

# Performance Guarantees for Random Fourier Features – Limitations and Merits

Zoltán Szabó

Joint work with Bharath K. Sriperumbudur (PSU)

ML@SITraN, University of Sheffield  
June 25, 2015

- Given:

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

- $\hat{k}(\mathbf{x}, \mathbf{y})$ : Monte-Carlo estimator of  $k(\mathbf{x}, \mathbf{y})$  using  $(\boldsymbol{\omega}_j)_{j=1}^m \stackrel{i.i.d.}{\sim} \Lambda$  [Rahimi and Recht, 2007].
- Motivation:
  - Primal form – fast linear solvers.
  - Kernel function approximation: out-of-sample extension.
  - Online applications.

- Uniform ( $r = \infty$ ):

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

- $L^r$  ( $1 \leq r < \infty$ ):

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

# Approximation of kernel derivatives

- One could also consider  $\partial^{\mathbf{p}, \mathbf{q}} k$ .
- Motivation [Zhou, 2008, Shi et al., 2010, Rosasco et al., 2010, Rosasco et al., 2013, Ying et al., 2012, Sriperumbudur et al., 2014]:
  - semi-supervised learning with gradient information,
  - nonlinear variable selection,
  - fitting of infD exp. family distributions.
- Many of the presented results hold for derivatives ( $[\mathbf{p}; \mathbf{q}] \neq \mathbf{0}$ ).

- Large deviation inequalities

$$\Lambda^m \left( \left\| k - \hat{k} \right\|_{\mathcal{S}} \leq \epsilon \right) \geq f_1(\epsilon, d, m, |\mathcal{S}|),$$

$$\Lambda^m \left( \left\| k - \hat{k} \right\|_{L^r} \leq \epsilon \right) \geq f_2(\epsilon, d, m, |\mathcal{S}|).$$

- Scaling of  $|\mathcal{S}|$  and  $m$  ensuring a.s. convergence?

Notations:  $X_n = \mathcal{O}_p(r_n)$  ( $\mathcal{O}_{a.s.}(r_n)$ ) denotes  $\frac{X_n}{r_n}$  boundedness in probability (almost surely).

- [Rahimi and Recht, 2007]:

$$\left\| \hat{k}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y}) \right\|_{\mathcal{S}} = \mathcal{O}_p \left( |\mathcal{S}| \sqrt{\frac{\log m}{m}} \right).$$

- [Sutherland and Schneider, 2015]: better constants.

- Uniform guarantee (empirical process theory),
- Two  $L^r$  guarantees (uniform consequence, direct).
- Kernel derivatives.

- 1 Empirical process form:

$$\left\| k - \hat{k} \right\|_{\mathcal{S}} = \sup_{g \in \mathcal{G}} |\Lambda g - \Lambda_m g| = \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

- 2  $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$  concentrates by its bounded difference property:

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

- 3  $\mathcal{G}$  is a uniformly bounded, separable Carathéodory family  $\Rightarrow$

$$\mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} \preceq \mathbb{E}_{\omega_{1:m}} \mathcal{R}(\mathcal{G}, \omega_{1:m}).$$



- 4 Using Dudley's entropy integral:

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

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

$$\begin{aligned} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} &\leq \sqrt{\log \left[ \frac{C(\omega_{1:m})}{r} + 1 \right]} \Rightarrow \\ \mathbb{E}_{\omega_{1:m}} \mathcal{R}(\mathcal{G}, \omega_{1:m}) &\lesssim \frac{1}{\sqrt{m}}. \end{aligned}$$

- 6 Putting together:

$$\|k - \hat{k}\|_S \lesssim \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{m}} = \mathcal{O} \left( \sqrt{\frac{\log |\mathcal{S}|}{m}} \right).$$

## Step-1: empirical process form

- Notation:  $\Lambda g = \int g(\omega) d\Lambda(\omega)$ ,  $\Lambda_m g = \int g(\omega) d\Lambda_m(\omega) = \frac{1}{m} \sum_{j=1}^m g(\omega_j)$ .

