

Optimal Rates for Random Fourier Feature Approximations

Zoltán Szabó*
Gatsby Unit, UCL

Joint work with Bharath K. Sriperumbudur*
Department of Statistics, PSU
(*equal contribution)

Department of Computing Science
University of Alberta
November 24, 2015

- Kernels and kernel derivatives.
- Random Fourier features (RFFs).
- Guarantees on RFF approximation: uniform, L^r .

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel on \mathcal{X} , if
 - $\exists \varphi : \mathcal{X} \rightarrow H$ (ilbert space) feature map,
 - $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_H$ ($\forall a, b \in \mathcal{X}$).

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel on \mathcal{X} , if
 - $\exists \varphi : \mathcal{X} \rightarrow H$ (ilbert space) feature map,
 - $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_H$ ($\forall a, b \in \mathcal{X}$).
- Kernel examples: $\mathcal{X} = \mathbb{R}^d$ ($p > 0, \theta > 0$)
 - $k(a, b) = (\langle a, b \rangle + \theta)^p$: polynomial,
 - $k(a, b) = e^{-\|a-b\|_2^2 / (2\theta^2)}$: Gaussian,
 - $k(a, b) = e^{-\theta \|a-b\|_2}$: Laplacian.

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel on \mathcal{X} , if
 - $\exists \varphi : \mathcal{X} \rightarrow H$ (hilbert space) feature map,
 - $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_H$ ($\forall a, b \in \mathcal{X}$).
- Kernel examples: $\mathcal{X} = \mathbb{R}^d$ ($p > 0, \theta > 0$)
 - $k(a, b) = (\langle a, b \rangle + \theta)^p$: polynomial,
 - $k(a, b) = e^{-\|a-b\|_2^2 / (2\theta^2)}$: Gaussian,
 - $k(a, b) = e^{-\theta \|a-b\|_2}$: Laplacian.
- In the $H = H(k)$ RKHS ($\exists!$): $\varphi(u) = k(\cdot, u)$.

Kernel: example domains (\mathcal{X})

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



Kernel: application example – ridge regression

- Given: $\{(x_i, y_i)\}_{i=1}^{\ell}$, $H = H(k)$.
- Task: find $f \in H$ s.t. $f(x_i) \approx y_i$,

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} [f(x_i) - y_i]^2 + \lambda \|f\|_H^2 \rightarrow \min_{f \in H} \quad (\lambda > 0).$$

Kernel: application example – ridge regression

- Given: $\{(x_i, y_i)\}_{i=1}^{\ell}$, $H = H(k)$.
- Task: find $f \in H$ s.t. $f(x_i) \approx y_i$,

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} [f(x_i) - y_i]^2 + \lambda \|f\|_H^2 \rightarrow \min_{f \in H} \quad (\lambda > 0).$$

- Analytical solution, $\mathcal{O}(\ell^3)$ – expensive:

$$f(x) = [k(x_1, x), \dots, k(x_{\ell}, x)](\mathbf{G} + \lambda \ell I)^{-1} [y_1; \dots; y_{\ell}],$$
$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^{\ell}.$$

Kernel: application example – ridge regression

- Given: $\{(x_i, y_i)\}_{i=1}^{\ell}$, $H = H(k)$.
- Task: find $f \in H$ s.t. $f(x_i) \approx y_i$,

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} [f(x_i) - y_i]^2 + \lambda \|f\|_H^2 \rightarrow \min_{f \in H} \quad (\lambda > 0).$$

- Analytical solution, $\mathcal{O}(\ell^3)$ – expensive:

$$f(x) = [k(x_1, x), \dots, k(x_{\ell}, x)](\mathbf{G} + \lambda \ell I)^{-1} [y_1; \dots; y_{\ell}],$$
$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^{\ell}.$$

- **Idea:** $\hat{\mathbf{G}}$, matrix-inversion lemma, fast primal solvers \rightarrow RFF.

