THALES

XAI – The Story So Far

May 2nd, 2019

Freddy Lecue
Chief Al Scientist, CortAlx, Thales, Montreal – Canada
Inria, Sophia Antipolis - France

@freddylecue https://tinyurl.com/freddylecue



Context







 $\label{eq:Gary Chavez} \mbox{ added a photo you might } \mbox{ ...} \\ \mbox{be in.}$

about a minute ago · 🔐









Markets we serve











Aerospace

Space

Ground Transportation

Defence

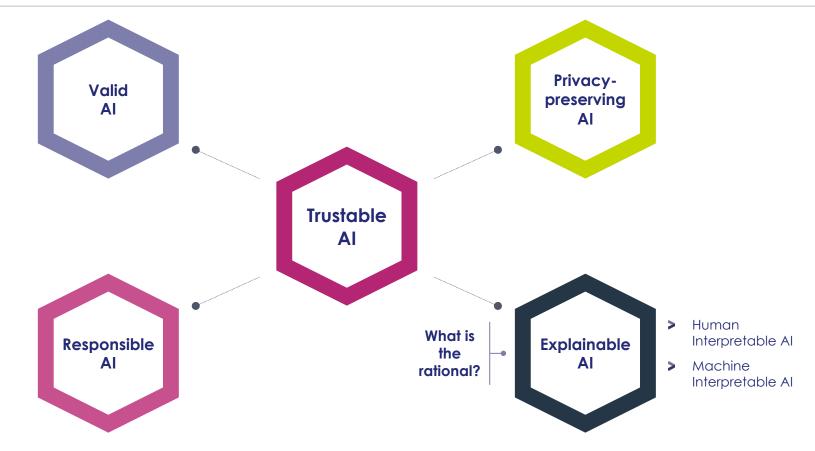
Security

Trusted Partner For A Safer World



Trustable Al







Motivations



Criminal Justice

- People wrongly denied parole
- Recidivism prediction
- Unfair Police dispatch

Opinion

The New Hork Times

OP-ED CONTRIBUTOR

When a Computer **Program Keeps You in Jail**

By Rebecca Wexler

June 13, 2017









nytimes.com/2017/06/13/opinion/how-computers-are-harming-criminal-justice.html

ACLU

GET UPDATES

DONATE

Q

STATEMENT OF CONCERN ABOUT PREDICTIVE POLICING BY ACLU AND 16 CIVIL RIGHTS PRIVACY, RACIAL JUSTICE, AND TECHNOLOGY **ORGANIZATIONS**



How We Analyzed the **COMPAS Recidivism Algorithm**

by Jeff Larson, Surva Mattu, Lauren Kirchner and Julia Angwin May 23, 2016

propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm



Motivation (2)

Finance:

- Credit scoring, loan approval
- > Insurance quotes







community.fico.com/s/explainable-machine-learning-challenge



Motivation (3)





Healthcare

- Applying ML methods in medical care is problematic.
- ➤ Al as 3rd-party actor in physician-patient relationship
- > Responsibility, confidentiality?
- > Learning must be done with available data.

Cannot randomize cares given to patients!

Must validate models before use.



Researchers say use of artificial intelligence in medicine raises ethical questions

In a perspective piece, Stanford researchers discuss the ethical implications of using machine-learning tools in making health care decisions for patients.

Patricia Hannon ,https://med.stanford.edu/news/all-news/2018/03/researchers-say-use-of-ai-in-medicine-raises-ethical-questions.html

Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission

Rich Caruana Microsoft Research rcaruana@microsoft.com

Paul Koch Microsoft Research paulkoch@microsoft.com Yin Lou LinkedIn Corporation ylou@linkedin.com

Marc Sturm NewYork-Presbyterian Hospital mas9161@nyp.org Johannes Gehrke Microsoft johannes@microsoft.com

Noémie Elhadad Columbia University noemie.elhadad@columbia.edu

Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, Noemie Elhadad: Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. KDD 2015: 1721-1730

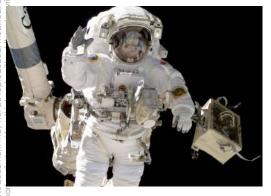
Motivation (4)

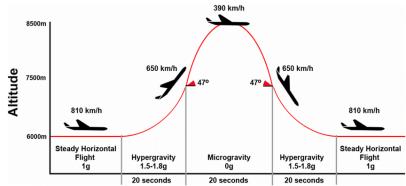
Critical Systems





https://www.sncf.com/sncv1/ressources/presskit_train_a utonome_september_2019_v2.pdf

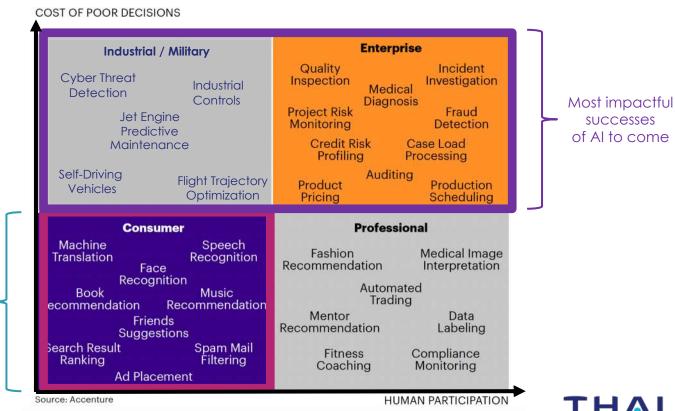






Trustable AI and eXplainable AI: a Reality Need

The need for explainable AI rises with the potential cost of poor decisions



Most prominent

successes

of AI to date

Definitions



Explanation in AI aims to create a suite of techniques that produce more explainable models, while maintaining a high level of searching, learning, planning, reasoning performance: optimization, accuracy, precision; and enable human users to understand, appropriately trust, and effectively manage the emerging generation of AI systems.



Oxford Dictionary of English

explanation | Eksplə'neIf(ə)n |

noun

a statement or account that makes something clear: the birth rate is central to any explanation of population trends.

Models, Outputs of the Intelligent System

interpret | In'terprit |

verb (interprets, interpreting, interpreted) [with object]

1 explain the meaning of (information or actions): the evidence is difficult to interpret.

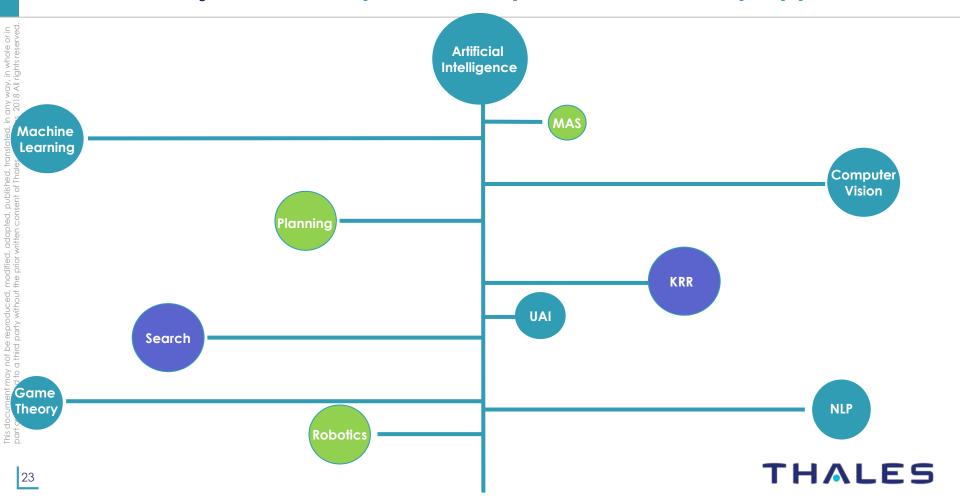
Models, Outputs of the Intelligent System



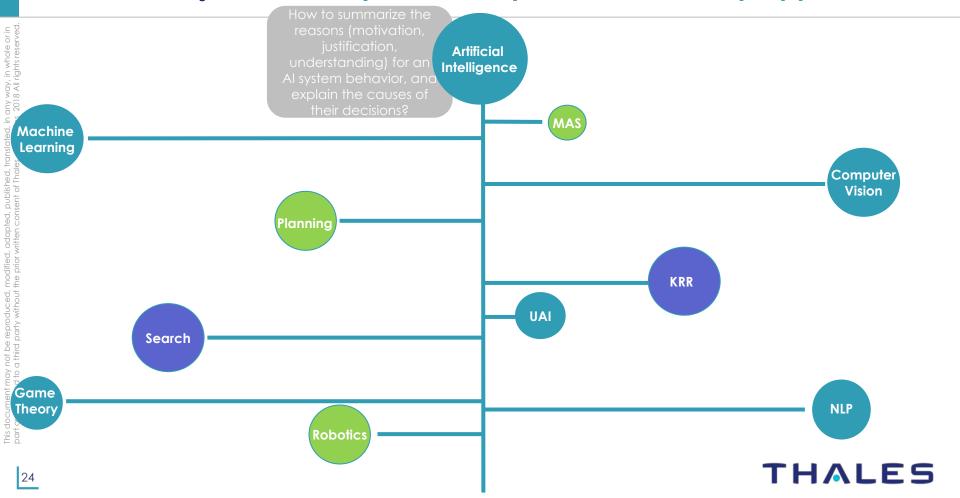
XAI in AI



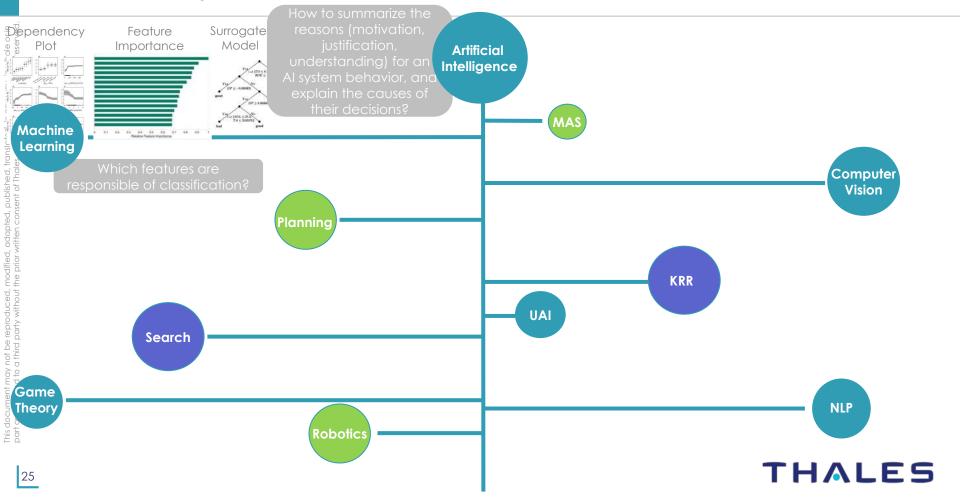
XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches



XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches



XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

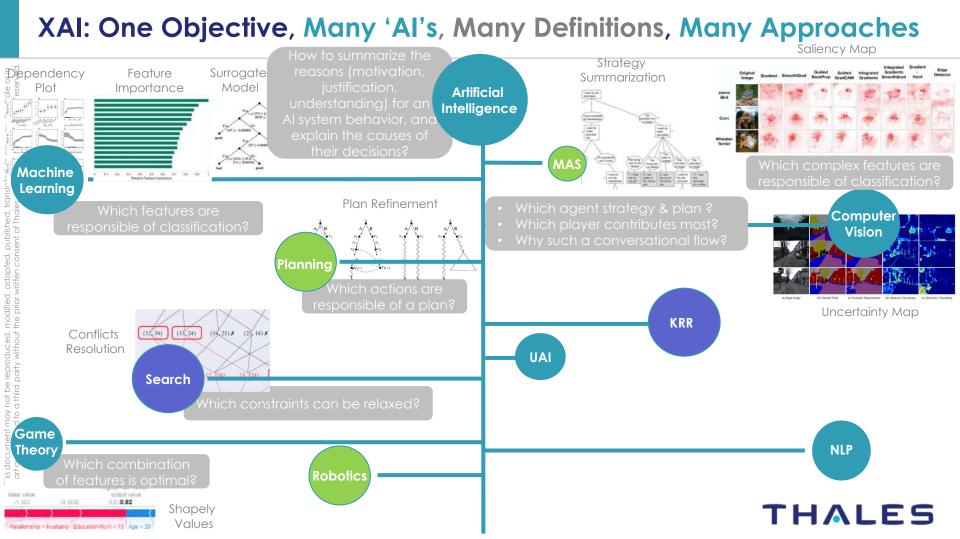


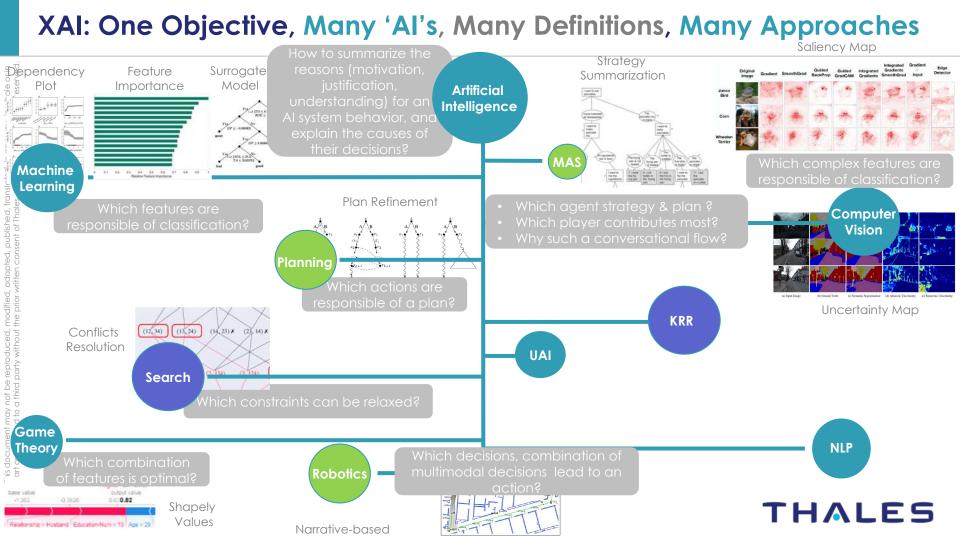
XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches Saliency Map Dependency Feature Surrogate Plot Importance Model **Artificial** Intelligence Al system behavior, and **Machine** Learning Computer Vision **Planning** Uncertainty Map KRR **UAI** Search Game NLP Theory Robotics THALES

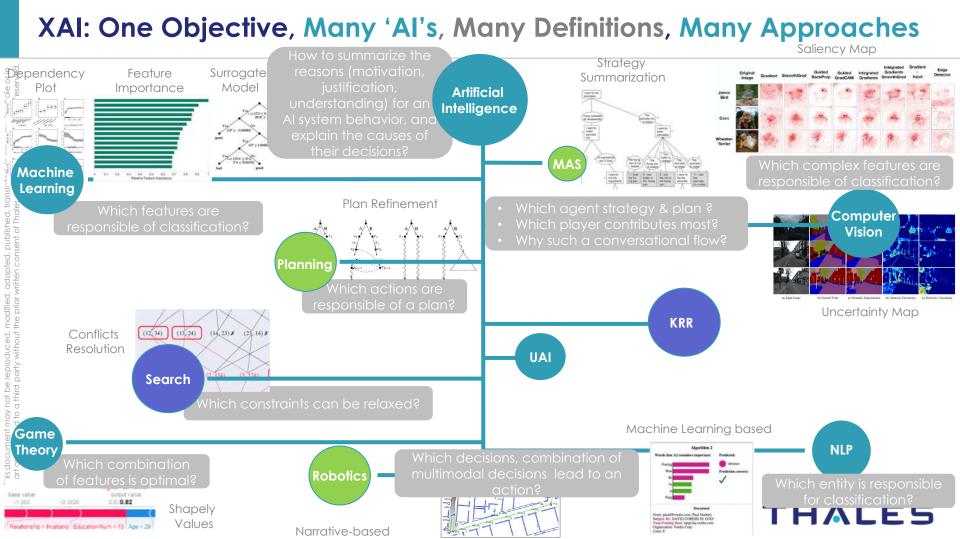
XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches Saliency Map Strategy Dependency Feature Surrogate Summarization Plot Importance Model **Artificial** Intelligence Al system behavior, and Machine Learnina Computer Vision **Planning** Uncertainty Map **KRR** UAI Search Game NLP Theory Robotics THALES

XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches Saliency Map Strategy Dependency Feature Surrogate Summarization Plot Importance Model **Artificial** Intelligence Al system behavior, and Machine Learnina Plan Refinement Computer Vision **Planning** Uncertainty Map **KRR** UAI Search Game NLP Theory **Robotics** THALES

XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches Saliency Map Strategy Dependency Feature Surrogate Summarization Plot Importance Model **Artificial** Intelliaence Al system behavior, and Machine Learnina Plan Refinement Computer Vision **Planning** Uncertainty Map KRR Conflicts (12, 34) (13, 24) (14, 23) x (23, 14) x Resolution UAI Search Game NLP Theory Robotics THALES







XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches Saliency Map Strategy Dependency Feature Surrogate Summarization Plot Importance Model **Artificial** Intelliaence Machine Learnina Plan Refinement Computer Vision Plannina Diagnosis Abduction Uncertainty Map **KRR** Conflicts (12, 34) (13, 24) (14, 23) x (23, 14) x Resolution UAI Search + (all p THING) = Machine Learning based Game **NLP** Theory Robotics Which entity is responsible Shapely Values Relationship = Husband Education-Num = 13 Age = 29 Narrative-based

XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches Saliency Map Strategy Dependency Feature Surrogate Summarization Plot Importance Model **Artificial** Intelliaence Al system behavior, and Machine Learnina Plan Refinement Computer Vision Plannina Diagnosis Abduction Uncertainty Map **KRR** Conflicts (12, 34) (13, 24) (14, 23) x (23, 14) x Resolution UAI Search) = THP to explanation
Machine Learning based Game **NLP** Theory Robotics Which entity is responsible Shapely Values Relationship = Husband Education-Num = 13 Age = 29 Narrative-based

Deep Dive



Machine Learning (except Artificial Neural Network)

Interpretable Models:

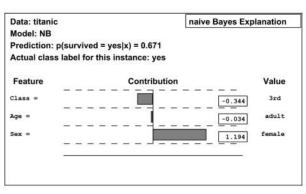
- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs



Machine Learning (except Artificial Neural Network)

Interpretable Models:

- Linear regression,
- · Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs



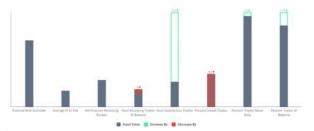
Naive Bayes model

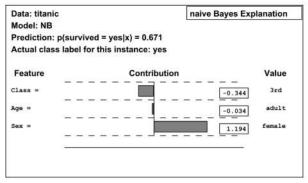


Machine Learning (except Artificial Neural Network)

Interpretable Models:

- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs





Naive Bayes model

Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. Artificial Intelligence in Medicine, 23:89–109, 2001.

Counterfactual What-if

Brent D. Mittelstadt, Chris Russell, Sandra Wachter: Explaining Explanations in Al. FAT 2019: 279-288

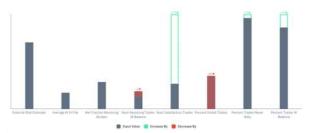
Rory Mc Grath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. CoRR abs/1811.05245 (2018)

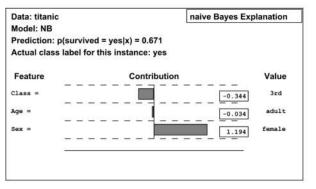


Machine Learning (except Artificial Neural Network)

Interpretable Models:

- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs





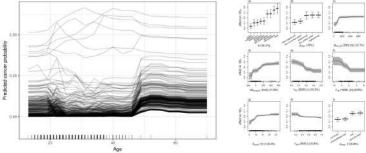
Naive Bayes model

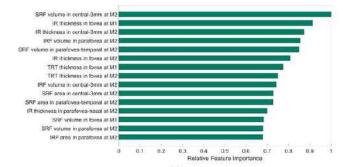
Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. Artificial Intelligence in Medicine, 23:89–109, 2001.

