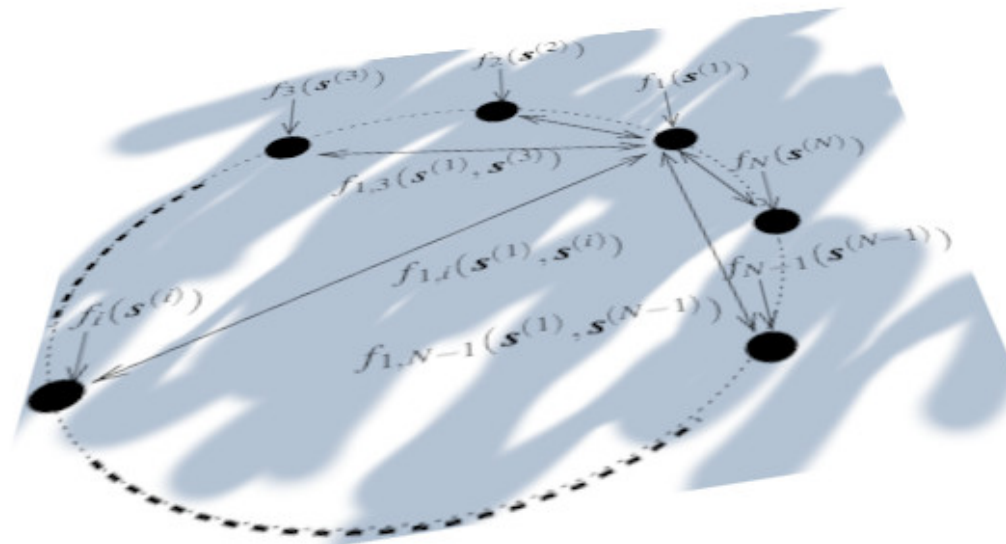


# Analyzing Complex Models by Orthogonal Input-Output Decompositions



**Pavel Loskot**

*pavelloskot@intl.zju.edu.cn*



**ZJU-UIUC INSTITUTE**

Zhejiang University-University of Illinois at Urbana-Champaign Institute  
浙江大学伊利诺伊大学厄巴纳香槟校区联合学院

**The First International Conference on Systems Explainability**  
**EXPLAINABILITY 2024**

**November 17, 2024 to November 21, 2024 - Valencia, Spain**

## ABOUT ME



**Pavel Loskot** joined the ZJU-UIUC Institute as Associate Professor in January 2021. He received his PhD degree in Wireless Communications from the University of Alberta in Canada, and the MSc and BSc degrees in Radioelectronics and Biomedical Electronics, respectively, from the Czech Technical University of Prague. He is the Senior Member of the IEEE, Fellow of the HEA in the UK, and the Recognized Research Supervisor of the UKCGE.

In the past 25 years, he was involved in numerous industrial and academic collaborative projects in the Czech Republic, Finland, Canada, the UK, Turkey, and China. These projects concerned mainly wireless and optical telecommunication networks, but also genetic regulatory circuits, air transport services, and renewable energy systems. This experience allowed him to truly understand the interdisciplinary workings, and crossing the disciplines boundaries.

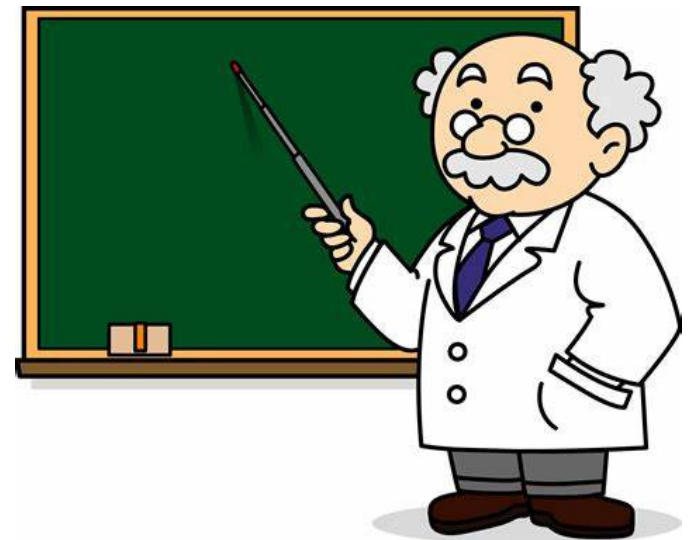
His current research focuses mathematical and probabilistic modeling, statistical signal processing and classical machine learning for multi-sensor data in biomedicine, computational molecular biology, and wireless communications.

