

April 30, 2020

kernel:  $k(x, y) = k(y, x)$   $\in \mathbb{R}^d$  measure similarity

Feature maps:  $\psi_1, \dots, \psi_m$   $\psi_j: \mathbb{R}^d \rightarrow \mathbb{R}$

$\Psi: \mathbb{R}^d \rightarrow \mathbb{R}^m$   $\Psi(x) = \begin{pmatrix} \psi_1(x) \\ \vdots \\ \psi_m(x) \end{pmatrix}$   $m > n = \#$  of data points

Data:  $X = \{x_1, \dots, x_n\}$   $x_i \in \mathbb{R}^d$

$\star \min_c \|\mathbf{c}\|_2^2$  s.t.  $\Psi \mathbf{c} = \mathbf{f}_X$

$\hat{f} = \sum_{j=1}^m c_j \psi_j$   $\hat{f}_X \approx f$

$\Psi = \begin{bmatrix} \psi_1(x_1) & \dots & \psi_m(x_1) \\ \vdots & & \vdots \\ \psi_1(x_n) & \dots & \psi_m(x_n) \end{bmatrix}$

$\mathbf{f}_X = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}$

if  $\psi_1, \dots, \psi_m$  are orthonormal basis for some inner product space  $H$

$\star \min \|\hat{f}\|_H^2$  s.t.  $\hat{f}_X = f_X$

$$k_y \in H \quad k_y = \sum_j \underbrace{\psi_j(y)}_{\text{coeff}} \underbrace{\psi_j}_{\text{function}}$$

$$\begin{aligned} k_y(x) &= \sum_j \psi_j(y) \psi_j(x) \\ &= \langle \psi(y), \psi(x) \rangle_{\mathcal{L}_2} \\ &= k(y, x) = k(x, y) \end{aligned} \quad \begin{aligned} &= \langle k_y, k_x \rangle_H \\ &= \langle \sum_i \psi_i(y) \psi_i, \sum_j \psi_j(x) \psi_j \rangle_H \\ &= \sum_{i,j} \psi_i(y) \psi_j(x) \langle \psi_i, \psi_j \rangle_H \delta_{ij} \\ &= \sum_j \psi_j(y) \psi_j(x) \end{aligned}$$

$$k_y(x) = k(y, x) = \langle \psi(y), \psi(x) \rangle_{\mathcal{L}_2} = \langle k_y, k_x \rangle_H$$

RKHS:  $g(x) = \langle g, k_x \rangle_H$  for any  $g \in H$

Note: need to be careful in infinite dimensions

Note!  $x_1, \dots, x_n$   $(K_{xx})_{ij} = k(x_i, x_j)$

$$K_{xx} = \underline{\Psi} \underline{\Psi}^T \quad \approx \begin{bmatrix} \Psi \\ \Psi \\ \vdots \\ \Psi \end{bmatrix}^m$$

positive semi-def.

RKHS  $\Rightarrow$  kernels  $k_y = \sum \psi_j(y) \psi_j$

Can we just start with a kernel?

We need an inner product:  $\langle \cdot, \cdot \rangle_H$

$x_1, \dots, x_n$   $k_{x_i}$   $k_{x_i}(y) = k(x_i, y)$

given  $k(x, y)$

Suppose  $f = \sum_i a_i k_{x_i}$ ,  $g = \sum_j b_j k_{x_j}$

$$\begin{aligned} \text{Propose: } \langle f, g \rangle_H &= \langle \sum_i a_i k_{x_i}, \sum_j b_j k_{x_j} \rangle_H \\ &= \sum_{i,j} a_i b_j \langle k_{x_i}, k_{x_j} \rangle_H \\ &= \sum_{i,j} a_i b_j k(x_i, x_j) \\ &= a^T K_{xx} b \end{aligned}$$

proposal

native space RKHS

Works if  $f, g$  finite lin. combos of  $k_{x_i}$   
Need to "complete" space with Cauchy sequences.

$$\langle f, f \rangle_H = f^T K_{xx} f \geq 0 \quad (\text{need})$$

## Moore-Aronszajn theorem

If  $K_{xx}$  is positive semidef for any finite collection of points  $X$ , then there exists a corresp. unique RKHS

Difficult to characterize functions in native space

$$\min \|\hat{f}\|_H^2 \text{ s.t. } \hat{f}_X = f_X$$

$$\left[ \hat{f}(z) = K_{zX} K_{XX}^{-1} f_X \quad (\text{last lecture}) \right]$$

$\downarrow$   
 $[k(z, x_1) \dots k(z, x_n)]$

Can derive error bounds  $|\hat{f}(z) - f(z)|$  in terms of  $\|f\|_H^2, \|\hat{f}\|_H^2$

