

April 23, 2019

function approx.

true $f: \mathbb{R}^n \rightarrow \mathbb{R}$

observe $f(x)$ $x \in X$, f_x

$\hat{f}: \mathbb{R}^n \rightarrow \mathbb{R}$

Today: kernel methods

$k(x, y)$ = similarity between
 x & y $x, y \in \mathbb{R}^n$

exp: $k(x, y) = s^2 \exp\left(-\frac{\|x - y\|^2}{\lambda}\right)$

cubic splines ($n=1$):

$$k(x, y) = |x - y|^3$$

① Feature maps

kernel is inner product
b/w feature vectors

② Quadratic forms on function
spaces

③ Covariance of random
process is given by
a kernel

⇒ Gaussian Processes

kernel on some set S

$$k: S \times S \Rightarrow \mathbb{R}$$

$$k(x, y) = k(y, x)$$

linear regression

$$f(x) \approx \hat{f}(x) = x^T c$$

$$\min_c \|X^T c - f_X\|_2^2$$

rows are data
points $x \in X$

true $f(x)$
at $x \in X$

$$\hat{f}(x) = \psi(x)^T c$$

$$= \sum_{j=1}^N c_j \psi_j(x)$$

ψ is the feature map

$$\min_c \|\Psi^T c - f_x\|_2^2$$

$$\Psi_{ij} = \psi_i(x_j)$$

if N (# of features) \leq
 m (# of data points)

\Rightarrow standard OLS

$$N > m$$

$$\boxed{\Psi^T} c = f_x$$

$$\min \|c\|_2^2$$

$$\text{s.t. } \Psi^T c = f_x$$

$$c = \underline{\Psi} (\underline{\Psi}^T \underline{\Psi})^{-1} f_x$$

$$\hat{f}(z) = \Psi(z)^T c$$

$$= \Psi(z)^T \underline{\Psi} (\underline{\Psi}^T \underline{\Psi})^{-1} f_x$$

$$k(x, \gamma) = \Psi(x)^T \Psi(\gamma)$$

$$\hat{f}(z) = k_z X K_{XX}^{-1} f_x$$

$$[K_{XX}]_{ij} = k(x_i, x_j)$$

$$[k_{zX}]_i = k(z, x_i)$$

$$x_i, x_j \in X$$

Only need $k(x, y) = \psi(x)^\top \psi(y)$
called "kernel trick"

If we start from kernel,
e.g. $k(x, y) = \exp(-\|x - y\|^2)$,
is there an associated
feature map?

Yes! (Mercer's theorem)

if K_{xx} is positive semi-def.
for any finite set X , then
there is an associated
feature map

$$\hat{f}(x) = \sum_i c_i \psi_i(x)$$

$$\min \|c\|^2$$

$$\text{s.t. } \underline{\psi}^T c = f_x$$

implicitly, defined an inner
product on functions with
orthonormal basis $\{\psi_1, \dots, \psi_N\}$

call this H

$$\min \|\hat{f}\|_{\mathcal{H}}^2$$

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

point evaluation

$$\hat{f}(z)$$

$$k_z(x) \triangleq \sum_i \psi_i(z) \psi_i(x)$$

$$\langle \hat{f}, k_z \rangle_{\mathcal{H}}$$

$$= \left\langle \sum_i c_i \psi_i, \sum_j \psi_j(z) \psi_j \right\rangle$$

$$= \sum_{i,j} c_i \psi_j(z) \langle \psi_i, \psi_j \rangle_{\mathcal{H}}$$

$$(\langle \psi_i, \psi_j \rangle_{\mathcal{H}} = \delta_{ij})$$

$$= \sum c_i \psi_i(z)$$

$$= \hat{f}(z)$$

More generally,

$$g(z) = \langle g, k_z \rangle_{\mathcal{H}}$$

\Rightarrow It is a reproducing
Kernel Hilbert space (RKHS)

Can we get the inner
product w/o explicitly
computing the feature map?

Main idea:

$$\left\langle \sum_i c_i k_{x_i}, \sum_j d_j k_{x_j} \right\rangle_{\mathcal{H}}$$

$$k_{x_i}(z) = k(x_i, z)$$

$$\begin{aligned} &\rightarrow \sum_{i,j} c_i d_j \langle k_{x_i}, k_{x_j} \rangle_{\mathcal{H}} \\ &= c^T K_{xx} d \end{aligned}$$

(need to extend to
infinite sums)

