

Bachelor's Thesis

fANOVA for Interpretable Machine Learning

Department of Statistics
Ludwig-Maximilians-Universität München

Juliet Fleischer

Munich, Month Dayth, Year



Submitted in partial fulfillment of the requirements for the degree of B. Sc.
Supervised by Prof. Dr. Thomas Nagler

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Contents

1	Introduction	1
2	Background Knowledge	5
2.1	Basic Setup	5
2.2	Conditional expectation	5
3	Foundations	8
3.1	Early Work on fANOVA	8
3.2	Modern Work on fANOVA	8
3.3	Formal Introduction to fANOVA	9
4	Generalization	14
4.1	Motivating Example	14
4.2	Formal Introduction to Generalized fANOVA	14
5	Example	15
5.1	Multivariate Normal Inputs	15
6	Simulation Study	17
6.1	Software implementations	17
7	Conclusion	18
8	Mathematical Statements	18
A	Appendix	V
B	Electronic appendix	VI

1 Introduction

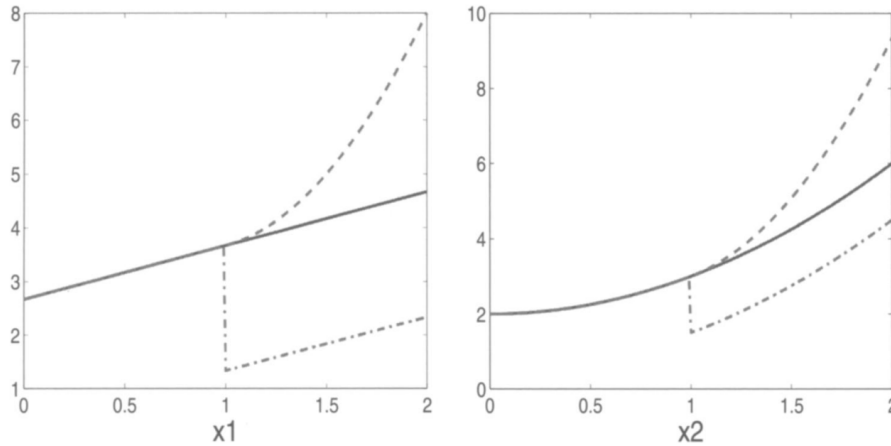


Figure 2. A comparison of Functional ANOVA effects between $F(x_1, x_2)$ (solid) and a learned approximation, $\hat{F}(x_1, x_2)$ (dashed). These two functions are indistinguishable on the data in Figure 1. The left-hand plot provides the effect for x_1 , the right for x_2 . Dotted lines provide the conditional expectations $E(\hat{F}(x_1, x_2)|x_1)$ and $E(\hat{F}(x_1, x_2)|x_2)$ with respect to the measure in Figure 1.

Questions

- I am confused: I was not clear about the distinction between deterministic input and random variables in the definitions and the examples; and now things are mixed up when it comes to notation, as well as the setup. In general I am now confused if we deal with random variables or deterministic inputs and on which basis fANOVA is really constructed?
- I am still confused if setting the zero-mean constraint for $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}[0, 1]$ is essentially saying that we centre the distribution and now assume $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}[-1, 1]$. So can we, instead of explicitly stating the zero-mean constraint just assume $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}[-1, 1]$? And following the same principle we would shift other distributions by altering their parameters, not by explicitly stating the zero-mean constraint? Then for the standard normal distribution we wouldn't need to do anything, it is already centred around 0. For other distributions we would need to change, and for some it doesn't make inhaltlichen Sinn e.g. Poisson distribution?
-
- How is the generalized fANOVA really computed? Could you compute it by hand?
- Effect of non-linear main effects in classical fANOVA (see "General_fANOVA_handnote")

- fANOVA decomposition via the integral, how would the zero mean constraint look here? (see “General fANOVA handnotes”)
- Can you reconstruct the function from only the fANOVA terms? I think it can be reconstructed only if variables are independent, have zero-mean, are orthogonal?
- Is it possible to perform fANOVA for non-square-integrable functions? I think in general yes but the variance decomposition doesn’t work then or might have problems.
- fANOVA decomposition for discrete variables possible? Does it make sense even?
- Connection between the (conditional) expected value, (partial) integral, projections (section ??)?
- In the hierarchical orthogonality condition (4.2) formulated in Hooker (2007) for the generalized fANOVA framework, shouldn’t we explicitly exclude the case that $v = u$, because then, we would require that the inner product of the fANOVA component is zero wouldn’t we (section 4)?
- I am a bit confused by Figure 2 in Hooker (2007) (see section 1), especially by the dotted line for conditional expectation. What should it tell us? I think what the dashed line (learned approximation) shows is that the model estimates a non-linear effect for x_1 , even though the true effect is linear. The reasons are the problematic data points in the top right region.
- Why is it a problem, when explainability methods also place large emphasis on regions of low probability mass when dependencies between variables exist - because in the end explainability is about explaining the model, not the data generating process; and after all it is how the model works in these regions. [But as the Hooker example illustrates, how the model works and what it estimates in these regions is wrong and then it’s better to not report any model behaviour or come closer to the DGP than to give wrong estimations?]
-
- Use of AI tools?
- Do we need to restrict ourselves to the unit hypercube? Or does fANOVA decomposition work in general, but maybe with some constraints? Originally it was constructed for models on the unit hypercube $[0, 1]$, but other papers also use models

from R^d Generally no restriction, so next step could be to generalize, to \mathbb{R}^n , other measures, dependent variables