## Step-1: empirical process form

- Notation:  $\Lambda g = \int g(\boldsymbol{\omega}) d\Lambda(\boldsymbol{\omega})$ ,  $\Lambda_m g = \int g(\boldsymbol{\omega}) d\Lambda_m(\boldsymbol{\omega}) = \frac{1}{m} \sum_{j=1}^m g(\boldsymbol{\omega}_j)$ .
- Reformulation of the objective:

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

where

$$\begin{aligned}\mathcal{G} &= \{g_{\mathbf{z}} : \mathbf{z} \in \mathcal{S}_{\Delta}\}, \\ \mathcal{S}_{\Delta} &= \mathcal{S} - \mathcal{S} = \{\mathbf{x} - \mathbf{y} : \mathbf{x}, \mathbf{y} \in \mathcal{S}\}, \\ g_{\mathbf{z}} : \boldsymbol{\omega} &\mapsto \cos(\boldsymbol{\omega}^T \mathbf{z}).\end{aligned}$$

## Step-2: bounded difference property of $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$

**McDiarmid inequality:** Let  $\omega_1, \dots, \omega_m \in D$  be independent r.v.-s, and  $f : D^m \rightarrow \mathbb{R}$  satisfy the bounded diff. property ( $\forall r$ ):

$$\sup_{u_1, \dots, u_m, u'_r \in D} |f(u_1, \dots, u_m) - f(u_1, \dots, u_{r-1}, u'_r, u_{r+1}, \dots, u_m)| \leq c_r.$$

Then for  $\forall \beta > 0$

$$\mathbb{P}(f(\omega_1, \dots, \omega_m) - \mathbb{E}[f(\omega_1, \dots, \omega_m)] \geq \beta) \leq e^{-\frac{2\beta^2}{\sum_{r=1}^m c_r^2}}.$$

## Step-2: bounded difference property of $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$

Our choice:  $f(\omega_1, \dots, \omega_m) := \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ .

$$\begin{aligned} & |f(\omega_1, \dots, \omega_{r-1}, \omega_r, \omega_{r+1}, \dots, \omega_m) - f(\omega_1, \dots, \omega_{r-1}, \omega'_r, \omega_{r+1}, \dots, \omega_m)| = \\ & = \left| \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) \right| - \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) + \frac{1}{m} [g(\omega_r) - g(\omega'_r)] \right| \right| \end{aligned}$$

## Step-2: bounded difference property of $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$

Our choice:  $f(\omega_1, \dots, \omega_m) := \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ .

$$\begin{aligned} & |f(\omega_1, \dots, \omega_{r-1}, \omega_r, \omega_{r+1}, \dots, \omega_m) - f(\omega_1, \dots, \omega_{r-1}, \omega'_r, \omega_{r+1}, \dots, \omega_m)| = \\ & = \left| \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1} g(\omega_j) \right| - \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1} g(\omega_j) + \frac{1}{m} [g(\omega_r) - g(\omega'_r)] \right| \right| \\ & \stackrel{(*)}{\leq} \frac{1}{m} \sup_{g \in \mathcal{G}} |g(\omega_r) - g(\omega'_r)| \end{aligned}$$

## Step-2: bounded difference property of $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$

Our choice:  $f(\omega_1, \dots, \omega_m) := \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ .

$$\begin{aligned} & |f(\omega_1, \dots, \omega_{r-1}, \omega_r, \omega_{r+1}, \dots, \omega_m) - f(\omega_1, \dots, \omega_{r-1}, \omega'_r, \omega_{r+1}, \dots, \omega_m)| = \\ & = \left| \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) \right| - \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) + \frac{1}{m} [g(\omega_r) - g(\omega'_r)] \right| \right| \\ & \stackrel{(*)}{\leq} \frac{1}{m} \sup_{g \in \mathcal{G}} |g(\omega_r) - g(\omega'_r)| \leq \frac{1}{m} \sup_{g \in \mathcal{G}} (|g(\omega_r)| + |g(\omega'_r)|) \end{aligned}$$

## Step-2: bounded difference property of $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$

Our choice:  $f(\omega_1, \dots, \omega_m) := \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ .

$$\begin{aligned} & |f(\omega_1, \dots, \omega_{r-1}, \omega_r, \omega_{r+1}, \dots, \omega_m) - f(\omega_1, \dots, \omega_{r-1}, \omega'_r, \omega_{r+1}, \dots, \omega_m)| = \\ & = \left| \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) \right| - \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) + \frac{1}{m} [g(\omega_r) - g(\omega'_r)] \right| \right| \\ & \stackrel{(*)}{\leq} \frac{1}{m} \sup_{g \in \mathcal{G}} |g(\omega_r) - g(\omega'_r)| \leq \frac{1}{m} \sup_{g \in \mathcal{G}} (|g(\omega_r)| + |g(\omega'_r)|) \\ & \leq \frac{1}{m} \left[ \sup_{g \in \mathcal{G}} |g(\omega_r)| + \sup_{g \in \mathcal{G}} |g(\omega'_r)| \right] \end{aligned}$$



