

# Explainable Approximation in High Dimensions: Fourier-Based Algorithms Meet Kernel Methods

Basics

## ANOVA Decomposition

Let  $f \in L_2(\mathbb{T}^d)$ . Based on the **ANOVA (analysis of variance) decomposition**

$$f(\mathbf{x}) = \sum_{u \subseteq \{1, \dots, d\}} f_u(\mathbf{x}_u)$$

the importance of single dimensions as well as of groups of dimensions can be studied.

$$f_\emptyset = \int_{\mathbb{T}^d} f(\mathbf{x}) d\mathbf{x},$$

$$f_u(\mathbf{x}_u) = \int_{\mathbb{T}^{d-|u|}} f(\mathbf{x}) d\mathbf{x}_{\bar{u}} - \sum_{v \subsetneq u} f_v(\mathbf{x}_v),$$

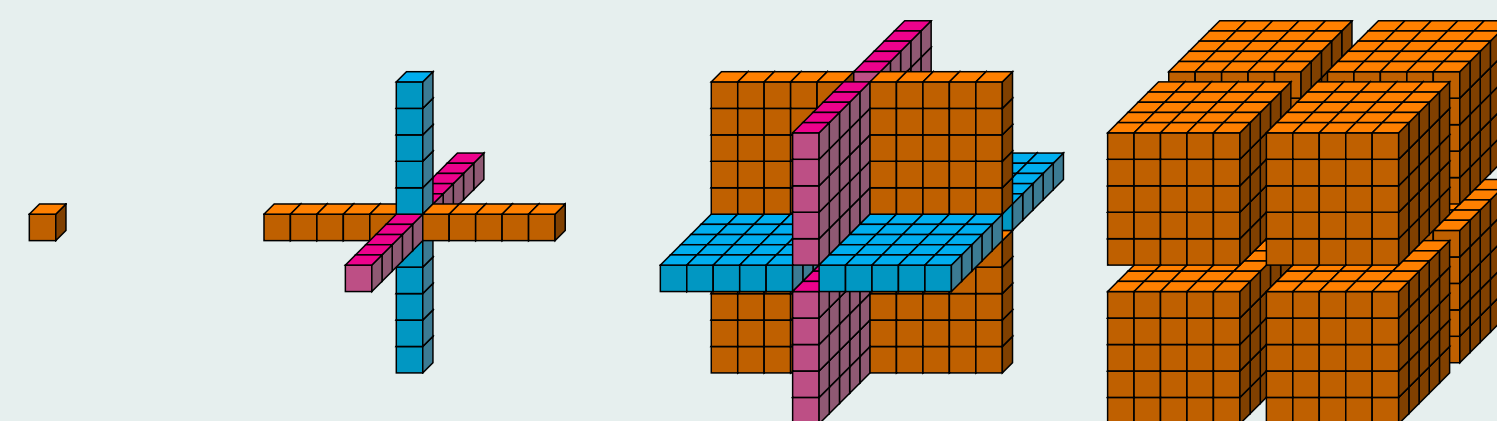
where  $\bar{u} = \{1, \dots, d\} \setminus u$ .

## Decomposition in Fourier Domain

For a subset  $u \subset \{1, \dots, d\} =: \mathcal{D}$  of dimensions we have [Potts, Schmischke 2021]

$$f_u(\mathbf{x}) = f_u(\mathbf{x}_u) = \sum_{\mathbf{k} \in \mathbb{Z}^d: \text{supp}(\mathbf{k})=u} c_{\mathbf{k}}(f) e^{2\pi i \mathbf{k}^\top \mathbf{x}_u}$$

Decomposition of the index set, 3D-visualization:



## Variations and GSI

Consider a trigonometric polynomial

$$f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathcal{K}} \hat{f}_{\mathbf{k}} e^{2\pi i \mathbf{k}^\top \mathbf{x}_j}.$$

We easily see that  $\sigma^2(f) = \sum_{\emptyset \neq u \subseteq \mathcal{D}} \sigma^2(f_u)$ .

Study the importance of subsets  $u$  in terms of the global sensitivity indices (GSI)

$$\rho_u(f) := \frac{\sigma^2(f_u)}{\sigma^2(f)} = \frac{\sum_{\text{supp}(\mathbf{k})=u} |\hat{f}_{\mathbf{k}}|^2}{\sum_{\mathbf{k} \in \mathcal{K} \setminus \{\mathbf{0}\}} |\hat{f}_{\mathbf{k}}|^2} \in [0, 1].$$

Methods

## Goal

Given:  $N$  data points  $\mathbf{x}_j \in \mathbb{T}^d$  or  $[0, 1]^d$  and corresponding values  $y_j \in \mathbb{R}$ , find a model or rather function  $f$  with

$$y_j \approx f(\mathbf{x}_j).$$

We can work efficiently with trigonometric models (nonuniform FFT, short: NFFT), but only in low dimensions (due to the curse of dimensionality).

