Kernelized algorithms (Lecture 7)

Zoltán Szabó

December 13, 2016

One-page summary

- Until now:
 - Regularized least-squares problems.

$$J(f) = \sum_{i=1}^{n} [y_i - f(x_i)]^2 + \lambda \|Lf\|_2^2 \to \min_f,$$

$$f(x) = \sum_{j=1}^{B} c_j \phi_j(x).$$

Basis expansion (finite).

One-page summary

- Until now:
 - Regularized least-squares problems.

$$J(f) = \sum_{i=1}^{n} [y_i - f(x_i)]^2 + \lambda \|Lf\|_2^2 \to \min_f,$$

$$f(x) = \sum_{j=1}^{B} c_j \phi_j(x).$$

- Basis expansion (finite).
- Today:
 - kernel → kernel ridge regression, kernel PCA.

Kernel

Definition (Inner product space)

 \mathcal{F} : vector space over \mathbb{R} . $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \to \mathbb{R}$ is an inner product on \mathcal{F} if for $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{F}$

Definition (Inner product space)

 \mathfrak{F} : vector space over \mathbb{R} . $\langle \cdot, \cdot \rangle : \mathfrak{F} \times \mathfrak{F} \to \mathbb{R}$ is an inner product on \mathfrak{F} if for $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathfrak{F}$

Definition (Inner product space)

 \mathcal{F} : vector space over \mathbb{R} . $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \to \mathbb{R}$ is an inner product on \mathcal{F} if for $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{F}$

Definition (Inner product space)

 \mathcal{F} : vector space over \mathbb{R} . $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \to \mathbb{R}$ is an inner product on \mathcal{F} if for $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{F}$

- $\langle f,g\rangle=\langle g,f\rangle$ (symmetry),

Notes: 1, $2 \Rightarrow$ bilinearity.

- Norm induced by the inner product: $||f|| = \sqrt{\langle f, f \rangle}$.
- CBS: $|\langle f, g \rangle| \leq ||f|| ||g|| \Rightarrow \cos(f, g)$.

Definition (Inner product space)

 \mathcal{F} : vector space over \mathbb{R} . $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \to \mathbb{R}$ is an inner product on \mathcal{F} if for $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{F}$

- $\langle f,g\rangle=\langle g,f\rangle$ (symmetry),

Notes: 1, $2 \Rightarrow$ bilinearity.

- Norm induced by the inner product: $||f|| = \sqrt{\langle f, f \rangle}$.
- CBS: $|\langle f, g \rangle| \leq ||f|| ||g|| \Rightarrow \cos(f, g)$.

Definition (Hilbert space)

Nice ('complete') inner product space.



Kernel: inner product of features

Let \mathcal{X} be a set. A $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ function is called kernel if

- ullet there exists a Hilbert space \mathcal{H} , and
- $\phi: \mathfrak{X} \to \mathfrak{H}$ feature map such that

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}.$$

Kernel: examples $(\mathfrak{X} = \mathbb{R}^d)$

•
$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathbb{R}^D}, \ \phi(x) = [\phi_1(x); \dots; \phi_B(x)].$$

Kernel: examples $(X = \mathbb{R}^d)$

- $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathbb{R}^D}, \ \phi(x) = [\phi_1(x); \dots; \phi_B(x)].$
- Polynomial, Gaussian, Laplacian kernel $(\theta > 0, d \in \mathbb{Z}^+)$:

$$k(x, x') = (\langle x, x' \rangle + \theta)^d,$$
 $k(x, x') = e^{-\frac{\|x - x'\|_2^2}{2\theta^2}},$ $k(x, x') = e^{-\theta \|x - x'\|_2}.$

Kernel: example domains (\mathfrak{X})

- Euclidean space: $\mathfrak{X} = \mathbb{R}^d$.
- Graphs, texts, time series, dynamical systems, probability distributions.





