

Lecture 10: Feature map vs Kernel

Perspectives

So far our discussion of Kernel methods & in particular, representer theorems has been focused on the so called "Kernel Perspective"

$$\begin{aligned} \xrightarrow{\text{equir}} \min_{f \in \mathcal{H}} \quad & l(f(x_1), \dots, f(x_n)) + R(\|f\|) \\ \xrightarrow{\text{ }} \min_{z \in \mathbb{R}^n} \quad & l(z) + R\left(\sqrt{z^\top K(X, X) z}\right) \end{aligned}$$

The second formula involving the kernel matrix $K(X, X)$, hence the name "Kernel Perspective".

Conceptually this method is very elegant but it has two potential issues.

(i) $K(X, X)$ can be ill-conditioned

(ii) $K(X, X) \in \mathbb{R}^{n \times n}$ & it is typically dense

the loss l involves more points ie $|X|=n$ groups, the cost of computing $K(X, X)^{-1} z$ grows rapidly!

①

Many potential workarounds exist for this issue some keywords are

- Nystrom or low rank approx
- Sparse GPs
- Sparse Cholesky factorizations

We won't have time to cover these but focus our attention on the so called "feature map perspective" which is both efficient & mathematically enlightening.

10.1 The feature map perspective

Recall that we showed in the RKHS thm.

That every PPS Kernel K is associated with a feature map $\varphi: \mathcal{X} \rightarrow \mathcal{F}$

$$K(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{F}}$$

a Hilbert Space

We showed (in HW2) that indeed the RKHS H associated with K is unique! But the feature maps φ are not unique.

(2)

To see this consider the Canonical feature map

$$\varphi(x) = K(x, \cdot), \quad \varphi: X \rightarrow H$$

as one example. At the same time, by Mercer's theorem we have that

$$K(x, y) = \sum_{j=1}^{\infty} \lambda_j \psi_j(x) \psi_j(y) \quad \text{in } L^2(X, \mu)$$

$$\tilde{\psi}_j = \sqrt{\lambda_j} \psi_j$$

$$= \sum_{j=1}^{\infty} \tilde{\psi}_j(x) \tilde{\psi}_j(y)$$

$$= \left\langle \{\tilde{\psi}_j(x)\}, \{\tilde{\psi}_j(y)\} \right\rangle_{L^2}$$

Thus, we could also pick

$$\varphi: X \rightarrow L^2(X, \mu)$$

$$\varphi(x) = \sum_{j=1}^{\infty} \tilde{\psi}_j(x) \tilde{\psi}_j$$

Indeed observe that, by Mercer's theorem

$$K(x, x) = \sum_{j=1}^{\infty} \lambda_j \tilde{\psi}_j(x)^2 < +\infty \quad \text{and so,}$$

$$\begin{aligned} \left\| \sum_{j=1}^{\infty} \sqrt{\lambda_j} \tilde{\psi}_j(x) \tilde{\psi}_j \right\|_{L^2(X, \mu)}^2 &= \left\langle \sum_{j=1}^{\infty} \sqrt{\lambda_j} \tilde{\psi}_j(x) \tilde{\psi}_j, \sum_{k=1}^{\infty} \sqrt{\lambda_k} \tilde{\psi}_k(x) \tilde{\psi}_k \right\rangle_{L^2(X, \mu)} \\ &= \sum_{j=1}^{\infty} \lambda_j \tilde{\psi}_j(x)^2 < +\infty \end{aligned}$$

③

Thus, $\sum_{j=1}^{\infty} \tilde{\psi}_j(x) \psi_j \in L^2(x, \mu)$. Now consider

$$f = \sum_{j=1}^{\infty} \alpha_j \psi_j$$

$$g = \sum_{j=1}^{\infty} \beta_j \psi_j$$

$$\{\alpha_j\}_{j=1}^{\infty}, \{\beta_j\}_{j=1}^{\infty} \in l^2$$

& inner product $\langle f, g \rangle := \sum_{j=1}^{\infty} \frac{1}{\lambda_j} \alpha_j \beta_j$.

Now observe,

$$\begin{aligned} \langle f, K(x, \cdot) \rangle &= \left\langle \sum_{j=1}^{\infty} \alpha_j \psi_j, \frac{1}{\lambda_j} \sum_{k=1}^{\infty} \lambda_k \psi_k(x) \psi_k \right\rangle_{L^2(x, \mu)} \\ &= \sum_{j=1}^{\infty} \alpha_j \psi_j(x) = f(x) \end{aligned}$$

(Reproducing Property)

$$\langle f, f \rangle = \sum_{j=1}^{\infty} \frac{\alpha_j^2}{\lambda_j} < +\infty$$

$\left\{ \frac{\alpha_j}{\sqrt{\lambda_j}} \right\} \in l^2$

Based on these calculations (skipping detailed proof) we can show that

Then: Suppose Mercer's theorem holds. Then the RKHS of a PDS kernel K can be identified as

$$H = \left\{ f : X \rightarrow \mathbb{R} \mid f = \sum_{j=1}^{\infty} \tilde{\alpha}_j \tilde{\psi}_j, \sum_{j=1}^{\infty} \tilde{\alpha}_j^2 < +\infty \right\}$$

(4)

Note: The $\tilde{\psi}_j$ in our notation are not the eigenfunctions of K & have been scaled by the eigenvalues $\lambda_j = \sqrt{\lambda_j} \psi_j$.