## ONB Approach

Ansatz:  $y_j \approx f(\mathbf{x}_j) := \sum_{\mathbf{k} \in \mathcal{K}} \hat{f}_{\mathbf{k}} e^{2\pi i \mathbf{k}^\top \mathbf{x}_j}$ , where we assume a

low superposition dimension:  $\mathbf{k} \in \mathcal{K} \iff |\text{supp}(\mathbf{k})| \leq d_s$ .

1. Solve  $\min_{\mathbf{f}} \|\mathbf{y} - \Phi \mathbf{f}\|_2^2 + \lambda \mathbf{f}^* W \mathbf{f}$  iterative (LSQR), where  $\Phi = [e^{2\pi i \mathbf{k}^\top \mathbf{x}_j}]_{\mathbf{k}, j}$  (fast mult. via NFFTs),  $W = \text{diag}(\hat{\omega}_{\mathbf{k}})$  (decay of  $\hat{f}_{\mathbf{k}} \sim$  smoothness of  $f$ )
2. Compute GSI to determine active subsets  $u$ . Re-compute the approximation by only keeping the active subsets. [Potts, Schmischke 2020, 2021]

## Kernel-Based Approach

Search  $f \in \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot), \mathbf{x} \in \mathbb{R}^d\}}$  (RKHS). By the representer theorem, we may solve

$$\min_{f \in \mathcal{H}} \sum_{j=1}^N (y_j - f(\mathbf{x}_j))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\iff \min_{\alpha \in \mathbb{C}^N} \|\mathbf{y} - K\alpha\|_2^2 + \lambda \alpha^* K \alpha,$$

where  $K \in \mathbb{R}^{N \times N}$  is the kernel matrix with entries  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ . Kernel ridge regression (KRR),  $y_j \approx \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}_j)$ .

Results

## Connection: Regularized Least Squares

Consider the case  $\lambda \neq 0$  with kernels of the form  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\mathbf{k} \in \mathcal{K}} \hat{\omega}_{\mathbf{k}}^{-1} e^{2\pi i \mathbf{k}^\top (\mathbf{x}_i - \mathbf{x}_j)}$ .

Then, both approaches are mathematically equivalent [Shawe-Taylor, Cristianini 2004]

$$\Phi^* \Phi \mathbf{f} + \lambda W \mathbf{f} = \Phi^* \mathbf{y} \iff \mathbf{f} = \lambda^{-1} W^{-1} \Phi^* (\mathbf{y} - \Phi \mathbf{f}) = W^{-1} \Phi^* \alpha$$

with

$$\alpha = \lambda^{-1} (\mathbf{y} - \Phi \mathbf{f}) \iff \lambda \alpha = \mathbf{y} - \Phi W^{-1} \Phi^* \alpha \iff \underbrace{(\Phi W^{-1} \Phi^* + \lambda I_N)}_{\mathcal{K}} \alpha = \mathbf{y}.$$

## Comments and Questions

- ▶ Other orthonormal systems analogously, e.g. half-period cosine basis in non-periodic case,  $\mathbf{x}_j \in [0, 1]^d$ .
- ▶ The restriction to small superposition dimensions  $d_s$  is motivated by the sparsity of effects: In practice most phenomena can be described by a few low-dim. interactions. The same applies to sufficiently smooth functions.
- ▶ In the case of  $\lambda = 0$  we need to work **over-/under-** determined.
- ▶ In the dual setting we compute the  $\alpha_j$  in addition. Benefit?

## Results - Least Squares Approach

Tested on real and synthetic data sets (regression and classification).

