

Weighting Inputs by Sensitivity in Random Feature Expansions

John Darges

Department of Mathematics
Emory University

Joint work with Laura Weidensager
Chemnitz University of Technology

Partially supported by NSF Grant DMS-2038118

May 22, 2025

High-dimensional Approximation

Approximate a function $z = f(\mathbf{x})$, where \mathbf{x} is a random vector taking values in \mathbb{R}^n whose law is μ

- Approximate using a set of interpolation or sample points
- Budget for sample points is limited
- Input dimension is high enough that ordered basis methods (orthogonal polynomials, Fourier series) are not feasible
- Make no other assumptions about structure of f except it is in $L_2(\mu)$ (*this doesn't mean f has no structure*)
- Methods of choice often are kernels or neural networks

Kernel Methods and Kernel Operator

- Let $\kappa : X \times X \rightarrow \mathbb{C}$ be a reproducing kernel with associated RKHS \mathcal{H}_κ .
- Kernel approximation aims to find f^* that solves $\min_{\hat{f} \in \mathcal{H}_\kappa} \|f - \hat{f}\|_2$
- Note that the integral operator $T_\kappa : L_2(\mu) \rightarrow L_2(\mu)$ has \mathcal{H}_κ as its image

$$(T_\kappa g)(\mathbf{x}) = \int_{\mathbb{R}^n} g(\mathbf{y}) \kappa(\mathbf{x}, \mathbf{y}) d\mu(\mathbf{y})$$

- There exists g^* such that $T_\kappa g^* = f^*$, we can describe an approximation

$$f^* \approx \sum_{l=1}^M \frac{g^*(\mathbf{x}_k)}{M} \kappa(\cdot, \mathbf{x}_k), \quad \mathbf{x}_k \sim \mu$$

- We don't know f^*, g^* – we find coefficients through least squares (kernel ridge regression)

Feature Map Representation

- Many favorite kernels come from a feature map ϕ and a measure τ , admitting a (non-unique) feature map representation

$$\kappa(\mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^n} \phi(\mathbf{x}^\top \boldsymbol{\omega}) \overline{\phi(\mathbf{y}^\top \boldsymbol{\omega})} d\tau(\boldsymbol{\omega})$$

- Bochner: shift invariant kernels are Fourier transforms of Borel measures
- Feature map induces an integral operator $T_\phi : L_2(\tau) \rightarrow L_2(\mu)$ whose image is \mathcal{H}_κ ¹

$$(T_\phi \alpha)(\mathbf{x}) = \int_{\mathbb{R}^n} \alpha(\boldsymbol{\omega}) \phi(\mathbf{x}^\top \boldsymbol{\omega}) d\tau(\boldsymbol{\omega})$$

- There exists α^* such that $T_\phi \alpha^* = f^*$, leading to an approximation

$$f^* \approx \sum_{k=1}^M \frac{\alpha^*(\boldsymbol{\omega}_k)}{K} \phi(\boldsymbol{\omega}_k^\top \cdot), \quad \boldsymbol{\omega}_k \sim \tau$$

¹F. Bach. On the Equivalence between Kernel Quadrature Rules and Random Feature Expansions. 2017.

Random Feature Expansions

- In practice, many RKHS dense in $L_2(\mu)$ – some f^* arbitrarily close to f exists²
- Class $\{\phi(\boldsymbol{\omega}^\top \cdot) : \boldsymbol{\omega} \sim \tau\}$ is dense in \mathcal{H}_κ , random feature expansion (RFE) is an expansion over random choices of basis functions from this class³
- Common choice is Fourier features $\{e^{i\boldsymbol{\omega}^\top \cdot}\}$ or $\{\cos(\boldsymbol{\omega}^\top \cdot)\}$
- Entries of $\boldsymbol{\omega}$ often identical and independent, e.g., each $\omega^j \sim \mathcal{N}(0, 1)$ so τ is a product of identical Gaussian measures
- Coefficients β_k determined by solving linear least squares problem

$$\min_{\beta} \sum_{i=1}^N \left(f(\mathbf{x}_i) - \sum_{k=1}^K \beta_k \phi(\boldsymbol{\omega}_k^\top \mathbf{x}_i) \right)^2, \quad \boldsymbol{\omega}_k \sim \tau$$

- Equivalent to single layer neural networks with random weights⁴

²C. Michelli, Y. Xu, H. Zhang. Universal Kernels. 2008.