10.2 Linear methods & Kernels

The above, spectral characterization of RKHS's has important implications in practice. Let us demonstrate the ideas with an example in Regression:

Given $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ & data

$$\mathbb{R}^d \ni \underline{y} = f^*(X) + \underline{\xi}, \quad \underline{\xi} \stackrel{iid}{\sim} N(0, \sigma^2 I)$$

Approximate f^*

The standard kernel approach via repn gives

$$f^* = \underset{f \in H}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|f(X) - \underline{y}\|_2^2 + \lambda \|f\|^2$$

$$\Rightarrow f^* = K(x, X) K(X, X)^{-1} \underline{z}^*$$

(5) When $\underline{z}^* = \underset{\underline{z} \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\underline{z} - \underline{y}\|_2^2 + \lambda \underline{z}^T K(X, X)^{-1} \underline{z}$

Now suppose n is large, then we need to invert $K(X|X) \in \mathbb{R}^{n \times n}$ which is also dense!

Linear solve will cost us $\mathcal{O}(n^3)$!

Note: Standard linear solvers are ineff. because they do not use the structure of the kernel matrix.

There are many ideas for improving this performance & this is an active area of research. e.g. see

- Snelson & Ghahramani, "Local & Global Sparse Gaussian Process Approximations" (2007)
- Cao et al, "Variational sparse Cholesky approximation for latent Gaussian processes via double Kullback-Liebler minimization" (2023).

We can also, reformulate the above problem using the spectral approach to RKHSs as

$$f^* = \underset{f \in H}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|f(X) - y\|_2^2 + \lambda \sum_{j=1}^{\infty} \tilde{\alpha}_j^2$$

$$\text{s.t. } f = \sum_{j=1}^{\infty} \tilde{\alpha}_j \tilde{\psi}_j$$

Or equivalently, $\hat{f}^* = \sum_{j=1}^{\infty} \tilde{\alpha}_j^* \tilde{\psi}_j$ when

$$\tilde{\alpha}^* = \underset{\tilde{\alpha} \in \mathbb{R}^{\infty}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \left\| \sum_{j=1}^{\infty} \tilde{\alpha}_j \tilde{\psi}_j(x) - \underline{y} \right\|_2^2 + \lambda \|\tilde{\alpha}\|_2^2$$

This is an instance of a so-called "feature map" formulation of our regression problem. We can further approximate \hat{f}^* simply by truncating the expansion:

$$\tilde{f} = \sum_{j=1}^m \tilde{\alpha}_j \tilde{\psi}_j \quad \text{& solve}$$

$$\tilde{\alpha}^* = \underset{\tilde{\alpha} \in \mathbb{R}^m}{\operatorname{argmin}} \frac{1}{2\sigma^2} \left\| \sum_{j=1}^m \tilde{\alpha}_j \tilde{\psi}_j(x) - \underline{y} \right\|_2^2 + \lambda \|\tilde{\alpha}\|_2^2$$

Indeed, define

$$A_{ij} = \tilde{\psi}_j(x_i)$$

then we have

$$\tilde{\alpha}^* = \underset{\tilde{\alpha}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|A\tilde{\alpha} - \underline{y}\|_2^2 + \lambda \|\tilde{\alpha}\|_2^2$$

which is nothing more than a penalized least squares problem!

(7)

The feature map perspective allows us to solve a problem for $\tilde{\alpha}^* \in \mathbb{R}^m$ after a small truncation error versus solving for $\tilde{z}^* \in \mathbb{R}^n$. If the kernel K is "nice" ie, its eigenvalues λ_j decay rapidly then we can take $m \ll n$.

A few observations:

(i) The above calculations show that all penalized least squares problems, & more broadly, all linear methods with ansatz of the form $f = \sum_j \alpha_j \psi_j$ with ℓ^2 -penalty on $\{\alpha_j\}$ coincide with kernel methods!

(ii) The function up to m -terms is the same as a rank m -approx. of the kernel K ie,

$$\tilde{K}(x, y) = \sum_{j=1}^m \tilde{\psi}_j(x) \tilde{\psi}_j(y) \approx K(x, y)$$

8