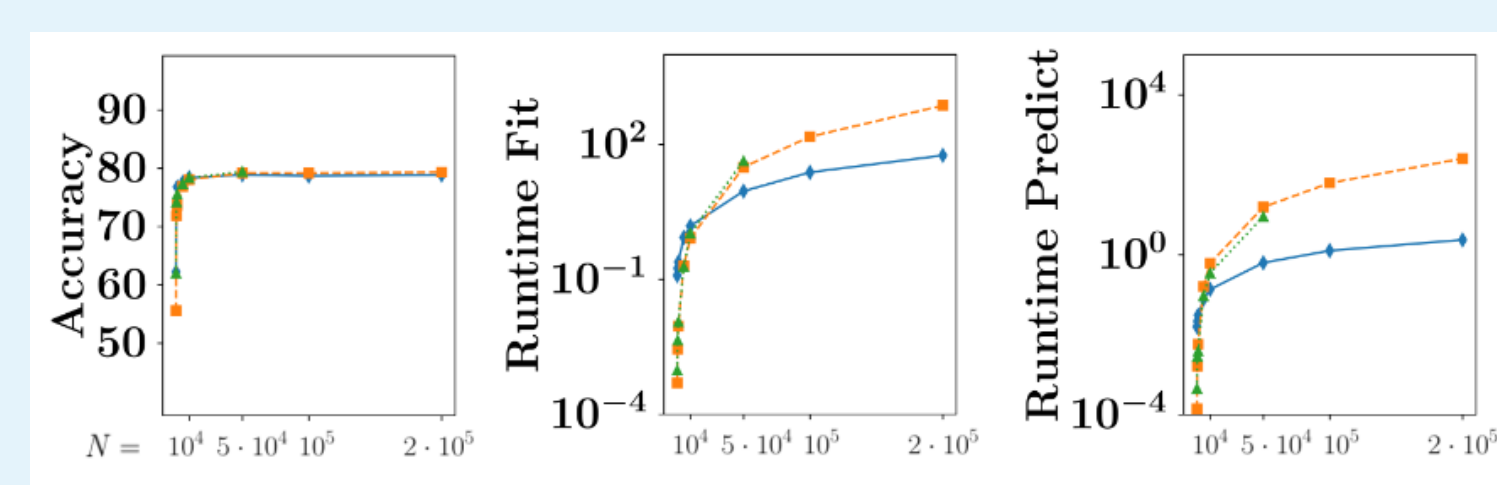
data set	$d$	$N$	error (type)	ref. method	ANOVA-LSQR
Forest Fires	12	517	12.71 (MAD)	SVM	12.65
Energy Eff. Housing	8	768	1.79 (RMSE)	Grad. Boost. Mach.	1.49
Energy Eff. Cooling	8	768	0.48 (RMSE)	Random Forest	0.44
Airfoil Self-Noise	5	1503	0.028 (rel. $\ell_2$ )	Sparse Rand. Features	0.016
California Housing	8	20640	0.115 (RMSE)	Local Learning Reg. NN	0.109
Ailerons	40	13750	0.0460 (RMSE)	Local Learning Reg. NN	0.0457

[Schmischke 2022, PhD thesis], publicly available julia software.

## Results - Kernel Approach

We consider real data sets for binary classification and apply Gaussian ANOVA-like kernels (which we approximate by trig. polynomials)

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{n_{\text{kernels}}} \sum_{|u| \leq d_s} e^{-\|\mathbf{x}_i - \mathbf{x}_j\|_u^2 / \sigma^2}.$$



We compared our NFFT-based KRR with standard sklearn algorithms (KRR and SVM) in python. Example: Results for SUSY data set with  $d = 18$  features. [Nestler, Stoll, Wagner 2022]

Ongoing & Future

## Support Vector Regression

$$\min_{\mathbf{f}} \frac{1}{2} \|\hat{\mathbf{f}}\|_2^2 + C \sum_{j=1}^N |\xi_j|^2 \quad \text{s.t. } |y_j - (\Phi \hat{\mathbf{f}})_j| \leq \epsilon + |\xi_j|$$

## Conditionally p.d. Kernels

$$f(\mathbf{x}) = \sum_{j=1}^N \alpha_j \kappa(\mathbf{x}_j, \mathbf{x}) + \sum_l \beta_l p_l(\mathbf{x}), \quad \begin{pmatrix} K & P \\ P^\top & 0 \end{pmatrix} \cdot \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}$$

## Gaussian Process Regression

Use ANOVA-type kernels  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  as covariance functions in Gaussian processes.

## Software

- ▶ Julia Package ANOVAapprox by M. Schmischke:
  - ▶ <https://github.com/NFFT/ANOVAapprox.jl>
  - ▶ Bases: "per" (Fourier system), "cos" (half-per. cosine), "cheb" (Chebyshev polynomials), "wav1"... "wav4" (Wavelets).
  - ▶ Solvers: "lsqr", "fista" (different regularization, group lasso), "krr" (solves the dual problem with CG,  $K = \Phi W^{-1} \Phi^*$ ).
  - ▶ Interpretability: GSI and attribute ranking.
- ▶ Python code NFFT4ANOVA by T. Wagner:
  - ▶ <https://github.com/wagnertheresa/NFFT4ANOVA>
  - ▶ KRR, where the kernel is a sum of equally weighted low-dimensional Gaussian kernels.

## References

- 1 Franziska Nestler, Martin Stoll and Theresa Wagner. **Learning in High-Dimensional Feature Spaces Using ANOVA-Based Fast Matrix-Vector Multiplication**. arXiv: 2111.10140, 2021. (accepted in Found. Data Sci.)
- 2 Daniel Potts and Michael Schmischke. **Interpretable Approximation of High-Dimensional Data**. SIAM J. Math. Data Sci., 3 (4), 1301–1323, 2021.
- 3 Daniel Potts and Michael Schmischke. **Approximations of high-dimensional periodic functions with Fourier-based methods**. SIAM J. Numer. Anal. 59, 2393–2429, 2021.
- 4 Felix Bartel, Michael Schmischke and Daniel Potts. **Grouped Transformations and Regularization in High-Dimensional Explainable ANOVA Approximation**. SIAM Journal on Scientific Computing, 2021 (accepted).
- 5 Laura Lippert, Daniel Potts and Tino Ullrich. **Fast Hyperbolic Wavelet Regression meets ANOVA**. ArXiv: 2108.13197, 2021.