## Step-2: bounded difference property of $\|\Lambda - \Lambda_m\|_{\mathcal{G}}$

Our choice:  $f(\omega_1, \dots, \omega_m) := \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ .

$$\begin{aligned} & |f(\omega_1, \dots, \omega_{r-1}, \omega_r, \omega_{r+1}, \dots, \omega_m) - f(\omega_1, \dots, \omega_{r-1}, \omega'_r, \omega_{r+1}, \dots, \omega_m)| = \\ & = \left| \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) \right| - \sup_{g \in \mathcal{G}} \left| \Lambda g - \frac{1}{m} \sum_{j=1}^m g(\omega_j) + \frac{1}{m} [g(\omega_r) - g(\omega'_r)] \right| \right| \\ & \stackrel{(*)}{\leq} \frac{1}{m} \sup_{g \in \mathcal{G}} |g(\omega_r) - g(\omega'_r)| \leq \frac{1}{m} \sup_{g \in \mathcal{G}} (|g(\omega_r)| + |g(\omega'_r)|) \\ & \leq \frac{1}{m} \left[ \sup_{g \in \mathcal{G}} |g(\omega_r)| + \sup_{g \in \mathcal{G}} |g(\omega'_r)| \right] \leq \frac{1+1}{m} = \frac{2}{m}. \end{aligned}$$

## Step-2: (\*) = reverse triangle inequality with sup

- Lemma:  $\mathcal{G}$ : set of functions,  $a, b : \mathcal{G} \rightarrow \mathbb{R}$  maps; then

$$\left| \sup_{g \in \mathcal{G}} |a(g)| - \sup_{g \in \mathcal{G}} |a(g) + b(g)| \right|$$

## Step-2: (\*) = reverse triangle inequality with sup

- Lemma:  $\mathcal{G}$ : set of functions,  $a, b : \mathcal{G} \rightarrow \mathbb{R}$  maps; then

$$\left| \sup_{g \in \mathcal{G}} |a(g)| - \sup_{g \in \mathcal{G}} |a(g) + b(g)| \right| \leq \sup_{g \in \mathcal{G}} |b(g)|.$$

## Step-2: (\*) = reverse triangle inequality with sup

- Lemma:  $\mathcal{G}$ : set of functions,  $a, b : \mathcal{G} \rightarrow \mathbb{R}$  maps; then

$$\left| \sup_{g \in \mathcal{G}} |a(g)| - \sup_{g \in \mathcal{G}} |a(g) + b(g)| \right| \leq \sup_{g \in \mathcal{G}} |b(g)|.$$

- Proof: combine

$$\sup_{g \in \mathcal{G}} |a(g) + b(g)| \leq \sup_{g \in \mathcal{G}} (|a(g)| + |b(g)|) \leq \sup_{g \in \mathcal{G}} |a(g)| + \sup_{g \in \mathcal{G}} |b(g)|,$$

## Step-2: (\*) = reverse triangle inequality with sup

- Lemma:  $\mathcal{G}$ : set of functions,  $a, b : \mathcal{G} \rightarrow \mathbb{R}$  maps; then

$$\left| \sup_{g \in \mathcal{G}} |a(g)| - \sup_{g \in \mathcal{G}} |a(g) + b(g)| \right| \leq \sup_{g \in \mathcal{G}} |b(g)|.$$

- Proof: combine

$$\sup_{g \in \mathcal{G}} |a(g) + b(g)| \leq \sup_{g \in \mathcal{G}} (|a(g)| + |b(g)|) \leq \sup_{g \in \mathcal{G}} |a(g)| + \sup_{g \in \mathcal{G}} |b(g)|,$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} |a(g)| &= \sup_{g \in \mathcal{G}} |a(g) + b(g) - b(g)| \\ &\leq \sup_{g \in \mathcal{G}} |a(g) + b(g)| + \sup_{g \in \mathcal{G}} |b(g)|. \end{aligned}$$

## Step-2: (\*) = reverse triangle inequality with sup

- Lemma:  $\mathcal{G}$ : set of functions,  $a, b : \mathcal{G} \rightarrow \mathbb{R}$  maps; then

$$\left| \sup_{g \in \mathcal{G}} |a(g)| - \sup_{g \in \mathcal{G}} |a(g) + b(g)| \right| \leq \sup_{g \in \mathcal{G}} |b(g)|.$$