Kernel: properties

Construction from old kernels ($c \ge 0$):

• k_1 , k_2 , k: kernel on $\mathfrak{X} \Rightarrow k_1 + k_2$, ck: kernel on \mathfrak{X} .

Kernel: properties

Construction from old kernels ($c \ge 0$):

- k_1 , k_2 , k: kernel on $\mathfrak{X} \Rightarrow k_1 + k_2$, ck: kernel on \mathfrak{X} .
- ullet Composition: Let \tilde{k} kernel on $\tilde{\mathfrak{X}}$, $M:\mathfrak{X} \to \tilde{\mathfrak{X}}$ mapping. Then

$$k(x, x') = \tilde{k}(M(x), M(x'))$$

is kernel on \mathfrak{X} .

Kernel: properties - continued

Product:

• k_i kernel on \mathfrak{X}_i (i=1,2). Then

$$(k_1 \times k_2) ((x,y),(x',y')) = k_1(x,x')k_2(y,y')$$

is kernel on $\mathfrak{X}_1 \times \mathfrak{X}_2$.

Kernel: properties – continued

Product:

• k_i kernel on \mathfrak{X}_i (i=1,2). Then

$$(k_1 \times k_2) ((x,y),(x',y')) = k_1(x,x')k_2(y,y')$$

is kernel on $\mathfrak{X}_1 \times \mathfrak{X}_2$.

 $\bullet \ \, \mathsf{For} \,\, \mathfrak{X} := \mathfrak{X}_1 = \mathfrak{X}_2 \colon$

$$k(x,x')=k_1(x,x')k_2(x,x')$$

is kernel on X.

Kernel: properties – Taylor series construction

Let

$$f(z) = \sum_{j=0}^{\infty} b_j z^j \quad (|z| < r)$$

with $r \in (0, \infty]$, $b_i \ge 0$ $(\forall j)$. Then

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

is kernel with ||x|| < r.

Kernel: properties – Taylor series construction

Let

$$f(z) = \sum_{j=0}^{\infty} b_j z^j \quad (|z| < r)$$

with $r \in (0, \infty]$, $b_i \ge 0$ $(\forall j)$. Then

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

is kernel with ||x|| < r.

• Example (exponential kernel, $b_j = \frac{1}{j!}$):

$$f(z) = e^z = \sum_{j=0}^{\infty} \frac{z^j}{j!},$$
 $k(x, x') = e^{\langle x, x' \rangle}.$

Kernel: properties – usage example

Sum, multiplication with a non-negative scalar, product ⇒ polynomial kernel.

Kernel: properties – usage example

- Sum, multiplication with a non-negative scalar, product ⇒ polynomial kernel.
- Gaussian kernel:

$$k(x, x') = e^{-\theta \|x - x'\|_2^2} = e^{-\theta \sum_i (x_i - x_i')^2} = \prod_i e^{-\theta (x_i - x_i')^2},$$

Kernel: properties – usage example

- Sum, multiplication with a non-negative scalar, product ⇒ polynomial kernel.
- Gaussian kernel:

$$k(x, x') = e^{-\theta \|x - x'\|_2^2} = e^{-\theta \sum_i (x_i - x_i')^2} = \prod_i e^{-\theta (x_i - x_i')^2},$$

$$k(x, x') = e^{-\theta (x - x')^2} = e^{-\theta (x^2 + x'^2 - 2xx')} = \left[e^{-\theta x^2} e^{-\theta x'^2} \right] \times e^{2\theta xx'}$$

by the product rule, composition rule $[M(x) = e^{-\theta x^2}]$ and Taylor-rule the result follows.

Kernel: reproducing view

Reproducing view \Rightarrow elements, kernel trick.

- Let \mathcal{H} be a Hilbert space of $\mathcal{X} \to \mathbb{R}$ functions.
- $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if for $\forall x \in \mathcal{X}$
 - $k(\cdot, x) \in \mathcal{H}$ ('generators'),

Kernel: reproducing view

Reproducing view ⇒ elements, kernel trick.

