

Consequences of Mercer's Theorem + NTK

Last time

Theorem (Mercer)

Suppose X is compact and $k: X \times X \rightarrow \mathbb{R}$ is positive-definite and continuous. Let ν be a finite Borel measure supported on all of X . Then

Define T_k as

$$(T_k f)(x) = \int_X k(x, y) f(y) d\nu(y)$$

Claim. T_k has eigenbasis e_1, e_2, \dots
 μ_1, μ_2, \dots

$$\text{s.t. } k(x, y) = \sum_{m=1}^{\infty} \mu_m e_m(x) e_m(y).$$

Corollary,

For any f ,

$$\|f\|_{\mathcal{H}}^2 = \sum_{m=1}^{\infty} c_m^2 / \mu_m, \text{ where}$$

$$c_m = \int f(x) e_m(x) d\nu(x).$$

Slightly weird b/c \mathcal{X} RHS depends on ν ,
but LHS does not.

Intuition

$$\|f\|_{\mathcal{H}}^2 = \sup_{S=\{x_1, \dots, x_n\}} f_S^T K_S^{-1} f_S$$

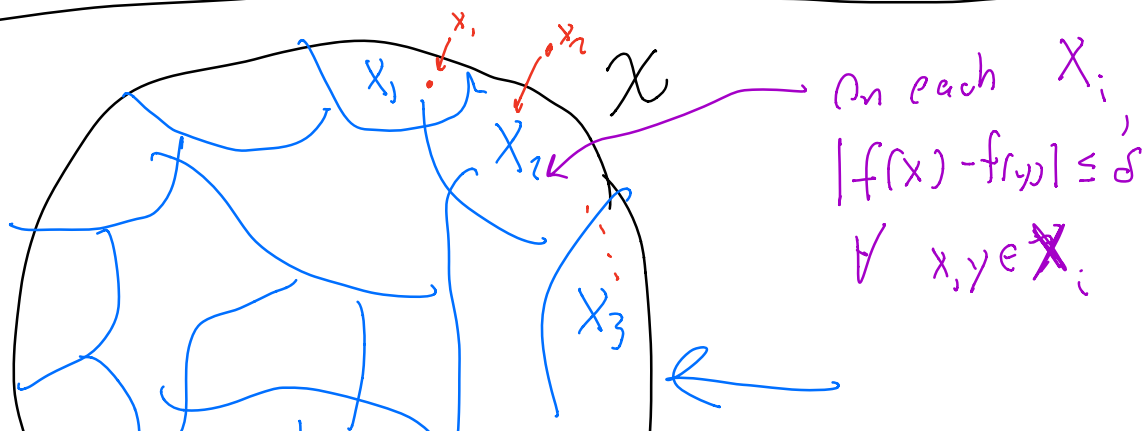
$$f_S = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}$$

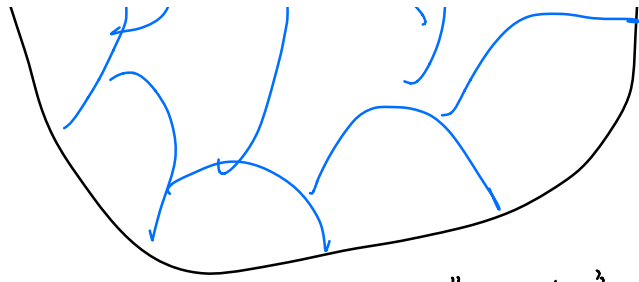
Sample $\overbrace{x_1, \dots, x_n}^S \sim \nu$

Since $\text{supp}(\nu) = \mathcal{X}$, "eventually" these samples
contain everything.

As $n \rightarrow \infty$, $f_S^T K_S^{-1} f_S \rightarrow \|f\|_{\mathcal{H}}^2$

K_S, f_S clearly do depend on ν Happens independently of ν





Think of x_i as "samples", but chosen this way for ease of analysis.

as $\delta \rightarrow 0$

①

$$S = \{x_1, \dots, x_n\}$$

$$f_S^T K_S^{-1} f_S \rightarrow \|f\|_H^2$$

as S gets large.

②

$$f_S^T K_S^{-1} f_S \rightarrow \sum_{m=1}^{\infty} \frac{1}{\mu_m} \left(\int_{\mathcal{X}} f(x) e_m(x) d\nu(x) \right)^2$$

Why ② is true.