# Gaussian Processes (GPs)

Multivariate normal: distribution for vectors

$$y \sim N(\mu, \Sigma) \quad p(y) = \frac{\exp(-\frac{1}{2}(y-\mu)^T \Sigma^{-1}(y-\mu))}{\sqrt{(2\pi)^d \det(\Sigma)}}$$

GPs: "distribution" for functions  $\{f: \mathbb{R}^d \rightarrow \mathbb{R}\}$

① mean field  $\mu: \mathbb{R}^d \rightarrow \mathbb{R}$

② kernel function

Definition:  $f \sim GP(\mu, k)$  if for any  $x_1, \dots, x_n$

$$f_x \sim N(\mu_x, K_{xx})$$

$$\downarrow$$
$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$$

$$\downarrow$$
$$\begin{pmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{pmatrix}$$

$$\downarrow$$
$$(K_{xx})_{ij} = k(x_i, x_j)$$

Bayesians: start with prior, observe data, get posterior

Multivariate normal

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \sim N(0, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix})$$

$$p(y_2 | y_1) \sim N(\underline{\Sigma_{21} \Sigma_{11}^{-1} y_1}, \underline{\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}})$$

GP: prior  $\mu(z) = 0$  observe  $X = x_1, \dots, x_n$ ,  $f_x = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$

$\Rightarrow$  conditioning on  $X \dots$  gives new GP

$$f_y | x, f_x \sim N(\hat{\mu}_y, \hat{K}_{yy})$$

$$\hat{\mu}(z) = k_{zx} \underline{K_{xx}^{-1} f_x}$$

$[k(z, x_1) \dots k(z, x_n)]$

$\hat{\mu} = \underset{\text{s.t. } \hat{f}_x = f_x}{\text{argmin}} \| \hat{f} \|_H^2$

$$\hat{K}(z, z') = k(z, z') - k_{zx} K_{xx}^{-1} k_{xz}$$

$$k_{xz} = k_{zx}^T$$

$x_i$ : demo features,  $f(x_i)$  election results

$$f_{Y|X} \sim N(\hat{\mu}_Y, \hat{K}_{YY})$$

now we get quantifiable uncertainty

Another: time series:  $x_i$ : timestamps  $f(x_i)$  observed values

Example: noisy observations

$f \sim GP(0, k)$  observe:  $x_i, f(x_i) + \epsilon_i, \epsilon_i \sim N(0, \sigma^2)$

Given  $\sigma, X, f_x + \epsilon$ , posterior is

$$f_{Y|X, \sigma} \sim N(\tilde{\mu}_Y, \tilde{K}_{YY}) \quad \begin{pmatrix} \tilde{\mu}(y_1) \\ \vdots \\ \tilde{\mu}(y_n) \end{pmatrix}$$

$$\tilde{\mu}(z) = k_{zx} (K_{xx} + \sigma^2 I)^{-1} f_x$$

$$\tilde{k}(z, z') = k(z, z') - k_{zx} (K_{xx} + \sigma^2 I)^{-1} k_{xz}$$

$K_{xx} + \sigma^2 I$  instead of  $K_{xx}$

$$\hat{m} = \operatorname{argmin} \|\hat{f}\|_H^2$$

s.t.  $\hat{f}_x = f_x$

$$\tilde{m} = \operatorname{argmin}_{\hat{f}} \sigma^2 \|\hat{f}\|_H^2 + \|\hat{f}_x - f_x\|_2^2$$

observe:  $y = f_x + \epsilon$ ,  $\epsilon \sim N(0, \sigma^2 I)$

$$y \sim N(0, K_{xx} + \sigma^2 I)$$

$$\Sigma = K_{xx} + \sigma^2 I$$

$$\max_{\sigma} \mathcal{L}(\sigma | y) = p_0(y) = \frac{\exp(-\frac{1}{2} y^T \Sigma^{-1} y)}{\sqrt{(2\pi)^d \det(\Sigma)}}$$

$$\log \mathcal{L} \Rightarrow \underbrace{-\frac{1}{2} y^T \Sigma^{-1} y}_{\text{fidelity}} - \frac{d}{2} \log(2\pi) - \underbrace{\frac{1}{2} \log \det(\Sigma)}_{\text{complexity}}$$



$$\frac{\partial \log \det(\Sigma)}{\partial \sigma}$$

,

$$\text{trace}(\Sigma^{-1} 2\sigma \mathbf{I})$$

$$\frac{\partial \mathbf{y}^T \Sigma^{-1} \mathbf{y}}{\partial \sigma} = -\mathbf{y}^T \Sigma^{-1} (2\sigma \mathbf{I}) \Sigma^{-1} \mathbf{y}$$

$$\Sigma = \mathbf{K}_{xx} + \sigma^2 \mathbf{I}$$

~~Annoying computations!!~~