- Let \mathcal{H} be a Hilbert space of $\mathcal{X} \to \mathbb{R}$ functions.
- $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if for $\forall x \in \mathcal{X}, f \in \mathcal{H}$
 - $k(\cdot, x) \in \mathcal{H}$ ('generators'),
 - $\langle f, k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$ (reproducing property).

Kernel: reproducing view

Reproducing view ⇒ elements, kernel trick.

- Let \mathcal{H} be a Hilbert space of $\mathcal{X} \to \mathbb{R}$ functions.
- $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if for $\forall x \in \mathcal{X}, f \in \mathcal{H}$
 - $k(\cdot, x) \in \mathcal{H}$ ('generators'),
 - $\langle f, k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$ (reproducing property).

Specifically: $\forall x, y \in \mathcal{X}$,

$$k(x,y) = \langle k(\cdot,x), k(\cdot,y) \rangle_{\mathcal{H}}.$$

• Let $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a symmetric function.

- Let $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a symmetric function.
- **G** := $[k(x_i, x_j)]_{i,j=1}^n$: Gram matrix.

- Let $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a symmetric function.
- **G** := $[k(x_i, x_j)]_{i,j=1}^n$: Gram matrix.
- k is called positive definite, if

$$\mathbf{a}^T \mathbf{G} \mathbf{a} \geqslant 0$$

for
$$\forall n \geqslant 1$$
, $\forall \mathbf{a} \in \mathbb{R}^n$, $\forall (x_1, \dots, x_n) \in \mathcal{X}^n$.

- Let $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a symmetric function.
- **G** := $[k(x_i, x_j)]_{i,j=1}^n$: Gram matrix.
- k is called positive definite, if

$$\mathbf{a}^T \mathbf{G} \mathbf{a} \geqslant 0$$

for
$$\forall n \geq 1$$
, $\forall \mathbf{a} \in \mathbb{R}^n$, $\forall (x_1, \dots, x_n) \in \mathcal{X}^n$.

The different views are equivalent!

Convergence in RKHS norm

Convergence in RKHS \Rightarrow uniform convergence! (k: bounded).

Indeed:

$$\begin{aligned} |f(x)| &\stackrel{k: \text{r.k.}}{=} |\langle f, k(\cdot, x) \rangle_{\mathfrak{H}}| \stackrel{CBS}{\leqslant} ||k(\cdot, x)||_{\mathfrak{H}} ||f||_{\mathfrak{H}} \\ &\stackrel{k: \text{r.k.}}{=} \sqrt{k(x, x)} ||f||_{\mathfrak{H}}. \end{aligned}$$

Convergence in RKHS norm

Convergence in RKHS \Rightarrow uniform convergence! (k: bounded).

Indeed:

$$|f(x)| \stackrel{k: \text{r.k.}}{=} |\langle f, k(\cdot, x) \rangle_{\mathcal{H}}| \stackrel{CBS}{\leqslant} ||k(\cdot, x)||_{\mathcal{H}} ||f||_{\mathcal{H}}$$

$$\stackrel{k: \text{r.k.}}{=} \sqrt{k(x, x)} ||f||_{\mathcal{H}}.$$

Kernel k is called bounded if

$$\sup_{x\in\mathcal{X}}k(x,x)<\infty.$$

Kernel algorithms

- Given: $\{(x_i, y_i)\}_{i=1}^n$, $\mathcal{H} = \mathcal{H}(k)$.
- Task $(\lambda > 0)$:

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} [y_i - f(x_i)]^2 + \lambda \|f\|_{\mathcal{H}}^2 \to \min_{f \in \mathcal{H}}.$$

- Given: $\{(x_i, y_i)\}_{i=1}^n$, $\mathcal{H} = \mathcal{H}(k)$.
- Task $(\lambda > 0)$:

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} [y_i - f(x_i)]^2 + \lambda \|f\|_{\mathcal{H}}^2 \to \min_{f \in \mathcal{H}}.$$

Analytical solution:

$$f(x) = [k(x_1, x), \dots, k(x_n, x)](\mathbf{G} + \lambda nI)^{-1}[y_1; \dots; y_n],$$

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n.$$

- Given: $\{(x_i, y_i)\}_{i=1}^n$, $\mathcal{H} = \mathcal{H}(k)$.
- Task $(\lambda > 0)$:

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} [y_i - f(x_i)]^2 + \lambda \|f\|_{\mathcal{H}}^2 \to \min_{f \in \mathcal{H}}.$$

Analytical solution:

$$f(x) = [k(x_1, x), \dots, k(x_n, x)](\mathbf{G} + \lambda nI)^{-1}[y_1; \dots; y_n],$$

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n.$$

Question

How do we get this solution?

Assume for a moment (representer theorem):

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

Kernel ridge regression

Assume for a moment (representer theorem):

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

Multiplying the objective by n, using the reproducing property:

$$J(f) = \sum_{j=1}^{n} [y_j - \langle f, k(\cdot, x_j) \rangle_{\mathfrak{H}}]^2 + \lambda n \|f\|_{\mathfrak{H}}^2$$
$$= \|\mathbf{y} - \mathbf{G}\mathbf{a}\|_2^2 + (\lambda n)\mathbf{a}^T \mathbf{G}\mathbf{a}$$
$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{G}\mathbf{a} + \mathbf{a}^T [\mathbf{G}^2 + (\lambda n)\mathbf{G}]\mathbf{a}.$$

Kernel ridge regression

Assume for a moment (representer theorem):

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

Multiplying the objective by n, using the reproducing property:

$$J(f) = \sum_{j=1}^{n} [y_j - \langle f, k(\cdot, x_j) \rangle_{\mathfrak{H}}]^2 + \lambda n \|f\|_{\mathfrak{H}}^2$$
$$= \|\mathbf{y} - \mathbf{G}\mathbf{a}\|_2^2 + (\lambda n)\mathbf{a}^T \mathbf{G}\mathbf{a}$$
$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{G}\mathbf{a} + \mathbf{a}^T [\mathbf{G}^2 + (\lambda n)\mathbf{G}]\mathbf{a}.$$

Solving
$$\mathbf{0} = \frac{\partial J}{\partial \mathbf{a}}$$
, one gets $\mathbf{a}^* = (\mathbf{G} + \lambda n \mathbf{I})^{-1} \mathbf{y}$

Kernel ridge regression

Assume for a moment (representer theorem):

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

Multiplying the objective by n, using the reproducing property:

$$J(f) = \sum_{j=1}^{n} [y_j - \langle f, k(\cdot, x_j) \rangle_{\mathfrak{H}}]^2 + \lambda n \|f\|_{\mathfrak{H}}^2$$

= $\|\mathbf{y} - \mathbf{G}\mathbf{a}\|_2^2 + (\lambda n)\mathbf{a}^T \mathbf{G}\mathbf{a}$
= $\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{G}\mathbf{a} + \mathbf{a}^T [\mathbf{G}^2 + (\lambda n)\mathbf{G}]\mathbf{a}$.

Solving
$$\mathbf{0}=rac{\partial J}{\partial \mathbf{a}}$$
, one gets $\mathbf{a}^*=(\mathbf{G}+\lambda n\mathbf{I})^{-1}\mathbf{y}$ by

$$\frac{\partial \mathbf{a}^T \mathbf{B} \mathbf{a}}{\partial \mathbf{a}} = \left(\mathbf{B} + \mathbf{B}^T \right) \mathbf{a}, \frac{\partial \mathbf{c}^T \mathbf{a}}{\partial \mathbf{a}} = \mathbf{c}.$$