- Proof: combine

$$\sup_{g \in \mathcal{G}} |a(g) + b(g)| \leq \sup_{g \in \mathcal{G}} (|a(g)| + |b(g)|) \leq \sup_{g \in \mathcal{G}} |a(g)| + \sup_{g \in \mathcal{G}} |b(g)|,$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} |a(g)| &= \sup_{g \in \mathcal{G}} |a(g) + b(g) - b(g)| \\ &\leq \sup_{g \in \mathcal{G}} |a(g) + b(g)| + \sup_{g \in \mathcal{G}} |b(g)|. \end{aligned}$$

$$\Rightarrow \pm \left[ \sup_{g \in \mathcal{G}} |a(g)| - \sup_{g \in \mathcal{G}} |a(g) + b(g)| \right] \leq \sup_{g \in \mathcal{G}} |b(g)|.$$

- Our choice:  $a(g) = \Lambda g - \frac{1}{m} \sum_{j=1} g(\omega_j)$ ,  $b(g) = \frac{1}{m} [g(\omega_r) - g(\omega'_r)]$ .

Applying McDiarmid to  $f$  ( $c_r = \frac{2}{m}$ ): with probability  $1 - e^{-\tau}$

$$\|\Lambda - \Lambda_m\|_{\mathcal{G}} \leq \underbrace{\mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}}}_{\text{Step-3: bounding this term}} + \frac{\sqrt{2\tau}}{\sqrt{m}}.$$

## Step-3: bounding $\mathbb{E}_{\omega_1, \dots, \omega_m} \|\Lambda - \Lambda_m\|_{\mathcal{G}}$

$\mathcal{G} = \{g_{\mathbf{z}} : \mathbf{z} \in \mathcal{S}_{\Delta}\}$  is a separable Carathéodory family, i.e.

- 1  $\omega \mapsto \cos(\omega^T \mathbf{z})$ : measurable for  $\forall \mathbf{z} \in \mathcal{S}_{\Delta}$ .



## Step-3: bounding $\mathbb{E}_{\omega_1, \dots, \omega_m} \|\Lambda - \Lambda_m\|_{\mathcal{G}}$

$\mathcal{G} = \{g_{\mathbf{z}} : \mathbf{z} \in \mathcal{S}_{\Delta}\}$  is a separable Carathéodory family, i.e.

- 1  $\omega \mapsto \cos(\omega^T \mathbf{z})$ : **measurable** for  $\forall \mathbf{z} \in \mathcal{S}_{\Delta}$ .
- 2  $\mathbf{z} \mapsto \cos(\omega^T \mathbf{z})$ : **continuous** for  $\forall \omega$ .

## Step-3: bounding $\mathbb{E}_{\omega_1, \dots, \omega_m} \|\Lambda - \Lambda_m\|_{\mathcal{G}}$

$\mathcal{G} = \{g_{\mathbf{z}} : \mathbf{z} \in \mathcal{S}_{\Delta}\}$  is a separable Carathéodory family, i.e.

- 1  $\omega \mapsto \cos(\omega^T \mathbf{z})$ : **measurable** for  $\forall \mathbf{z} \in \mathcal{S}_{\Delta}$ .
- 2  $\mathbf{z} \mapsto \cos(\omega^T \mathbf{z})$ : **continuous** for  $\forall \omega$ .
- 3  $\mathbb{R}^d$  is separable,  $\mathcal{S}_{\Delta} \subseteq \mathbb{R}^d \Rightarrow \mathcal{S}_{\Delta}$ : **separable**.

## Step-3: bounding $\mathbb{E}_{\omega_1, \dots, \omega_m} \|\Lambda - \Lambda_m\|_{\mathcal{G}}$

$\mathcal{G} = \{g_{\mathbf{z}} : \mathbf{z} \in \mathcal{S}_{\Delta}\}$  is a separable Carathéodory family, i.e.

①  $\omega \mapsto \cos(\omega^T \mathbf{z})$ : **measurable** for  $\forall \mathbf{z} \in \mathcal{S}_{\Delta}$ .

②  $\mathbf{z} \mapsto \cos(\omega^T \mathbf{z})$ : **continuous** for  $\forall \omega$ .