## OBJECTIVES

- explore strategies for extracting information about complex models to achieve their explainability
- consider specifically the methods for factorizations, decompositions, and expansions of univariate and multivariate functions
- investigate Sobol's decomposition of multivariate functions, which did not receive sufficient attention it likely deserves

## OUTLINE

- Mathematical modeling
- Decompositions of multivariate functions
- Orthogonalizing inputs and outputs
- Numerical example
- Observations and future work



# MATHEMATICAL MODELING

## Modeling economics

- fast and cheap explorations  
→ replacing lab/field experiments
- used everywhere nowadays
- improving information gains

## Model formats

- mathematical expressions  
→ reproducibility
- algorithms, computer simulations
- datasets

## Model interpretability

- understand how outputs are obtained from inputs  
→ understand the model structure
- study how parameters affect model properties  
→ sensitivity analysis
- various model-related tasks:  
→ calibration, optimization, selection, validation, simplification

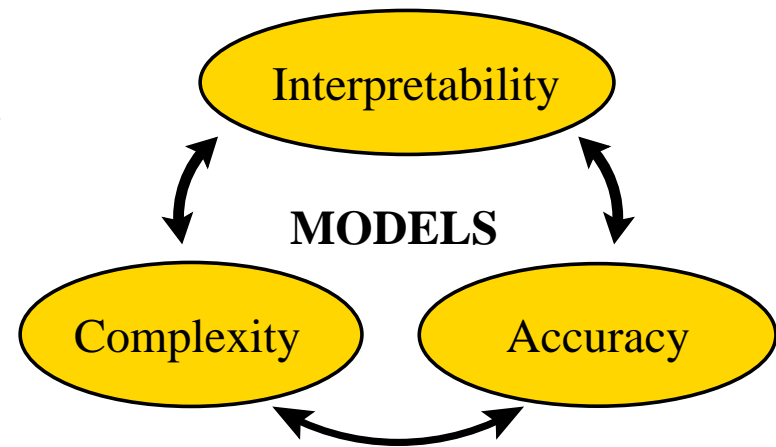
## Interpretability strategies

- local and global sensitivity
- surrogate and meta models
- output (variance) expansions

## MATHEMATICAL MODELING (CONT.)

### Model explainability

- narrower objective than interpretability
- which inputs more important for outputs  
→ factor screening, attribution problem
- model-agnostic methods  
→ permutation importance  
→ dependency plots  
→ SHAP explanations



### Orthogonal decomposition

$$M(x) = a_1 M_1(x) + a_2 M_2(x) + \dots + M_N(x)$$

- generalized Fourier series
- if components  $M_i(x)$  are mutually orthogonal  
→ straightforward to find coefficients  $a_i$   
→ analysis of  $M_i(x)$  can be separated  
→ faster learning (convergence)  
→ clear high-level structure i.e. explainability
- methods  
→ Gram-Schmidt process, matrix factorizations (SVD, PCA, QR, ... )

## DECOMPOSITIONS OF MULTIVARIATE FUNCTIONS

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) = \mathbf{f}(x_1, \dots, x_I) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_O(\mathbf{x}) \end{bmatrix} \in \mathcal{R}^O, \mathbf{x} \in \mathcal{R}^I$$

### Objectives

- uncover latent structure
- reducing computational complexity  
→ divide & conquer
- provide explainability
- approximations  
→ optimization and analysis

