# Please join the Biometric Colloquium

## MATTHIAS SCHMID

Institute for Medical Biometry, Informatics and Epidemiology, University of Bonn, Bonn, Deutschland

## ACHIEVING EXPLAINABLE MACHINE LEARNING BY FUNCTIONAL DECOMPOSITION OF BLACK-BOX MODELS INTO EXPLAINABLE PREDICTOR EFFECTS

### May 8th, 2025 at 9:00-10:00 am

Seminarraum Center for Medical Data Science (previously CeMSIIS),
Spitalgasse 23, Room 88.03.513
Medical University of Vienna, 1090 Wien
Host: Georg Heinze

## Abstract:

The high accuracy of supervised machine learning models is usually achieved by optimizing complex, uninterpretable "black-box" architectures, hindering their applicability in fields where an understanding of the model and its inner workings is paramount to ensuring user acceptance and fairness. Functional decomposition is a well explored tool that improves interpretability by splitting the prediction function into a sum of main and interaction effects. Existing implementations are often computationally infeasible. We present a novel implementation by fitting a neural additive model with DNN-based submodels using the model predictions as outcome variable. We enable identifiability and interpretability by orthogonalizing submodels against higher-order terms. Information is shifted into the explainable and visualizable lower-order effects, ensuring that these effects capture as much of the model variance as possible. By having minimal prerequisites on DNN architecture and model fitting, the method can be widely applied without constraining the learning algorithm and model predictive performance. It also yields a variance decomposition of the predictions, giving an intuitive quantification of the degree of explainable model variance. We illustrate the use of our method by applying it to an ecological dataset, yielding insights into the effects of geological features on the prediction of stream biological condition in the U.S. Chesapeake Bay watershed.