Counterfactual What-if

Brent D. Mittelstadt, Chris Russell, Sandra Wachter: Explaining Explanations in Al. FAT 2019: 279-288

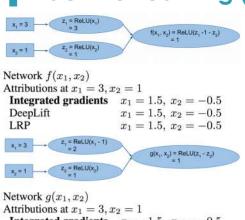
Rory Mc Grath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. CoRR abs/1811.05245 (2018)





Feature Importance Translation Partial Dependence Plot Individual Conditional Expectation Sensitivity Analysis

Machine Learning (only Artificial Neural Network)



Integrated gradients $x_1 = 1.5, x_2 = -0.5$ DeepLift $x_1 = 2, x_2 = -1$ LRP $x_1 = 2, x_2 = -1$

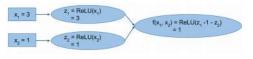
Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319-3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation 49ifferences. ICML 2017: 3145-3153



Machine Learning (only Artificial Neural Network)

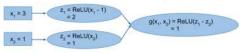


Network $f(x_1, x_2)$

Attributions at $x_1 = 3, x_2 = 1$

Integrated gradients $x_1 = 1.5, x_2 = -0.5$

DeepLift $x_1 = 1.5, x_2 = -0.5$ LRP $x_1 = 1.5, x_2 = -0.5$



Network $q(x_1, x_2)$

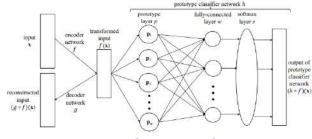
Attributions at $x_1 = 3, x_2 = 1$

Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features
Through Propagating Activation

Differences, ICML 2017: 3145-3153

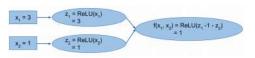


Auto-encoder

Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537



Machine Learning (only Artificial Neural Network)

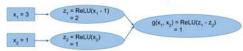


Network $f(x_1, x_2)$

Attributions at $x_1 = 3, x_2 = 1$

Integrated gradients $x_1 = 1.5, x_2 = -0.5$

DeepLift $x_1 = 1.5, x_2 = -0.5$ LRP $x_1 = 1.5, x_2 = -0.5$



Network $q(x_1, x_2)$

Attributions at $x_1 = 3, x_2 = 1$

Integrated gradients $x_1 = 1.5, x_2 = -0.5$ DeepLift $x_1 = 2, x_2 = -1$

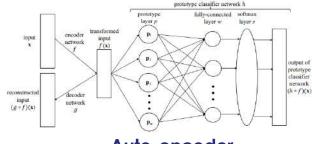
LRP $x_1 = 2, x_2 = -1$

Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

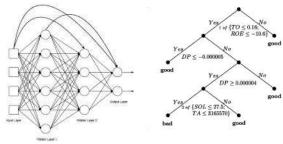
Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation

Bifferences. ICML 2017: 3145-3153



Auto-encoder

Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537

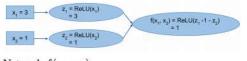


Surogate Model

Mark Craven, Jude W. Shavlik: Extracting Tree-Structured Representations of Trained Networks. NIPS 1995: 24-30

THALES

Machine Learning (only Artificial Neural Network)

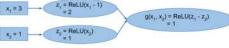


Network $f(x_1, x_2)$

Attributions at $x_1 = 3, x_2 = 1$

Integrated gradients $x_1 = 1.5, x_2 = -0.5$ DeepLift $x_1 = 1.5, x_2 = -0.5$

LRP $x_1 = 1.5, x_2 = -0.5$



Network $g(x_1, x_2)$

LRP

Attributions at $x_1 = 3, x_2 = 1$

Integrated gradients $x_1 = 1.5, x_2 = -0.5$ DeepLift $x_1 = 2, x_2 = -1$

 $x_1 = 2, x_2 = -1$

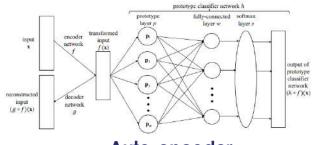
Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation 49ifferences. ICML 2017: 3145-3153

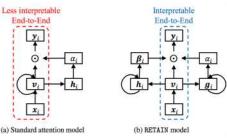
Attention Mechanism

D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. International Conference on Learning Representations, 2015

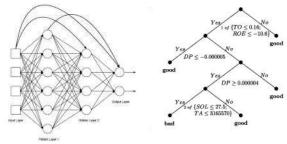


Auto-encoder

Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537



Edward Choi, Mohammad Taha Bahadori, Jimeng Sun, Joshua Kulas, Andy Schuetz, Walter F. Stewart: RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time Attention Mechanism. NIPS 2016: 3504-3512



Surogate Model

Mark Craven, Jude W. Shavlik: Extracting Tree-Structured Representations of Trained Networks. NIPS 1995: 24-30

THALES

Computer Vision



Interpretable Units

res5c unit 1243
res5c unit 1379
inception_4e unit 92

Airplane

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017: 3319-3327



Computer Vision Train

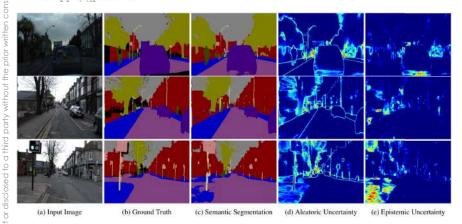


Interpretable Units

Airplane



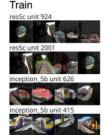
David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations, CVPR 2017: 3319-3327



Uncertainty Map



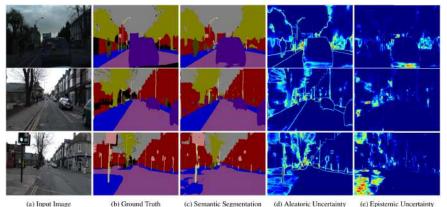
Computer Vision



Interpretable Units

Airplane res5c unit 1243 res5c unit 1379 nception 4e unit 92

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations, CVPR 2017: 3319-3327



Uncertainty Map

Description: This is a large bird with a white neck and a black back in the water Class Definition: The Western Grebe is a waterbird with a yellow pointy beak, white neck and bell-



Explanation: This is a Western Grebe because this bird has a long white neck, pointy yellow beak

Laysan Albatross

Description: This is a large flying bird with black wings and a white belly. Class Definition: The Laysan Albatross is a large seabird with a hooked yellow beak, black back

Visual Explanation: This is a Laysan Albatross because this bird has a large wingspan, hooked vellow beak, and white belly

Laysan Albatross Description: This is a large bird with a white neck and a black back in the water. Class Definition: The Laysan Albatross is a large seabird with a hooked yellow beak, black back and white belly.

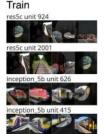
Visual Explanation: This is a Laysan Albatross because this bird has a hooked yellow beak white neck and black back

Visual Explanation

Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, Trevor Darrell: Generating Visual Explanations. ECCV (4) 2016: 3-19



Computer Vision

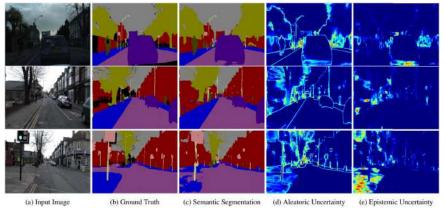


Interpretable Units

Airplane res5c unit 1243 res5c unit 1379

nception 4e unit 92

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations, CVPR 2017: 3319-3327



Uncertainty Map

Alex Kendall, Yarin Gal: What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? NIPS 2017: 5580-5590

Description: This is a large bird with a white neck and a black back in the water.

Class Definition: The Western Grebe is a waterbird with a yellow pointy beak, white neck and belland black back

Explanation: This is a Western Grebe because this bird has a long white neck, pointy yellow beak

Laysan Albatross Description: This is a large flying bird with black wings and a white belly.

Class Definition: The Laysan Albatross is a large seabird with a hooked yellow beak, black back

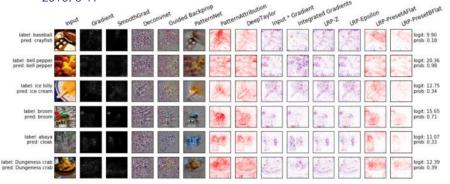
Visual Explanation: This is a Laysan Albatross because this bird has a large wingspan, hooked vellow beak, and white belly,

Laysan Albatross Description: This is a large bird with a white neck and a black back in the water. Class Definition: The Laysan Albatross is a large seabird with a hooked yellow beak, black back

and white belly. Visual Explanation: This is a Laysan Albatross because this bird has a hooked yellow beak white neck and black back

Visual Explanation

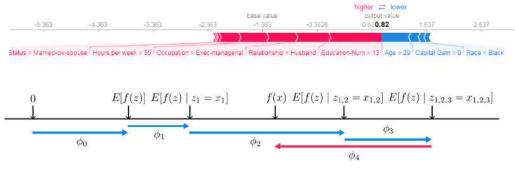
Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, Trevor Darrell: Generating Visual Explanations. ECCV (4) 2016: 3-19



Saliency Map

Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, Been Kim: Sanity Checks for Saliency Marps Neurl 2018:

Game Theory

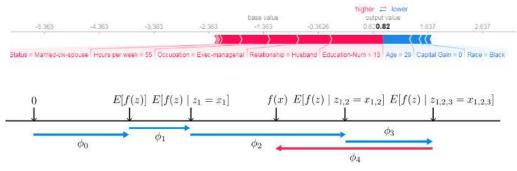


Shapley Additive Explanation

Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions. NIPS 2017: 4768-4777

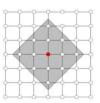


Game Theory

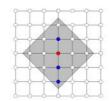


Shapley Additive Explanation

Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions. NIPS 2017: 4768-4777





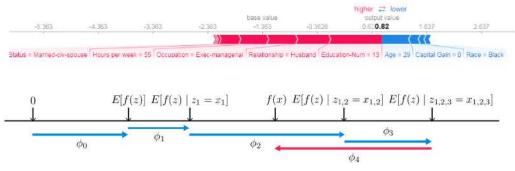


L-Shapley and C-Shapley (with graph structure)

Jianbo Chen, Le Song, Martin J. Wainwright, Michael I. Jordan: L-Shapley and C-Shapley: Efficient Model Interpretation for Structured Data. ICLR 2019

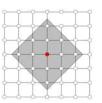


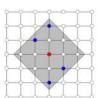
Game Theory

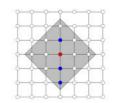


Shapley Additive Explanation

Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions. NIPS 2017: 4768-4777







L-Shapley and C-Shapley (with graph structure)

Jianbo Chen, Le Song, Martin J. Wainwright, Michael I. Jordan: L-Shapley and C-Shapley: Efficient Model Interpretation for Structured Data. ICLR 2019

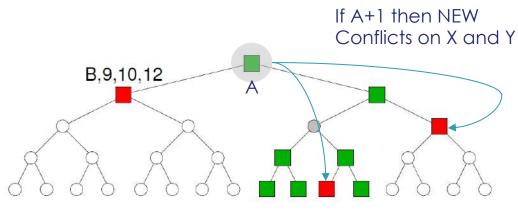
~ instancewise feature importance (causal influence)

Erik Štrumbelj and Igor Kononenko. An efficient explanation of individual classifications using game theory. Journal of Machine Learning Research, 11:1–18, 2010.

