

Convergence Theory for Vector-Valued Random Features

Nicholas H. Nelsen

Caltech

Email: nnelsen@caltech.edu

Web: nicholashnelsen.com

Supported By: **Amazon, NSF, ONR, AFOSR, DOD**

BIRS Scientific Machine Learning Workshop
Banff, Alberta, Canada

June 2023

N.H.N. and Andrew M. Stuart

The Random Feature Model for Input-Output Maps between Banach Spaces

SIAM Journal on Scientific Computing, Vol. 43, No. 5, pp. A3212–A3243, 2021.

Random feature map (\mathcal{X} and \mathcal{Y} infinite-dimensional)

$$\varphi: \mathcal{X} \times \Theta \rightarrow \mathcal{Y} \quad \text{and} \quad \text{a probability measure } \mu$$

Parametric model (looks like Monte Carlo)

$$\Psi_{\text{RFM}}(u; \alpha) := \frac{1}{M} \sum_{m=1}^M \alpha_m \varphi(u; \theta_m), \quad \theta_m \stackrel{\text{iid}}{\sim} \mu$$

Random Feature Ridge Regression

Data

$$u_n \stackrel{\text{iid}}{\sim} \nu \quad \text{and} \quad y_n = \Psi^\dagger(u_n) + \text{Noise} \quad \text{for} \quad n = 1, \dots, N$$

Convex problem (RF-RR)

$$\hat{\alpha}^{(N,M,\lambda)} := \arg \min_{\alpha \in \mathbb{R}^M} \left\{ \frac{1}{N} \sum_{n=1}^N \|y_n - \Psi_{\text{RFM}}(u_n; \alpha)\|_{\mathcal{Y}}^2 + \lambda \left(\frac{1}{M} \sum_{m=1}^M |\alpha_m|^2 \right) \right\}$$

Samuel Lanthaler and **N.H.N.**

Error Bounds for Learning with Vector-Valued Random Features

Submitted 2023 ([arXiv:2305.17170](https://arxiv.org/abs/2305.17170) [stat.ML](https://arxiv.org/abs/2305.17170))

Well-Specified Error Analysis (for any input and output dimension)

Trained RFM $\Psi_{\text{RFM}}(\cdot; \hat{\alpha}^{(N, M, \lambda)})$, where $\hat{\alpha}^{(N, M, \lambda)} \in \mathbb{R}^M$ solves RF-RR

Assumptions (no spectral assumptions on K are needed due to matrix-free analysis)

- ▶ $y_n = \Psi^\dagger(u_n) + \eta_n$, where $\eta_n \stackrel{\text{iid}}{\sim} \eta$ is subexponential on \mathcal{Y}
- ▶ Ψ^\dagger and φ are bounded a.s.

Theorem (Squared error: well-specified convergence rate)

If Ψ^\dagger belongs to the RKHS of the RF pair (φ, μ) (relaxations too), then

$$\mathbb{E}^{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \hat{\alpha}^{(N, M, \lambda)}) \right\|_{\mathcal{Y}}^2 \lesssim \lambda + \frac{1}{M} + \frac{1}{\sqrt{N}} \quad \text{with high probability.}$$

Comparison to Existing Well-Specified Results

Recall notation

- ▶ Regularization parameter: λ
- ▶ Training data sample size: N
- ▶ Number of random features: M

Paper	Approach	λ	$\dim(\mathcal{Y})$	M	Squared Error
Rahimi & Recht '08	“kitchen sinks”	—	1	N	$\omega(N^{-1/2})$
Rudi & Rosasco '17	matrix concn.	$N^{-1/2}$	1	$\sqrt{N} \log(N)$	$O(N^{-1/2})$
Li et al. '21	matrix concn.	$N^{-1/2}$	1	$\sqrt{N} \log(O(N))$	$O(N^{-1/2})$
This Talk (SOTA)	“kitchen sinks”	$N^{-1/2}$	∞	\sqrt{N}	$O(N^{-1/2})$

Core Proof Idea

Loss

$$L(u; \alpha) := \|\Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \alpha)\|_y^2$$

