Explainable Artificial Intelligence

Interrogating the AI systems

Nick Bassiliades, Professor School of Informatics Aristotle University of Thessaloniki, Greece <u>nbassili@csd.auth.gr</u>

Talk @ Rules: Logic and Applications, 17/12/2019





Intelligent Systems



A few words about the speaker

Nick Bassiliades (Νικόλαος Βασιλειάδης)

http://intelligence.csd.auth.gr/people/bassiliades

- Professor, School of Informatics, Aristotle University of Thessaloniki, Greece
- Scientific specialization: Knowledge Systems
 - Knowledge Representation & Reasoning
 - Rule-based systems, Logic programming, Defeasible Reasoning, Knowledge-based / expert systems
 - Semantic Web
 - Ontologies, Linked Open Data, Semantic Web Services
 - Multi-agent systems
 - Reputation/trust, knowledge-based interaction (argumentation, negotiation, brokering)
 - Applications on e-Learning, e-Government, e-Commerce, Electric Vehicles

#explain(model #explain(mode) #explain(model #explain(model #explain(model #explain(model #explain(model): plain(model) #explain(model) #explain(model) #explain(model) #explain(model #explain(model) #explain(model #explain(model

Credits

Ioannis Mollas, PhD student

- Explainable AI (mostly works on explainable Machine Learning)
- Credits for the most part of the talk, including the fancy looks
- Funded by project AI4EU (A European AI On Demand Platform and Ecosystem)
 - EU Horizon 2020 research and innovation programme under grant agreement No 825619
 - <u>https://www.ai4eu.eu/</u>

A few words about my institution

- Aristotle University of Thessaloniki, Greece
 - Largest (?) University in Greece and South-East Europe
 - Since 1925, 41 Schools, ~1.8K faculty, ~45K students
- School of Informatics
 - Since 1992, 30 faculty, 3 departments, 6 research labs, ~1500 undergraduate students, ~180 MSc students, ~110 PhD students, ~160 PhD graduates, >5000 pubs
- Intelligent Systems Laboratory (<u>http://intelligence.csd.auth.gr</u>)
 - 4 faculty, 18 PhD students, 4 post-graduate affiliates, 19 PhD graduates
 - Research on Artificial Intelligence, Machine Learning / Data Mining, Knowledge Representation & Reasoning / Semantic Web, Planning, Multi-Agent Systems
 - >470 publications, >40 projects





Intelligent Systems



SCHOOL OF INFORMATICS



This project has received funding from the European Union's Horizon 2020 research and innovation program under grant agreement No 825619.

Our Team







as Nick Bassiliades Professor



Grigorios Tsoumakas Assistant Professor



Ioannis Mollas PhD Student







AI4EU

About AI4EU

- European Union's landmark Artificial Intelligence project
- Al ecosystem containing knowledge, algorithms, tools and resources
- 80 partners, 21 countries
- 3 years project
- €20m budget

Goals

Create a leading collaborative AI European platform

- Open and sustainable Al On-Demand-Platform
- Unite stakeholders via high-profile conferences and virtual events
- Develop a Strategic Agenda for European Al
- Establish an Ethics Observatory to ensure development of human-centred AI
- Roll out of €3m in Cascade Funding





Filling AI Technological Gaps

This WP will reinforce European *excellence* and *leading* position worldwide in major AI research and application domains, through research and innovation efforts to fill important technology gaps.

Objectives: 1. Develop new AI tools and techniques (to be included in the AI4EU Platform)

2. Consolidate and strengthen excellence in AI in EU

WP7



AI4EU

Task 7.1 Explainable AI

Explainable Interpretable Comprehensible Intelligible, Justifiable Understandable

Definition: An AI system should allow humans to understand the reasons behind its recommendations or decisions. It should be possible to know the data, rationale and arguments that lead to a result, to question them and to correct them

Tasks:

- Community creation and aggregation activities
 - Summer School covering explainable AI
 - Workshops and tutorials
 - ♦.
- ✓ Research activities
 - Survey for explainable AI
 - Methodology and software components

***** ...







Back to XAI

TOD

FIRISHSEFOLLING

andle

Stop: 9. SUGCESS: B. EFFOT.

body: (value: y

depen

SCFOL

Darw

XAI: Explainable Artificial Intelligence

Explainable AI (XAI) refers to methods and techniques in the application of <u>artificial intelligence</u> technology (AI) such that the results of the solution can be understood by human experts.