Anupam Datta, Shayak Sen, and Yair Zick. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In Security and Privacy (SP), 2016 IEEE Symposium on, pp. 598–617. IEEE, 2016.



Search and Constraint Satisfaction



Conflicts resolution

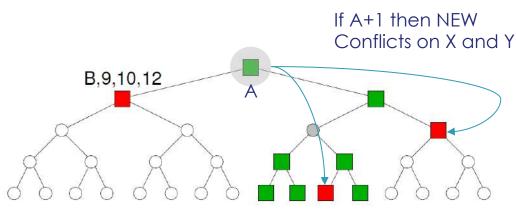
Barry O'Sullivan, Alexandre Papadopoulos, Boi Faltings, Pearl Pu: Representative Explanations for Over-Constrained Problems. AAAI 2007: 323-328

Robustness Computation

Hebrard, E., Hnich, B., & Walsh, T. (2004, July). Robust solutions for constraint satisfaction and optimization. In ECAI (Vol. 16, p. 186).



Search and Constraint Satisfaction

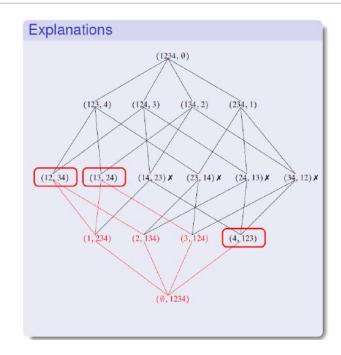


Conflicts resolution

Barry O'Sullivan, Alexandre Papadopoulos, Boi Faltings, Pearl Pu: Representative Explanations for Over-Constrained Problems. AAAI 2007: 323-328

Robustness Computation

Hebrard, E., Hnich, B., & Walsh, T. (2004, July). Robust solutions for constraint satisfaction and optimization. In ECAI (Vol. 16, p. 186).



Constraints relaxation

Ulrich Junker: QUICKXPLAIN: Preferred Explanations and Relaxations for Over-Constrained Problems. AAAI 2004: 167-172



Knowledge Representation and Reasoning

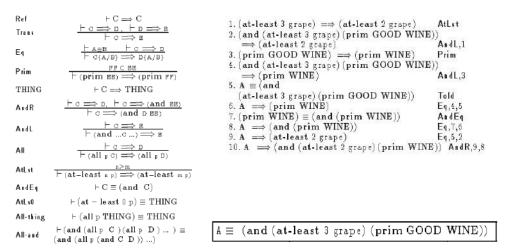
```
Ref
                                                       1. (at-least 3 grape) =⇒ (at-least 2 grape)
                                                                                                             AtLst
Trans
                                                       2. (and (at-least 3 grape) (prim GOOD WINE))
                                                                                                             AndL,1
                                                          ⇒ (at-least 2 grape)
                                                       3. (prim GOOD WINE) ==> (prim WINE)
                                                                                                             Prim
                  C(A/B) ==> D(A/B)
                                                      4. (and (at-least 3 grape) (prim GOOD WINE))
Prim
                                                          ⇒ (prim WINE)
                                                                                                             AndL.3
              \vdash (prim EE) \Longrightarrow (prim FF)
                                                       5. A = (and
                    \vdash C \Longrightarrow THING
THING
                                                                                                             Told
                                                          (at-least 3 grape) (prim GOOD WINE))
             \vdash c \Longrightarrow D, \vdash c \Longrightarrow (and EE)
                                                       6. A =⇒ (prim WINE)
                                                                                                             Eq.4.5
AndR
                  F c ⇒ (and DEE)
                                                       7. (prim WINE) = (and (prim WINE))
                                                                                                             AndEq
                                                      8. A =⇒ (and (prim WINE))
                                                                                                             Eq.7,6
AndL
                                                       9. A =⇒ [at-least 2 grape]
                                                                                                             Eq.5.2
                                                      10. A ⇒ (and (at-least 2 grape) (prim WINE)) AndR,9,8
                \vdash (all \mid g \mid G) \Longrightarrow (all \mid g \mid B)
AtLst
          F(at-least np) ⇒ (at-least mp)
                    \vdash C \equiv (and C)
AndEa
Atl s0
              \vdash (at - least 0 p) \equiv THING
              \vdash (all p THING) \equiv THING
All-thing
                                                   A \equiv (and (at\text{-least 3 grape}) (prim GOOD WINE))
          \vdash (and (all p C)(all p D)...) \equiv
          (and (all p (and C D )) ...)
```

Explaining Reasoning (through Justification) e.g., Subsumption

Deborah L. McGuinness, Alexander Borgida: Explaining Subsumption in Description Logics. IJCAI (1) 1995: 816-821

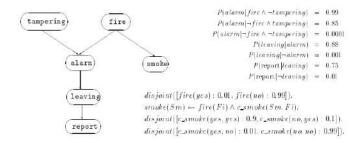


Knowledge Representation and Reasoning



Explaining Reasoning (through Justification) e.g., Subsumption

Deborah L. McGuinness, Alexander Borgida: Explaining Subsumption in Description Logics. IJCAI (1) 1995: 816-821

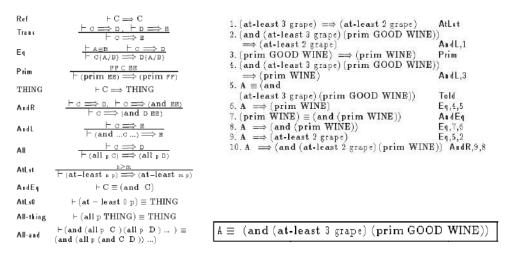


Abduction Reasoning (in Bayesian Network)

David Poole: Probabilistic Horn Abduction and Bayesian Networks. Artif. Intell. 64(1): 81-129 (1993)

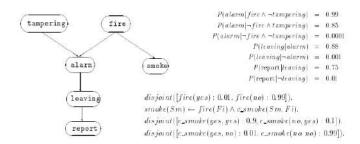


Knowledge Representation and Reasoning



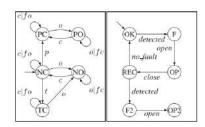
Explaining Reasoning (through Justification) e.g., Subsumption

Deborah L. McGuinness, Alexander Borgida: Explaining Subsumption in Description Logics. IJCAI (1) 1995: 816-821



Abduction Reasoning (in Bayesian Network)

David Poole: Probabilistic Horn Abduction and Bayesian Networks. Artif. Intell. 64(1): 81-129 (1993)



Diagnosis Inference

Alban Grastien, Patrik Haslum, Sylvie Thiébaux: Conflict-Based Diagnosis of Discrete Event Systems: Theory and Practice. KR 2012

Multi-agent Systems

MAS INFRASTRUCTURE	INDIVIDUAL AGENT INFRASTRUCTURE		
MAS INTEROPERATION Translation Services Interoperation Services	INTEROPERATION Interoperation Modules		
CAPABILITY TO AGENT MAPPING Middle Agents	CAPABILITY TO AGENT MAPPING Middle Agents Components		
NAME TO LOCATION MAPPING ANS	NAME TO LOCATION MAPPING ANS Component		
SECURITY Certificate Authority Cryptographic Services	SECURITY Security Module private/public Keys		
PERFORMANCE SERVICES MAS Monitoring Reputation Services	PERFORMANCE SERVICES Performance Services Modules		
MULTIAGENT MANAGEMENT SERVICES Logging, Acivity Visualization, Launching	MANAGEMENT SERVICES Logging and Visualization Components		
ACL INFRASTRUCTURE Public Ontology Protocols Servers	ACL INFRASTRUCTURE ACL Parser Private Ontology Protocol Engine		
COMMUNICATION INFRASTRUCTURE Discovery Message Transfer	COMMUNICATION MODULES Discovery Component Message Tranfer Module		

Explanation of Agent Conflicts & Harmful Interactions

Katia P. Sycara, Massimo Paolucci, Martin Van Velsen, Joseph A. Giampapa: The RETSINA MAS Infrastructure. Autonomous Agents and Multi-Agent Systems 7(1-2): 29-48 (2003)

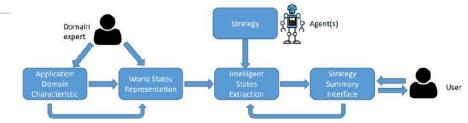


Multi-agent Systems

MAS INFRASTRUCTURE	INDIVIDUAL AGENT INFRASTRUCTURE		
MAS INTEROPERATION Translation Services Interoperation Services	INTEROPERATION Interoperation Modules CAPABILITY TO AGENT MAPPING Middle Agents Components NAME TO LOCATION MAPPING ANS Component		
CAPABILITY TO AGENT MAPPING Middle Agents			
NAME TO LOCATION MAPPING ANS			
SECURITY Certificate Authority Cryptographic Services	SECURITY Security Module private/public Keys		
PERFORMANCE SERVICES MAS Monitoring Reputation Services	PERFORMANCE SERVICES Performance Services Modules		
MULTIAGENT MANAGEMENT SERVICES Logging, Acivity Visualization, Launching	MANAGEMENT SERVICES Logging and Visualization Components		
ACL INFRASTRUCTURE Public Ontology Protocols Servers	ACL INFRASTRUCTURE ACL Parser Private Ontology Protocol Engine		
COMMUNICATION INFRASTRUCTURE Discovery Message Transfer	COMMUNICATION MODULES Discovery Component Message Tranfer Module		