### Function factorization

$$\mathbf{f}(\mathbf{x}) = \prod_{i=1}^n \mathbf{f}_i(\mathbf{s}_i), \mathbf{s}_i \subseteq \{x_1, \dots, x_I\} \quad (\text{product factors})$$

$$\mathbf{f}(\mathbf{x}) = \sum_{i=1}^n \mathbf{f}_i(\mathbf{s}_i), \mathbf{s}_i \subseteq \{x_1, \dots, x_I\} \quad (\text{sum factors})$$

$$\mathbf{f}(\mathbf{x})A(\mathbf{x}) = \begin{cases} \mathbf{f}(\mathbf{x}), & \mathbf{x} \in \mathcal{A} \\ 0, & \mathbf{x} \notin \mathcal{A} \end{cases} \quad (\text{adding constraints})$$

## DECOMPOSITIONS OF MULTIVARIATE FUNCTIONS (CONT.)

### Stochastic function

$$\mathbf{y} \approx \sum_{i=1}^n a_i \|\mathbf{x} - \mathbf{E}[\mathbf{x}]\|_1^i \quad [\text{Loskot, 2021}]$$

→ multivariate Taylor expansion

### Universal approximation theorem

$$\mathbf{y} \approx \dots \sigma \circ (\mathbf{A}_i, \mathbf{b}_i) \circ \dots \sigma(\mathbf{A}_1 \mathbf{x} + \mathbf{b}_1) \quad [\text{Hornik, 1989}]$$

### Kolmogorov-Arnold approximation

$$f(\mathbf{x}) = \sum_{i=0}^{2n} \Phi \left( \sum_{j=1}^n \phi_{i,j}(x_j) \right) \quad [\text{Lorentz, 1962}]$$

### Sobol's decomposition

$$f(\mathbf{x}) = f_0 + \sum_{i=1}^I f_i(\mathbf{x}_i) + \sum_{\substack{i,j=1 \\ i \neq j}}^I f_{i,j}(\mathbf{x}_i, \mathbf{x}_j) \cdots + \sum_{i=1}^I f_{\{1:I\} \setminus i}(\mathbf{x})$$

# ORTHOGONAL INPUT-OUTPUT DECOMPOSITIONS

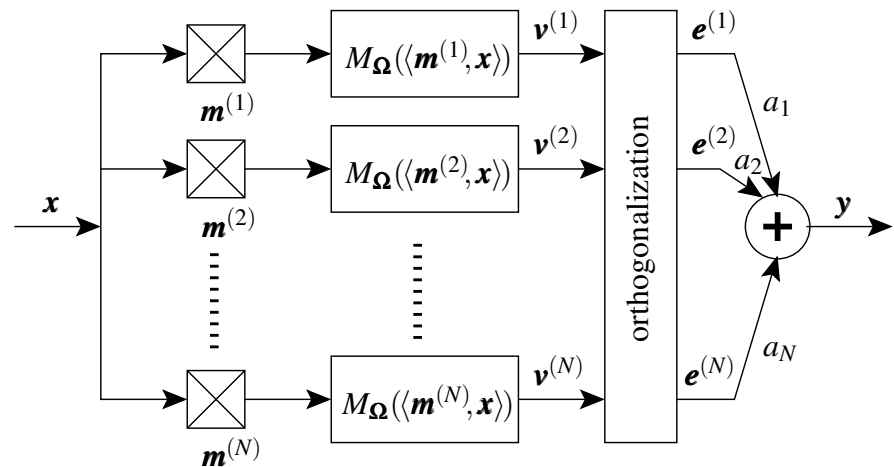
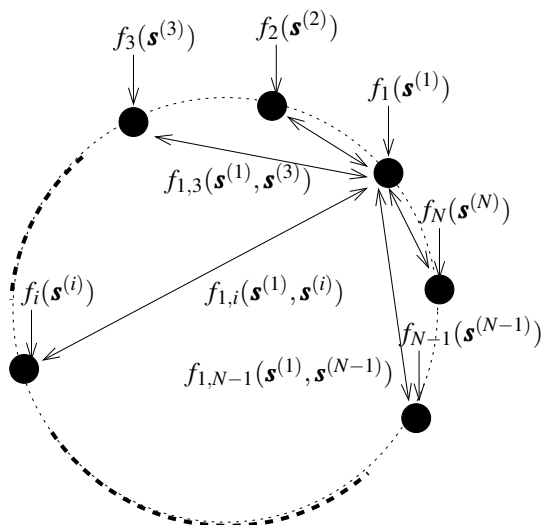
## Main task