³A. Rahimi, B. Recht. Uniform approximation of functions with random bases. 2008.

⁴W. Cao, X. Wang, Z. Ming, J. Gao. A review on neural networks with random weights. 2018.

Choices in Kernels and Choosing Kernels

- Even if kernel is universal, choice must be catered to the function or data
- If kernel chosen from a parametric class, must make choice of parameters, e.g., shape parameter
- Class of kernels might be unsuitable for f
- Choosing kernel with a feature map representation comes down to choosing the measure τ
- One should choose kernel based on smoothness and **multivariable interaction structure** of f
- Even if inputs are high-dimensional, f may be dominated by only a few inputs or a few lower-order interaction terms


Decompositions of Multivariable Functions

- Multivariable functions can be decomposed (in different ways) into terms depending on groups of inputs⁵

$$f = \sum_{\mathbf{u} \subseteq [n]} f_{\mathbf{u}}, \quad f_{\mathbf{u}} = T_{\mathbf{u}}f - \sum_{\mathbf{v} \subset \mathbf{u}} f_{\mathbf{v}}$$

- Terms $f_{\mathbf{u}}$ depends only on subset of inputs $x_j, j \in \mathbf{u}$
- $T_{\mathbf{u}}$ is a linear transformation that determines the type of decomposition
- Anchored: $(T_{\mathbf{u}}f)(\mathbf{x}_{\mathbf{u}}) = f(\mathbf{x}_{\mathbf{u}}, \mathbf{x}_{\mathbf{u}^c} = 0)$ fixes values of inputs not in indices \mathbf{u}
- Hoeffding-Sobol' (aka ANOVA)⁶: $(T_{\mathbf{u}}f)(\mathbf{x}_{\mathbf{u}}) = \int f d\mathbf{x}_{\mathbf{u}^c}$ integrates out inputs not in indices \mathbf{u}
- Allow us to project f into 2^n orthogonal components

⁵F. Kuo, I. Sloan, G. Wasilkowski, H. Woźniakowski. On Decompositions of Multivariable Functions. 2010.

⁶I.M. Sobol'. Multidimensional Quadrature Formulas and Haar Functions. 1969. 

Multivariable Structure Matters for the Sampling Distribution

- Suppose f is an additive function, meaning it has no higher order interaction terms
- Then $f = \sum_{i=1}^n f_i$ and $f_{\mathbf{u}} = 0$ if $|\mathbf{u}| > 1$
- Let $f^N = \sum_{k=1}^K \beta_k e^{i\omega_k^\top \cdot}$ be a random Fourier approximation to f where ω_k is sampled from standard normal distribution
- The error in the approximation is

$$\|f - f^N\| = \sum_{\mathbf{u} \subseteq [n]} \|f_{\mathbf{u}} - f_{\mathbf{u}}^N\| \quad (1)$$

$$= \sum_{i=1}^n \|f_i - f_i^N\| + \sum_{|\mathbf{u}| \geq 2} \|f_{\mathbf{u}}^N\| \quad (2)$$

- This happens because we almost surely sample basis functions from the class $\{e^{i\omega^\top \cdot}\}$ containing every order of interaction


Sparse Sampling

- Recent work proposed imposing interaction structure on RFEs through sparse weights
- Sparse features $\phi(\omega^\top \mathbf{x})$ have weights where most entries in ω are 0
- Robust sparse reconstruction method proposed in⁷
- Sparse random features proposed in⁸ assume target model only has low order interactions
- Weight sparsity in⁹ determined by hyperparameter tuning
- Sparse RFEs in¹⁰ sample weights based on ANOVA decomposition of target model

⁷E. Saha, H. Schaeffer, G. Tran. HARFE: hard-ridge random feature expansion. 2023.

⁸A. Hashemi, H. Schaeffer, R. Shi, U. Topcu, G. Tran, R. Ward. Generalization bounds for sparse random feature expansions. 2023.

⁹J.D., A. Alexanderian, P.A. Gremaud. Extreme Learning Machines for Variance-Based Global Sensitivity Analysis. 2024.

¹⁰L. Weidensager, D. Potts ANOVA-boosting for Random Fourier Features. 2024. 