Explanation of Agent Conflicts & Harmful Interactions

Katia P. Sycara, Massimo Paolucci, Martin Van Velsen, Joseph A. Giampapa: The RETSINA MAS Infrastructure. Autonomous Agents and Multi-Agent Systems 7(1-2): 29-48 (2003)



Agent Strategy Summarization

Ofra Amir, Finale Doshi-Velez, David Sarne: Agent Strategy Summarization. AAMAS 2018: 1203-1207

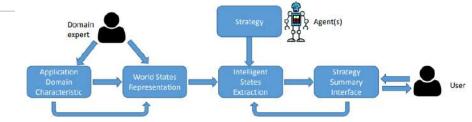


Multi-agent Systems

MAS INFRASTRUCTURE INDIVIDUAL AGENT INFRASTRUCTURE MAS INTEROPERATION INTEROPERATION Translation Services Interoperation Services Interoperation Modules CAPABILITY TO AGENT MAPPING CAPABILITY TO AGENT MAPPING Middle Agents Middle Agents Components NAME TO LOCATION MAPPING NAME TO LOCATION MAPPING ANS Component SECURITY SECURITY Certificate Authority Cryptographic Services Security Module private/public Keys PERFORMANCE SERVICES PERFORMANCE SERVICES Reputation Services Performance Services Modules MULTIAGENT MANAGEMENT SERVICES MANAGEMENT SERVICES Logging, Acivity Visualization, Launching Logging and Visualization Components **ACL INFRASTRUCTURE** ACL INFRASTRUCTURE Public Ontology ACL Parser Private Ontology Protocol Engine COMMUNICATION INFRASTRUCTURE COMMUNICATION MODULES Discovery Component Message Tranfer Module Discovery Message Transfer **OPERATING ENVIRONMENT** Machines, OS, Network Multicast Transport Laver: TCP/IP, Wireless, Infrared, SSL

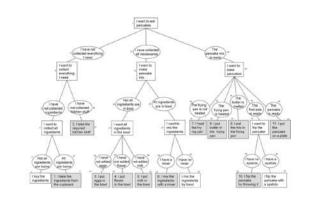
Explanation of Agent Conflicts & Harmful Interactions

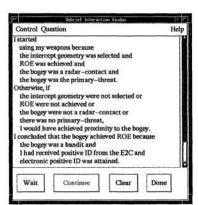
Katia P. Sycara, Massimo Paolucci, Martin Van Velsen, Joseph A. Giampapa: The RETSINA MAS Infrastructure. Autonomous Agents and Multi-Agent Systems 7(1-2): 29-48 (2003)



Agent Strategy Summarization

Ofra Amir, Finale Doshi-Velez, David Sarne: Agent Strategy Summarization. AAMAS 2018: 1203-1207

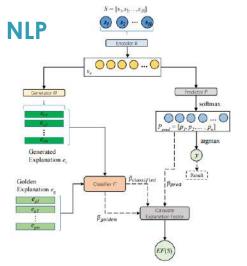




Explainable Agents

Joost Broekens, Maaike Harbers, Koen V. Hindriks, Karel van den Bosch, Catholijn M. Jonker, John-Jules Ch. Meyer: Do You Get It? User-Evaluated Explainable BDI Agents. MATES 2010: 28-39 W. Lewis Johnson: Agents that Learn to Explain Themselves. AAAI 1994: 1257-1263





Fine-grained explanations are in the form of:

- texts in a realworld dataset;
- Numerical scores

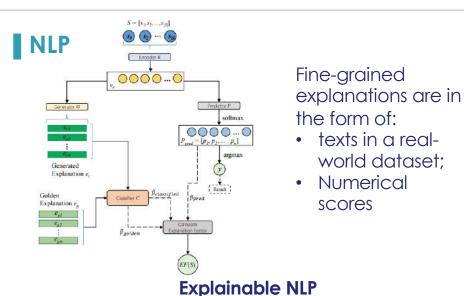
Explainable NLP

Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

Hendrik Strobelt, Sebastian Gehrmann, Hanspeter Pfister, Alexander M. Rush: LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks. IEEE Trans. Vis. Comput. Graph. 24(1): 667-676 (2018)

Hendrik Strobelt, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, Alexander M. Rush: Seq2seq-Vis: A Visual Debugging Tool for Sequence-to-Sequence Models. IEEE Trans. Vis. Comput. Graph. 25(1): 353-363 (2019)

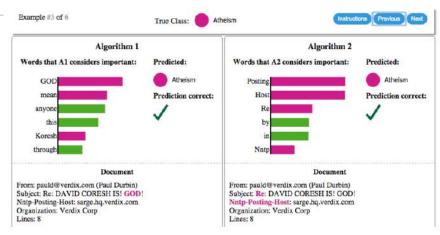




Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

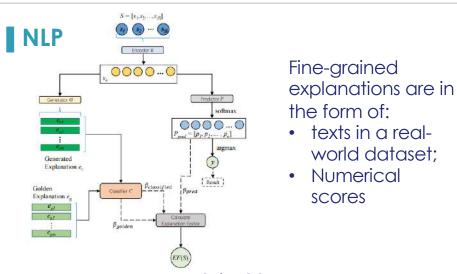
Hendrik Strobelt, Sebastian Gehrmann, Hanspeter Pfister, Alexander M. Rush: LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks. IEEE Trans. Vis. Comput. Graph. 24(1): 667-676 (2018)

Hendrik Strobelt, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, Alexander M. Rush: Seq2seq-Vis: A Visual Debugging Tool for Sequence-to-Sequence Models. IEEE Trans. Vis. Comput. Graph. 25(1): 353-363 (2019)



LIME for NLP

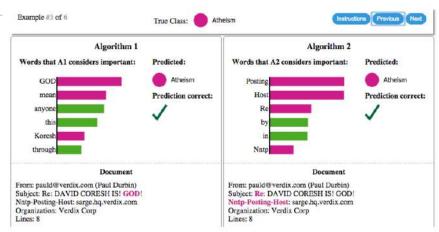
Marco Túlio Ribeiro, Sameer Singh, Carlos Guestrin: "Why Should I Trust You?": Explaining the Predictions of Any Classifier. KDD 2016: 1135-1144



Explainable NLP

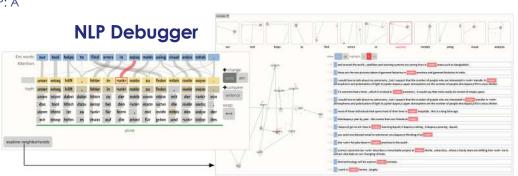
Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

Hendrik Strobelt, Sebastian Gehrmann, Hanspeter Pfister, Alexander M. Rush: LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks. IEEE Trans. Vis. Comput. Graph. 24(1): 667-676 (2018) Hendrik Strobelt, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, Alexander M. Rush: Seq2seq-Vis: A Visual Debugging Tool for Sequence-to-Sequence Models. IEEE Trans. Vis. Comput. Graph. 25(1): 353-363 (2019)



LIME for NLP

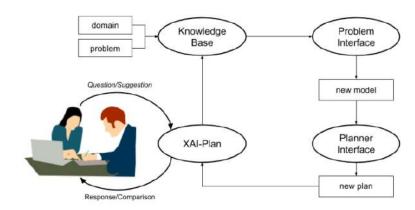
Marco Túlio Ribeiro, Sameer Singh, Carlos Guestrin: "Why Should I Trust You?": Explaining the Predictions of Any Classifier. KDD 2016: 1135-1144



Planning and Scheduling

Explanation Type		R2	R3	R4
Plan Patch Explanation / VAL	×	1	×	1
Model Patch Explanation		X	/	1
Minimally Complete Explanation		1	×	?
Minimally Monotonic Explanation		1	1	?
(Approximate) Minimally Complete Explanation		1	X	1

Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for Al Planner Decisions. CoRR abs/1810.06338 (2018)



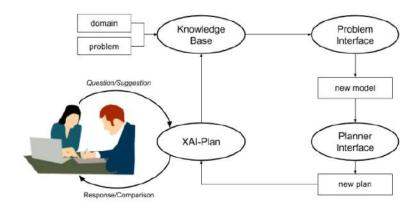
XAI Plan

THALES

Planning and Scheduling

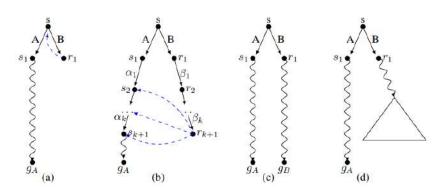
Explanation Type		R2	R3	R4
Plan Patch Explanation / VAL	×	1	×	1
Model Patch Explanation		X	1	1
Minimally Complete Explanation		1	×	?
Minimally Monotonic Explanation		1	1	?
(Approximate) Minimally Complete Explanation		1	X	1

Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for Al Planner Decisions. CoRR abs/1810.06338 (2018)



XAI Plan

Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for Al Planner Decisions. CoRR abs/1810.06338 (2018)



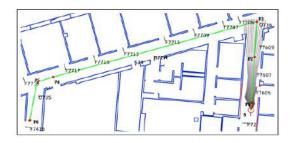
Human-in-the-loop Planning

Maria Fox, Derek Long, Daniele Magazzeni: Explainable Planning. CoRR abs/1709.10256 (2017)

(Manual) Plan Comparison



Robotics



			ction, A		
Specificity, S		Level 1	Level 2	Level 3	Level 4
	General Picture	Start and finish point of the complete route	Total distance and time taken for the complete route	Total distance and time taken for the complete route	Starting and ending land- mark of complete route
	Summary	Start and finish point for subroute on each floor of each building	Total distance and time taken for subroute on each floor of each build- ing	Total distance and angles for subroute on each floor of each building	Starting and ending land- mark for subroute on each floor of each build- ing
	Detailed	Start and finish points of complete route plus time taken for each edge of route	Angle turned at each point plus the total dis- tance and time taken for each edge of route	Turn direction at each point plus total distance for each edge of route	All landmarks encoun- tered on the route

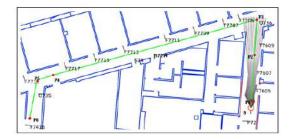
Narration of Autonomous Robot Experience

Stephanie Rosenthal, Sai P Selvaraj, and Manuela Veloso. Verbalization: Narration of autonomous robot experience. In IJCAI, pages 862–868. AAAI Press, 2016.