- Still unclear: Are the terms fully orthogonal or hierarchically? See subsection on Orthogonality of the fANOVA terms (especially the example) I think in the original fANOVA decomposition the terms are orthogonal but in the generalized fANOVA (Hooker, 2007) they are hierarchically orthogonal. *fully orthogonal when independence assumption, probably partially when no independence*
- x_1, \dots, x_k are simply the standardized features, right? *Yes*
- **My current understanding:** we need independence of x_1, \dots, x_k so that fANOVA decomposition is unique (and orthogonality holds). We need zero-mean constraint for the orthogonality of the components. We need orthogonality for the variance decomposition. *zero-mean \rightarrow orthogonality \rightarrow uniqueness; Lemma 1 in Hooker 2007 ist verallgemeinert ds zero-mean constraint*
- Next step might be to investigate the (mathematical) parallels of fANOVA decomposition and other IML methods (PDP, ALE, SHAP), e.g. there is definitely a strong relationship between Partial dependence (PD) and fANOVA terms, and PD is itself again related to other IML methods; Also look how are other IML models studied and study fANOVA in a similar way (e.g. other IML methods are defined, checked for certain properties, examined under different conditions (dependent features, independent features) etc.) (see dissertation by Christoph Molnar for this); Also I would be very interested in investigating the game theory paper further (Fumagalli et al., 2025) but still a bit unsure if it is too complex.
- Why does a fANOVA decomposition of a simple GAM not lead to the “true” coefficients? <https://christophm.github.io/interpretable-ml-book/decomposition.html> talks about this a bit in the subchapter “Statistical regression models” *It should actually lead to the GAM; at least under all the constraint like zero-mean constraint and orthogonality*
-
- In Hooker (2004) they work with $F(x)$ and $f(x)$, but in Sobol (2001) they only work with $f(x)$. I think this is only notation? *Only notation.*
- Does orthogonality in fANOVA context mean that all terms are orthogonal to each other? Or that a term is orthogonal to all lower-order terms (“Hierarchical orthogonality”)? *The terms are hierarchically orthogonal, so each term is orthogonal to*

all lower-order terms, but not to the same-order terms! So f_1 is not necessarily orthogonal to f_2 but it is orthogonal to f_{12} , f_0 .

- Do the projections here serve as approximations? (linalg skript 2024 5.7.4 Projektionen als beste Annäherung) *Yes, they can be interpreted as sort of approximation.*
- Which sub-space are we exactly projecting onto? Are the projections orthogonal by construction (orthogonal projections) or only when the zero-mean constraint is set? *The subspace we project onto depends on the component. For f_0 we project onto the subspace of constant functions, for f_1 we project onto the subspace of all functions that involve x_1 and have an expected value of 0 (zero-mean constraint to ensure orthogonality). It depends on the formulation of the fANOVA decomposition if you need to explicitly set the zero-mean constraint for orthogonality or if it is met by construction.*
- How “far” should I go back, formally introduce L^2 space, etc. or assume that the reader is familiar with it? *Yes, space, the inner product on this space should be formally introduced.*

2 Background Knowledge

2.1 Basic Setup

Let $(\Omega, \mathcal{F}, \nu)$ be a measure space, where Ω is a sample space, \mathcal{F} is a σ -algebra on Ω and $\nu : \mathcal{F} \rightarrow [0, 1]$ is a probability measure. \mathcal{B}^N is the Borel σ -algebra on \mathbb{R}^N , $N \in \mathbb{N}$. $\mathbf{X} = (X_1, \dots, X_N) : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^N, \mathcal{B}^N)$ denotes a \mathbb{R}^N -valued random vector.

We assume that the probability distribution of \mathbf{X} is continuous and completely defined by the joint probability density function $f_{\mathbf{X}} : \mathbb{R}^N \rightarrow \mathbb{R}_0^+$.

Let u denote a subset of indices $\{1, \dots, N\}$, and $-u := \{1, \dots, N\} \setminus u$ its complement. $\mathbf{X}_u = (X_1, \dots, X_{|u|})$, $u \neq \emptyset$, $1 \leq i_1 < \dots < i_{|u|} \leq N$ is a subvector of \mathbf{X} and $\mathbf{X}_{-u} = \mathbf{X}_{\{1, \dots, N\} \setminus u}$ is the complement of \mathbf{X}_u .

The marginal density function is $f_u(\mathbf{x}_u) := \int f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}_{-u}$ for a given set $\emptyset \neq u \subseteq \{1, \dots, N\}$. $f(\mathbf{X}) := f(X_1, \dots, X_N)$ is a mathematical model with random variables as inputs. We write a vector space of square-integrable functions as

$$\mathcal{L}^2(\Omega, \mathcal{F}, \nu) = \{f : \Omega \rightarrow \mathbb{R} \text{ s.t. } \mathbb{E}[f^2(\mathbf{X})] < \infty\}$$

$\mathcal{L}^2(\Omega, \mathcal{F}, \nu)$ is a Hilbert space with the inner product defined as:

$$\langle f, g \rangle = \int f(x)g(x) d\nu(x) = \mathbb{E}[fg], \quad \forall f, g \in \mathcal{L}^2.$$

The norm is then defined as:

$$\|f\| = \sqrt{\langle f, f \rangle} = \sqrt{\int f^2(x) d\nu(x)} = \mathbb{E}[f^2], \quad \forall f \in \mathcal{L}^2.$$

Which resource should I cite for these “general” definitions? e.g. <https://apachepersonal.miun.se/andrli/Bok.pdf>?