Weighted ANOVA Kernel

- To account for multivariable structure, use a weighted ANOVA kernel

$$\kappa_W(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{u} \subseteq [n]} \gamma_{\mathbf{u}} \kappa_{\mathbf{u}}(\mathbf{x}_{\mathbf{u}}, \mathbf{y}_{\mathbf{u}})$$

- If we use a kernel $\kappa(\mathbf{x}, \mathbf{y}) = \int e^{i\mathbf{x}^\top \boldsymbol{\omega}} \overline{e^{i\mathbf{y}^\top \boldsymbol{\omega}}} d\tau(\boldsymbol{\omega})$ and replace with a measure that places weights on inputs

$$\tau = \bigotimes_{i=1}^n \nu \quad \rightarrow \quad \tau_W = \bigotimes_{i=1}^n (\rho_i \nu + (1 - \rho_i) \delta_0), \quad \rho_i \in (0, 1)$$

- This results in a weighted kernel derived from the original

$$\kappa_W(\mathbf{x}, \mathbf{y}) = \int e^{i\mathbf{x}^\top \boldsymbol{\omega}} \overline{e^{i\mathbf{y}^\top \boldsymbol{\omega}}} d\tau_W(\boldsymbol{\omega}) = \sum_{\mathbf{u} \subseteq [n]} \gamma_{\mathbf{u}} \int e^{i\mathbf{x}_{\mathbf{u}}^\top \boldsymbol{\omega}_{\mathbf{u}}} \overline{e^{i\mathbf{y}_{\mathbf{u}}^\top \boldsymbol{\omega}_{\mathbf{u}}}} d\tau(\boldsymbol{\omega})$$

- Kernel weights are $\gamma_{\mathbf{u}} = \prod_{i \in \mathbf{u}} \rho_i \prod_{j \notin \mathbf{u}} (1 - \rho_j)$

Weighted Random Features

- Weighted ANOVA kernel is too cumbersome to work with – use feature representation instead

$$T_\phi \alpha = \int \alpha(\boldsymbol{\omega}) e^{i\mathbf{x}^\top \boldsymbol{\omega}} d\tau_W(\boldsymbol{\omega}) = \sum_{\mathbf{u} \subseteq [n]} \int \alpha(\boldsymbol{\omega}_u) e^{i\mathbf{x}_u^\top \boldsymbol{\omega}_u} d\tau(\boldsymbol{\omega})$$

- Random weighted feature expansion is

$$f^N = \sum_{k=1}^K \hat{\beta}_k e^{i\mathbf{x}^\top \boldsymbol{\omega}_k}, \quad \boldsymbol{\omega}_k \sim \tau_W$$

- This means the entry in each $\boldsymbol{\omega}$ is $\omega^i = YZ_i$, $Y \sim \nu$, $Z_i \sim \text{Bern}(1 - \rho_i)$

Proposed Approach

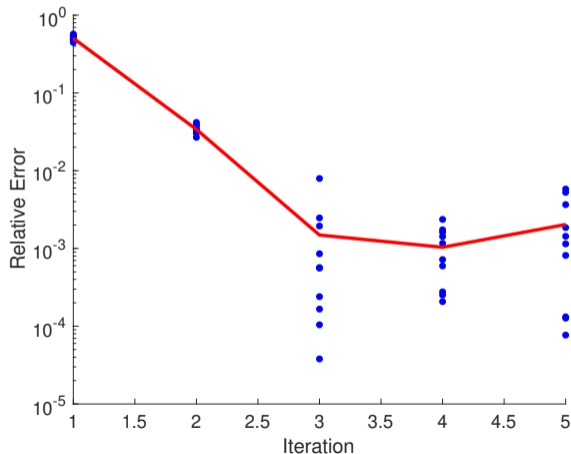
- Introduce random sparsity to feature weights based on sensitivity of inputs they correspond to
- If f is not sensitive to x_i , all its corresponding weights should be sparse (set to 0)
- Sample entries in weight vector ω^i corresponding to θ_i as random variable $\omega^i = Z_i X$, where $X \sim \mathcal{N}(0, 1)$, $Z_i \sim \text{Bern}(1 - \rho_i)$,
- ρ_i gives sensitivity of θ_i (**total Sobol' index**, derivative-based measures, Shapley value, etc.)
- **Begin by training a rough approximation to determine estimates for ρ_i , then build new approximation with updated sampling distribution**

- 1.
2. Set sampling distribution to standard normal distribution, $\omega^i \sim \mathcal{N}(0, 1)$ for all i
3. Sample feature weights and compute RFE approximation
4. Estimate input sensitivities ρ_i using the approximation
5. Update sampling distribution with random sparsity, $\omega^i = Z_i X$, where $X \sim \mathcal{N}(0, 1)$, $Z_i \sim \text{Bern}(1 - \rho_i)$
6. Iteratively repeat 2–5

How does this work in practice?