Daniel J Brooks et al. 2010. Towards State Summarization for Autonomous Robots.. In AAAI Fall Symposium: Dialog with Robots, Vol. 61. 62.



Robotics



		Abstraction, A					
		Level 1	Level 2	Level 3	Level 4		
Specificity, S	General Picture	Start and finish point of the complete route	Total distance and time taken for the complete route	Total distance and time taken for the complete route	Starting and ending land- mark of complete route		
	Summary	Start and finish point for subroute on each floor of each building	Total distance and time taken for subroute on each floor of each build- ing	Total distance and angles for subroute on each floor of each building	Starting and ending land- mark for subroute on each floor of each build- ing		
	Detailed Narrative	Start and finish points of complete route plus time taken for each edge of route	Angle turned at each point plus the total dis- tance and time taken for each edge of route	Turn direction at each point plus total distance for each edge of route	All landmarks encoun- tered on the route		

Narration of Autonomous Robot Experience

Stephanie Rosenthal, Sai P Selvaraj, and Manuela Veloso. Verbalization: Narration of autonomous robot experience. In IJCAI, pages 862–868. AAAI Press, 2016.

Daniel J Brooks et al. 2010. Towards State Summarization for Autonomous Robots.. In AAAI Fall Symposium: Dialog with Robots, Vol. 61. 62.

Robot: I have decided to turn left. **Human:** Why did you do that?

Robot: I believe that the correct action is to turn left

BECAUSE:

I'm being asked to go forward

AND This area in front of me was 20 cm higher than me *highlights area*

AND the area to the left has maximum protrusions of less than 5 cm *highlights area*

AND I'm tilted to the right by more than 5 degrees.

Here is a display of the path through the tree that lead to this decision. *displays tree*

Human: How confident are you in this decision?

Robot: The distribution of actions that reached this leaf node is shown in this histogram. *displays histogram* This action is predicted to be correct 67% of the time.

Human: Where did the threshold for the area in front come from?

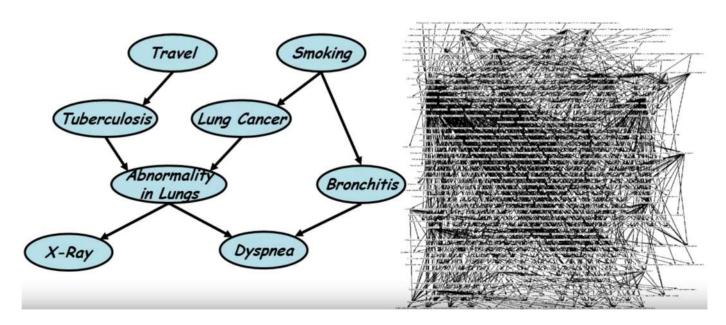
Robot: Here is the histogram of all training examples that reached this leaf. 80% of examples where this area was above 20 cm predicted the appropriate action to be "drive forward".

From Decision Tree to human-friendly information

Raymond Ka-Man Sheh: "Why Did You Do That?" Explainable Intelligent Robots. AAAI Workshops 2017



Reasoning under uncertainty



Probabilistic Graphical Models

Daphne Koller, Nir Friedman: Probabilistic Graphical Models - Principles and Techniques. MIT Press 2009, ISBN 978-0-262-01319-2, pp. I-XXXV, 1-1231



Evaluation



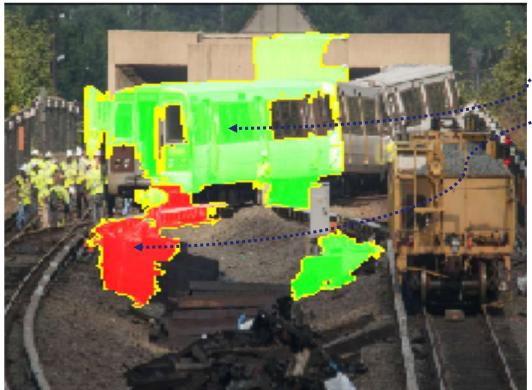
XAI: One Objective, Many Metrics





On the role of Knowledge Graphs in Explainable Machine Learning

Knowledge Graph Embeddings in Machine Learning

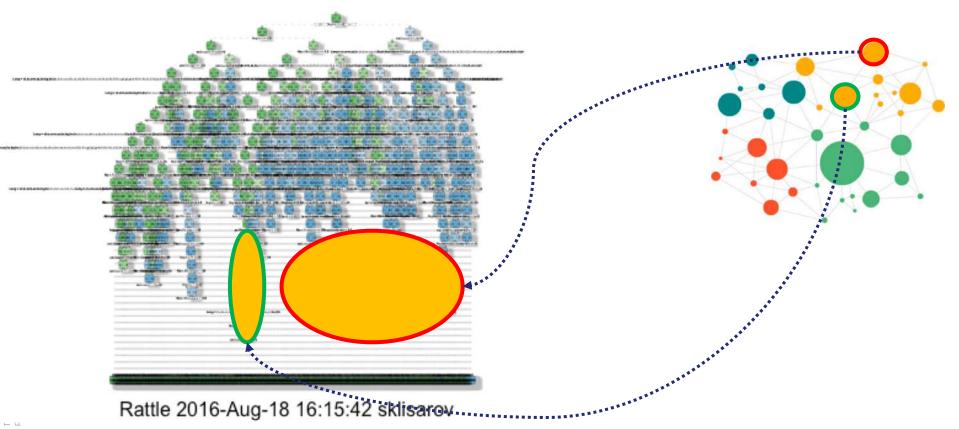




https://stats.stackexchange.com/questions/23058 1/decision-tree-too-large-to-interpret



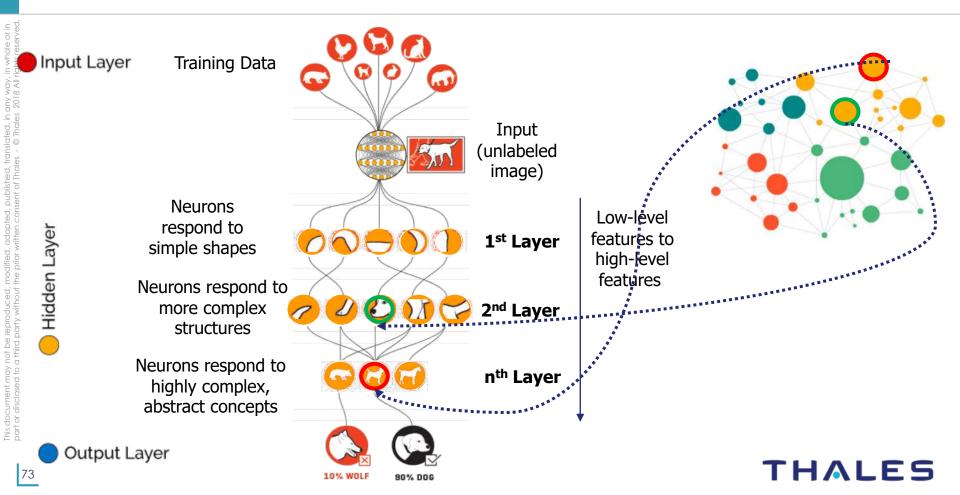
Knowledge Graph for Decision Trees



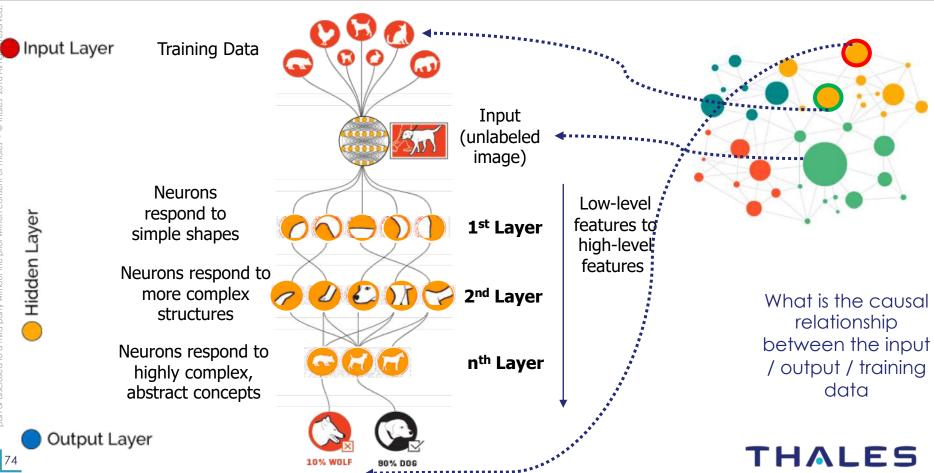
https://stats.stackexchange.com/questions/23058 1/decision-tree-too-large-to-interpret



Knowledge Graph for Deep Neural Network (1)



Knowledge Graph for Deep Neural Network (2)



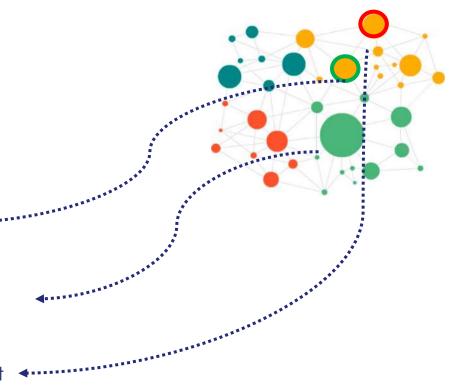
Knowledge Graph for Personalized XAI



Description 1: This is an orange train accident

Description 2: This is an train accident between two speed merchant trains of characteristics X43-B and Y33-C in a dry environment

Description 3: This is a public transportation accident



"How to explain transfer learning with appropriate knowledge representation?