Error analysis

Approx $f \in \mathcal{H}$ using X

$$\min \|\hat{f}\|_{\mathcal{H}}^2$$

$$\text{s.t. } \hat{f}_X \approx f_X$$

$$|\hat{f}(y) - f(y)| < ?$$

for $y \notin X$

Define \tilde{f} by

$$\min \|\tilde{f}\|_{\mathcal{H}}^2$$

$$\text{s.t. } \tilde{f}_X \approx f_X$$

$$\tilde{f}(y) = f(y)$$

$$\|\hat{f}\|_{\mathcal{H}}^2 \leq \|\tilde{f}\|_{\mathcal{H}}^2 \leq \|f\|_{\mathcal{H}}^2$$

$$X' \triangleq X \cup \{y\}$$

$$e \triangleq \tilde{f} - \hat{f}$$

$$e_{X'} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ e(y) \end{bmatrix}$$

$$e(y) = \tilde{f}(y) - \hat{f}(y)$$

$$\langle e, \hat{f} \rangle_{\mathcal{H}} = 0$$

$$\|e\|_{\mathcal{H}}^2 + \|\hat{f}\|_{\mathcal{H}}^2 = \|\tilde{f}\|_{\mathcal{H}}^2$$

$$\|e\|_{\mathcal{H}}^2 \leq \|f\|_{\mathcal{H}}^2 - \|\hat{f}\|_{\mathcal{H}}^2$$

$$\|e\|_{\mathcal{H}}^2 = d^{\top} K_{X'X'} d$$

$$(d = K_{x'x'}^{-1} e_{x'}) \text{ (can show)}$$

$$= e_{x'}^T K_{x'x'}^{-1} K_{x'x'} K_{x'x'}^{-1} e_{x'}$$

$$= e_{x'}^T K_{x'x'}^{-1} e_{x'}$$

$$= (K_{x'x'}^{-1})_{yy} e(y)^2$$

Putting the algebra together...

$$e(y)^2 \leq \frac{1}{(K_{x'x'}^{-1})_{yy}} (\|F\|_K^2 - \|\hat{f}\|_K^2)$$

$$R(y,y) = k_{yx} K_{xx}^{-1} k_{xy}$$

no x'

$$k_{yx}(x_i) = k(y, x_i)$$

$$x_i \in X \quad "$$

$$k_{xy}(x_i) = k(x_i, y)$$

Gaussian Processes

GP on a set $S \subseteq \mathbb{R}^d$

with mean field μ

covariance kernel k

$$f \sim \text{GP}(\mu, k)$$

for any $X \subseteq S$ finite

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

multivariate normal,

$$Y \sim N(0, K)$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}$$

Theorem

$$Y_2 | Y_1 = y_1 \sim N \left(\begin{matrix} K_{12} K_{22}^{-1} y_1 \\ K_{11} - K_{12} K_{22}^{-1} K_{21} \end{matrix} \right)$$

idea with GPs:

- observe set of points X

- condition on X
- get a new GP

Observe f_X

$$\hat{f} \sim \text{GP}(\hat{\mu}, \hat{k})$$

$$\hat{\mu}(z) = k_{zX} K_{XX}^{-1} f_X + \mu$$

→ same as least squares featurization

$$\hat{k}(y, z) = k(y, z) - k_{yX} K_{XX}^{-1} k_{Xz}$$

$$k_{yX}(x) = k(y, x) \quad \left. \vphantom{k_{yX}(x)} \right\} x \in X$$

$$k_{zX}(x) = k(z, x)$$

$$\hat{k}(y, y) = k(y, y) - k_{yX} K_{XX}^{-1} k_{Xy}$$

→ same as in error bound

Two ways to think about error in kernel methods for function approx.

$$\textcircled{1} \quad \min \|\hat{f}\|_K^2$$

$$\text{s.t. } \hat{f}_X = f_X$$

worst-case error bound

$\textcircled{2}$ Bayesian inference

Assume $f \sim \mathcal{GP}(\mu, k)$ (prior)

Observe X, f_X

Update to get \hat{f} (posterior)

MLE with GPs

k parameterized by some θ

$$\mathcal{L}(\theta) = \log(p(X; \theta))$$

$$p(y) \approx \frac{1}{\sqrt{\det(2\pi K)}}$$

$$\exp\left(-\frac{1}{2}(y - \mu)^T K^{-1}(y - \mu)\right)$$

$$\log p(y)$$

$$\approx -\frac{1}{2}(y - \mu)^T K^{-1}(y - \mu) \quad \text{fidelity}$$

$$-\frac{1}{2} \log \det K \quad \text{complexity}$$

$$-\frac{n}{2} \log(2\pi) \quad \text{constant}$$

$$\frac{d}{d\theta} \log(p(y))$$

$$\frac{d}{d\theta} [\text{fidelity}]$$

$$\approx \frac{1}{2} (y - \mu)^T C^T \left(\frac{d}{d\theta} K \right) C$$

$$C \approx K^{-1} (y - \mu)$$

requires solving ^{linear} system

$$\frac{d}{d\theta} [\text{complexity}]$$

$$\approx \text{tr} \left(K^{-1} \frac{d}{d\theta} K \right)$$

also expensive

Efficient optimization
active area of research
(at Cornell)