Representer theorem

Let $r:[0,\infty)\to\mathbb{R}$ be monotonically increasing. Then $\exists f\in\mathcal{H}(k)$ minimizer of

$$J(f) = c(x_1, y_1, f(x_1), \dots, x_n, y_n, f(x_n)) + r(\|f\|_{\mathcal{H}})$$

admitting the form

$$f = \sum_{i=1}^n a_i k(\cdot, x_i), \quad a_i \in \mathbb{R}.$$

Representer theorem - proof

Decompose f to $span(\{k(\cdot,x_i)\}_{i=1}^n)$ and its orthogonal complement:

$$f = f_D + f_\perp,$$
 $f_D = \sum_{i=1}^n a_i k(\cdot, x_i).$

Representer theorem - proof

Decompose f to $span(\{k(\cdot,x_i)\}_{i=1}^n)$ and its orthogonal complement:

$$f = f_D + f_\perp,$$
 $f_D = \sum_{i=1}^n a_i k(\cdot, x_i).$

Objective terms:

$$\frac{f(x_i)}{f(x_i)} = \langle f, k(\cdot, x_i) \rangle_{\mathcal{H}} = \langle f_D + f_{\perp}, k(\cdot, x_i) \rangle_{\mathcal{H}} \\
= \langle f_D, k(\cdot, x_i) \rangle_{\mathcal{H}} + \underbrace{\langle f_{\perp}, k(\cdot, x_i) \rangle_{\mathcal{H}}}_{=0},$$

Representer theorem - proof

Decompose f to $span\left(\{k(\cdot,x_i)\}_{i=1}^n\right)$ and its orthogonal complement:

$$f = f_D + f_\perp,$$
 $f_D = \sum_{i=1}^n a_i k(\cdot, x_i).$

Objective terms:

$$\begin{aligned} f(\mathbf{x}_{i}) &= \langle f, k(\cdot, \mathbf{x}_{i}) \rangle_{\mathcal{H}} = \langle f_{D} + f_{\perp}, k(\cdot, \mathbf{x}_{i}) \rangle_{\mathcal{H}} \\ &= \langle f_{D}, k(\cdot, \mathbf{x}_{i}) \rangle_{\mathcal{H}} + \underbrace{\langle f_{\perp}, k(\cdot, \mathbf{x}_{i}) \rangle_{\mathcal{H}}}_{=0}, \\ r(\|f\|_{\mathcal{H}}) &= r(\|f_{D} + f_{\perp}\|_{H}) \geqslant r(\|f_{D}\|_{\mathcal{H}}). \end{aligned}$$

Kernel principal component analysis

Kernel PCA

Let
$$\mathcal{H} = \mathcal{H}(\mathbf{k})$$
, $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{\mathcal{H}}$.

Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} \left\langle f, \phi(x_i) - \frac{1}{n} \sum_{j=1}^{n} \phi(x_j) \right\rangle^2 = var(f) \to \max_{f: \|f\|_{\mathcal{H}} \leqslant 1}.$$

Kernel PCA

Let
$$\mathcal{H} = \mathcal{H}(k)$$
, $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{\mathcal{H}}$.

Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} \left\langle f, \phi(x_i) - \frac{1}{n} \sum_{j=1}^{n} \phi(x_j) \right\rangle^2 = var(f) \to \max_{f: \|f\|_{\mathcal{H}} \leqslant 1}.$$

• The solution can be searched in the form $(f \leftrightarrow \mathbf{a})$:

$$f = \sum_{i=1}^{n} a_i \tilde{\phi}(x_i)$$

since component $\perp span(\{\phi(x_i)\}_{i=1}^n)$ has no contribution.



Kernel PCA

Let
$$\mathcal{H} = \mathcal{H}(k)$$
, $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{\mathcal{H}}$.

Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} \left\langle f, \phi(x_i) - \frac{1}{n} \sum_{j=1}^{n} \phi(x_j) \right\rangle^2 = var(f) \to \max_{f: \|f\|_{\mathcal{H}} \leqslant 1}.$$

• The solution can be searched in the form $(f \leftrightarrow \mathbf{a})$:

$$f = \sum_{i=1}^{n} a_i \tilde{\phi}(x_i)$$

since component $\perp span(\{\phi(x_i)\}_{i=1}^n)$ has no contribution.

• We will get an eigenvalue problem for a.



(Empirical) covariance operator

$$C := \frac{1}{n} \sum_{i=1}^{n} \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i).$$

 $c \otimes d$ is the analogue of cd^T :

$$(c \otimes d)(e) = c \langle d, e \rangle_{\mathcal{H}}.$$

Similarly to the finite-dimensional case:

$$Cf_j = \lambda_j f_j$$
.

Challenge

How do we solve this eigenvalue problem?

Computation of Cf_j

Assume *j* is fixed ($Cf = \lambda f$):

$$Cf = \left[\frac{1}{n}\sum_{i=1}^{n}\tilde{\phi}(x_{i})\otimes\tilde{\phi}(x_{i})\right]f$$

$$\stackrel{\otimes}{=} \frac{1}{n}\sum_{i=1}^{n}\tilde{\phi}(x_{i})\left\langle\tilde{\phi}(x_{i}),\sum_{j=1}^{n}a_{j}\tilde{\phi}(x_{j})\right\rangle_{\mathfrak{H}} = \frac{1}{n}\sum_{i=1}^{n}\tilde{\phi}(x_{i})\sum_{j=1}^{n}a_{j}\tilde{k}(x_{i},x_{j})$$

with
$$\tilde{\mathbf{G}} = \mathbf{H}\mathbf{G}\mathbf{H} = \left[\tilde{k}(x_i, x_j)\right]_{i,j=1}^n$$
, $\mathbf{H} = \mathbf{I}_n - \frac{1}{n}$.

Eigenvalue problem

- We want to solve $Cf = \lambda f$, $Cf = \frac{1}{n} \sum_{i=1}^{n} \tilde{\phi}(x_i) \sum_{j=1}^{n} a_j \tilde{k}(x_i, x_j)$.
- Idea: multiple by $\tilde{\phi}(x_r)$

$$\left\langle \tilde{\phi}(x_r), \lambda f \right\rangle_{\mathfrak{H}} = \left\langle \tilde{\phi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\phi}(x_j) \right\rangle_{\mathfrak{H}} = \lambda \underbrace{\sum_{j=1}^n a_j \tilde{G}_{rj}}_{(\tilde{\mathsf{G}}\mathsf{a})_{rj}},$$

Eigenvalue problem

- We want to solve $Cf = \lambda f$, $Cf = \frac{1}{n} \sum_{i=1}^{n} \tilde{\phi}(x_i) \sum_{j=1}^{n} a_j \tilde{k}(x_i, x_j)$.
- Idea: multiple by $\tilde{\phi}(x_r)$

$$\begin{split} \left\langle \tilde{\phi}(x_r), \lambda f \right\rangle_{\mathfrak{H}} &= \left\langle \tilde{\phi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\phi}(x_j) \right\rangle_{\mathfrak{H}} = \lambda \sum_{j=1}^n a_j \tilde{G}_{rj}, \\ \left\langle \tilde{\phi}(x_r), Cf \right\rangle_{\mathfrak{H}} &= \left\langle \tilde{\phi}(x_r), \frac{1}{n} \sum_{i=1}^n \tilde{\phi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j) \right\rangle_{\mathfrak{H}} \\ &= \frac{1}{n} \sum_{i=1}^n \tilde{G}_{ri} \sum_{j=1}^n a_j \tilde{G}_{ij} = \frac{1}{n} (\tilde{\mathbf{G}}^2 \mathbf{a})_{rj}. \end{split}$$