Kernels: more generally

- Requirement: inner product on the inputs ($k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$).
- Loss function ($\lambda > 0$):

$$J(f) = \sum_{i=1}^{\ell} V(y_i, f(x_i)) + \lambda \|f\|_{H(k)}^2 \rightarrow \min_{f \in H(k)}.$$

Kernels: more generally

- Requirement: inner product on the inputs ($k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$).
- Loss function ($\lambda > 0$):

$$J(f) = \sum_{i=1}^{\ell} V(y_i, f(x_i)) + \lambda \|f\|_{H(k)}^2 \rightarrow \min_{f \in H(k)} .$$

- By the representer theorem [$f(\cdot) = \sum_{i=1}^{\ell} \alpha_i k(\cdot, x_i)$]:

$$J(\alpha) = \sum_{i=1}^{\ell} V(y_i, (\mathbf{G}\alpha)_i) + \lambda \alpha^T \mathbf{G}\alpha \rightarrow \min_{\alpha \in \mathbb{R}^{\ell}} .$$

- $\Rightarrow k(x_i, x_j)$ matters.

Motivation:

- fitting ∞ -D exp. family distributions [Sriperumbudur et al., 2014],
- $k \leftrightarrow$ sufficient statistics,
- rich family,

Kernel derivatives: application example

Motivation:

- fitting ∞ -D exp. family distributions [Sriperumbudur et al., 2014],
- $k \leftrightarrow$ sufficient statistics,
- rich family,
- fitting = linear equation:
 - coefficient matrix: $(d\ell) \times (d\ell)$, $d = \dim(x)$,
 - entries: **kernel values and derivatives.**

- Objective:

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} V(y_i, \{\partial^{\mathbf{p}} f(x_i)\}_{\mathbf{p} \in J_i}) + \lambda \|f\|_{H(k)}^2 \rightarrow \min_{f \in H(k)} .$$

Kernel derivatives: more generally

- Objective:

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} V(y_i, \{\partial^{\mathbf{p}} f(x_i)\}_{\mathbf{p} \in J_i}) + \lambda \|f\|_{H(k)}^2 \rightarrow \min_{f \in H(k)} .$$

- [Zhou, 2008, Shi et al., 2010, Rosasco et al., 2010, Rosasco et al., 2013, Ying et al., 2012]:
 - semi-supervised learning with gradient information,
 - nonlinear variable selection.
- Kernel HMC [Strathmann et al., 2015].

- $\mathcal{X} = \mathbb{R}^d$. k : continuous, shift-invariant [$k(\mathbf{x}, \mathbf{y}) = \tilde{k}(\mathbf{x} - \mathbf{y})$].
- By Bochner's theorem:

$$k(\mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^d} e^{i\boldsymbol{\omega}^T(\mathbf{x}-\mathbf{y})} d\Lambda(\boldsymbol{\omega}) = \int_{\mathbb{R}^d} \cos(\boldsymbol{\omega}^T(\mathbf{x} - \mathbf{y})) d\Lambda(\boldsymbol{\omega}).$$

- $\mathcal{X} = \mathbb{R}^d$. k : continuous, shift-invariant [$k(\mathbf{x}, \mathbf{y}) = \tilde{k}(\mathbf{x} - \mathbf{y})$].
- By Bochner's theorem:

$$k(\mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^d} e^{i\boldsymbol{\omega}^T(\mathbf{x}-\mathbf{y})} d\Lambda(\boldsymbol{\omega}) = \int_{\mathbb{R}^d} \cos(\boldsymbol{\omega}^T(\mathbf{x} - \mathbf{y})) d\Lambda(\boldsymbol{\omega}).$$

- RFF trick [Rahimi and Recht, 2007] (MC): $\boldsymbol{\omega}_{1:m} := (\boldsymbol{\omega}_j)_{j=1}^m \stackrel{i.i.d.}{\sim} \Lambda$,

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\boldsymbol{\omega}_j^T(\mathbf{x} - \mathbf{y})) = \int_{\mathbb{R}^d} \cos(\boldsymbol{\omega}^T(\mathbf{x} - \mathbf{y})) d\Lambda_m(\boldsymbol{\omega}).$$

- Hoeffding inequality + union bound:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_p\left(\underbrace{|\mathcal{S}|}_{\text{linear}} \sqrt{\frac{\log m}{m}}\right).$$

- Hoeffding inequality + union bound:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_p\left(\underbrace{|\mathcal{S}|}_{\text{linear}} \sqrt{\frac{\log m}{m}}\right).$$

- Characteristic function point of view [Csörgő and Totik, 1983] (asymptotic!):
 - 1 $|\mathcal{S}_m| = e^{o(m)}$ is the optimal rate for a.s. convergence,
 - 2 For faster growing $|\mathcal{S}_m|$: even convergence in probability fails.