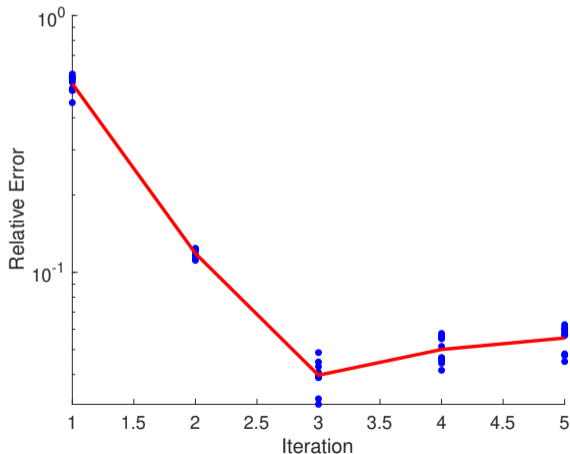
Example: Griewank Function

$$F(\mathbf{x}) = \sum_{k=1}^{50} \frac{x_k^2}{4000} + \prod_{k=1}^{50} \left(\cos\left(\frac{x}{\sqrt{k}}\right) + 1 \right), \quad \mathbf{x} \sim \mathcal{U}([-600, 600]^{50}).$$



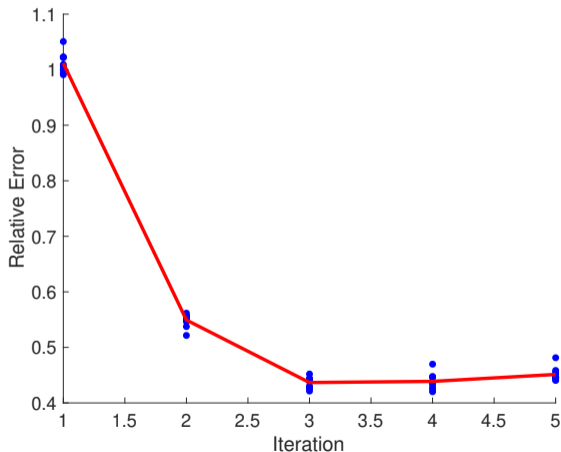
Example: Dixon-Price Function

$$f(\mathbf{x}) = (x_1 - 1)^2 + \sum_{k=2}^{50} k(2x_k - x_{k-1})^2, \quad \mathbf{x} \sim \mathcal{U}([-10, 10]^{50}).$$



Example: Sobol' g-function

$$F(\mathbf{x}) = \prod_{k=1}^{10} \frac{|4x_k - 2| + a_k}{1 + a_k}, \quad a_k = \frac{k-2}{2}, \quad \mathbf{x} \sim \mathcal{U}([0, 1]^{10})$$



- Does altering the sampling distribution based on weight sensitivity provide benefits for certain classes of functions?
- Does $\|f - f_W^N\|$ converge faster than $\|f - f^N\|$ for certain classes of f dominated by a subset of inputs or interaction terms?
- This analysis relies on showing $\|f^* - f_W^N\|$ converges faster than $\|f^* - f^N\|$, recall

$$\|f^* - f^N\| = \left\| \int \alpha^*(\omega) e^{i\omega^\top \mathbf{x}} d\tau(\omega) - \sum_{k=1}^K \frac{\alpha^*(\omega_k)}{K} e^{i\omega_k^\top \mathbf{x}} \right\|$$

-

Summary and Goals

- Proposed an algorithm for random feature expansions that modifies weight sampling by taking relative input sensitivities into account
- Incorporate input weighting into the state-of-the-art RFE algorithms
- Identify ideal applications and adversarial examples
- Are there better approaches to input sensitivity based on the choice of kernel?
- Should input weighting be implemented in different fashion, e.g., by scaling inputs or scaling variances in the sampling distribution?