2.2 Conditional expectation

In general, we define the conditional expectation of a vector of random variables $\mathbf{X} = (X_1, X_2)$ as follows:

$$\mathbb{E}[g(X_1, X_2) \mid X_1 = x_1] = \int g(x_1, s_2) p_{X_2|X_1}(s_2 \mid x_1) ds_2.$$

Only when X_1 and X_2 are independent can we write

$$\mathbb{E}[g(X_1, X_2) \mid X_1 = x_1] = \int g(x_1, s_2) p_{X_2|X_1}(s_2 \mid x_1) ds_2 = \int g(x_1, s_2) p_{X_2}(s_2) ds_2 = \mathbb{E}_{X_2}[g(x_1, X_2)].$$

Extended to n random variables it looks as follows. Without loss of generality, we condition on $X_1 = x_1$:

$$\begin{aligned} \mathbb{E}[g(X_1, \dots, X_n) \mid X_1 = x_1] &= \int g(x_1, s_2, \dots, s_n) p_{X_2, \dots, X_n|X_1}(s_2, \dots, s_n \mid x_1) ds_2 \dots ds_n \\ &= \int g(x_1, s_2, \dots, s_n) p_{X_2}(s_2, \dots, s_n) ds_2 \dots ds_n \\ &= \mathbb{E}_{X_2, \dots, X_n}[g(x_1, X_2, \dots, X_n)] \end{aligned}$$

Orthogonal projection

Let $\mathcal{G} \subset \mathcal{L}^2$ denote a linear subspace. The projection of f onto \mathcal{G} is defined by the function $\Pi_{\mathcal{G}}f$ which minimizes the distance to f in \mathcal{L}^2 :

$$\Pi_{\mathcal{G}}f = \arg \min_{g \in \mathcal{G}} \|f - g\|^2 d\nu = \arg \min_{g \in \mathcal{G}} \mathbb{E}[(f - g)^2].$$

I think this is closely related to Hilbert projection theorem?

Definition of \mathcal{L}^2 space and projection modified from <https://tnagler.github.io/mathstat-lmu-2024.pdf>.

Properties of the Multivariate Normal Distribution

Let $\mathbf{X} = (X_1, \dots, X_d)^\top \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be a d -dimensional multivariate normal (MVN) random vector, where $\boldsymbol{\mu} \in \mathbb{R}^d$ is the mean vector and $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ is the symmetric positive semi-definite covariance matrix.

The marginal distribution of X_i is generally given by an univariate normal distribution:

$$X_i \sim \mathcal{N}(\mu_i, \Sigma_{ii}) \quad \text{for all } i = 1, \dots, d.$$

If we condition on a subset of the variables, we can also make statements about the conditional distribution. For this we partition the random vector \mathbf{X} into two parts, \mathbf{X}_A and \mathbf{X}_B , where \mathbf{X}_A contains the variables we condition on and \mathbf{X}_B contains the remaining variables. The joint distribution of \mathbf{X} can be expressed as:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_A \\ \mathbf{X}_B \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{AA} & \boldsymbol{\Sigma}_{AB} \\ \boldsymbol{\Sigma}_{BA} & \boldsymbol{\Sigma}_{BB} \end{pmatrix} \right).$$

The conditional distribution of \mathbf{X}_B given $\mathbf{X}_A = \mathbf{x}_A$ is

$$\mathbf{X}_B \mid \mathbf{X}_A = \mathbf{x}_A \sim \mathcal{N}(\boldsymbol{\mu}_B + \boldsymbol{\Sigma}_{BA}\boldsymbol{\Sigma}_{AA}^{-1}(\mathbf{x}_A - \boldsymbol{\mu}_A), \boldsymbol{\Sigma}_{BB} - \boldsymbol{\Sigma}_{BA}\boldsymbol{\Sigma}_{AA}^{-1}\boldsymbol{\Sigma}_{AB}).$$

For normally distributed random variables, we also know that $\text{Cov}(X_i, X_j) = 0$, implies $X_i \perp X_j$. Lastly, for any real vector $\mathbf{a} \in \mathbb{R}^d$, the linear combination $\mathbf{a}^\top \mathbf{X}$ is normally distributed:

$$\mathbf{a}^\top \mathbf{X} \sim \mathcal{N}(\mathbf{a}^\top \boldsymbol{\mu}, \mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{a}).$$

[Find official resource for these properties; review how this can be shown.](#)

3 Foundations

3.1 Early Work on fANOVA

The main idea of the fANOVA decomposition is to decompose a statistical model into the sum of the main effects and interaction effects of its input variables. The underlying principle of fANOVA decomposition dates back to Hoeffding (1948). In his famous paper he introduced U-statistics, along with the “Hoeffding decomposition”, which allows to write a symmetric function of the data as a sum of orthogonal components. Sobol (1993) used the same principle and applied it to deterministic mathematical models. He built on the originally called “decomposition into summands of different dimension” in Sobol (2001), where he introduces Sobol indices and renames the method to the “ANOVA-representation”. Sobol indices are now commonly used in sensitivity analysis. Efron and Stein (1981) use the idea of the decomposition to prove their famous lemma on jackknife variances. Stone (1994) mainly uses fANOVA decomposition to base smooth regression models with interactions on it and his paper is the building block for a broader body of work of fANOVA-based models [example citations needed](#).

3.2 Modern Work on fANOVA

The fANOVA decomposition has a long history with roots in mathematical statistics and non-parametric estimation theory. In more recent years, the method has been rediscovered by the machine-learning community, especially in the context of interpretable machine learning (IML) and explainable AI (XAI). Hooker (2004) introduces the fANOVA decomposition with the goal of providing a global explanation method for black-box models. Since the assumptions of independent variables in classical fANOVA is often too restrictive in practice, Hooker (2007) generalizes the method to dependent variables. A recent paper by Il Idrissi et al. (2025) can be seen as another approach to generalize the principle of fANOVA decomposition to dependent inputs.