- analyze explainability of complex model,  $f(x, \Omega) \equiv M(x; \Omega) \equiv M_{\Omega}(x)$

## Key paper contribution

$$f(x) \approx f_0 + \sum_{i=1}^N f_i(s^{(i)}) + \sum_{\substack{i,j=1 \\ i \neq j}}^N f_{i,j}(s^{(i)}, s^{(j)})$$

orthogonal input projections (masking)  $\left\{ \begin{array}{l} s^{(i)} \cap s^{(j \neq i)} = \emptyset \quad (\text{disjoint}) \\ \cup_i s^{(i)} = x \quad (\text{full coverage}) \end{array} \right.$





## NUMERICAL EXAMPLE

### Problem

- classifying hand-written digits in MNIST dataset
- train a MLP with two hidden layers  
→ training accuracy 97.56%, testing accuracy 94.31%

### Explore

1. training dataset vs. randomly generated inputs
2. orthogonal (disjoint) masking vs. random masking of inputs

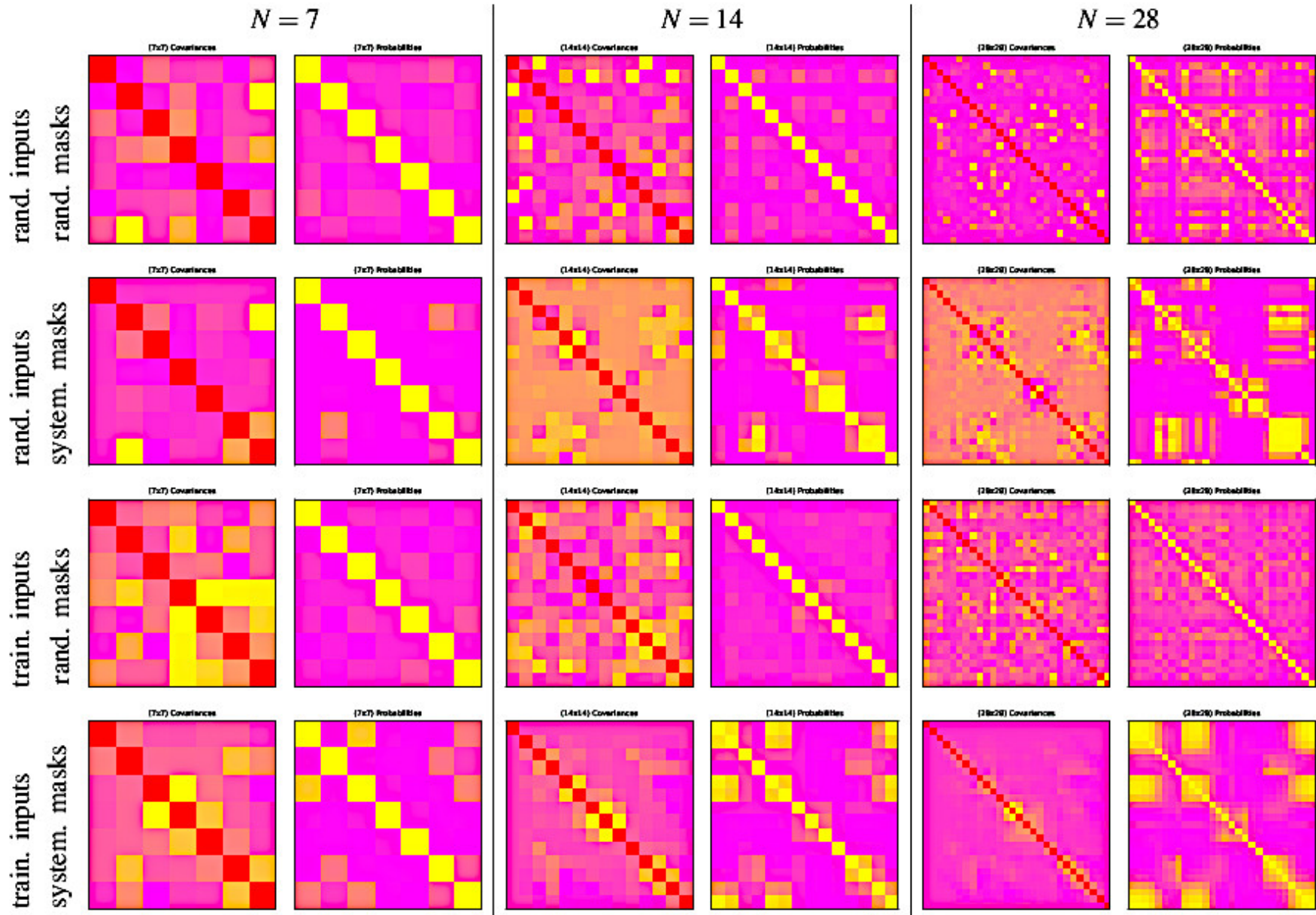
$$\overline{\text{MSE}}_0 = \frac{1}{K} \sum_{k=1}^K \|\mathbf{y}(k)\|^2$$