$\nu(x_1)$ small but $\nu(x_2)$ is large

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

$$K_S = \begin{bmatrix} k(x_1, x_1) & \dots & \vdots \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix}$$

$$f_{\sqrt{}} = \begin{bmatrix} \sqrt{v(x_1)} f(x_1) \\ \vdots \\ \sqrt{v(x_n)} f(x_n) \end{bmatrix} \quad (K_{\sqrt{}})_{ij} = \sqrt{v(x_i)v(x_j)} \cdot k(x_i, x_j)$$

$$f_{\sqrt{}}^T f_{\sqrt{}} = \sum_i v(x_i) f(x_i)^2$$

expand this $\approx \int_{\mathcal{X}} f(x)^2 d\nu(x)$
 ↓ using Mercer

$$f_{\sqrt{}}^T K_{\sqrt{}}^{-1} f_{\sqrt{}} = f_S^T K_S^{-1} f_S$$

$$f_{\sqrt{}} = \text{diag}(v(x_i))^{1/2} f_S$$

$$K_{\sqrt{}} = \text{diag}(v(x_i))^{1/2} K_S \text{diag}(v(x_i))^{1/2}$$

$$f_{\sqrt{}}^T K_{\sqrt{}}^{-1} f_{\sqrt{}} = f_{\sqrt{}}^T \left(\sum_{m=1}^{\infty} \frac{1}{\mu_m} \begin{bmatrix} v(x_1)^{1/2} e_m(x_1) \\ \vdots \\ v(x_n)^{1/2} e_m(x_n) \end{bmatrix} \begin{bmatrix} v(x_1)^{1/2} e_m(x_1) \\ \vdots \\ v(x_n)^{1/2} e_m(x_n) \end{bmatrix}^T \right) f_{\sqrt{}}$$

$$D^{\frac{1}{2}} K S^{-1} D^{\frac{1}{2}} = \sum \mu_m e_m e_m^T$$

$D = \text{diag}(\sqrt{X_1}, \dots, \sqrt{X_n})$

μ_m, e_m : $k(x, y) = \sum_{m=1}^{\infty} \mu_m e_m(x) e_m(y)$

And e_m is an eigenfunction of \hat{T}_K
w/ eigenvalue μ_m

Note.

$$\begin{bmatrix} \sqrt{X_1}^{\frac{1}{2}} e_m(x_1) \\ \vdots \\ \sqrt{X_n}^{\frac{1}{2}} e_m(x_n) \end{bmatrix} = h_m$$

$$\langle h_m, h_l \rangle \approx 0 \quad \text{if } m \neq l$$

$$\approx 1 \quad \text{if } m = l$$

$$\langle h_m, h_l \rangle = \sum_{i=1}^n \sqrt{X_i} e_m(x_i) e_l(x_i)$$

$$\approx \int e_m(x) e_l(x) dx(x)$$

$$= \delta_{m=l}$$

$$f_V^T K_V^{-1} f_V \approx f_V^T \left(\sum_{m=1}^{\infty} \frac{1}{\mu_m} h_m h_m^T \right) f_V$$

$$= \sum_{m=1}^{\infty} \frac{1}{\mu_m} \langle f_V, h_m \rangle$$

$$= \sum_{k=1}^n v(x_i) f(x_i) e_m(x_i)$$

$$\approx \int f(x) e_m(x) v(x) dx = C_m$$

$$\sum_{m=1}^{\infty} \frac{C_m^2}{\mu_m}$$

$$\sup_x K(x, x) < \infty \quad \leftarrow \leq M$$

$$\|f\|_{\mathcal{H}} < \infty$$

$$\Leftrightarrow \|f\|_{\infty} < \infty$$

$$\|f\|_{\mathcal{H}}^2 = \sup_S f^T K_S^{-1} f_S$$

$$\geq \sup_x f(x) K(x, x)^{-1} f(x)$$

$$\geq \frac{1}{M} \sup_x f(x)^2$$

$$\Rightarrow |f(x)| \leq \|f\|_{\mathcal{H}} \cdot \sqrt{M}$$

$V = P^{\dagger}$: eigenvalues $\mu_m \Rightarrow$ eigenvalues λ_m of S

If $K(a, \omega)$ is t times differentiable,
 this is enough to bound μ_m .

Theorem (Kühn, 1987).

