

## An Open Source Repository of Explainable Artificial Intelligence Projects

### **Sohee Cho**

### Explainable Artificial Intelligence Center KAIST

### **Explainable AI Program in Korea**

#### Goal

#### Human-level Learning and Inference to overcome the limitations of Deep Neural Networks



- It is hard to know the decision, so called Blackbox model
- It does not work well when we do not have enough training data



- **Explainable learners** which can provide the reasons of decisions
- Learning explainable models even with **data deficient environment**



Institute of Information & Communication Technology Promotion (IITP) under Ministry of Science and ICT (MSICT) as part of Innovative Growth Engine Project

Period

> July 2017 ~ December 2021 (54 months)

### **Project Organization**







2019 ICCV Workshop on

Interpretating and Explaining Visual Artificial Intelligence Models



## https://Openxai.org

## http://xai.unist.ac.kr/opensource /relatedproject/





- Institutions: TU Berlin, Fraunhofer Heinrich Hertz Institut
  - Authors: Maximilian Alber, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele, Kristof T. Schütt, Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, Pieter-Jan Kindermans
- **Publication :** iNNvestigate neural networks!
- Source Code: https://github.com/albermax/innvestigate
- Further Info: https://innvestigate.readthedocs.io/en/latest/





#### **Network Dissection**



- Publication: Network Dissection: Quantifying Interpretability of Deep Visual Representations
- **Source Code :** https://github.com/CSAILVision/NetDissect
- Further Info: http://netdissect.csail.mit.edu/co



#### XCAD



Publication: ICADx: Interpretable computer aided diagnosis of breast masses

Source Code: https://github.com/xairc/XCAD

#### **Relational Automatic Statistician**









Institutions: VAIT Authors: Secing Tex Kim, Hakmin Lee, Hak Gu Kim, Yong Man Bo Pelotication: CADs interpretable computer aded diagnoss of breaz masses Source Code : https://gthub.com/uurch/CAD



stiftutions : Frauntofer inennin hertz Institute, TU Berll Authors : Sebastian Lapuschin, Alexander Binder, Klaus Robert Muller, Wojdech Sairkek Walicatian : (Understanding and Companing Deep Narual Networks seurce Code : https://github.com/sebastian lapusch/innum/instanding age gender deep learning mode

+ 1 2 3 4 5 -

**Publication :**Discovering Latent Covariance Structures for Multiple Time Series**Source Code :**https://github.com/OpenXAIProject/Automatic-Stock-Report

#### Workshop on Visual XAI

Institutions: UNIST

Authors: Anh Tong, Jaesik Choi



	Project Title	Institutions	Authors	Publication_title	Sourcecode
1	Principles of Explanatory Debugging to Personalize Interactive Machine Learning	Oregon State, City University London	T. Kulesza, M. Burnett, W-K. Wong and S. Stumpf	Principles of Explanatory Debugging to Personalize Interactive Machine Learning, IUI, 2015	https://github.com/fflewddur/IMLPla yground
2	Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model	MIT, U of Washington, Columbia	B. Letham, C. Rudin, T. McCormick and D. Madigan	Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model, Annals of Applied Statistics, 2015	https://github.com/nlarusstone/corel s
3	Explaining Recurrent Neural Network Predictions in Sentiment Analysis	Fraunhofer, TU Berlin, Korea University, Max	L. Arras, G. Montavon, K-R. M체ller and W. Samek	Explaining Recurrent Neural Network Predictions in Sentiment Analysis, EMNLP, 2017	https://github.com/ArrasL/LRP_for_L STM
4	Why Should I Trust You?: Explaining the Predictions of Any Classifier""	U of Washington	M. T. Ribeiro, S. Singh, S. and C. Guestrin	Why Should I Trust You?: Explaining the Predictions of Any Classifier, KDD, 2016	https://github.com/marcotcr/lime
5	Multimedia Event Detection and Recounting	SRI-International Sarnoff, U. of Massachusett	H. Cheng et. al.	SRI-Sarnoff AURORA at TRECVID 2014 - Multimedia Event Detection and Recounting	https://www.nist.gov/itl/iad/mig/tool s
6	Examples are not Enough, Learn to Criticize! Criticism for Interpretability	Allen Institute, UT Austin, UIUC	B. Kim, R. Khanna, S. Koyejo	Examples are not Enough, Learn to Criticize! Criticism for Interpretability, NIPS, 2016	https://github.com/BeenKim/MMD- critic
7	Learning AND-OR Templates for Object Recognition and Detection	UC Los Angeles	Z. Si and S. Zhu	Learning AND-OR Templates for Object Recognition and Detection, TPAMI, 2013	http://www.stat.ucla.edu/~zzsi/AOT. html
8	Human-level concept learning through probabilistic program introduction	New York University, U of Toronto, MIT	B. H. Lake, R. Salakhutdinov, and J. B. Tenenbaum	Human-level concept learning through probabilistic program introduction, Science, 2015	https://github.com/brendenlake/BPL
9	Generating Visual Explanations	UC Berkeley, Max Planck Institute for Informa	L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell	Generating Visual Explanations, ECCV 2016	https://github.com/LisaAnne/ECCV2 016/tree/ECCV2016/examples/ECCV 2016
10	The Automatic Statistician	Cambridge, MIT	J. R. Lloyd, D. Duvenaud, R. Grosse, J. B. Tenenbaum and Z. Ghahramani	Automatic Construction and Natural- Language Description of Nonparametric Regression Models, AAAI, 2014	https://github.com/jamesrobertlloyd /gpss-research