[#]explain(mode). explain(model): xplain(model): xplain(model): this) equir: value (unction) return this ender to a light the second s plain(model): anno (function) (d. startScrolling (st.), this if anno (function) (d. startScrolling (st.), this if anno (function) (d. startScrolling (st.), this if the start scrolling (st.), the start scrolling (st.), the start anno (function) (d. startScrolling (st.), the start scrolling vin(model): art.scrollop=this.dragStatPosition (mod mo

odel odel de de P de de 10

Outline

□ Why XAI?

- □ Interpretable Machine Learning
 - ✓ Transparent and Black Box Models
 - ✓ Explainable ML
 - ✓ Types of Data
 - \checkmark Opening the Black Box
 - ✓ Implementations

Where we need Explainable Al



«At the individual level, designers have both ethical and legal responsibilities to provide such justification for decisions that could result in death, financial loss, or denial of parole.», Don Monroe 2018^[3]

Interpretable Machine Learning



"A way to present the result of a black box model in human readable terms"^[1]

Each community addresses this meaning from a different perspective

Recognized Interpretable Models

Or Transparent Boxes:

IF Proline>=990.0 THEN Wine=1

intensity<=3.52 THEN Wine=2

IF Flavanoids<=1.41 AND

IF Proline>=680.0 AND

IF Color intensity <= 3.85 AND Color

Proline>=470.0 THEN Wine=3

Alcohol>=12.93 THEN Wine=1

IF Hue>=0.69 THEN Wine=2

Accepted

Rejected

IF TRUE THEN Wine=2

#explain(model): #explain(model):

- Decision Trees
- Rule Based Systems
- Decision Tables
- Linear Models
- Decision Sets
- ✤ K-NN



	Rule 1	Rule 2	Rule 3	Rule 4
Conditions				
Valid User Name	F	т	т	т
Valid Password	-	F	т	т
Adequate balance in the account	-	-	F	т
Actions				
Login accepted	F	F	т	т
Amount transferred	-	-	F	т

Black Box Predictor... Means?



A black box predictor is the result of a machine learning algorithm, whose internals are:

- known to his creator or
- unknown to his creator

and uninterpretable to other people.

Some known black box models are: Support Vector Machines, Neural Networks, Tree Ensemble, Deep Neural Networks and Non-Linear Models

plain (model) .rexplain(model): #explain(model): #explain(model): #explain(model): #explain(model): #explain(model): #explain(model): #explain(model): # exting the context scropl Top this drag Start Position scroll Position scrol #explain(model): #explain(model): יxplain(model): (ain(model)

Dimensions of Interpretability

1. Is it Global or Local?

2. Has Time Limitation?

3. Which is the Nature of User Expertise?

4. Which is the Shape of Explanation?

Desired Features of Explainable Model

Accuracy: Measuring the accuracy and other metrics of the model

Interpretability: Measuring Comprehensibility of model Fidelity: Measuring the imitation of the black box <u>Responsive-</u> <u>ness</u>: Measuring the time for explaining an instance

> <u>Other</u>: Fairness, privacy, usability, reliability, robustness and scalability

> > 20



Shapes of Explanations

plain(mode plain(model plain(model plain(model plain(model nodel olain de model plaint plain odel olain model olain(model plain alain(model

Comprehensibility Measure^[2,6]

- ✓ Human Centered Evaluation
- ✓ Number of regions (parts of the model)
- ✓ Number of conditions (in a rule)
- ✓ Number of non-zero weights (in linear models)
- \checkmark Depth of the tree (in a decision tree)

Simplification Methods like pruning are used in order to enhance the comprehensibility and to avoid overfitting. *But do we want to avoid overfitting when explaining?* plain(mode plain(m ode plain(model plain(model olain(mode plain plain(mode plain(r olain model plain(model

Comprehensibility Measure^[2,6]

If f1 > 20 and f2 <= 123 and f4 > 0.5 and f4 <= 1 and f5 <= 0 and f6 > 9 and f7 <= 10 and f8 <= 1 and f8 > 0.7 and f9 <= 19 and f10 > 10 and f10 <= 100 and f11 <= 1.23 and f12 > 12.1452 and f12 <= 15 and then "class A"

VS

You classified as "class A" because $0.5 < f5 \le 1$ and $f9 \le 19$.

