

# Optimal Uniform and $L^p$ Rates for Random Fourier Features\*

Bharath K. Sriperumbudur<sup>†,1</sup>

<sup>1</sup>Department of Statistics, Pennsylvania State University

Zoltán Szabó<sup>†,2</sup>

<sup>2</sup>Gatsby Unit, University College London

## Brief Summary

- **Kernel methods** [1]:  
▪ Pro: flexible modelling toolkit.  
▪ Contra: computationally intensive, poor scalability.
- **Randomized algorithms**:  
▪ Low-D feature representation → fast linear methods.  
▪ Random Fourier features (RFF) [2]:  
▪ simple, popular, practically efficient, but theoretically not well-understood.
- **Contribution**: detailed theoretical analysis,  
▪  $L^\infty$ -optimal performance guarantees (RFF dimension, growing set size),  
▪  $L^r$  ( $1 \leq r < \infty$ )-guarantees,  
▪ RFF approximation for kernel derivatives + analysis.

## RFF Idea (Kernel Approximation)

- $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ : continuous, bounded, translation-invariant kernel.  
Bochner's theorem ⇒  

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

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T(\mathbf{x}-\mathbf{y})).$$

Here:  $(\omega_j)_{j=1}^m \stackrel{i.i.d.}{\sim} \Lambda$ ,  $\hat{k}(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathbb{R}^{2m}}$  with

$$\phi(\mathbf{x}) = \frac{1}{\sqrt{m}} [\cos(\omega_1^T \mathbf{x}); \dots; \cos(\omega_m^T \mathbf{x}); \sin(\omega_1^T \mathbf{x}); \dots; \sin(\omega_m^T \mathbf{x})].$$

## Existing RFF Guarantees

- [2]:  $\hat{k}$  is consistent (compact convergence),  

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S} \times \mathcal{S})} := \sup_{(x,y) \in \mathcal{S} \times \mathcal{S}} |k(x, y) - \hat{k}(x, y)| = \mathcal{O}_p(|\mathcal{S}| \sqrt{m^{-1} \log m}).$$
- [3]: 3 RFF variants, better constants.

## Theorem-1: $k$ approximation, $L^\infty(\mathcal{S} \times \mathcal{S})$

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

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

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

### Remark-1:

- A.s. convergence on compact sets:  $\hat{k} \xrightarrow{m \rightarrow \infty} k$  at rate  $\sqrt{\frac{\log |\mathcal{S}|}{m}}$  (B-C. lemma).
- Growing diameter ( $\mathcal{S}_m$ ):  
▪  $\frac{\log |\mathcal{S}_m|}{m} \xrightarrow{m \rightarrow \infty} 0$  is enough (i.e.,  $|\mathcal{S}_m| = e^{o(m)}$ ) ⇔ old result:  $|\mathcal{S}_m| = o(\sqrt{m/\log(m)})$ .
- Specifically: asymptotic optimality [4, Theorem 2].

## Theorem-2: $k$ approximation, $L^r(\mathcal{S} \times \mathcal{S})$ , $1 \leq r < \infty$

For any  $\tau > 0$ , compact  $\mathcal{S} \subset \mathbb{R}^d$

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

### Remark-2:

- Consequence of Theorem-1.
- $L^r(\mathcal{S} \times \mathcal{S})$ -consistency:  $\|k - \hat{k}\|_{L^r(\mathcal{S} \times \mathcal{S})} = \mathcal{O}_{a.s.} \left( \underbrace{m^{-1/2} |\mathcal{S}|^{2d/r} \sqrt{\log |\mathcal{S}|}}_{=: (*)_1; \text{ if } \xrightarrow{m \rightarrow \infty} 0} \right)$ , i.e.
- Growing diameter:  $(*)_1 \rightarrow 0 \Rightarrow |\mathcal{S}_m| = \tilde{o}(m^{\frac{r}{2d}})$ ;  $L^\infty$ -case:  $|\mathcal{S}_m| = e^{m^{\delta < 1}}$ .

