

## Department of Data Science

### Data Science Seminar Series

## Deep Neural Networks Explainability: Algorithms and Applications



### **Mengnan Du**

**Ph.D. Candidate**

**Texas A&M University**

**Date:** Tuesday, April 26<sup>th</sup>, 2022

**Time:** 11:00 AM – 12:00 PM EDT

**Location:** Zoom Virtual Room

**Web Link:** [Zoom Meeting Room Link](#)

Deep neural networks (DNN) have achieved extremely high prediction accuracy in a wide range of fields such as computer vision, natural language processing, and recommender systems. Despite the superior performance, DNN models are often regarded as black boxes and criticized for the lack of interpretability, since these models cannot provide meaningful explanations on how a certain prediction is made. Without the explanations to enhance the transparency of DNN models, it would become difficult to build up trust and credibility among end-users. In this talk, I will present our efforts to tackle the black-box problem and to make powerful DNN models more interpretable and trustworthy. First, I will introduce post-hoc interpretation approaches for predictions made by two standard DNN architectures, including Convolution Neural Network (CNN) and Recurrent Neural Network (RNN). Second, I will introduce the usage of explainability as a debugging tool to improve the generalization ability and fairness of DNN models.

Mengnan Du obtained the Ph.D. degree in Computer Science at Texas A&M University, under the supervision of Dr. Xia Ben Hu. His research is on the broad area of trustworthy machine learning, including model explainability, fairness, and robustness. He has had around 40 papers published in prestigious venues such as NeurIPS, AAAI, KDD, WWW, NAACL, ICLR, and TPAMI. He received over 1,400 citations with an H-index of 13. Three of his papers were selected for the Best Paper (Candidate) at WWW 2019, ICDM 2019, and INFORMS 2019, respectively. His paper on Explainable AI was highlighted on the cover page of Communications of the ACM, January 2020 issue. He served as the Registration Chair of WSDM'22, and is the program committee member of conferences including NeurIPS, ICML, ICLR, KDD, AAAI, ACL, EMNLP, etc. More detail can be found at <https://mengnandu.com/>.