③  $\mathbb{R}^d$  is separable,  $\mathcal{S}_{\Delta} \subseteq \mathbb{R}^d \Rightarrow \mathcal{S}_{\Delta}$ : **separable**.

Thus, by [Steinwart and Christmann, 2008, Prop. 7.10]

$$\mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} \leq 2 \mathbb{E}_{\omega_{1:m}} \left[ \underbrace{\mathcal{R}(\mathcal{G}, \omega_{1:m})}_{:= \mathbb{E}_{\epsilon} \sup_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{j=1}^m \epsilon_j g(\omega_j) \right|} \right]$$

using the **uniformly boundedness** of  $\mathcal{G}$  ( $\sup_{g \in \mathcal{G}} \|g\|_{\infty} \leq 1$ ).

## Step-4: bounding $\mathcal{R}$

$$\mathcal{R}(\mathcal{G}, (\omega_j)_{j=1}^m) \leq \frac{8\sqrt{2}}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr,$$

## Step-4: bounding $\mathcal{R}$

$$\mathcal{R}(\mathcal{G}, (\omega_j)_{j=1}^m) \leq \frac{8\sqrt{2}}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr,$$

where

- $L^2(\Lambda_m) = L^2(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d), \Lambda_m)$ ,  $\|g\|_{L^2(\Lambda_m)} = \sqrt{\frac{1}{m} \sum_{j=1}^m g^2(\omega_j)}$ ,

## Step-4: bounding $\mathcal{R}$

$$\mathcal{R}(\mathcal{G}, (\omega_j)_{j=1}^m) \leq \frac{8\sqrt{2}}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr,$$

where

- $L^2(\Lambda_m) = L^2(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d), \Lambda_m)$ ,  $\|g\|_{L^2(\Lambda_m)} = \sqrt{\frac{1}{m} \sum_{j=1}^m g^2(\omega_j)}$ ,
- $|\mathcal{G}|_{L^2(\Lambda_m)} = \sup_{g_1, g_2 \in \mathcal{G}} \|g_1 - g_2\|_{L^2(\Lambda_m)}$ ,

## Step-4: bounding $\mathcal{R}$

$$\mathcal{R}(\mathcal{G}, (\omega_j)_{j=1}^m) \leq \frac{8\sqrt{2}}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr,$$

where

- $L^2(\Lambda_m) = L^2(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d), \Lambda_m)$ ,  $\|g\|_{L^2(\Lambda_m)} = \sqrt{\frac{1}{m} \sum_{j=1}^m g^2(\omega_j)}$ ,
- $|\mathcal{G}|_{L^2(\Lambda_m)} = \sup_{g_1, g_2 \in \mathcal{G}} \|g_1 - g_2\|_{L^2(\Lambda_m)}$ ,
- $\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)$ :  $r$ -covering number.
  - $r$ -net:  $S \subseteq \mathcal{G}$ , for  $\forall g \in \mathcal{G} \exists s \in S$  such that  $\|g - s\|_{L^2(\Lambda_m)} \leq r$ .
  - $\mathcal{N}$ : size of the smallest  $r$ -net of  $\mathcal{G}$ .

## Step-5: bound on $|\mathcal{G}|_{L^2(\Lambda_m)}$

$$\begin{aligned} |\mathcal{G}|_{L^2(\Lambda_m)} &= \sup_{g_1, g_2 \in \mathcal{G}} \|g_1 - g_2\|_{L^2(\Lambda_m)} \leq \sup_{g_1, g_2 \in \mathcal{G}} \left( \|g_1\|_{L^2(\Lambda_m)} + \|g_2\|_{L^2(\Lambda_m)} \right) \\ &\leq \sup_{g_1 \in \mathcal{G}} \|g_1\|_{L^2(\Lambda_m)} + \sup_{g_1 \in \mathcal{G}} \|g_2\|_{L^2(\Lambda_m)} \stackrel{*}{\leq} 2 \times 1, \end{aligned}$$