Proceedings of the Sixteenth International Conference on Principles of Knowledge Representation and Reasoning (KR 2018)

Knowledge-Based Transfer Learning Explanation

Jiaoyan Chen

Department of Computer Science University of Oxford, UK

Jeff Z. Pan

Department of Computer Science University of Aberdeen, UK

Freddy Lecue

INRIA, France Accenture Labs, Ireland

Ian Horrocks

Department of Computer Science University of Oxford, UK

Huajun Chen

College of Computer Science, Zhejiang University, China Alibaba-Zhejian University Frontier Technology Research Center

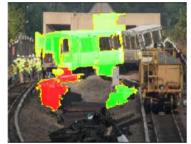


Applications



Obstacle Identification Certification (Trust) - Transportation

















Challenge: Public transportation is getting more and more self-driving vehicles. Even if trains are getting more and more autonomous, the human stays in the loop for critical decision, for instance in case of obstacles. In case of obstacles trains are required to provide recommendation of action i.e., go on or go back to station. In such a case the human is required to validate the recommendation through an explanation exposed by the train or machine.

Al Technology: Integration of Al related technologies i.e., Machine Learning (Deep Learning / CNNs), and semantic segmentation.

XAI Technology: Deep learning and Epistemic uncertainty



Explainable On-Time Performance - Transportation

in whole or in ights reserved.

KLM /	Transavia	Flight	Delay	Prediction
-------	-----------	--------	-------	------------

PLANE INFO ARRIVAL		TURNAROUND		DEPARTURE								
Status / Aircraft	Flight	ETA	Status	Delay Code	Gate	Slot	Progress	Milestones	Flight	ETA	Status	Delay Code
o urtwet ~	4567	16:30	Scheduled		345345	1			5678	19:00	Scheduled	
O idalaw v	4567	18:30	Delayed	ABC, DEF, GHI	345345	1			5678	19.00	Delayed	ABC, DEF, GH
O maide v	4567	18.00	Scheduled	ABC, DEF, GHI	345345	1			5678	19:00	Scheduled	ABC, DEF, GH
Ø kshdks ♥	4567	196	Cancelled	ABC, DEF, OHI	1.00				5678	34	Cancelled	ABC, DEE, OH
vallasen 0	4567	18.35	Delayed	ABC, DEF, GHI	345345	1			5678	19:00	Delayed	ABC, DEF, GH
O odlobs ~	4567	1630	Delayed	ABC, DEF, GHI	345345	35			5678	19:00	Scheduled	ABC, DEF, GH
aedbsc v	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1			5678	19:00	Scheduled	ABC, DEF, GH
O sedies V	4567	36.30	Scheduled	ABC, DEF, GHI	345345	1			5678	19.00	Scheduled	ABC, DEF, GH
O nedles V	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1			5678	19:00	Scheduled	ABC, DEF, GH
O aediss V	4567	18:30	Scheduled	ABC, DEF, OHI	345345	1			5678	19.00	Scheduled	ABC, DEF, GH
O neding ~	A567	18:30	Scheduled	ABC, DEF, GHI	345345	1			5678	19:00	Scheduled	ABC, DEF, GH
O seifies V	4567	16:30	Scheduled	ABC, DEE, GHE	345345	1			5678	19:00	Scheduled	ABC, DEE, GH
O sedbec v	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1.			5678	19:00	Scheduled	ABC, DEF, GH
O sedbec ~	4567	18:30	Scheduled	ABC, DEF, GHI	345345	3			5678	19:00	Scheduled	ABC, DEF, OH
aedbsc v	4567	18.30	Scheduled	ABC, DEF, OHI	345345	1	- 80		5678	19:00	Scheduled	ABC, DEF, GH

Challenge: Globally 323,454 flights are delayed every year. Airline-caused delays totaled 20.2 million minutes last year, generating huge cost for the company. Existing in-house technique reaches 53% accuracy for **predicting flight delay**, does not provide any time estimation (in **minutes** as opposed to True/False) and is unable to capture the underlying reasons (explanation).

Al Technology: Integration of Al related technologies i.e., Machine Learning (Deep Learning / Recurrent neural Network), Reasoning (through semantics-augmented case-based reasoning) and Natural Language Processing for building a robust model which can (1) predict flight delays in minutes, (2) explain delays by comparing with historical cases.

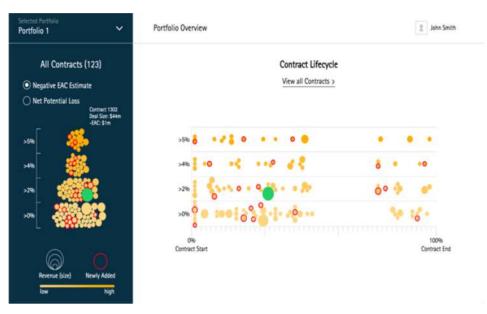
XAI Technology: Knowledge graph embedded Sequence Learning using LSTMs

aoyan Chen, Freddy Lécué, Jeff Z. Pan, Ian Horrocks, Huajun Chen: Knowledge-Based Transfer Learning Explanation. KR 2018: 349-358

Nicholas McCarthy, Mohammad Karzand, Freddy Lecue: Amsterdam to Dublin Eventually Delayed? LSTM and Transfer Learning for Predicting Delays of Low Cost Airlines: AAAI 2019



Explainable Risk Management - Finance



Jiewen Wu, Freddy Lécué, Christophe Guéret, Jer Hayes, Sara van de Moosdijk, Gemma Gallagher, Peter McCanney, Eugene Eichelberger: Personalizing Actions in Context for Risk Management Using Semantic Web Technologies. International Semantic Web Conference (2) 2017: 367-383

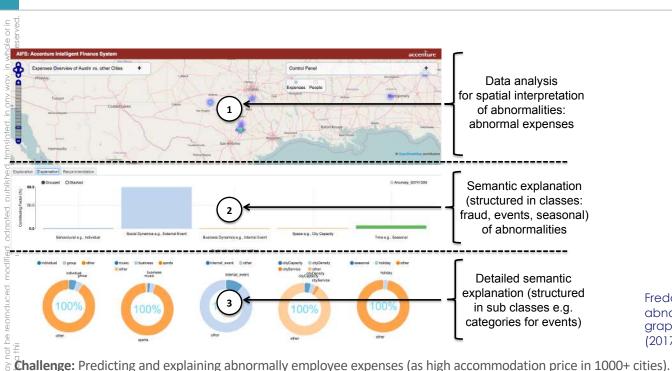
Challenge: Accenture is managing every year more than 80,000 opportunities and 35,000 contracts with an expected revenue of \$34.1 billion. Revenue expectation does not meet estimation due to the complexity and risks of critical contracts. This is, in part, due to the (1) large volume of projects to assess and control, and (2) the existing non-systematic assessment process.

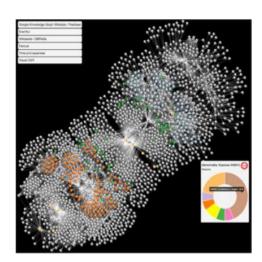
Al Technology: Integration of Al technologies i.e., Machine Learning, Reasoning, Natural Language Processing for building a robust model which can (1) predict revenue loss, (2) recommend corrective actions, and (3) explain why such actions might have a positive impact.

XAI Technology: Knowledge graph embedded Random Forrest



Explainable anomaly detection – Finance (Compliance)





Freddy Lécué, Jiewen Wu: Explaining and predicting abnormal expenses at large scale using knowledge graph based reasoning. J. Web Sem. 44: 89-103 (2017)

The mention of the state of the

Technology: Various techniques have been matured over the last two decades to achieve excellent results. However most methods address the problem from a statistic and pure data-centric angle, which in turn limit any interpretation. We elaborated a web application running live with real data from (i) travel and expenses from Accenture, (ii) external data from third party such as Google Knowledge Graph, DBPedia (relational DataBase version of Wikipedia) and social events from Eventful, for explaining abnormalities.

81

KAI Technology: Knowledge graph embedded Ensemble Learning

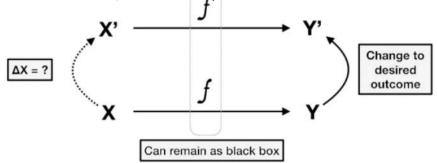
Counterfactual Explanations for Credit Decisions (1) - Finance

- Local, post-hoc, contrastive explanations of black-box classifiers
- Required minimum change in input vector to flip the decision of the classifier.
- Interactive Contrastive Explanations

Challenge: We predict loan applications with off-the-shelf, interchangeable black-box estimators, and we explain their predictions with counterfactual explanations. In counterfactual explanations the model itself remains a black box; it is only through changing inputs and outputs that an explanation is obtained.

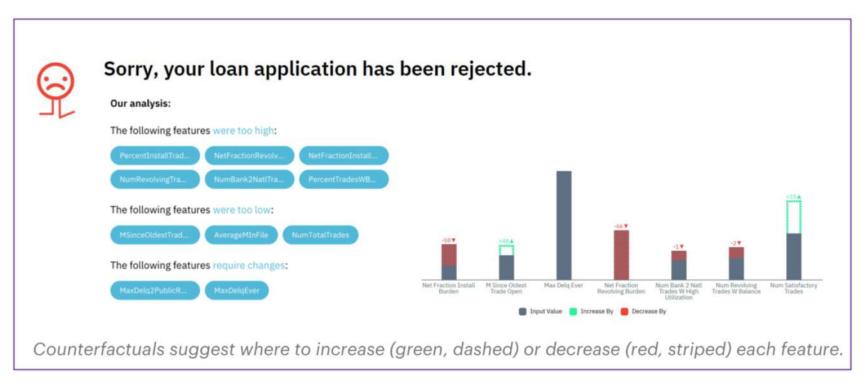
Al Technology: Supervised learning, binary classification.

XAI Technology: Post-hoc explanation, Local explanation, Counterfactuals, Interactive explanations



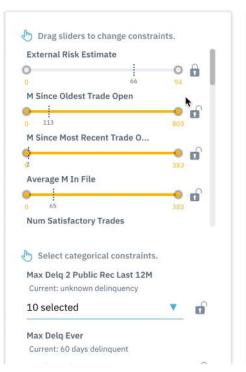
Rory Mc Grath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. FEAP-Al4fin workshop, NeurIPS, 2018.