There are specific domains of statistics, such as geostatistics, that explicitly build models on fANOVA framework (see Muehlenstaedt et al. (2012) for fANOVA Kriging models). And recent work discovered interesting mathematical parallels between fANOVA and other IML methods, such as PDP Friedman (2001), or Shapley values (Fumagalli et al. (2025), Herren, Owen preprint).

Liu and Owen (2006) use of fANOVA and sensitivity analysis for functions arising in computational finance. Owen (2013) formal intro to fANOVA decomposition and generalization of Sobol indices. Owen has generally a lot of work related to fANOVA decomposition, either lecture notes explaining the decomposition, methods based on it Owen (2003), or

deeper into sensitivity analysis and fANOVA Owen (2013).

fANOVA and U-statistics, fANOVA and sensitivity analysis, fANOVA and GAMs (with interactions)

3.3 Formal Introduction to fANOVA

This chapter is based on the formal introductions by Rahman (2014), Sobol (1993, 2001), Hooker (2004), Owen (2013), Muehlenstaedt et al. (2012). We show both formulations of the fANOVA, via the integral and via the expected value and in general prefer the expected value formulation as it is more intuitive in a probabilistic setting. Originally, Sobol (1993) presented the fANOVA decomposition with independent input variables with support bounded to the unit interval, i.e. he considered the measure space $([0, 1]^n, \mathcal{B}([0, 1]^n), \nu)$. Later work shows that this restriction is not necessary, and we can work with the Borel σ -algebra on the n -dimensional real number line, i.e. $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \nu)$, and with a general measure ν defined on it (see e.g. Rahman (2014)). Since we assume independence of the input variables, their joint distribution is given by the product over the marginal distributions, i.e. $f_{\mathbf{X}}(\mathbf{x}) d\nu(\mathbf{x}) = \prod_{i=1}^N f_{X_i}(x_i) d\nu(x_i)$. $f_{X_i} : \mathbb{R} \rightarrow \mathbb{R}_0^+$ is the marginal probability density function of X_i defined on $(\Omega_i, \mathcal{F}_i, \nu_i)$.

Definition 3.1. *Let $y(\mathbf{X})$ be a mathematical model with realizations of independent random variables x_1, \dots, x_N as input. We can represent such a model y as the hierarchical sum of specific basis functions with increasing dimensionality:*

$$y(\mathbf{X}) = \sum_{u \subseteq \{1, \dots, N\}} y_u(\mathbf{X}_u), \quad (1)$$

If $|u| = 0$ it describes the constant term, if $|u| = 1$ it describes the main effects, if $|u| > 1$ it describes the interaction effects of the variables in u . The expansion consists of 2^N terms.

The fANOVA terms should be constructed in such a way that they have two specific properties crucial for identifiability and interpretation.

Proposition 3.1. *The zero-mean property states that all effects, except for the constant terms, are centred around zero. Mathematically this means that the effects integrate to zero w.r.t. their own variables:*

$$\int y_u(\mathbf{x}_u) f_{\mathbf{X}_u}(\mathbf{x}_u) d\nu(\mathbf{x}_u) = \mathbb{E}[y_u(\mathbf{X}_u)] = 0 \quad (2)$$

Proposition 3.2. *The second property is the orthogonality of the fANOVA terms. If two sets of indices are not completely equivalent, i.e. $\emptyset \neq u \subseteq \{1, \dots, N\}, \emptyset \neq v \subseteq$*

$\{1, \dots, N\}$, and $u \neq v$, then it holds that their fANOVA terms are orthogonal to each other:

$$\int y_u(\mathbf{x}_u) y_v(\mathbf{x}_v) f_{\mathbf{X}}(\mathbf{x}) d\nu(\mathbf{x}) = \mathbb{E}[y_u(\mathbf{X}_u) y_v(\mathbf{X}_v)] = 0 \quad (3)$$

This means that fANOVA terms are “fully orthogonal” to each other, meaning not only terms of different order are orthogonal to each other but also terms of the same order are. Rahman (2014) derives these two properties (Equation 2, Equation 3) from a more general condition, he calls the “strong annihilating conditions”.

The strong annihilating conditions require that the fANOVA terms integrate to zero w.r.t the individual variables contained in u and weighted by the individual marginal probability density functions:

$$\int y_u(\mathbf{x}_u) f_{X_i}(x_i) d\nu(x_i) = 0, \quad \text{for } i \in u \neq \emptyset. \quad (4)$$

We can reassure ourselves that the properties in fact follow from the strong annihilating conditions. For the zero-mean constraint we can write:

$$\begin{aligned} \mathbb{E}[y_u(\mathbf{X}_u)] &= \int_{\mathbb{R}^{|u|}} y_u(\mathbf{x}_u) f_{\mathbf{X}_u}(\mathbf{x}_u) d\nu(\mathbf{x}_u) \\ &= \int_{\mathbb{R}^{|u|}} y_u(\mathbf{x}_u) \prod_{i \in u} f_{X_i}(x_i) d\nu(\mathbf{x}_u) \\ &= \int_{\mathbb{R}^{|u|-1}} \int_{\mathbb{R}} y_u(\mathbf{x}_u) f_{X_i}(x_i) dx_u \prod_{j \in u, j \neq i} f_{X_j}(x_j) = 0 \end{aligned}$$