## Step-5: bound on $|\mathcal{G}|_{L^2(\Lambda_m)}$

$$\begin{aligned} |\mathcal{G}|_{L^2(\Lambda_m)} &= \sup_{g_1, g_2 \in \mathcal{G}} \|g_1 - g_2\|_{L^2(\Lambda_m)} \leq \sup_{g_1, g_2 \in \mathcal{G}} \left( \|g_1\|_{L^2(\Lambda_m)} + \|g_2\|_{L^2(\Lambda_m)} \right) \\ &\leq \sup_{g_1 \in \mathcal{G}} \|g_1\|_{L^2(\Lambda_m)} + \sup_{g_1 \in \mathcal{G}} \|g_2\|_{L^2(\Lambda_m)} \stackrel{*}{\leq} 2 \times 1, \end{aligned}$$

$$\begin{aligned} \sup_{g \in \mathcal{G}} \|g\|_{L^2(\Lambda_m)} &= \sup_{z \in \mathcal{S}_\Delta} \sqrt{\frac{1}{m} \sum_{j=1}^m g_z^2(\omega_j)} \\ &= \sup_{z \in \mathcal{S}_\Delta} \sqrt{\frac{1}{m} \sum_{j=1}^m \cos^2(\omega_j^T z)} \leq \sup_{z \in \mathcal{S}_\Delta} \sqrt{\frac{1}{m} \sum_{j=1}^m 1} = 1. \end{aligned}$$

## Step-5: bound on $\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)$

Let  $g_{z_1}, g_{z_2} \in \mathcal{G}$ . We want to bound  $\|g_{z_1} - g_{z_2}\|_{L^2(\Lambda_m)}$ . One term:

$$\begin{aligned} & \left| \cos(\omega^T \mathbf{z}_1) - \cos(\omega^T \mathbf{z}_2) \right| \\ &= \left\| \nabla_{\mathbf{z}} \cos(\omega^T \mathbf{z}_c) \right\|_2 \|\mathbf{z}_1 - \mathbf{z}_2\|_2 \\ &= \left\| -\sin(\omega^T \mathbf{z}_c) \omega \right\|_2 \|\mathbf{z}_1 - \mathbf{z}_2\|_2 \\ &\leq \|\omega\|_2 \|\mathbf{z}_1 - \mathbf{z}_2\|_2, \end{aligned}$$

where  $\mathbf{z}_c \in (\mathbf{z}_1, \mathbf{z}_2)$ .

## Step-5: bound on $\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)$

- Smooth parameterization:

$$\begin{aligned}\|g_{\mathbf{z}_1} - g_{\mathbf{z}_2}\|_{L^2(\Lambda_m)} &\leq \sqrt{\frac{1}{m} \sum_{j=1}^m (\|\omega_j\|_2 \|\mathbf{z}_1 - \mathbf{z}_2\|_2)^2} \\ &= \|\mathbf{z}_1 - \mathbf{z}_2\|_2 \underbrace{\sqrt{\frac{1}{m} \sum_{j=1}^m \|\omega_j\|_2^2}}_{=:A}.\end{aligned}$$

- $r$ -net on  $(\mathcal{S}_\Delta, \|\cdot\|_2) \Rightarrow r' = Ar$ -net on  $(\mathcal{G}, L^2(\Lambda_m))$ .
- In other words,  $\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r) \leq \mathcal{N}(\mathcal{S}_\Delta, \|\cdot\|_2, \frac{r}{A})$ .

## Step-5: bound on $\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)$

- Note that  $\mathcal{S}_\Delta \subseteq B_{\|\cdot\|_2} \left( \mathbf{t}, \frac{|\mathcal{S}_\Delta|}{2} \right)$  for some  $\mathbf{t} \in \mathbb{R}^d$ .
- $\mathcal{N}(B_{\|\cdot\|_2}(\mathbf{s}, R), \|\cdot\|_2, \epsilon) \leq \left( \frac{2R}{\epsilon} + 1 \right)^d$  for  $\forall \mathbf{s} \in \mathbb{R}^d$ .
- Thus

$$\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r) \leq \left( \frac{2|\mathcal{S}|A}{r} + 1 \right)^d$$

by  $|\mathcal{S}_\Delta| \leq 2|\mathcal{S}|$  and the compactness of  $\mathcal{S}_\Delta$ .