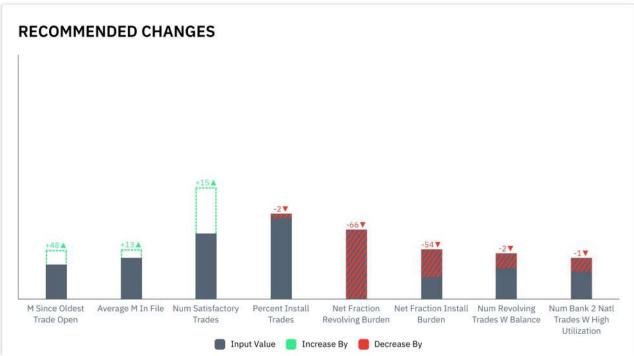
Counterfactual Explanations for Credit Decisions (2) - Finance



Rory Mc Grath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. FEAP-Al4fin workshop, NeurlPS, 2018.

Counterfactual Explanations for Credit Decisions (3) - Finance





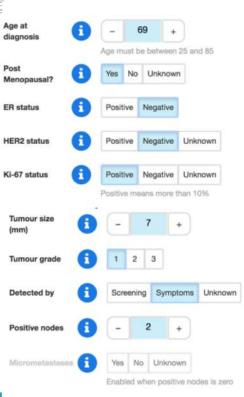
Rory Mc Grath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. FEAP-Al4fin workshop, NeurlPS, 2018.

Breast Cancer Survival Rate Prediction - Health

least 10 years.

Show ranges?









shows the percentage of women who survive at least 5 10 15 years after surgery, based on the information you have provided.

Treatment	Additional Benefit	Overall Survival %
Surgery only	-	72%
+ Hormone therapy	0%	72%
If death from breast o	ancer were excluded 83	2% would survive at

These results are for women who have already had surgery. This table

Challenge: Predict is an online tool that helps patients and clinicians see how different treatments for early invasive breast cancer might improve survival rates after surgery.

Al Technology: competing risk analysis

XAI Technology: Interactive explanations, Multiple representations.

David Spiegelhalter, Making Algorithms trustworthy, NeurlPS 2018 Keynote

predict.nhs.uk/tool



More on XAI



(Some) Tutorials, Workshops, Challenge

Tutorial:

- AAAI 2019 Tutorial on On Explainable AI: From Theory to Motivation, Applications and Limitations (#1) https://xaitutorial2019.github.io/
- ICIP 2018 / EMBC 2019 Interpretable Deep Learning: Towards Understanding & Explaining Deep Neural Networks (#2) http://interpretable-ml.org/icip2018tutorial/ http://interpretable-ml.org/embc2019tutorial/

Workshop:

- ISWC 2019 Workshop on Semantic Explainability (#1) http://www.semantic-explainability.com/
- IJCAI 2019 Workshop on Explainable Artificial Intelligence (#3) https://sites.google.com/view/xai2019/home
- IJCAI 2019 Workshop on Optimisation and Explanation in AI (#1) https://www.doc.ic.ac.uk/~kc2813/OXAI/
- ICAPS 2019 Workshop on Explainable Planning (#2)- https://kcl-planning.github.io/XAIP-Workshops/ICAPS_2019
- ICCV 2019 Workshop on Interpreting and Explaining Visual Artificial Intelligence Models (#1) http://xai.unist.ac.kr/workshop/2019/
- NeurIPS 2019 Workshop on Challenges and Opportunities for AI in Financial Services: the Impact of Fairness, Explainability, Accuracy, and Privacy https://sites.google.com/view/feap-ai4fin-2018/
- CD-MAKE 2019 Workshop on Explainable AI (#2) https://cd-make.net/special-sessions/make-explainable-ai/
- AAAI 2019 / CVPR 2019 Workshop on Network Interpretability for Deep Learning (#1 and #2) https://explainai.net/

Challenge:

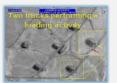
2018: FICO Explainable Machine Learning Challenge (#1) - https://community.fico.com/s/explainable-machine-learning-challenge

(Some) Software Resources

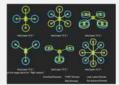
- DeepExplain: perturbation and gradient-based attribution methods for Deep Neural Networks interpretability. <u>aithub.com/marcoancona/DeepExplain</u>
- INNvestigate: A toolbox to iNNvestigate neural networks' predictions. github.com/albermax/innvestigate
- SHAP: SHapley Additive exPlanations. github.com/slundberg/shap
- GANDissect: Pytorch-based tools for visualizing and understanding the neurons of a GAN. https://github.com/CSAILVision/GANDissect
- ELI5: A library for debugging/inspecting machine learning classifiers and explaining their predictions. aithub.com/TeamHG-Memex/eli5
- Skater: Python Library for Model Interpretation/Explanations, github.com/datascienceinc/Skater
- Yellowbrick: Visual analysis and diagnostic tools to facilitate machine learning model selection. github.com/DistrictDataLabs/yellowbrick
- Lucid: A collection of infrastructure and tools for research in neural network interpretability. github.com/tensorflow/lucid
- LIME: Agnostic Model Explainer. https://github.com/marcotcr/lime
- Sklearn_explain: model individual score explanation for an already trained scikit-learn model. https://github.com/antoinecarme/sklearn_explain
- Heatmapping: Prediction decomposition in terms of contributions of individual input variables
- Deep Learning Investigator: Investigation of Saliency, Deconvnet, GuidedBackprop and more. https://github.com/albermax/innvestigate
- Google PAIR What-if: Model comparison, counterfactual, individual similarity. https://pair-code.github.io/what-if-tool/
- IBM AI Fairness: Set of fairness metrics for datasets and ML models, explanations for these metrics. https://github.com/IBM/aif360
- Blackbox auditing: Auditing Black-box Models for Indirect Influence. https://github.com/algofairness/BlackBoxAuditing
- Model describer: Basic statiscal metrics for explanation (visualisation for error, sensitivity). https://github.com/DataScienceSquad/model-describer

(Some) Initiatives: XAI in USA

Challenge Problem Areas



Data Analytics
Multimedia Data



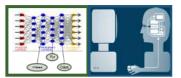
Autonomy
ArduPilot &
SITL Simulation

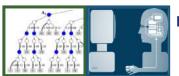
TA 1:

Explainable Learners

Teams that provide prototype systems with both components:

- Explainable Model
- Explanation Interface







Deep Learning Teams

Interpretable Model Teams

> Model Induction Teams

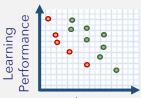
TA 2:

Psychological Model of Explanation



- Psych. Theory of Explanation
- Computationa I Model
- Consulting

Evaluation Framework



Explanation Effectiveness

Explanation Measures

- User Satisfaction
- Mental Model
- Task Performance
- Trust Assessment
- Correctability

Evaluator

TA1: Explainable Learners

Explainable learning systems that include both an explainable model and an explanation interface

TA2: Psychological Model of Explanation

Psychological theories of explanation and develop a computational model of explanation from those theories

(Some) Initiatives: XAI in Canada

DEEL (Dependable Explainable Learning) Project 2019-2024

Research institutions







Industrial partners









- Academic partners
 - Science and technology to develop new methods towards Trustable and Explainable Al







System Robustness

- To biased data
- Of algorithm
- To change
- To attacks

Certificability

- Structural warranties
- Risk auto evaluation
- External audit

Explicability & Interpretability

Privacy by design

- Differential privacy
- Homomorphic coding
- Collaborative learning
- To attacks



(Some) Initiatives: XAI in EU





















































































































W























































Explainable AI is motivated by real-world applications in AI

Not a new problem – a reformulation of past research challenges in Al

Multi-disciplinary: multiple AI fields, HCI, social sciences (multiple definitions)

In AI (in general): many interesting / complementary approaches

- Creating awareness! Success stories!
- Foster multi-disciplinary collaborations in XAI research.
- Help shaping industry standards, legislation.
- More work on transparent design.
- Investigate symbolic and sub-symbolic reasoning.

Evaluation:

- We need benchmark Shall we start a task force?
- > We need an XAI challenge Anyone interested?
- Rigorous, agreed upon, human-based evaluation protocols



Research and Technology Applied AI (Artificial Intelligence) Scientist

Wherever safety and Security are Critical, Thales c build smarter solutions. Everywhere.

protecting the national security interests of count

Established in 1972, Thales Canada has over 1,800 Toronto and Vancouver working in Defence, Avior

This is a unique opportunity to play a key role on Technology (TRT) in Canada (Quebec and Montre applied R&T experts at five locations worldwide. intelligence technologies. Our passion is imagining cutting edge AI technologies. Not only will you joi network, but this TRT is also co-located within Co-Intelligence expertise) i.e., the new flagship progr to work.

Job Description

An AI (Artificial Intelligence) Research and Techno developing innovative prototypes to demonstrate intelligence. To be successful in this role, one mos what's new, and a strong ability to learn new tech hand-on technical skills and be familiar with latest will contribute as technical subject matter experts and its business units. In addition to the implementarion Qualifications individual will also be involved in the initial project thinking, and team work is also critical for this role

As a Research and Technology Applied AI Scientist paced projects.

Professional Skill Requirements

· Good foundation in mathematics, statistic

MAY 2ND, 2019

Freddy Lecue Chief Al Scientist, CortAlx, Thales, Montreal – Canada

@freddylecue https://tinvurl.com/freddylecue Freddy.lecue.e@thalesdiaital.io

- · Strong knowledge of Machine Learning foundations
- Strong development skills with Machine Learning frameworks e.g., Scikit-learn, Tensoflow, PyTorch, Theano
- Knowledge of mainstream Deep Learning architectures (MLP, CNN, RNN, etc).
- Strong Python programming skills
- Working knowledge of Linux OS
- Eagerness to contribute in a team-oriented environment
- Demonstrated leadership abilities in school, civil or business organisations
- Ability to work creatively and analytically in a problem-solving environment
- Proven verbal and written communication skills in English (talks, presentations, publications, etc.)

Basic Qualifications

- Master's degree in computer science, engineering or mathematics fields
- Prior experience in artificial intelligence, machine learning, natural language processing, or advanced analytics

- Minimum 3 years of analytic experience Python with interest in artificial intelligence with working structured and unstructured data (SQL, Cassandra, MongoDB, Hive, etc.)
- · A track record of outstanding AI software development with Github (or similar) evidence
- Demonstrated abilities in designing large scale AI systems
- Demonstrated interes in Explainable AI and or relational learning
- Work experience with programming languages such as C, C++, Java, scripting languages (Perl/Python/Ruby) or similar
- Hands-on experience with data visualization, analytics tools/languages
- · Demonstrated teamwork and collaboration in professional settings
- · Ability to establish credibility with clients and other team members