One can follow the same reasoning for the orthogonality condition:

$$\begin{aligned} \mathbb{E}[y_u(\mathbf{X}_u) y_v(\mathbf{X}_v)] &= \int_{\mathbb{R}^N} y_u(\mathbf{x}_u) y_v(\mathbf{x}_v) f_{\mathbf{X}}(\mathbf{x}) d\nu(\mathbf{x}) \\ &= \int_{\mathbb{R}^N} y_u(\mathbf{x}_u) y_v(\mathbf{x}_v) \prod_{i=1}^N f_{X_i}(x_i) d\nu(x_i) \\ &= \int_{\mathbb{R}^{N-1}} \int_{\mathbb{R}} y_u(\mathbf{x}_u) y_v(\mathbf{x}_v) f_{X_i}(x_i) dx_u \prod_{j \in \{1, \dots, N\}, j \neq i} f_{X_j}(x_j) = 0 \end{aligned}$$

Construction of fANOVA terms

The individual fANOVA term for the variables with indices in u are constructed from integrating the original function $y(\mathbf{X})$ w.r.t all variables except for the ones in u , and subtracting the lower order terms. Intuitively the integral is averaging the original function over all other variables except the ones of interest, which makes sense as we are then left with a function of the variables of interest only. Subtracting lower order terms

corresponds to accounting for effects that are already explained by other variables or interactions so that we obtain the isolated effects.

Since $u = \emptyset$ for the constant term, we integrate w.r.t all variables:

$$y_{\emptyset} = \int y(\mathbf{x}) \prod_{i=1}^N f_{X_i}(x_i) d\nu(x_i) = \mathbb{E}[y(\mathbf{X})]. \quad (5)$$

For all other effects $u \neq \emptyset$ we can write:

$$y_u(\mathbf{X}_u) = \int y(\mathbf{X}_u, \mathbf{x}_u) \prod_{i=1, i \notin u}^N f_{X_i}(x_i) d\nu(x_i) - \sum_{v \subsetneq u} y_v(\mathbf{X}_v), \quad (6)$$

The basis components offer a clear interpretation of the model, decomposing it into main effects, two-way interaction effects, and so on. This is why fANOVA decomposition has received increasing attention in the IML and XAI literature, holding the potential for a global explanation method of black box models.

Second-moment statistics

We already established that $\mathbb{E}[y(\mathbf{X})] = y_{\emptyset}$. For the variance of $y(\mathbf{X})$, we find that the total variance can be decomposed into the sum of the fANOVA term variances. The variance decomposition is a major result in Sobol (1993) and forms the basis for the Sobol indices in sensitivity analysis. We sketch the variance decomposition here and note that it is only possible under independence assumption.

If $y \in \mathcal{L}^2$, then $y_{i_1, \dots, i_n} \in \mathcal{L}^2$ [proof? reference?; Sobol 1993 says it is easy to show using Schwarz inequality and the definition of the single fANOVA terms](#). Therefore, we define the variance of f as follows:

$$\begin{aligned} \sigma^2 &:= \int y^2(\mathbf{X}) f_{\mathbf{X}}(\mathbf{x}) d\nu(\mathbf{x}) - (y_{\emptyset})^2 \\ &= \int y^2(\mathbf{X}) f_{\mathbf{X}}(\mathbf{x}) d\nu(\mathbf{x}) - \left(\int y(\mathbf{X}) f_{\mathbf{X}}(\mathbf{x}) d\nu(\mathbf{x}) \right)^2 \\ &= \mathbb{E}[y^2(\mathbf{X})] - \mathbb{E}[y(\mathbf{X})]^2 \end{aligned}$$

The variance of the fANOVA components is then defined as

$$\begin{aligned} \sigma_{x_{i_1}, \dots, x_{i_n}}^2 &= \int \cdots \int y_{i_1, \dots, i_n}^2 f_{\mathbf{X}}(\mathbf{x}) d\nu(x_1) \cdots d\nu(x_n) - \left(\int \cdots \int y_{i_1, \dots, i_n} f_{\mathbf{X}}(\mathbf{x}) d\nu(x_1) \cdots d\nu(x_n) \right)^2 \\ &= \mathbb{E}[y_{i_1, \dots, i_n}^2] - \mathbb{E}[y_{i_1, \dots, i_n}]^2 \end{aligned}$$

Because of the orthogonality property, the second term vanished and we get:

$$\begin{aligned}\sigma_{x_{i_1}, \dots, x_{i_n}}^2 &= \int \cdots \int y_{i_1, \dots, i_n}^2 f_{\mathbf{X}}(\mathbf{x}) d\nu(x_1) \cdots d\nu(x_n) \\ &= \mathbb{E}[y_{i_1, \dots, i_n}^2]\end{aligned}$$

With the definition of the total variance σ^2 and the component-wise variance $\sigma_{x_{i_1}, \dots, x_{i_n}}^2$ we can now see that the total variance can be decomposed into the sum of the component-wise variances.

Alternatively we can formulate this via the expected value. We write the sum over u for the sum over $\emptyset \neq u \subseteq \{1, \dots, N\}$ and the sum over $u \neq v$ for the sum over $\emptyset \neq u \subseteq \{1, \dots, N\}, \emptyset \neq v \subseteq \{1, \dots, N\}, u \neq v$.

