is the sum of  $D$  i.i.d. r.v., each bounded by  $\frac{2}{D}$ . So, by Hoeffding and a union bound,

$$P\left(\exists j \in \{1, \dots, B\} \mid |S(z_j)| \geq \frac{\epsilon}{2}\right) \leq \sum_{i=1}^B P\left(|S(z_i)| \geq \frac{\epsilon}{2}\right)$$
$$\leq 2B P\left(|S(z_1)| \geq \frac{\epsilon}{2}\right)$$
$$\leq 2B \exp(-D\epsilon^2/8).$$

Putting (a), (b) together gives the desired result:

$$P\left(\sup_{z \in \mathcal{Z}} |S(z)| \geq \epsilon\right) \leq \underbrace{2N(2R, r)}_B \exp(-D\epsilon^2/8) + \left(\frac{4rZ_\mu}{\epsilon}\right)^2.$$