Error decomposition

(Rahimi & Recht '08)

$$\underbrace{\mathbb{E}^{u \sim \nu} [L(u; \hat{\alpha}^{(N, M, \lambda)})]}_{\text{squared error}} = \underbrace{\frac{1}{N} \sum_{n=1}^N L(u_n; \hat{\alpha}^{(N, M, \lambda)})}_{\text{empirical approximation error (Monte Carlo)}} + \underbrace{\left[\mathbb{E}^{u \sim \nu} [L(u; \hat{\alpha}^{(N, M, \lambda)})] - \frac{1}{N} \sum_{n=1}^N L(u_n; \hat{\alpha}^{(N, M, \lambda)}) \right]}_{\text{generalization gap (linearity and empirical processes)}}$$

Unifying Sources of Error

Approximation, finite data, noise, and discretization (optimization error is zero)

Corollary (Stability to discretization error)

Let the data be discretized as

$$y_n = \Psi_h^\dagger(u_n) + \eta_n,$$

where $h > 0$ denotes a discretization parameter corresponding to bounded discretized operator Ψ_h^\dagger . If bounded Ψ^\dagger belongs to the RKHS of (φ, μ) , then $\lambda \asymp N^{-1/2} \asymp M^{-1}$ guarantees that

$$\mathbb{E}^{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \hat{\alpha}^{(N, M, \lambda)}) \right\|_{\mathcal{Y}}^2 \lesssim \frac{1}{\sqrt{N}} + \varepsilon_h^2 \quad \text{with high probability,}$$

where the *discretization error* is

$$\varepsilon_h := \operatorname{ess\,sup}_{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_h^\dagger(u) \right\|_{\mathcal{Y}}.$$

Theorem (Strong statistical consistency)

If the number of features $M = \tilde{\Omega}(\sqrt{N})$ and the penalty strength $\lambda = \tilde{\Omega}(1/\sqrt{N})$, then

$$\lim_{N \rightarrow \infty} \mathbb{E}^{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \hat{\alpha}^{(N, M_N, \lambda_N)}) \right\|_{\mathcal{Y}}^2 = 0 \quad \text{with probability one.}$$

(Above, Ψ^\dagger is just bounded)

Theorem (Squared error: slow rates)

If Ψ^\dagger *does not* belong to the RKHS of (φ, μ) but satisfies a “regularity source condition,” and $M \asymp \sqrt{N}$ and $\lambda \asymp 1/\sqrt{N}$, then there exists $0 < r \leq 1/2$ such that

$$\mathbb{E}^{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \hat{\alpha}^{(N, M, \lambda)}) \right\|_{\mathcal{Y}}^2 \lesssim N^{-r} \quad \text{with high probability.}$$

Theorem (Strong statistical consistency)

If the number of features $M = \tilde{\Omega}(\sqrt{N})$ and the penalty strength $\lambda = \tilde{\Omega}(1/\sqrt{N})$, then

$$\lim_{N \rightarrow \infty} \mathbb{E}^{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \hat{\alpha}^{(N, M_N, \lambda_N)}) \right\|_{\mathcal{Y}}^2 = 0 \quad \text{with probability one.}$$

(Above, Ψ^\dagger is just bounded)

Theorem (Squared error: slow rates)

If Ψ^\dagger **does not** belong to the RKHS of (φ, μ) but satisfies a “regularity source condition,” and $M \asymp \sqrt{N}$ and $\lambda \asymp 1/\sqrt{N}$, then there exists $0 < r \leq 1/2$ such that

$$\mathbb{E}^{u \sim \nu} \left\| \Psi^\dagger(u) - \Psi_{\text{RFM}}(u; \hat{\alpha}^{(N, M, \lambda)}) \right\|_{\mathcal{Y}}^2 \lesssim N^{-r} \quad \text{with high probability.}$$

Complete theory for vector-valued RF-RR algorithm

- ▶ Statistical consistency of RF for supervised learning and UQ
- ▶ SOTA rates in any dimension (matrix-free analysis)
- ▶ Includes error due to model misspecification and discretization