

Tülay Adalı, Rodrigo Capobianco Guido, Tin Kam Ho,  
Klaus-Robert Müller, and Stephen Strother

## Interpretability, Reproducibility, and Replicability

Most of the work we do in signal processing these days is data driven. The shift from the more traditional and model-driven approaches to those that are data driven has also underlined the importance of explainability of our solutions. Because most traditional signal processing approaches start with a number of modeling assumptions, they are comprehensible by the very nature of their construction. However, this is not necessarily the case when we choose to rely more heavily on the data and minimize modeling assumptions.

Explainability is critical not only for the simple reason that one would like to have confidence over solutions, but also because one would like to obtain further insights into the problem from learned models. This includes interpretability and completeness so that one can not only “audit” them but also ask appropriate questions to probe for insights beyond the initial solution, and address additional concerns such as safety, fairness, and reliability. Interpretability, i.e., the ability to attach a physical meaning to a solution, along with reproducibility and replicability are three key aspects of explainability. Following definitions from the National Academies of Sciences, Engineering, and Medicine, *reproducibility* refers to obtaining consistent results using the same data and code—i.e., method,—as the original study, and *replicability* is obtaining

consistent results across studies aimed at answering the same scientific question using new data or other computational methods.

In this special issue of *IEEE Signal Processing Magazine*, we have nine articles that demonstrate the multifaceted nature of explainability, and span the related concepts of interpretability, reproducibility, and replicability. They successfully exhibit the rich nature of these concepts while also highlighting the fact that they take on slightly different meanings in various contexts, and the considerations might be slightly different as well. These articles also emphasize the fact that explainability is a key theme requiring attention across different application domains and types of solutions, i.e., well beyond neural networks, where they have been mostly emphasized to date.

The first two articles of the special issue study the questions of reproducibility, replicability, and interpretability for two important classes of machine learning solutions: matrix and tensor decompositions (MTDs) and graph data science. The first article, “Reproducibility in Matrix and Tensor Decompositions,” by Adalı et al., addresses reproducibility in MTD solutions that have been growing in importance, where, in addition to the discovery of structure in data, the resulting decomposition is also directly interpretable. With an applied focus where there is no ground truth, the authors study the intricate relationships of interpretability, model match, and uniqueness. They make use of two widely used

methods with relaxed uniqueness guarantees, independent component analysis, and the canonical-polyadic decomposition, and provide examples to solidify these concepts and demonstrate the tradeoffs. Finally, a reproducibility checklist for MTDs is provided similar to those developed for supervised learning. In “Explainability in Graph Data Science,” Aviyente and Karaaslanli explain methods and metrics from network science to quantify three different aspects of explainability, i.e., interpretability, replicability, and reproducibility, in the context of community detection. Specifically, the strategies described by the authors can be used to address some common issues and provide guidelines to reduce the opacity of community-detection algorithms and their outputs. In addition, they can be extended to other community-detection and data-clustering algorithms as well as different learning tasks on graphs.

The second group of articles take a broader view of explainability, with a focus on neural networks. Letzgun et al. discuss an area of explainable artificial intelligence (XAI) that has thus far received comparatively little attention, namely, XAI for regression models. Their review, “Toward Explainable Artificial Intelligence for Regression Models,” provides a novel theory showing that there are important conceptual differences between XAI for regression and classification, not only algorithmically but also with respect to the choice of reference for the explanation. “Explanatory Paradigms in

Neural Networks,” by AlRegib and Prabhushankar characterizes a complete explanation as an additive combination of observed correlations and counterfactuals, and contrastive explanations. It then discusses how existing explanation methods can be analyzed within this framework, and how well they are suited to different evaluation strategies under a proposed taxonomy.

The next two articles consider generative adversarial networks, which have been growing in importance. In “Robust Explainability,” Nielsen et al. present a timely and comprehensive tutorial on gradient-based attribution/saliency methods, their relationship to adversarial robustness, and the practical importance of robust explainability of computer vision classification models based on these techniques, together with many of the associated terms that appear in explainability literature. They provide a useful list of best practices to consider when choosing an explainability method and conclude with future directions of research in the area of robust explainability. They augment their article with a website containing links to all the explainability methods discussed, the article’s figures, and code for generating the figures. In the second article of this group, “Explaining Artificial Intelligence Generation and Creativity,” Das and Varshney review different motivations, algorithms, and methods intended to explain the principles of AI algorithms or the possible artifacts they produce by using a generative point of view, with creativity as the focus. In particular, they observe that discussions of interpretable AI, especially in settings of decisions and predictions, frequently start with the misconception that there is a fundamental tradeoff between interpretability and accuracy, however, as reviewed, numerous examples show the opposite.