$$\begin{aligned}\sigma^2 &:= \mathbb{E}[(y(\mathbf{X}) - \mu)^2] = \mathbb{E}[(y_\emptyset + \sum_u y_u(\mathbf{X}_u) - y_\emptyset)^2] \\ &= \mathbb{E}[(\sum_u y_u(\mathbf{X}_u))^2] \\ &= \mathbb{E}[\sum_u y_u^2(\mathbf{X}_u)] + 2\mathbb{E}[\sum_{u \neq v} y_u(\mathbf{X}_u) y_v(\mathbf{X}_v)] \\ &= \sum_u \mathbb{E}[y_u^2(\mathbf{X}_u)]\end{aligned}$$

fANOVA as projection

Referring to the general connection between the expected value and orthogonal projections presented in section 2, the fANOVA terms can also be understood from a viewpoint of projections. This will also help to understand the generalization of fANOVA in section 4. y_\emptyset is the projections of the original function y onto the space of all constant functions $\mathcal{G}_0 = \{g(x) = a; a \in \mathbb{R}\}$. It is an unconditional expected value and the best approximation of y given a constant function:

$$\begin{aligned}\Pi_{\mathcal{G}_0} y &= \arg \min_{g \in \mathcal{G}_0} \|y(\mathbf{X}) - g\|^2 \\ &= \arg \min_{g \in \mathcal{G}_0} \mathbb{E}[(y(\mathbf{X}) - g)^2] \\ &= \mathbb{E}[y(\mathbf{X})]\end{aligned}$$

The main effect $f_i(x_i)$ is the projection of y onto the subspace of all functions that only depend on x_i and have an expected value of zero while accounting for the lower-order

effects. The subspace we project onto is $\mathcal{G}_i = \{g(x) = g_i(x_i); \int g(x)d\nu(x_i) = 0\}$.

$$\begin{aligned}\Pi_{\mathcal{G}_i}y - y_0 &= \arg \min_{g \in \mathcal{G}_i} \|y(x) - g(x_i)\|^2 - f_0 \\ &= \arg \min_{g \in \mathcal{G}_i} \mathbb{E}_{-x_i}[(f(x) - g(x_i))^2] - \mathbb{E}[y(x)] \\ &= \mathbb{E}_{-x_i}[y(X_1, \dots, x_i, \dots, X_n)] - \mathbb{E}[y(X)]\end{aligned}$$

The two-way interaction effect $y_{ij}(x_i, x_j)$ is the projection of y onto the subspace of all functions that depend on x_i and x_j and have an expected value of zero in each of it's single components, i.e. $\mathcal{G}_{i,j} = \{g(x) = g_{ij}(x_i, x_j); \int g(x)d\nu(x_i) = 0 \wedge \int g(x)d\nu(x_j) = 0\}$. Again, we account for lower-order effects by subtracting the constant term and all main effects:

$$\begin{aligned}\Pi_{\mathcal{G}_{ij}}y - y_0 - f_1(x_i) - \dots &= \arg \min_{g \in \mathcal{G}_{ij}} \|y(x) - g(x_i, x_j)\|^2 - f_0 - f_1(x_i) - \dots \\ &= \arg \min_{g \in \mathcal{G}_{ij}} \mathbb{E}_{-x_i, -x_j}[(y(x) - g(x_i, x_j))^2] - \mathbb{E}[y(x)] - \mathbb{E}_{-x_i}[y(x)] \\ &= \mathbb{E}_{-x_i, -x_j}[y(X_1, \dots, x_i, x_j, \dots, X_n)] - \mathbb{E}[y(x)] - \mathbb{E}_{-x_i}[y(X)]\end{aligned}$$

I think Hilbert space theorem tells us that the orthogonal projection minimizes the squared difference in a Hilbert space? So the projection is the solution to the minimization problem that wants to minimize the squared differences between two elements of the vector space. This would be the first equality. The last equality that the solution is equal to the (conditional) expected value also has to be shown, still have to look which theorem this is proven by.

In general, we can write:

$$y_u(x) = \Pi_{\mathcal{G}_u}y - \sum_{v \subsetneq u} f_v(x) \tag{7}$$

We project y onto the subspace spanned by the own terms of the fANOVA component to be defined, while accounting for all lower-order terms.

4 Generalization

4.1 Motivating Example

4.2 Formal Introduction to Generalized fANOVA

In practice assuming independent input variables is often not realistic. When this assumption is violated, it can lead to misleading results (Hooker, 2007). The need for a generalization to dependent variables is evident. We base this chapter mainly on the generalization of Rahman (2014), while there exists other work from Hooker (2007) or Chastaing et al. (2012).

Letting go of the independence assumption means that we no longer work with a product-type probability measure. $f_{\mathbf{X}} : \mathbb{R} \rightarrow \mathbb{R}_0^+$ denotes an arbitrary probability density function and $f_{\mathbf{X}_u} : \mathbb{R}^u \rightarrow \mathbb{R}_0^+$ the marginal probability density function of the subset of variables $u \subseteq d$.

The definition of the fANOVA decomposition, and the two main properties - zero mean and orthogonality - can be stated in the same way as for the classical case, with a slight change for the orthogonality. We build now on the weak annihilating conditions instead of the strong ones, which results in hierarchical orthogonality instead of full orthogonality.

Definition 4.1. *Generalized fANOVA decomposition.*

Weak annihilating conditions.

From this we derive

Proposition 4.1.

Proposition 4.2.

5 Example

5.1 Multivariate Normal Inputs

Consider the function $g = a + X_1 + 2X_2 + X_1X_2$. We assume that $\mathbf{X} = (X_1, X_2)^T$ follows a MVN distribution:

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 \\ \rho_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}\right)$$

where we assume that $\sigma_1 = \sigma_2 = 1$, while we will vary the values for ρ_{12} to distinguish between varying degrees of dependence between the two variables. From the properties of the MVN, we know that marginal distributions are standard normal:

$$X_i \sim \mathcal{N}(0, \sigma_i) \quad \text{for } i = 1, 2 \quad \text{and} \quad \sigma_i = 1$$

We also know that the conditional distributions are:

$$X_1 \mid X_2 = x_2 \sim \mathcal{N}(\rho_{12}x_2, 1 - \rho_{12}^2), \quad X_2 \mid X_1 = x_1 \sim \mathcal{N}(\rho_{12}x_1, 1 - \rho_{12}^2)$$