## Step-5: bound on $\mathcal{R}$

Combining the obtained

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{8\sqrt{2}}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr,$$

$$|\mathcal{G}|_{L^2(\Lambda_m)} \leq 2,$$

$$\log [\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)] \leq d \log \left( \frac{2|\mathcal{S}|A}{r} + 1 \right)$$

results

## Step-5: bound on $\mathcal{R}$

Combining the obtained

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{8\sqrt{2}}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr,$$

$$|\mathcal{G}|_{L^2(\Lambda_m)} \leq 2,$$

$$\log [\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)] \leq d \log \left( \frac{2|\mathcal{S}|A}{r} + 1 \right)$$

results, we have ( $r \leq 2$ )

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{8\sqrt{2d}}{\sqrt{m}} \int_0^2 \sqrt{\log \left( \frac{2|\mathcal{S}|A + 2}{r} \right)} dr.$$

## Step-5: bound on $\mathcal{R}$

Using  $|\mathcal{S}|A + 1 \leq (|\mathcal{S}| + 1)(A + 1)$

$$\begin{aligned}\mathcal{R}(\mathcal{G}, \omega_{1:m}) &\leq \frac{8\sqrt{2d}}{\sqrt{m}} \int_0^2 \sqrt{\log\left(\frac{2|\mathcal{S}|A + 2}{r}\right)} dr \\ &\leq \frac{8\sqrt{2d}}{\sqrt{m}} \left[ \int_0^2 \sqrt{\log\frac{2(|\mathcal{S}| + 1)}{r}} dr + 2\sqrt{\log(A + 1)} \right] \\ &= \frac{16\sqrt{2d}}{\sqrt{m}} \left[ \int_0^1 \sqrt{\log\frac{|\mathcal{S}| + 1}{r}} dr + \sqrt{\log(A + 1)} \right].\end{aligned}$$

Applying  $\int_0^1 \sqrt{\log\frac{a}{r}} dr \leq \sqrt{\log a} + \frac{1}{2\sqrt{\log a}}$  ( $a > 1$ )

## Step-5: bound on $\mathcal{R}$

we get

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{16\sqrt{2d}}{\sqrt{m}} \left[ \sqrt{\log(|\mathcal{S}| + 1)} + \frac{1}{2\sqrt{\log(|\mathcal{S}| + 1)}} + \sqrt{\log(A + 1)} \right]. \quad (1)$$



## Step-5: bound on $\mathcal{R}$

we get

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{16\sqrt{2d}}{\sqrt{m}} \left[ \sqrt{\log(|\mathcal{S}| + 1)} + \frac{1}{2\sqrt{\log(|\mathcal{S}| + 1)}} + \sqrt{\log(A + 1)} \right]. \quad (1)$$

By the Jensen inequality

## Step-5: bound on $\mathcal{R}$

we get

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{16\sqrt{2d}}{\sqrt{m}} \left[ \sqrt{\log(|\mathcal{S}| + 1)} + \frac{1}{2\sqrt{\log(|\mathcal{S}| + 1)}} + \sqrt{\log(A + 1)} \right]. \quad (1)$$

By the Jensen inequality

$$\mathbb{E}_{\omega_{1:m}} \sqrt{\log(A + 1)} \leq \sqrt{\mathbb{E}_{\omega_{1:m}} \log(A + 1)} \leq \sqrt{\log(\mathbb{E}_{\omega_{1:m}} A + 1)},$$

## Step-5: bound on $\mathcal{R}$

we get

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{16\sqrt{2d}}{\sqrt{m}} \left[ \sqrt{\log(|\mathcal{S}| + 1)} + \frac{1}{2\sqrt{\log(|\mathcal{S}| + 1)}} + \sqrt{\log(A + 1)} \right]. \quad (1)$$

By the Jensen inequality

$$\begin{aligned} \mathbb{E}_{\omega_{1:m}} \sqrt{\log(A + 1)} &\leq \sqrt{\mathbb{E}_{\omega_{1:m}} \log(A + 1)} \leq \sqrt{\log(\mathbb{E}_{\omega_{1:m}} A + 1)}, \\ \mathbb{E}_{\omega_{1:m}} A &\leq \sqrt{\frac{1}{m} \sum_{j=1}^m \mathbb{E}_{\omega_j} [\|\omega_j\|_2^2]} =: \sigma. \Rightarrow \end{aligned}$$

## Step-5: bound on $\mathcal{R}$

we get

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \leq \frac{16\sqrt{2d}}{\sqrt{m}} \left[ \sqrt{\log(|\mathcal{S}| + 1)} + \frac{1}{2\sqrt{\log(|\mathcal{S}| + 1)}} + \sqrt{\log(A + 1)} \right]. \quad (1)$$