Eigenvalue problem

- We want to solve $Cf = \lambda f$, $Cf = \frac{1}{n} \sum_{i=1}^{n} \tilde{\phi}(x_i) \sum_{j=1}^{n} a_j \tilde{k}(x_i, x_j)$.
- Idea: multiple by $\tilde{\phi}(x_r)$

$$\begin{split} \left\langle \tilde{\phi}(x_r), \lambda f \right\rangle_{\mathfrak{H}} &= \left\langle \tilde{\phi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\phi}(x_j) \right\rangle_{\mathfrak{H}} = \lambda \sum_{j=1}^n a_j \tilde{G}_{rj}, \\ \left\langle \tilde{\phi}(x_r), Cf \right\rangle_{\mathfrak{H}} &= \left\langle \tilde{\phi}(x_r), \frac{1}{n} \sum_{i=1}^n \tilde{\phi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j) \right\rangle_{\mathfrak{H}} \\ &= \frac{1}{n} \sum_{i=1}^n \tilde{G}_{ri} \underbrace{\sum_{j=1}^n a_j \tilde{G}_{ij}}_{(\tilde{\mathbf{G}}\mathbf{a})_{ij}} = \frac{1}{n} (\tilde{\mathbf{G}}^2 \mathbf{a})_{rj}. \end{split}$$

• Eigenvalue problem: $\tilde{\mathbf{G}}^2 \mathbf{a} = n\lambda \tilde{\mathbf{G}} \mathbf{a}$, i.e. $\tilde{\mathbf{G}} \mathbf{a} = (n\lambda)\mathbf{a}$.

Ortogonal eigenvectors in kernel PCA

Taking two (eigenvector, eigenvalue) pairs:

$$f_1 = \sum_{i=1}^n a_{1i} \tilde{\phi}(x_i),$$
 $\tilde{\mathbf{G}} \mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$ $f_2 = \sum_{i=1}^n a_{2i} \tilde{\phi}(x_i),$ $\tilde{\mathbf{G}} \mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathfrak{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\phi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\phi}(x_j) \right\rangle_{\mathfrak{H}} = \mathbf{a}_1^T \tilde{\mathbf{G}} \mathbf{a}_2 = \mathbf{a}_1^T \lambda_2 \mathbf{a}_2.$$

Orthogonality ⇒ projection is easy

Projection of a new x^* to the first d-PCs:

$$\Pi[\tilde{\phi}(x^*)] = \sum_{j=1}^{d} \left\langle \tilde{\phi}(x^*), f_j \right\rangle_{\mathcal{H}} f_j.$$

Orthogonality \Rightarrow projection is easy

Projection of a new x^* to the first d-PCs:

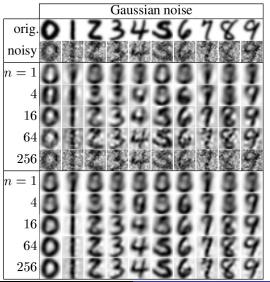
$$\Pi[\tilde{\phi}(x^*)] = \sum_{j=1}^{d} \left\langle \tilde{\phi}(x^*), f_j \right\rangle_{\mathfrak{H}} f_j.$$

For fixed $f = f_j$, using $f = \sum_{i=1}^n a_i \tilde{\phi}(x_i)$:

$$\left\langle \tilde{\phi}\left(x^{*}\right), f\right\rangle_{\mathfrak{H}} f = \sum_{i} a_{i} \tilde{k}\left(x_{i}, x^{*}\right) f = \sum_{i, j=1}^{n} a_{i} a_{j} \tilde{k}\left(x_{i}, x^{*}\right) \tilde{\phi}(x_{j}).$$

In denoising application

The pre-image problem to solve: $\widehat{x^*} = \arg\min_{x \in \mathcal{X}} \left\| \widetilde{\phi}(x) - \Pi[\widetilde{\phi}(x^*)] \right\|_{\mathcal{H}}^2$.



Summary

- We covered:
 - basic properties of kernels.
 - kernel ridge regression, representer theorem.
 - kernel PCA.
- References: [4-6].