Worksnop on visual XAL



	Project Title	Institutions	Authors	Publication_title	Sourcecode
11	Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car	NVIDIA, New York University, Google	B. Mariusz et. al.	Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car	https://github.com/maxritter/SDC- End-to-end-driving
12	The LRP Toolbox for Artificial Neural Networks	Fraunhofer, Berlin Institute of Tech., Korea	S. Lapuschkin, A. Binder, G. Montavon, K-R. M利ller, W. Samek	The LRP Toolbox for Artificial Neural Networks, JMLR, 2016	https://github.com/VigneshSrinivasa n10/interprettensor
13	PatternNet and PatternAttribution	Google Brain, TU Berlin	Pieter-Jan Kindermans, Kristof T. Schutt & Maximilian Alber, K-R. M체 Iler, Dumitru Erhan & Been Kim, Sven Dahne	LEARNING HOW TO EXPLAIN NEURAL NETWORKS: PATTERNNET AND PATTERNATTRIBUTION	https://openreview.net/pdf?id=Hkn7 CBaTW
14	Network Dissection	MIT	David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba	Network Dissection: Quantifying Interpretability of Deep Visual Representations	https://github.com/CSAILVision/Net Dissect
15	iNNvestigate neural networks	TU Berlin, Fraunhofer Heinrich Hertz Institut	Maximilian Alber et al	iNNvestigate neural networks!	https://github.com/albermax/innvest igate
16	Understanding-age-gender- deep-learning-models	Fraunhofer Heinrich Hertz Institute, TU Berli	Sebastian Lapuschkin, Alexander Binder, Klaus- Robert Muller, Wojciech Samek	Understanding and Comparing Deep Neural Networks	https://github.com/sebastian- lapuschkin/understanding-age- gender-deep-learning-models
17	XCAD	KAIST	Seong Tae Kim, Hakmin Lee Hak Gu Kim, Yong Man Ro	ICADx: Interpretable computer aided diagnosis of breast masses	https://github.com/xairc/XCAD
18	Relational Automatic Statistician	UNIST	Anh Tong, Jaesik Choi	Discovering Latent Covariance Structures for Multiple Time Series	https://github.com/OpenXAIProject/ Automatic-Stock-Report
19	SHAP (SHapley Additive exPlanations)	University of Washington	Scott M. Lundberg, Gabriel G. Erion, Su-In Lee	Consistent Individualized Feature Attribution for Tree Ensembles	https://github.com/slundberg/shap

Worksnop on visual XAL



## https://github.com/OpenXAIProject

ICCV 2019 Seoul, Korea



Forked from kirarenctaon/xai



Edit

OpenXAIProject / PyConKorea2019-Tutorials						• Wa	ntch 🕶 13	★ Unstar	37	<b>%</b> Fork	14
<> Code	() Issues 0	11 Pull requests 0	Projects 0	🔳 Wiki	C Security	11 Insights	🌣 Settings				

#### Tutorials about Layer-wise Relevance Propagation (LRP) and SHapley Additive exPlanations (SHAP)

Manage topics

52 commits	52 commits 1 branch		♥ 0 releases ▲ 4 contributors		മ് Apache-2.0		
Branch: master - New pull request			Create new file	Upload files	Find file	Clone or download 🗸	
F SeongmanHeo Update README.md					Latest con	nmit 4fa4f78 on 19 Aug	
LRP		Add files via u	pload			3 months ago	
SHAP		presentation n	naterial			3 months ago	
		Initial commit				3 months ago	
PyConKorea2019-Introduction-pres	entation.pdf	Add files via u	pload			3 months ago	
README.md		Update READN	/IE.md			2 months ago	
Colab_setting.ipynb		Add files via u	pload			3 months ago	

I README.md

### 딥러닝 안에서 일어나는 과정을 설명하는 인공지능 기술 (PyCon Korea 2019 Tutorial)

Introduction

Code 1 Pull requests 0	Projects 0	'iki 🕕 Security 🔟 Insight	s 🔅 Settings	
nsorflow tutorial for variou	is Deep Neural Network	visualization techniques		Edit
anage topics				
🕝 <b>151</b> commits	2 branches	𝔝 0 releases	<b>L</b> contributor	MIT کړه
Branch: master 🔻 New pull req	uest		Create new file Upload files Fi	nd file Clone or download 🗸
his branch is even with 1202kl	os:master.			🕅 Pull request 🖹 Compare
1202kbs Fix typo			Latest co	ommit <del>f9ddafc</del> on 14 Nov 2018
assets		Update Figure		last year
models		Add Grad-CAM++		last year
.gitignore		Add CAM tutorial		2 years ago
1.1 Activation Maximization	ipynb	Update explanation		2 years ago
1.3 Performing AM in Code	Space.ipynb	Update explanation		2 years ago
2.1 Sensitivity Analysis.ipynl	)	Update tutorials		2 years ago
2.2 Simple Taylor Decomposition	sition.ipynb	Update tutorials		2 years ago
2.3 Layer-wise Relevance Pr	opagation (1).ipynb	Fix typo		last year
2.3 Layer-wise Relevance Pr	opagation (2).ipynb	Reorder tutorials		2 years ago
2.4 Deep Taylor Decomposit	tion (1).ipynb	Fix image paths		2 years ago
2.4 Deep Taylor Decomposi	tion (2).ipynb	Reorder tutorials		2 years ago
2.5 DeepLIFT.ipynb		Add DeepLIFT tutorial		2 years ago
		Fix type		2 years ago

Seoul, Korea



## https://www.youtube.com/channel/UCGx sfIsOry\_LdBaPSet2p7g









XAI open software project - Data-driven Open-domain Neural Conversation Models (Part 3)

Seoul, Korea



# Thank you!

### Let us know other related projects! sohee.cho@kaist.ac.kr