- ① Finite-sample L^∞ -guarantee $\xrightarrow{\text{specifically}}$

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\frac{\sqrt{\log |\mathcal{S}|}}{\sqrt{m}} \right)$$

$\Rightarrow \mathcal{S}$ can grow **exponentially** [$|\mathcal{S}_m| = e^{o(m)}$] – optimal!

- ① Finite-sample L^∞ -guarantee $\xrightarrow{\text{specifically}}$

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\frac{\sqrt{\log |\mathcal{S}|}}{\sqrt{m}} \right)$$

$\Rightarrow \mathcal{S}$ can grow **exponentially** $[|\mathcal{S}_m| = e^{o(m)}]$ – optimal!

- ② Finite sample L^r guarantees, $r \in [1, \infty)$.

Today: one-page summary

- ① Finite-sample L^∞ -guarantee $\xrightarrow{\text{specifically}}$

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\frac{\sqrt{\log |\mathcal{S}|}}{\sqrt{m}} \right)$$

$\Rightarrow \mathcal{S}$ can grow **exponentially** $[|\mathcal{S}_m| = e^{o(m)}]$ – optimal!

- ② Finite sample L^r guarantees, $r \in [1, \infty)$.
- ③ Derivatives: $\partial^{\mathbf{p}, \mathbf{q}} k$.

..., where

- Uniform ($r = \infty$), L^r ($1 \leq r < \infty$) norm:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} := \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})|,$$

$$\|k - \hat{k}\|_{L^r(\mathcal{S})} := \left(\int_{\mathcal{S}} \int_{\mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})|^r d\mathbf{x} d\mathbf{y} \right)^{\frac{1}{r}}.$$

..., where

- Uniform ($r = \infty$), L^r ($1 \leq r < \infty$) norm:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} := \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})|,$$

$$\|k - \hat{k}\|_{L^r(\mathcal{S})} := \left(\int_{\mathcal{S}} \int_{\mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})|^r d\mathbf{x} d\mathbf{y} \right)^{\frac{1}{r}}.$$

- Kernel derivatives:

$$\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{x}, \mathbf{y}) = \frac{\partial^{|\mathbf{p}|+|\mathbf{q}|} k(\mathbf{x}, \mathbf{y})}{\partial x_1^{p_1} \dots \partial x_d^{p_d} \partial y_1^{q_1} \dots \partial y_d^{q_d}}, \quad |\mathbf{p}| = \sum_{j=1}^d |p_j|.$$

$\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- 1 Empirical process form $[\mathbb{P}g := \int g d\mathbb{P}]$:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{g \in \mathcal{G}} |\Lambda g - \Lambda_m g| = \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

$\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- ① Empirical process form $[\mathbb{P}g := \int g d\mathbb{P}]$:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{g \in \mathcal{G}} |\Lambda g - \Lambda_m g| = \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

- ② $f(\omega_{1:m}) = \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ concentrates (bounded difference):

$$\|\Lambda - \Lambda_m\|_{\mathcal{G}} \lesssim \mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} + \frac{1}{\sqrt{m}}.$$

$\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- ① Empirical process form $[\mathbb{P}g := \int g d\mathbb{P}]$:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{g \in \mathcal{G}} |\Lambda g - \Lambda_m g| = \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

- ② $f(\omega_{1:m}) = \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ concentrates (bounded difference):

$$\|\Lambda - \Lambda_m\|_{\mathcal{G}} \preceq \mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} + \frac{1}{\sqrt{m}}.$$