$$\overline{\text{MSE}}_1 = \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K \|\mathbf{v}^{(i)}(k) - \mathbf{y}(k)\|^2$$

$$\min \text{MSE}_1 = \min_{1 \leq i \leq N} \frac{1}{K} \sum_{k=1}^K \|\mathbf{y}^{(i)}(k) - \mathbf{y}(k)\|^2$$

$$\overline{\text{MSE}}_2 = \frac{1}{K} \sum_{k=1}^K \|\tilde{\mathbf{y}}(k) - \mathbf{y}(k)\|^2$$

# NUMERICAL EXAMPLE (CONT.)



→ correlations and probabilities (decisions with masks  $m^{(i)}$  and  $m^{(j)}$  the same as with combined mask  $m^{(i)} + m^{(j)}$ )

## OBSERVATIONS AND CONCLUSION

### MSE values

- input masking substantially reduces accuracy in exchange for interpretability
- combining outputs corresponding to masks inputs restores some accuracy
- larger sensitivity for training vs. random inputs than for orthogonal vs. random masks

### Correlations vs. probabilities

- increasing  $N$  increases resolution and variance of calculated values
- some orthogonal inputs have dependent decisions (square patterns)
- probability patterns require sufficiently large  $N$  to be visible  
→ i.e. explainability requires the minimum resolution

### Summary

- Sobol-based decomposition of multivariate functions  
→ graph representation of multivariate functions
- input masking  
→ projections into orthogonal input subspaces
- output orthogonalization via decorrelation
- SVD of correlation matrix
- easy to obtain linear combining coefficients

## FUTURE WORK

### Some ideas

- explainability methods must be much simpler than the model investigated
- optimizing input projections (granularity vs. complexity vs. accuracy)  
→ different explainability objectives
- averaging out some inputs instead of using default values (zero)
- other strategies for obtaining component models from the base model
- constructing complex models from orthogonal component models  
→ machine learning architectures
- model decomposition as structural causal model  
→ causal inferences
- Sobol's decomposition  
→ assuming higher-order graphs  
→ approximation guarantees

*The End*

*Thank you!*

*pavelloskot@intl.zju.edu.cn*