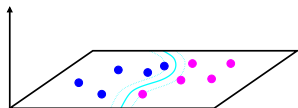


LEARNING FROM DISTRIBUTIONS VIA SUPPORT MEASURE MACHINES

K. Muandet, K. Fukumizu, F. Dinuzzo, B. Schölkopf

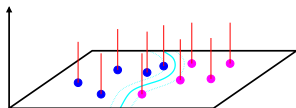


MAX-PLANCK-GESELLSCHAFT



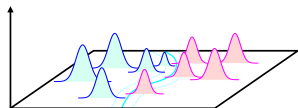
$$\mathcal{X} \rightarrow \mathcal{H}$$

$$x \mapsto k(x, \cdot)$$



$$\Delta \rightarrow \mathcal{H}$$

$$\delta_x \mapsto \int k(x, \cdot) d\delta_x$$



$$\mathcal{P} \rightarrow \mathcal{H}$$

$$\mathbb{P} \mapsto \int k(x, \cdot) d\mathbb{P}(x)$$

Potential Applications:

- ▶ Uncertain/noisy data (astronomical/biological data)
- ▶ Groups of samples (group anomaly, preference learning)
- ▶ Changing environments (domain adaptation/generalization)
- ▶ Large-scale machine learning (data squashing)

Given a sample $(\mathbb{P}_1, y_1), (\mathbb{P}_2, y_2), \dots, (\mathbb{P}_m, y_m)$, any solution f to

$$\ell(\mathbb{P}_1, y_1, \mathbb{E}_{\mathbb{P}_1}[f], \dots, \mathbb{P}_m, y_m, \mathbb{E}_{\mathbb{P}_m}[f]) + \Omega(\|f\|_{\mathcal{H}}) \quad (1)$$

admits a form $f = \sum_{i=1}^m \alpha_i \mathbb{E}_{\mathbb{P}_i}[k(x, \cdot)]$ for some $\alpha_i \in \mathbb{R}$, $i = 1, \dots, m$.

Our framework (1) is different from

1. $\mathbb{E}_{\mathbb{P}_1} \mathbb{E}_{\mathbb{P}_2} \dots \mathbb{E}_{\mathbb{P}_m} \ell(\{\mathbf{x}_i, y_i, f(\mathbf{x}_i)\}_{i=1}^m) + \Omega(\|f\|_{\mathcal{H}})$
2. $\ell(\{M_i, y_i, f(M_i)\}_{i=1}^m) + \Omega(\|f\|_{\mathcal{H}})$, $M_i = \mathbb{E}_{\mathbb{P}_i}[\mathbf{x}]$

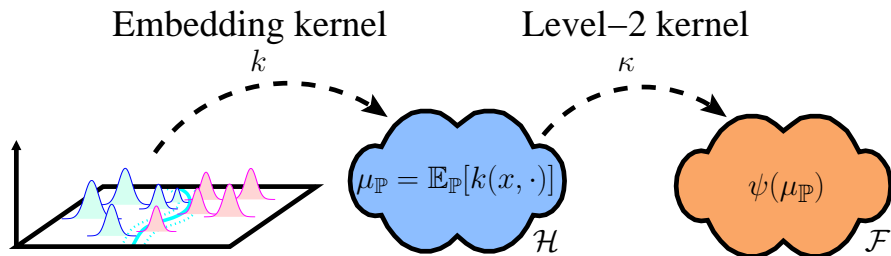
**intracable
information loss**

Risk Deviation Bound

Given a distribution \mathbb{P} with variance σ^2 , a Lipschitz continuous function f with constant C_f , a loss function ℓ with constant C_ℓ , it follows for any $y \in \mathbb{R}$ that

$$|\mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[\ell(y, f(\mathbf{x}))] - \ell(y, \mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[f(\mathbf{x})])| \leq 2C_\ell C_f \sigma$$

Information preserving + computationally efficient.

**Feature maps**

$$\begin{aligned} \mathcal{K}(\delta_x, \delta_y) &= \langle k(x, \cdot), k(y, \cdot) \rangle_{\mathcal{H}} \\ &= k(x, y) \end{aligned}$$

The SVM is recovered
as a special case.

Linear kernels

$$\begin{aligned} \mathcal{K}(\mathbb{P}, \mathbb{Q}) &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} \\ &= \mathbb{E}_{x \sim \mathbb{P}, z \sim \mathbb{Q}}[k(x, z)] \end{aligned}$$

It defines the feature for
distributions.

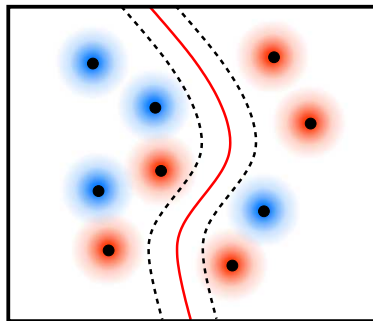
Nonlinear kernels

$$\begin{aligned} \mathcal{K}(\mathbb{P}, \mathbb{Q}) &= \kappa(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}}) \\ &= \langle \psi(\mu_{\mathbb{P}}), \psi(\mu_{\mathbb{Q}}) \rangle_{\mathcal{F}} \end{aligned}$$

It allows for nonlinear
learning algorithms.

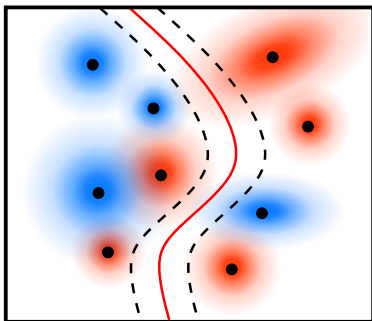
$$K(\mathbb{P}, \mathbb{Q}) = \left\langle \int k(\tilde{x}, \cdot) g(x, \tilde{x}) d\tilde{x}, \int k(\tilde{z}, \cdot) g(z, \tilde{z}) d\tilde{z} \right\rangle_{\mathcal{H}} = k_g(x, z)$$

Standard Support Vector Machine



$$f = \sum_{i=1}^n k(x_i, \cdot)$$

Flexible Support Vector Machine



$$f = \sum_{i=1}^n k_i(x_i, \cdot)$$

The flexible SVM places different kernels on training samples.