- ③ \mathcal{G} is 'nice' (uniformly bounded, separable Carathéodory) \Rightarrow

$$\mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} \preceq \mathbb{E}_{\omega_{1:m}} \underbrace{\mathcal{R}(\mathcal{G}, \omega_{1:m})}_{\mathbb{E}_\epsilon \sup_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{j=1}^m \epsilon_j g(\omega_j) \right|}.$$

- Using Dudley's entropy bound:

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \lesssim \frac{1}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), u)} du.$$

- 4 Using Dudley's entropy bound:

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \lesssim \frac{1}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), u)} du.$$

- 5 \mathcal{G} is smoothly parameterized by a compact set \Rightarrow

$$\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), u) \leq \left(\frac{4|\mathcal{S}|A}{u} + 1 \right)^d, \quad A(\omega_{1:m}) = \sqrt{\frac{1}{m} \sum_{j=1}^m \|\omega_j\|_2^2}.$$

- 4 Using Dudley's entropy bound:

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \lesssim \frac{1}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), u)} du.$$

- 5 \mathcal{G} is smoothly parameterized by a compact set \Rightarrow

$$\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), u) \leq \left(\frac{4|\mathcal{S}|A}{u} + 1 \right)^d, \quad A(\omega_{1:m}) = \sqrt{\frac{1}{m} \sum_{j=1}^m \|\omega_j\|_2^2}.$$

- 6 Putting together [$|\mathcal{G}|_{L^2(\Lambda_m)} \leq 2$, Jensen inequality] we get ...

Let k be continuous, $\sigma^2 := \int \|\omega\|^2 d\Lambda(\omega) < \infty$. Then for $\forall \tau > 0$ and compact set $\mathcal{S} \subset \mathbb{R}^d$

$$\Lambda^m \left(\|\hat{k} - k\|_{L^\infty(\mathcal{S})} \geq \frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau},$$

$$h(d, |\mathcal{S}|, \sigma) := 32\sqrt{2d \log(2|\mathcal{S}| + 1)} + 16\sqrt{\frac{2d}{\log(2|\mathcal{S}| + 1)}} + 32\sqrt{2d \log(\sigma + 1)}.$$

Consequence-1 (Borel-Cantelli lemma)

- A.s. convergence on compact sets: $\hat{k} \xrightarrow{m \rightarrow \infty} k$ at rate $\sqrt{\frac{\log |\mathcal{S}|}{m}}$.

Consequence-1 (Borel-Cantelli lemma)

- A.s. convergence on compact sets: $\hat{k} \xrightarrow{m \rightarrow \infty} k$ at rate $\sqrt{\frac{\log |\mathcal{S}|}{m}}$.
- Growing diameter:
 - $\frac{\log |\mathcal{S}_m|}{m} \xrightarrow{m \rightarrow \infty} 0$ is enough, i.e. $|\mathcal{S}_m| = e^{o(m)}$.

Consequence-1 (Borel-Cantelli lemma)

- A.s. convergence on compact sets: $\hat{k} \xrightarrow{m \rightarrow \infty} k$ at rate $\sqrt{\frac{\log |\mathcal{S}|}{m}}$.
- Growing diameter:
 - $\frac{\log |\mathcal{S}_m|}{m} \xrightarrow{m \rightarrow \infty} 0$ is enough, i.e. $|\mathcal{S}_m| = e^{o(m)}$.
- Specifically:
 - *asymptotic* optimality [Csörgő and Totik, 1983, Theorem 2] (if $k(\mathbf{z})$ vanishes at ∞).