Case 1: Independent Inputs

Computing the fANOVA decomposition of $g(x_1, x_2)$ by hand, we start with the constant term and make use of formulation via the expected value instead of the integral for notational simplicity:

$$f_0 = \mathbb{E}[g_1(X_1, X_2)] = \mathbb{E}[a + X_1 + 2X_2 + X_1X_2] = \mathbb{E}[a] + \mathbb{E}[X_1] + 2\mathbb{E}[X_2] + \mathbb{E}[X_1X_2]$$

Making use of the independence assumption of x_1 and x_2 , the last term can be written as the product of the expected values. Additionally, given the zero-mean constraint, all terms, except for the constant, vanish and we obtain:

$$f_0 = \mathbb{E}[a] + \mathbb{E}[X_1] + 2\mathbb{E}[X_2] + \mathbb{E}[X_1]\mathbb{E}[X_2] = a$$

Under zero-mean constraint and independence, the main effects and the interaction effect can be computed as follows:

$$\begin{aligned}
f_1(x_1) &= \mathbb{E}_{X_2}[g_1(x_1, X_2)] - f_0 \\
&= \mathbb{E}_{X_2}[a + x_1 + 2X_2 + x_1X_2] - a \\
&= x_1 + 2\mathbb{E}[X_2] + x_1\mathbb{E}[X_2] = x_1 \\
f_2(x_2) &= \mathbb{E}_{X_1}[g_1(X_1, x_2)] - f_0 \\
&= \mathbb{E}_{X_1}[a + X_1 + 2x_2 + X_1x_2] - a \\
&= \mathbb{E}_{X_1}[X_1] + 2x_2 + x_2\mathbb{E}_{X_1}[X_1] = 2x_2 \\
f_{12}(x_1, x_2) &= \mathbb{E}[g_1(x_1, x_2)] - f_0 - f_1(x_1) - f_2(x_2) \\
&= a + x_1 + 2x_2 + x_1x_2 - a - x_1 - 2x_2 = x_1x_2
\end{aligned}$$

It comes as no surprise that in this simple case the fANOVA decomposition does not provide any additional insights, as the isolated effects can be directly seen from the function. In case of independence the generalized fANOVA decomposition is of course equal to the classical one. We show this simple example nevertheless to illustrate at which step which assumption is used. This motivates the need for a generalization of the fANOVA decomposition to dependent variables.

Case 2: Dependent Inputs (weak)

Now assume $\rho_{12} = 0.2$, $\sigma_1 = \sigma_2 = 1$. The constant term f_0 is given by:

$$\begin{aligned}
f_0 &= \mathbb{E}[g(X_1, X_2)] = a + \mathbb{E}[X_1] + 2\mathbb{E}[X_2] + \mathbb{E}[X_1X_2] \\
&= a + \mathbb{E}[X_1X_2] = a + (\text{Cov}(X_1, X_2) + \mathbb{E}[X_1]\mathbb{E}[X_2]) \\
&= a + \rho_{12}
\end{aligned}$$

The main effects can be computed as:

$$\begin{aligned}
f_1(x_1) &= \mathbb{E}[g(X_1, X_2) | X_1 = x_1] - f_0 \\
&= \mathbb{E}[a + x_1 + 2X_2 + x_1X_2 | X_1 = x_1] - (a + \rho_{12}) \\
&= a + x_1 + 2\mathbb{E}[X_2 | X_1 = x_1] + x_1\mathbb{E}[X_2 | X_1 = x_1] - a - \rho_{12} \\
&= x_1 - \rho_{12} + \rho_{12}(2 + x_1) \\
f_2(x_2) &= \mathbb{E}[g(X_1, X_2) | X_2 = x_2] - f_0 \\
&= \mathbb{E}[a + X_1 + 2x_2 + X_1x_2 | X_2 = x_2] - (a + \rho_{12}) \\
&= a + 2x_2 + x_2\mathbb{E}[X_1 | X_2 = x_2] - a - \rho_{12} \\
&= 2x_2 + x_2^2\rho_{12} + \rho_{12}
\end{aligned}$$

Finally, the interaction effect $f_{12}(x_1, x_2)$ is given by:

$$\begin{aligned}
f_{12}(x_1, x_2) &= g(x_1, x_2) - f_0 - f_1(x_1) - f_2(x_2) \\
&= a + x_1 + 2x_2 + x_1x_2 - (a + \rho_{12}) - (x_1 - \rho_{12} + 2\rho_{12}x_1 + \rho_{12}x_1^2) - (2x_2 + x_2^2\rho_{12} + \rho_{12}) \\
&= x_1x_2 - 2\rho_{12}x_1 - \rho_{12}x_1^2 - \rho_{12}x_2^2 - \rho_{12}
\end{aligned}$$

Case 3: Dependent Inputs (strong)

6 Simulation Study

6.1 Software implementations

- Suitable but currently problems installing locally: fanova
- Context of kriging models; create own graphs (not super informative): fanovaGraph
- mlr3 function
- tntorch
- shapley values implementation python

7 Conclusion

8 Mathematical Statements

Square Integrability of $f_1(x_1)$

For now we want to show that the single fANOVA term $f_1(x_1)$ is square integrable, given that the original function $f(x) \in \mathcal{L}^2$. We need to show that:

$$\int |f_1(x_1)|^2 dx_1 < \infty$$

The single fANOVA term is defined as:

$$f_1(x_1) = \int f(x) dx_{-1} - f_0$$

We take the squared norm, and integrate w.r.t. x_1 to use the Cauchy-Schwarz inequality:

$$\begin{aligned} \int |f_1(x_1)|^2 dx_1 &= \int \left| \int f(x) dx_{-1} - f_0 \right|^2 dx_1 \\ &= \int \left| \left(\int f(x) dx_{-1} \right)^2 - 2 \int f(x) dx_{-1} f_0 + f_0^2 \right| dx_1 \end{aligned}$$