The last group of articles of the special issue consider explainability with an application focus, and across a wide array of data-driven solutions. The first two articles examine applications in health care, and the last one surveys time-series classification. The first article of this group, “Explainability of Methods for Critical Information

Extraction From Clinical Documents,” by Ho et al., reviews a collection of representative works that address several natural language understanding tasks in health care, and discusses their explainability. It showcases the complex dimensions of considerations in providing explainable methods for an essential application of AI. In “Interpreting Brain Biomarkers,” Jiang et al. review predictive methods and their applications for interpreting brain signatures in neuroimaging based on a survey of more than 300 articles. This way, better validation and assessment of the reliability and interpretability of biomarkers across multiple data sets and contexts can be achieved. Finally, in “Post Hoc Explainability for Time Series Classification,” Mochaourab et al. discuss the explainability advantages of methods for time-series classification that are based on representations well established in signal processing. The article highlights the relevance of such conventional transformations for the important concerns of understanding feature importance and providing counterfactual explanations.

We thank our contributors for their comprehensive and interesting articles, Robert Heath for his support for our proposal, and Christian Jutten for providing valuable guidance and support at every step of the process. We also extend thanks to our reviewers for their detailed and insightful comments, Rebecca Wollman for guidance and support along the way, and Sharon Turk for special care in putting together this issue of *IEEE Signal Processing Magazine*.

Data-driven solutions are becoming the dominant approach to solving practical problems across multiple fields, and explainability is a key aspect that will further enhance their utility. Signal processing is at the heart of data science and is where the connection with applications is natural. Hence, we are hoping that the insights, as well as critical perspectives provided by contributions to the special issue, will prove to be a useful reference and help identify some of the new and emerging directions in the area.

## Guest Editors



**Tülay Adalı** (adali@umbc.edu) received her Ph.D. degree in electrical engineering from North Carolina State University. She is

a distinguished university professor in the Department of Computer Science and Electrical Engineering, the University of Maryland, Baltimore County, Baltimore, Maryland, 21250, USA. She is a Fulbright Scholar and an IEEE Signal Processing Society (SPS) Distinguished Lecturer. She is a past vice president for technical directions for the SPS and is currently the chair-elect of IEEE Brain. She is the recipient of a Humboldt Research Award, an IEEE SPS Best Paper Award, the University System of Maryland Regents’ Award for Research, and a National Science Foundation CAREER Award. Her research interests include statistical signal processing and machine learning, and their applications, with an emphasis on applications in medical image analysis and fusion. She is a Fellow of IEEE and the American Institute for Medical and Biological Engineering.



**Rodrigo Capobianco Guido** (guido@ieee.org) received his Ph.D. degree in computational applied physics from the University of São

Paulo (USP) in 2003. After two postdoctoral programs in signal processing at USP, he obtained the title of associate professor (livre-docência) in signal processing, also from USP, in 2008. Currently, he is an associate professor at São Paulo State University (UNESP) at São José do Rio Preto, São Paulo, 15054-000, Brazil. He has been an area editor for *IEEE Signal Processing Magazine* and was recently included in Stanford University’s ranking of the world’s top 2% of scientists. He is a Senior Member of IEEE.



**Tin Kam Ho** (tkh@ieee.org) received her Ph.D. degree in computer science from the State University of New York at Buffalo in 1992. She is a senior artificial

intelligence scientist at IBM Watson Health, Yorktown Heights, New York, 10598-0218, USA, where she leads projects in semantic modeling of natural languages in clinical applications. Prior to 2014, she was with Bell Labs as the head of the Statistics and Learning Research department. She pioneered research in multiple classifier systems and ensemble learning, random decision forests, and data complexity analysis and also contributed to many applications of pattern recognition and computational modeling. She is a Fellow of IEEE and the International Association for Pattern Recognition.



**Klaus-Robert Müller** (klaus-robert.mueller@tu-berlin.de) received his Ph.D. degree in computer science from Technische

Universität Karlsruhe in 1992. He has been a professor of computer science with TU Berlin, 10587, Germany, since 2006. Since 2012, he has been a distinguished professor at Korea University, Seoul, 02841, South Korea. In 2020 and 2021, he was on sabbatical leave from TU Berlin and with the Brain Team, Google Research, as a principal researcher. He is currently directing the Berlin Institute for the Foundations of Learning and Data. He was elected a member of the German National Academy of Sciences, Leopoldina (2012); the Berlin Brandenburg Academy of Sciences (2017); an External Scientific Member of the Max Planck Society (2017); and the National Academy of Science and Engineering (2021). He is lead of the European Laboratory for Learning and Intelligent Systems unit Berlin. From 2019, he became an

Institute for Scientific Information Highly Cited Researcher in the cross-disciplinary area.



**Stephen Strother** (sstrother@research.baycrest.org) received his Ph.D. degree in electrical engineering from McGill University

in 1986. He is a member of the Rotman Research Institute, Baycrest Hospital, Toronto, M6A 2E1, and a professor of medical biophysics at the University of Toronto, Toronto, Ontario, M5G 1L7, Canada. His research interests include neuroinformatics and data science for neuroimaging and big clinical data sets through statistical and machine learning techniques, applying these techniques in cognitive neuroscience and brain disease and translating this work to non-academic settings. **SP**



420,000+ members in 160 countries. Embrace the largest, global, technical community.

People Driving Technological Innovation.

[ieee.org/membership](https://www.ieee.org/membership)

#IEEEmember