Consequence-2: L^r result for k ($1 \leq r$)

Idea:

- Note that

$$\begin{aligned}\|\hat{k} - k\|_{L^r(\mathcal{S})} &= \left(\int_{\mathcal{S}} \int_{\mathcal{S}} |\hat{k}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y})|^r \, d\mathbf{x} \, d\mathbf{y} \right)^{\frac{1}{r}} \\ &\leq \|\hat{k} - k\|_{L^\infty(\mathcal{S})} \text{vol}^{2/r}(\mathcal{S}).\end{aligned}$$

Consequence-2: L^r result for k ($1 \leq r$)

Idea:

- Note that

$$\begin{aligned}\|\hat{k} - k\|_{L^r(\mathcal{S})} &= \left(\int_{\mathcal{S}} \int_{\mathcal{S}} |\hat{k}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y})|^r \, d\mathbf{x} \, d\mathbf{y} \right)^{\frac{1}{r}} \\ &\leq \|\hat{k} - k\|_{L^\infty(\mathcal{S})} \text{vol}^{2/r}(\mathcal{S}).\end{aligned}$$

- $\text{vol}(\mathcal{S}) \leq \text{vol}(B)$, where $B := \left\{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 \leq \frac{|\mathcal{S}|}{2} \right\}$,

Consequence-2: L^r result for k ($1 \leq r$)

Idea:

- Note that

$$\begin{aligned}\|\hat{k} - k\|_{L^r(\mathcal{S})} &= \left(\int_{\mathcal{S}} \int_{\mathcal{S}} |\hat{k}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y})|^r \, d\mathbf{x} \, d\mathbf{y} \right)^{\frac{1}{r}} \\ &\leq \|\hat{k} - k\|_{L^\infty(\mathcal{S})} \text{vol}^{2/r}(\mathcal{S}).\end{aligned}$$

- $\text{vol}(\mathcal{S}) \leq \text{vol}(B)$, where $B := \left\{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 \leq \frac{|\mathcal{S}|}{2} \right\}$,
- $\text{vol}(B) = \frac{\pi^{d/2} |\mathcal{S}|^d}{2^d \Gamma(\frac{d}{2} + 1)}$, $\Gamma(t) = \int_0^\infty u^{t-1} e^{-u} \, du. \Rightarrow$

Under the previous assumptions, and $1 \leq r < \infty$:

$$\Lambda^m \left(\|\hat{k} - k\|_{L^r(\mathcal{S})} \geq \left(\frac{\pi^{d/2} |\mathcal{S}|^d}{2^d \Gamma(\frac{d}{2} + 1)} \right)^{2/r} \frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau}.$$

Under the previous assumptions, and $1 \leq r < \infty$:

$$\Lambda^m \left(\|\hat{k} - k\|_{L^r(\mathcal{S})} \geq \left(\frac{\pi^{d/2} |\mathcal{S}|^d}{2^d \Gamma(\frac{d}{2} + 1)} \right)^{2/r} \frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau}.$$

Hence,

$$\|\hat{k} - k\|_{L^r(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\underbrace{\frac{|\mathcal{S}|^{2d/r} \sqrt{\log |\mathcal{S}|}}{\sqrt{m}}}_{L^r(\mathcal{S})\text{-consistency if } \frac{m \rightarrow \infty}{\rightarrow} 0} \right).$$

L^r result for k

Under the previous assumptions, and $1 \leq r < \infty$:

$$\Lambda^m \left(\|\hat{k} - k\|_{L^r(\mathcal{S})} \geq \left(\frac{\pi^{d/2} |\mathcal{S}|^d}{2^d \Gamma(\frac{d}{2} + 1)} \right)^{2/r} \frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau}.$$

Hence,

$$\|\hat{k} - k\|_{L^r(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\underbrace{\frac{|\mathcal{S}|^{2d/r} \sqrt{\log |\mathcal{S}|}}{\sqrt{m}}}_{L^r(\mathcal{S})\text{-consistency if } \frac{m \rightarrow \infty}{\rightarrow} 0} \right).$$

Uniform guarantee: $|\mathcal{S}_m| = e^{m^{\delta < 1}}$; now: $\frac{|\mathcal{S}_m|^{2d/r}}{\sqrt{m}} \rightarrow 0 \Rightarrow |\mathcal{S}_m| = \tilde{o}(m^{\frac{r}{4d}})$.