## Theorem-3: $k$ approximation, $L^r(\mathcal{S} \times \mathcal{S})$ , $1 < r < \infty$

Applying a direct reasoning: for any  $\tau > 0$ , compact  $\mathcal{S} \subset \mathbb{R}^d$

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

### Remark-3:

- $C_r = \mathcal{O}(\sqrt{r})$ , universal constant.
- $L^r(\mathcal{S} \times \mathcal{S})$ -consistency: if  $2 \leq r$ , then  $\|k - \hat{k}\|_{L^r(\mathcal{S} \times \mathcal{S})} = \mathcal{O}_{a.s.} \left( \underbrace{m^{-1/2} |\mathcal{S}|^{2d/r}}_{\text{if } \xrightarrow{m \rightarrow \infty} 0} \right)$ ,  
the  $\sqrt{\log(\mathcal{S})}$  term disappeared ( $\tilde{o} \rightarrow o$ , see Remark-2).

## Kernel Derivative Approximation

$$\begin{aligned} \partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{x}, \mathbf{y}) &= \int_{\mathbb{R}^d} \omega^{\mathbf{p}}(-\omega)^{\mathbf{q}} h_{|\mathbf{p}+\mathbf{q}|}(\omega^T(\mathbf{x}-\mathbf{y})) d\Lambda(\omega), \quad h_n = \cos^{(n)}, n \in \mathbb{N} \\ \widehat{\partial^{\mathbf{p}, \mathbf{q}} k}(\mathbf{x}, \mathbf{y}) &= \frac{1}{m} \sum_{j=1}^m \omega_j^{\mathbf{p}}(-\omega_j)^{\mathbf{q}} h_{|\mathbf{p}+\mathbf{q}|}(\omega_j^T(\mathbf{x}-\mathbf{y})) = \langle \phi^{\mathbf{p}}(\mathbf{x}), \phi^{\mathbf{q}}(\mathbf{y}) \rangle_{\mathbb{R}^{2m}}. \end{aligned}$$

## Th.-4: $\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{x}, \mathbf{y})$ appr., $\text{supp}(\Lambda)$ : bounded, $L^\infty(\mathcal{S} \times \mathcal{S})$

Let  $\mathbf{p}, \mathbf{q} \in \mathbb{N}^d$ ,  $T_{\mathbf{p}, \mathbf{q}} := \sup_{\omega \in \text{supp}(\Lambda)} |\omega^{\mathbf{p}+\mathbf{q}}|$ ,  $C_{\mathbf{p}, \mathbf{q}} := \mathbb{E}_{\omega \sim \Lambda} [|\omega^{\mathbf{p}+\mathbf{q}}| \|\omega\|_2^2]$ . Assume:  $C_{2\mathbf{p}, 2\mathbf{q}} < \infty$ ;  $\text{supp}(\Lambda)$  is bounded if  $[\mathbf{p}; \mathbf{q}] \neq \mathbf{0}$ . Then for  $\forall \tau > 0$ , compact  $\mathcal{S} \subset \mathbb{R}^d$

$$\Lambda^m \left( \|\partial^{\mathbf{p}, \mathbf{q}} k - \widehat{\partial^{\mathbf{p}, \mathbf{q}} k}(\mathbf{x}, \mathbf{y})\|_{L^\infty(\mathcal{S} \times \mathcal{S})} \geq \frac{H(d, \mathbf{p}, \mathbf{q}, |\mathcal{S}|) + T_{\mathbf{p}, \mathbf{q}} \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau},$$

where

$$\begin{aligned} H(d, \mathbf{p}, \mathbf{q}, |\mathcal{S}|) &= \left[ \sqrt{U(\mathbf{p}, \mathbf{q}, |\mathcal{S}|)} + \frac{1}{2\sqrt{U(\mathbf{p}, \mathbf{q}, |\mathcal{S}|)}} + \sqrt{\log(\sqrt{C_{2\mathbf{p}, 2\mathbf{q}}} + 1)} \right], \\ U(\mathbf{p}, \mathbf{q}, |\mathcal{S}|) &= \log \left( \frac{2|\mathcal{S}|}{\sqrt{T_{\mathbf{p}, \mathbf{q}}}} + 1 \right). \end{aligned}$$