Break this into three terms:

$$(1) : \quad \int \left| \int f(x) dx_{-1} \right|^2 dx_1 \leq \int \left(\int 1^2 dx_{-1} \right) \left(\int |f(x)|^2 dx_{-1} \right) dx_1 = \int |f(x)|^2 dx < \infty$$

$$(2) : \quad 2 \int \left(\int f(x) dx_{-1} \right) f_0 dx_1 = 2f_0 \int \left(\int f(x) dx_{-1} \right) dx_1 = 2f_0^2 < \infty$$

$$(3) : \quad \int f_0^2 dx_1 = f_0^2 < \infty$$

Since each term (1)–(3) is finite, and $\int |f_1(x_1)|^2 dx_1$ is a linear combination of them: $\int |f_1(x_1)|^2 dx_1 < \infty$

A Appendix

B Electronic appendix

Data, code and figures are provided in electronic form.

References

- Chastaing, G., Gamboa, F. and Prieur, C. (2012). Generalized Hoeffding-Sobol decomposition for dependent variables - application to sensitivity analysis, *Electronic Journal of Statistics* **6**(none).
- Efron, B. and Stein, C. (1981). The Jackknife Estimate of Variance, **The Annals of Statistic**(Vol. 9, No. 3): pp. 586–596.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine., *The Annals of Statistics* **29**(5): 1189–1232. Publisher: Institute of Mathematical Statistics.
URL: <https://projecteuclid.org/journals/annals-of-statistics/volume-29/issue-5/Greedy-function-approximation-A-gradient-boosting-machine/10.1214/aos/1013203451.full>
- Fumagalli, F., Muschalik, M., Hüllermeier, E., Hammer, B. and Herbringer, J. (2025). Unifying Feature-Based Explanations with Functional ANOVA and Cooperative Game Theory. arXiv:2412.17152 [cs].
URL: <http://arxiv.org/abs/2412.17152>
- Hoeffding, W. (1948). A Class of Statistics with Asymptotically Normal Distribution, *The Annals of Mathematical Statistics* **19**(3): 293–325. Publisher: Institute of Mathematical Statistics.
URL: <https://www.jstor.org/stable/2235637>
- Hooker, G. (2004). Discovering additive structure in black box functions, *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, Seattle WA USA, pp. 575–580.
URL: <https://dl.acm.org/doi/10.1145/1014052.1014122>
- Hooker, G. (2007). Generalized Functional ANOVA Diagnostics for High-Dimensional Functions of Dependent Variables, *Journal of Computational and Graphical Statistics* **16**(3): 709–732.
URL: <http://www.tandfonline.com/doi/abs/10.1198/106186007X237892>
- Il Idrissi, M., Bousquet, N., Gamboa, F., Iooss, B. and Loubes, J.-M. (2025). Hoeffding decomposition of functions of random dependent variables, *Journal of Multivariate Analysis* **208**: 105444.
URL: <https://linkinghub.elsevier.com/retrieve/pii/S0047259X25000399>

- Liu, R. and Owen, A. B. (2006). Estimating Mean Dimensionality of Analysis of Variance Decompositions, *Journal of the American Statistical Association* **101**(474): 712–721.
URL: <https://www.tandfonline.com/doi/full/10.1198/016214505000001410>
- Muehlenstaedt, T., Roustant, O., Carraro, L. and Kuhnt, S. (2012). Data-driven Kriging models based on FANOVA-decomposition, *Statistics and Computing* **22**(3): 723–738.
URL: <http://link.springer.com/10.1007/s11222-011-9259-7>
- Owen, A. (2003). The dimension distribution and quadrature test functions, *Statistica Sinica* **13**: 1–17.
- Owen, A. B. (2013). Variance components and generalized sobol’ indices, *SIAM/ASA Journal on Uncertainty Quantification* **1**(1): 19–41. tex.eprint: <https://doi.org/10.1137/120876782>.
URL: <https://doi.org/10.1137/120876782>
- Rahman, S. (2014). A generalized anova dimensional decomposition for dependent probability measures, *SIAM/ASA Journal on Uncertainty Quantification* **2**(1): 670–697.
URL: <https://doi.org/10.1137/120904378>
- Sobol, I. (2001). Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates, *Mathematics and Computers in Simulation* **55**(1-3): 271–280.
URL: <https://linkinghub.elsevier.com/retrieve/pii/S0378475400002706>
- Sobol, I. M. (1993). Sensitivity estimates for nonlinear mathematical models, *Mathematical Modelling and Computational Experiments* **1**: 407–414.
- Stone, C. J. (1994). The Use of Polynomial Splines and Their Tensor Products in Multivariate Function Estimation, *The Annals of Statistics* **22**(1): 118–171. Publisher: Institute of Mathematical Statistics.
URL: <https://www.jstor.org/stable/2242446>

Declaration of authorship

I hereby declare that the report submitted is my own unaided work. All direct or indirect sources used are acknowledged as references. I am aware that the Thesis in digital form can be examined for the use of unauthorized aid and in order to determine whether the report as a whole or parts incorporated in it may be deemed as plagiarism. For the comparison of my work with existing sources I agree that it shall be entered in a database where it shall also remain after examination, to enable comparison with future Theses submitted. Further rights of reproduction and usage, however, are not granted here. This paper was not previously presented to another examination board and has not been published.

Location, date

Name