Direct L^r result for k (proof idea after discussion)

Under the previous assumptions, and $1 < r < \infty$:

$$\Lambda^m \left(\|\hat{k} - k\|_{L^r(\mathcal{S})} \geq \left(\frac{\pi^{d/2} |\mathcal{S}|^d}{2^d \Gamma(\frac{d}{2} + 1)} \right)^{2/r} \left(\frac{C'_r}{m^{1 - \max\{\frac{1}{2}, \frac{1}{r}\}}} + \frac{\sqrt{2\tau}}{\sqrt{m}} \right) \right) \leq e^{-\tau},$$

$C'_r = \mathcal{O}(\sqrt{r})$: universal constant; only r -dependent (not $|\mathcal{S}|$ or m -dep.).

Direct L^r result for k (proof idea after discussion)

Under the previous assumptions, and $1 < r < \infty$:

$$\Lambda^m \left(\|\hat{k} - k\|_{L^r(\mathcal{S})} \geq \left(\frac{\pi^{d/2} |\mathcal{S}|^d}{2^d \Gamma(\frac{d}{2} + 1)} \right)^{2/r} \left(\frac{C'_r}{m^{1 - \max\{\frac{1}{2}, \frac{1}{r}\}}} + \frac{\sqrt{2\tau}}{\sqrt{m}} \right) \right) \leq e^{-\tau},$$

$C'_r = \mathcal{O}(\sqrt{r})$: universal constant; only r -dependent (not $|\mathcal{S}|$ or m -dep.).

Note: if $2 \leq r$, then

- 1 $m^{1 - \max\{\frac{1}{2}, \frac{1}{r}\}} = \sqrt{m}$,
- 2 $\|\hat{k} - k\|_{L^r(\mathcal{S}_m)} \xrightarrow{a.s.} 0$ if $|\mathcal{S}_m| = o\left(m^{\frac{r}{4d}}\right)$ as $m \rightarrow \infty$.
- 3 In short, we got rid of $\sqrt{\log(|\mathcal{S}|)}$: $\tilde{o} \rightarrow o$.

Direct L^r result for k : proof idea

① $f(\omega_1, \dots, \omega_m) = \|k - \hat{k}\|_{L^r(\mathcal{S})}$ concentrates (bounded difference):

$$\|k - \hat{k}\|_{L^r(\mathcal{S})} \leq \mathbb{E}_{\omega_{1:m}} \|k - \hat{k}\|_{L^r(\mathcal{S})} + \text{vol}^{2/r}(\mathcal{S}) \sqrt{\frac{2\tau}{m}}.$$

Direct L^r result for k : proof idea

- ① $f(\omega_1, \dots, \omega_m) = \|k - \hat{k}\|_{L^r(\mathcal{S})}$ concentrates (bounded difference):

$$\|k - \hat{k}\|_{L^r(\mathcal{S})} \leq \mathbb{E}_{\omega_{1:m}} \|k - \hat{k}\|_{L^r(\mathcal{S})} + \text{vol}^{2/r}(\mathcal{S}) \sqrt{\frac{2\tau}{m}}.$$

- ② By $L^r \cong (L^{r'})^*$ ($\frac{1}{r} + \frac{1}{r'} = 1$), the separability of $L^{r'}(\mathcal{S})$ ($r > 1$) and symmetrization:

$$\mathbb{E}_{\omega_{1:m}} \|k - \hat{k}\|_{L^r(\mathcal{S})} \leq \frac{2}{m} \mathbb{E}_{\omega_{1:m}} \underbrace{\mathbb{E}_{\varepsilon} \left\| \sum_{i=1}^m \varepsilon_i \cos(\langle \omega_i, \cdot - \cdot \rangle) \right\|_{L^r(\mathcal{S})}}_{=:(*)}.$$

Direct L^r result for k : proof idea

- ① $f(\omega_1, \dots, \omega_m) = \|k - \hat{k}\|_{L^r(\mathcal{S})}$ concentrates (bounded difference):

$$\|k - \hat{k}\|_{L^r(\mathcal{S})} \leq \mathbb{E}_{\omega_{1:m}} \|k - \hat{k}\|_{L^r(\mathcal{S})} + \text{vol}^{2/r}(\mathcal{S}) \sqrt{\frac{2\tau}{m}}.$$