Suppose \mathcal{X} is compact d -dimensional
 smooth manifold and K has continuous
 t^{th} derivative. Then $\mu_m = O(m^{-t/d-1})$.

Consequences

$\nu = p^{\text{tr}} \Rightarrow T_K \rightarrow \frac{1}{n}K$ on $X_1, \dots, X_n - p^{\text{tr}}$
 as $n \rightarrow \infty$

$\rightarrow \textcircled{S}$

eigenvalues: power law decay $m^{-\alpha}$
 w/ $\alpha = \frac{t}{d} + 1$.

\Rightarrow error: $n^{\frac{1}{\alpha}-1}$ if σ large \Rightarrow $n^{-\frac{t}{t+d}}$

$$n^{1-\alpha} \text{ if } \sigma \text{ small} \Rightarrow n^{-f/d}$$



Classical Most people about
learning rates for these smoother classes.

Another application Random features method

Kernels \leftrightarrow random features

$$\varphi(x) = \frac{1}{\sqrt{m}} \begin{bmatrix} \varphi_1(x) \\ \vdots \\ \varphi_m(x) \end{bmatrix} \quad \varphi_j \sim P$$

$$\langle \varphi(x), \varphi(y) \rangle = \frac{1}{m} \sum_{j=1}^m \varphi_j(x) \varphi_j(y)$$
$$\rightarrow \mathbb{E}_{\varphi \sim P} [\varphi(x) \varphi(y)]$$

$$\Rightarrow K(x, y) = \mathbb{E}_{\varphi \sim P}[\varphi(x)\varphi(y)]$$

$$K(x, y) = \sum_{\ell} \mu_{\ell} e_{\ell}(x) e_{\ell}(y)$$

$$\varphi = \sqrt{Z} e_{\ell} \quad \text{with probability } \frac{\mu_{\ell}}{Z}$$

$$Z = \mu_1 + \mu_2 + \dots$$

$$\left(\varphi \sim P \Rightarrow K = \mathbb{E}_{\varphi}[\varphi(x)\varphi(y)] \right)$$

$$\left(K \Rightarrow (\mu_{\ell}, e_{\ell}) \Rightarrow P(\varphi = e_{\ell}) \propto \mu_{\ell} \right)$$

$$\text{Theorem (Bochner), } \underbrace{K(x, y) = K(x-y)}_r$$

$$k(x, y) = \int e^{-i\langle w, x-y \rangle} d\mu(w)$$

$$\mu \propto \hat{k}(w) = \int e^{i\langle w, \Delta \rangle} k(\Delta) d\Delta$$

Mercer: $k(x, y) = \sum_{i\langle w, x \rangle} \mu_m e_m(x) \overline{e_m(y)}$

$$e_m(x) = e^{-i\langle w, x \rangle}$$

$$\sum \mu_m e_m(x) \overline{e_m(y)} = \sum \mu_m e^{-i\langle w, x-y \rangle}$$

Example 1.

$$k(x, y) = e^{-\|x-y\|_2^2 / 2\sigma^2}$$

$$k(\Delta) = e^{-\|\Delta\|_2^2 / 2\sigma^2}$$

$$\hat{k}(w) \propto e^{-\sigma^2 \|w\|_2^2 / 2}$$

$$\hat{k}(w) \propto e^{-\sigma^2 \|w\|_2^2 / 2}$$

$$k(x, y) \propto \int e^{i\langle w, x-y \rangle} e^{-\sigma^2 \|w\|_2^2 / 2} dw$$

⇒ random features:

Feature j : sample $w_j \sim \mathcal{N}(0, \frac{1}{\sigma^2} I)$

$$\phi_j(x) = e^{i \langle w_j, x \rangle}$$

← has complex numbers i

$$\mathbb{E}_q[\phi(x) \bar{\phi}(y)] = k(x, y)$$

$$\phi_j(x) = \cos(\langle w_j, x \rangle + b), \quad b \sim \text{Uniform}([0, 2\pi])$$

Example 2. $k(x, y) = e^{-\|x-y\|_1 / \lambda}$

$$k(\Delta) = e^{-\|\Delta\|_1 / \lambda}$$

$$\hat{k}(w) \propto \prod_{j=1}^d \frac{1}{1 + \lambda^2 w_j^2}$$

⇒ Sample w from Cauchy dist^d

$$\cos(\langle w, x \rangle + b)$$

Rehimi & Recht, '07