By the Jensen inequality

$$\begin{aligned} \mathbb{E}_{\omega_{1:m}} \sqrt{\log(A + 1)} &\leq \sqrt{\mathbb{E}_{\omega_{1:m}} \log(A + 1)} \leq \sqrt{\log(\mathbb{E}_{\omega_{1:m}} A + 1)}, \\ \mathbb{E}_{\omega_{1:m}} A &\leq \sqrt{\frac{1}{m} \sum_{j=1}^m \mathbb{E}_{\omega_j} [\|\omega_j\|_2^2]} =: \sigma. \Rightarrow \\ \mathbb{E}_{\omega_{1:m}} \mathcal{R}(\mathcal{G}, \omega_{1:m}) &\leq (1), \text{ but with } A \rightarrow \sigma. \end{aligned}$$

## Step-6: putting together

Result:  $k$  continuous, shift-invariant kernel; for any  $\tau > 0$ ,  $\mathcal{S} \neq \emptyset$  compact set,

$$\Lambda^m \left( \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |\hat{k}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y})| \geq \frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau},$$

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

## Step-6: putting together

Result:  $k$  continuous, shift-invariant kernel; for any  $\tau > 0$ ,  $\mathcal{S} \neq \emptyset$  compact set,

$$\Lambda^m \left( \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |\hat{k}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y})| \geq \underbrace{\frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}}}_{:=\epsilon} \right) \leq e^{-\tau},$$

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

Equivalently

$$\Lambda^m \left( \left\| \hat{k} - k \right\|_{\mathcal{S}} \geq \epsilon \right) \leq e^{-\frac{[\epsilon\sqrt{m} - h(d, |\mathcal{S}|, \sigma)]^2}{2}}.$$

## Discussion (Borel-Cantelli lemma)

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

# Discussion (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.,  $|\mathcal{S}_m| = e^{o(m)}$ )  $\leftrightarrow$
  - Old:  $|\mathcal{S}_m| = o\left(\sqrt{m/\log m}\right)$ .



# Discussion (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.,  $|\mathcal{S}_m| = e^{o(m)}$ )  $\leftrightarrow$
  - Old:  $|\mathcal{S}_m| = o\left(\sqrt{m/\log m}\right)$ .
- Specifically:
  - *asymptotically* optimal result [Csörgő and Totik, 1983, Theorem 2] (if  $\psi$  vanishes at  $\infty$ ),
  - at faster rate  $\Rightarrow$  even conv. in prob. would fail.

# Direct consequence: $L^r$ guarantee ( $1 < r$ )

Idea:

- Note that

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

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

# $L^r$ large deviation inequality

Under the previous assumptions:

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

In other words,

$$\|\hat{k} - k\|_{L^r(\mathcal{S})} = O_{a.s.} \left( m^{-1/2} |\mathcal{S}|^{2d/r} \sqrt{\log |\mathcal{S}|} \right).$$

For  $2 \leq r$ : direct  $L^r$  proof  $\Rightarrow \sqrt{\log(|\mathcal{S}|)}$  factor can be discarded.

- If  $\text{supp}(\Lambda)$  is bounded
  - $k$ -proof can be extended ( $L^r$  as well), but
  - Gaussian kernel:(
- [Rahimi and Recht, 2007]'s proof:
  - Hoeffding inequality (boundedness!) + Lipschitzness,
- Bernstein + Lipschitzness: handles  $\partial^{\mathbf{p}, \mathbf{q}} k$  with
  - moment constraints on  $\Lambda$  (example: Gaussian kernel).
  - slightly worse rates.





- Kernel + derivative approximations.
- Performance: uniform,  $L^r$ .
- Detailed finite-sample analysis, optimal rates.
- Paper (submitted to NIPS):
  - RFF: <http://arxiv.org/abs/1506.02155>,
  - infD exp. fitting: <http://arxiv.org/abs/1506.02564>.



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 Sci. Math. (Szeged)*, 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>.
-  Steinwart, I. and Christmann, A. (2008).  
*Support Vector Machines*.  
Springer.
-  Sutherland, D. and Schneider, J. (2015).  
On the error of random fourier features.  
In *Conference on Uncertainty in Artificial Intelligence (UAI)*.
-  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.

- Ingredients:
  - $(X, \tau)$ : topological space with a countable basis.
  - $\mathcal{B} = \sigma(\tau)$ : sigma-algebra generated by  $\tau$ .
  - $\Lambda$ : measure on  $(X, \mathcal{B})$ .

Then

$$\text{supp}(\Lambda) = \overline{\cup\{A \in \tau : \Lambda(A) = 0\}},$$

i.e., the complement of the union of all open  $\Lambda$ -null sets.

- Our choice:  $X = \mathbb{R}^d$ .