### Remark-4:

- Theorem-4  $\xrightarrow{\text{spec. } \mathbf{p}=\mathbf{q}=0}$  Theorem-1,  $T_{\mathbf{p}, \mathbf{q}} = T_{2\mathbf{p}, 2\mathbf{q}} = 1$ . Else:  $\text{supp}(\Lambda)$ : bounded ⇒  $T_{\mathbf{p}, \mathbf{q}} < \infty$  and  $T_{2\mathbf{p}, 2\mathbf{q}} < \infty$ .
- Growth of  $|\mathcal{S}_m|$ : A la Remarks 1-2  
▪  $\|\partial^{\mathbf{p}, \mathbf{q}} k - \widehat{\partial^{\mathbf{p}, \mathbf{q}} k}(\mathbf{x}, \mathbf{y})\|_{L^\infty(\mathcal{S}_m \times \mathcal{S}_m)} \xrightarrow{a.s.} 0$  if  $|\mathcal{S}_m| = e^{o(m)}$ .  
▪  $\|\partial^{\mathbf{p}, \mathbf{q}} k - \widehat{\partial^{\mathbf{p}, \mathbf{q}} k}(\mathbf{x}, \mathbf{y})\|_{L^r(\mathcal{S}_m \times \mathcal{S}_m)} \xrightarrow{a.s.} 0$  if  $m^{-1/2} |\mathcal{S}_m|^{2d/r} \sqrt{\log |\mathcal{S}_m|} \xrightarrow{m \rightarrow \infty} 0$  ( $1 \leq r < \infty$ ).

## Theorem-5: $\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{x}, \mathbf{y})$ approximation, $\text{supp}(\Lambda)$ : unbounded, $L^\infty(\mathcal{S} \times \mathcal{S})$

Assume: (i)  $\mathbf{z} \mapsto \nabla_{\mathbf{z}} [\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{z})]$ : continuous; (ii)  $\mathcal{S} \subset \mathbb{R}^d$ : compact, (iii)  $E_{\mathbf{p}, \mathbf{q}} := \mathbb{E}_{\omega \sim \Lambda} |\omega^{\mathbf{p}+\mathbf{q}}| \|\omega\|_2 < \infty$ , (iv)  $\exists L > 0, \sigma > 0$  such that with  $\mathcal{S}_\Delta := \mathcal{S} - \mathcal{S}$

$$\mathbb{E}_{\omega \sim \Lambda} |f(\mathbf{z}; \omega)|^M \stackrel{(*)_2}{\leq} \frac{M! \sigma^2 L^{M-2}}{2} \quad (\forall M \geq 2, \forall \mathbf{z} \in \mathcal{S}_\Delta), \\ f(\mathbf{z}; \omega) = \partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{z}) - \omega^{\mathbf{p}}(-\omega)^{\mathbf{q}} h_{|\mathbf{p}+\mathbf{q}|}(\omega^T \mathbf{z}).$$

Then with  $F_d := d^{-\frac{d}{d+1}} + d^{\frac{1}{d+1}}$ ,  $D_{\mathbf{p}, \mathbf{q}, \mathcal{S}} := \sup_{\mathbf{z} \in \text{conv}(\mathcal{S}_\Delta)} \|\nabla_{\mathbf{z}} [\partial^{\mathbf{p}, \mathbf{q}} k(\mathbf{z})]\|_2$

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

### Remark-5:

- $F_d$ : monotonically decreasing in  $d$ ,  $F_1 = 2$ .
- $(*)_2$  holds if  $|f(\mathbf{z}; \omega)| \leq \frac{L}{2}$  and  $\mathbb{E}_{\omega \sim \Lambda} [|f(\mathbf{z}; \omega)|^2] \leq \sigma^2$  ( $\forall \mathbf{z} \in \mathcal{S}_\Delta$ ).
- $\|\partial^{\mathbf{p}, \mathbf{q}} k - \widehat{\partial^{\mathbf{p}, \mathbf{q}} k}(\mathbf{x}, \mathbf{y})\|_{L^\infty(\mathcal{S} \times \mathcal{S})} = \mathcal{O}_{a.s.} \left( |\mathcal{S}| \sqrt{\frac{\log m}{m}} \right)$ :  
▪ slightly worse than Theorem-4, but it handles unbounded functions.

## Future Research Directions

- Kernel derivatives: tighter guarantees, (ii) prediction using kernel (derivative) estimates, (iii) analysis of smart RFF approximations [5].

## References

- [1] Bernhard Schölkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- [2] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Neural Information Processing Systems (NIPS)*, pages 1177–1184, 2007.
- [3] Dougal J. Sutherland and Jeff Schneider. On the error of random Fourier features. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 862–871, 2015.
- [4] Sándor Csörgő and Vilmos Totik. On how long interval is the empirical characteristic function uniformly consistent? *Acta Scientiarum Mathematicarum*, 45:141–149, 1983.
- [5] Quoc Le, Tamás Sarlós, and Alexander Smola. Fastfood - computing Hilbert space expansions in loglinear time. *JMLR W&CP - International Conference on Machine Learning (ICML)*, 28:244–252, 2013.