- ② By $L^r \cong (L^{r'})^*$ ($\frac{1}{r} + \frac{1}{r'} = 1$), the separability of $L^{r'}(\mathcal{S})$ ($r > 1$) and symmetrization:

$$\mathbb{E}_{\omega_{1:m}} \|k - \hat{k}\|_{L^r(\mathcal{S})} \leq \frac{2}{m} \mathbb{E}_{\omega_{1:m}} \underbrace{\mathbb{E}_{\varepsilon} \left\| \sum_{i=1}^m \varepsilon_i \cos(\langle \omega_i, \cdot - \cdot \rangle) \right\|_{L^r(\mathcal{S})}}_{=:(*)}.$$

- ③ Since $L^r(\mathcal{S})$ is of type $\min(2, r)$ [' \diamond -rule'] $\exists C'_r$ such that

$$(*) \leq C'_r \left(\sum_{i=1}^m \left\| \cos(\langle \omega_i, \cdot - \cdot \rangle) \right\|_{L^r(\mathcal{S})}^{\min(2, r)} \right)^{\frac{1}{\min(2, r)}}.$$

Goal: $\widehat{k^{\mathbf{p}, \mathbf{q}}}$. If

① $\text{supp}(\Lambda)$ is bounded:

- $C_{k, \mathbf{p}, \mathbf{q}} := \mathbb{E}_{\omega \sim \Lambda} \left[|\omega^{\mathbf{p} + \mathbf{q}}| \|\omega\|_2^2 \right] < \infty$: $L^\infty, L^r \checkmark$, but
- Gaussian, Laplacian, inverse multiquadratic, Matern:(
- c_0 universality $\Leftrightarrow \text{supp}(\Lambda) = \mathbb{R}^d$, if $k(\mathbf{z}) \in C_0(\mathbb{R}^d)$.

Goal: $\widehat{k^{\mathbf{p}, \mathbf{q}}}$. If

① $\text{supp}(\Lambda)$ is bounded:

- $C_{k, \mathbf{p}, \mathbf{q}} := \mathbb{E}_{\boldsymbol{\omega} \sim \Lambda} \left[|\boldsymbol{\omega}^{\mathbf{p} + \mathbf{q}}| \|\boldsymbol{\omega}\|_2^2 \right] < \infty$: $L^\infty, L^r \checkmark$, but
- Gaussian, Laplacian, inverse multiquadratic, Matern:(
- c_0 universality $\Leftrightarrow \text{supp}(\Lambda) = \mathbb{R}^d$, if $k(\mathbf{z}) \in C_0(\mathbb{R}^d)$.

② $\text{supp}(\Lambda)$ is unbounded:

- \mathcal{G} : becomes unbounded.
- [Rahimi and Recht, 2007]: 'Hoeffding \rightarrow Bernstein', but

Kernel derivatives: unbounded $\text{supp}(\Lambda)$

Assumptions [$h_a = \cos^{(a)}$, $\mathcal{S}_\Delta = \mathcal{S} - \mathcal{S}$]:

- 1 $\mathbf{z} \mapsto \nabla_{\mathbf{z}} [\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{z})]$: continuous; $\mathcal{S} \subset \mathbb{R}^d$: compact,
 $E_{\mathbf{p}, \mathbf{q}} := \mathbb{E}_{\boldsymbol{\omega} \sim \Lambda} |\boldsymbol{\omega}^{\mathbf{p} + \mathbf{q}}| \|\boldsymbol{\omega}\|_2 < \infty$.
- 2 $\exists L > 0, \sigma > 0$

$$\mathbb{E}_{\boldsymbol{\omega} \sim \Lambda} |f(\mathbf{z}; \boldsymbol{\omega})|^M \leq \frac{M! \sigma^2 L^{M-2}}{2} \quad (\forall M \geq 2, \forall \mathbf{z} \in \mathcal{S}_\Delta),$$
$$f(\mathbf{z}; \boldsymbol{\omega}) = \partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{z}) - \boldsymbol{\omega}^{\mathbf{p}} (-\boldsymbol{\omega})^{\mathbf{q}} h_{|\mathbf{p} + \mathbf{q}|}(\boldsymbol{\omega}^T \mathbf{z}).$$

