

Kernel methods: ... "Using the kernel trick" ...
 (related to: Kernel regression, SVM, Kernel PCA, Kernel discriminant analysis)

• Kernel regression $\left\{ \begin{array}{l} \text{feature embedding into RKHS} \\ \text{optimization perspective (representer thm.)} \end{array} \right.$

Nonparametric regression = $(y_i, x_i)_{i=1}^n$ $\begin{array}{c} \downarrow \text{response} \quad \downarrow \text{covariates} \end{array}$ $x^* \rightarrow y^*$ predict.

Learning the regression f_n :
 $f(x) = E[Y | X=x]$

Recall = The prediction at x^* using linear ridge regression is $x^{*T} \hat{\beta}_R = x^{*T} (X^T X + \lambda I_d)^{-1} (X^T Y)$
 $= (X \cdot x^*)^T (X X^T + \lambda I_n)^{-1} Y$

$$X X^* = \begin{pmatrix} \langle x_1, x^* \rangle \\ \vdots \\ \langle x_n, x^* \rangle \end{pmatrix}$$

$$(X^T X)_{jk} = \langle x_j, x_k \rangle$$

$\langle x_j, x_k \rangle$ is a (bilinear) measure of similarity b/w x_j & $x_k \in \mathbb{R}^d$.
 c.s. $\theta = \frac{\langle x_i, x_j \rangle}{\|x_i\| \|x_j\|}$

(Idea) • Replace $\langle x_j, x_k \rangle$ with $k(x_j, x_k)$: nonlinear measure of similarity.

Define $K^* = \begin{pmatrix} k(x_1, x^*) \\ \vdots \\ k(x_n, x^*) \end{pmatrix}$ $(K)_{jk} = k(x_j, x_k)$

$$K^* \in \mathbb{R}^n, K \in \mathbb{R}^{n \times n}$$

$$\boxed{f(x^*) = K^{*T} (K + \lambda I_n)^{-1} Y}$$
 Kernel Ridge Regression (KRR) estimator.

Conditions on K : $\left\{ \begin{array}{l} \text{i) symmetric: } k(x, x') = k(x', x) \quad \forall x, x' \\ \text{ii) } \forall n \in \mathbb{N}, \forall x_1, \dots, x_n \in \mathcal{X} \\ K_n = (k(x_i, x_j))_{i,j} \text{ p.d. (n.n.d at least)} \end{array} \right.$
 $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

Hadamard
product:
 $(A \circ B)_{ij} = A_{ij} \cdot B_{ij}$

$K_n = X^T X$ satisfies i) & ii)

Claim: $k(x, x') = \langle x, x' \rangle^2$ is pd kernel $(J_n)_{ij} = \langle x_i, x_j \rangle^2$

Fact: if A and B are nnd,
then so is $A \circ B$ (Schur's theorem)

$$J_n = K_n \circ K_n$$

$$\text{Let } K_m(x_i, x_j) = \langle x_i, x_j \rangle^m \quad m \in \mathbb{N}^+$$

Then k_m is a pd kernel \Rightarrow polynomial kernel.

$$x_i \in \mathbb{R}^2$$

$$x_i = (x_{i1}, x_{i2})$$

$$\begin{aligned} \langle x_i, x_j \rangle^2 &= (x_{i1}x_{j1} + x_{i2}x_{j2})^2 \\ &= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2x_{i1}x_{j1}x_{i2}x_{j2} \\ &= \langle \psi(x_i), \psi(x_j) \rangle \end{aligned}$$

$$\text{where } \psi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$x \mapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$(J_n)_{ij} = \langle \psi(x_i), \psi(x_j) \rangle \Rightarrow J_n \text{ is nnd.}$$

KRR \Leftrightarrow ordinary ridge regression in the transformed
feature space produced by the feature map ψ

"Kernel Engineering"

Creating new kernels out of simple ones.

Ex: if k_1 and k_2 are kernels, then so is

any of the following:

$$i) \quad k(x, x') = k_1(x, x') + k_2(x, x')$$

$$ii) \quad k(x, x') = f(x) \cdot k_1(x, x') \cdot f(x')$$

$$iii) \quad k(x, x') = k_1(x, x') \cdot k_2(x, x')$$

$$iv) \quad k(x, x') = \exp\{k_1(x, x')\} \quad (\text{Taylor expansion})$$

Gaussian kernel

$$k(x, x') = \exp \left\{ -\frac{\|x - x'\|^2}{2\tau^2} \right\}$$

Ex: Show (using i) ~ iv) that k is p.d. kernel.

If $k(x, x') =$

$$\phi(\|x - x'\|)$$

then k is called an isotropic kernel.

$k(x, x') = \phi(\|x - x'\|)$
stationary kernel

Gaussian
Matern

Mercer's theorem (general way of getting a feature map)

If $k: X \times X \rightarrow \mathbb{R}$ is a p.d. f. and $X \subseteq \mathbb{R}^d$ is compact, then there is an orthonormal set of functions $\{\psi_j\}$ st.

$$k(x, x') = \sum_{j=1}^{\infty} \lambda_j \psi_j(\cdot) \cdot \psi_j(x')$$

$$\text{i.e. } x \mapsto (\lambda_1 \psi_1(x), \lambda_2 \psi_2(x), \dots)$$

$$k(x, x') = \langle \psi(x), \psi(x') \rangle$$

Bochner's

~~Rechner's~~ Thm (complete characterization of stationary kernels)

Let $k(x, x') = \phi(x - x')$. Then k is p.d. $\Leftrightarrow \hat{\phi}(\omega) \geq 0 \forall \omega$.

$$\text{where } \hat{\phi}(\omega) = (2\pi)^{-d/2} \int e^{i\langle \omega, t \rangle} \phi(t) dt.$$

proof: $\hat{\phi}(\omega) \geq 0 \forall \omega \Rightarrow k$ is p.d.

$$\text{Inverse-Fourier Transform: } \phi(t) = (2\pi)^{-d/2} \int e^{-i\langle \omega, t \rangle} \hat{\phi}(\omega) d\omega$$

(For any x_1, \dots, x_n , and any $a_1, \dots, a_n \in \mathbb{R}$.)
we want to show: $\sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j) \geq 0$

$$\sum_{i,j} a_i a_j k(x_i, x_j)$$

$$= (2\pi)^{-d/2} \sum_{i,j} a_i a_j \int e^{-i\langle \omega, x_i - x_j \rangle} \hat{\phi}(\omega) d\omega$$

$$= (2\pi)^{-d/2} \int \sum_{i,j} a_i a_j \cdot e^{-i\omega x_i} \cdot \overline{e^{-i\omega x_j}} \cdot \hat{\phi}(\omega) d\omega$$

$$= (2\pi)^{-d/2} \int \left| \sum_i a_i e^{-i\omega x_i} \right|^2 \hat{\phi}(\omega) d\omega$$

$$\geq 0$$

This gives a recipe for constructing kernel fns
taking Take any non-negative fun. (prob. density) and