Kernel derivatives: unbounded $\text{supp}(\Lambda)$

Then with $F_d := d^{-\frac{d}{d+1}} + d^{\frac{1}{d+1}}$

$$\begin{aligned} \Lambda^m \left(\|\partial^{\mathbf{p}, \mathbf{q}} k - \widehat{\partial^{\mathbf{p}, \mathbf{q}} k}\|_{L^\infty(\mathcal{S})} \geq \epsilon \right) &\leq \\ &\leq 2^{d-1} e^{-\frac{m\epsilon^2}{8\sigma^2 \left(1 + \frac{\epsilon L}{2\sigma^2}\right)}} + F_d 2^{\frac{4d-1}{d+1}} \left[\frac{|\mathcal{S}| (D_{\mathbf{p}, \mathbf{q}, \mathcal{S}} + E_{\mathbf{p}, \mathbf{q}})}{\epsilon} \right]^{\frac{d}{d+1}} e^{-\frac{m\epsilon^2}{8(d+1)\sigma^2 \left(1 + \frac{\epsilon L}{2\sigma^2}\right)}}, \end{aligned}$$

where $D_{\mathbf{p}, \mathbf{q}, \mathcal{S}} := \sup_{\mathbf{z} \in \text{conv}(\mathcal{S}_\Delta)} \|\nabla_{\mathbf{z}} [\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{z})]\|_2$.

Finite sample

- $L^\infty(\mathcal{S})$ guarantees $\xrightarrow{\text{spec.}} |\mathcal{S}_m| = e^{o(m)}$ – asymp. optimal!
- $L^r(\mathcal{S})$ results (\Leftarrow uniform, type of L^r).

Finite sample

- $L^\infty(\mathcal{S})$ guarantees $\xrightarrow{\text{spec.}} |\mathcal{S}_m| = e^{o(m)}$ – asymp. optimal!
- $L^r(\mathcal{S})$ results (\Leftarrow uniform, type of L^r).
- derivative approximation guarantees:
 - bounded spectral support: ✓
 - unbounded spectral support: trickier – to be continued;)

Thank you for the attention!



Acknowledgments: This work was supported by the Gatsby Charitable Foundation.



Csörgő, S. and Totik, V. (1983).

On how long interval is the empirical characteristic function uniformly consistent?

Acta Scientiarum Mathematicarum, 45:141–149.



Rahimi, A. and Recht, B. (2007).

Random features for large-scale kernel machines.

In *Neural Information Processing Systems (NIPS)*, pages 1177–1184.



Rosasco, L., Santoro, M., Mosci, S., Verri, A., and Villa, S. (2010).

A regularization approach to nonlinear variable selection.

JMLR W&CP – International Conference on Artificial Intelligence and Statistics (AISTATS), 9:653–660.



Rosasco, L., Villa, S., Mosci, S., Santoro, M., and Verri, A. (2013).

Nonparametric sparsity and regularization.

Journal of Machine Learning Research, 14:1665–1714.

-  Shi, L., Guo, X., and Zhou, D.-X. (2010).
Hermite learning with gradient data.
Journal of Computational and Applied Mathematics,
233:3046–3059.
-  Sriperumbudur, B. K., Fukumizu, K., Gretton, A., Hyvärinen, A., and Kumar, R. (2014).
Density estimation in infinite dimensional exponential families.
Technical report.
<http://arxiv.org/pdf/1312.3516.pdf>.
-  Strathmann, H., Sejdinovic, D., Livingstone, S., Szabó, Z., and Gretton, A. (2015).
Gradient-free Hamiltonian Monte Carlo with efficient kernel exponential families.
In Neural Information Processing Systems (NIPS).
-  Ying, Y., Wu, Q., and Campbell, C. (2012).
Learning the coordinate gradients.
Advances in Computational Mathematics, 37:355–378.



Zhou, D.-X. (2008).

Derivative reproducing properties for kernel methods in learning theory.

Journal of Computational and Applied Mathematics,
220:456–463.