Which one would you prefer?

Types of Data





Agnostic Explanator?

- Can explain indifferently any kind of Black Box in most of the papers
- Most of the times, It has a specific type of data (Tabular, Images, Texts, Other), but rarely It could be agnostic on type of data too

Some Advantages of Agnostic Explanators:

- They can explain models, even without dataset knowledge
- They can explain models remotely

Reverse Engineering on Black Box

→ Model Explanation:

Is aimed at understanding the overall logic behind the black box. Provides a global explanation, through a transparent box. Examples: Surrogate Method Explanators, Knowledge Distillation

• Outcome Explanation:

Is aimed at understanding the outcome of an instance. Provides a local explanation, through a transparent box. Examples: Lime, Anchors

→ <u>Model Inspection</u>:

Provides a visual or textual representation of some specific properties.

Examples: Permutation Importance, PDP, ICE, SHAP





Use of Scikit Learn for Explainable Models https://scikit-learn.org

 Image: Surrogate method

 Global or Local

Use of Orange for Explainable Models https://github.com/biola b/orange3

Synthetic Neighborhood: by removing words from the instance randomly

audits

trains

IP

Linear Model

creates

Ρ

New Instance

Dreodict Dreodict Caused 0.26 Rice 0.15 Genocide 0.13 certainty 0.09 scri 0.09 owlnet 0.08

Final explanation

is.Selement,c=this.options, .fallback),tip:function(){retu.

State of

<u>local</u> interpretable <u>model-agnostic</u> explanations https://github.com/marcotcr/lime

> Mathematic Formulation $L(f, g, \pi x) + \Omega(g)$

> > unction f(){v

30

Disadvantage on Sparse Data!

Why?

LIME can only generate 2ⁿ different neighbours, where n the number of non-zero values.

We proposed LioNets^[9]...

State of the Art

31

fallback), tip:function(){retu.

<u>local</u> interpretable <u>model-agnostic</u> explanations https://github.com/marcotcr/lime

> Mathematic Formulation L(f, g, πx) + $\Omega(g)$

LioNets Hypothesis

Create neighbours to the abstract space using a decoder trained with the neural network as encoder will lead to better, bigger and more representative neighbourhoods, thus better explanations.

LioNets Architecture

10, 93 (10, 22 () 103 () (2, 0) (0) 10

LioNets will try to interpret a neural network's prediction. Thus, it is a **model-specific** outcome explanator



Deep Neural Network Classifier

LioNets Architecture









SHAP/Shapley Value Examples^[4]

Can do anything!



Model Inspection via Feature Importance (left) and Partial Dependence Plots (right) (computed using SHAP values)

Reverse Engineering on Black Box

h 👝 v

Reverse Engineering on Black Box

in tion) vano (,b;this;\$ei 's' s \$), + 8, , + 8, , ≠ []€, fallt

e

turnab(...dead-u)ton

emove:fu

wall

ti

SHAP/Shapley Value Examples^[4]

Can do anything!



Outcome explanation

Knowledge Distillation^[8]

in tion) vano (,b;this;\$ei 's' s \$), + 8, , + 8, , ≠ []€, fallt

Using Transparent boxes to mimic black boxes' accuracies



#explain(model)



Existing XML Libraries

Bibliography

- 1. Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Dino Pedreschi, Fosca Giannotti. 2018. A Survey Of Methods For Explaining Black Box Models.
- 2. Amina Adadi and Mohammed Berrada. 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI).
- 3. Don Monroe. 2018. Al, Explain Yourself.
- 4. Christoph Molnar. 2018. Interpretable Machine Learning.
- 5. Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier.
- 6. Johan Huysmans, Karel Dejaeger, Christophe Mues, Jan Vanthienen, and Bart Baesens. 2011. An empirical evaluation of the comprehensibility of decision table, tree and rule based predictive models.
- 7. Pedro Domingos. 1998. Knowledge Discovery Via Multiple Models
- 8. Geoffrey Hinton, Oriol Vinyals, Jeff Dean. 2015. Distilling the Knowledge in a Neural Network
- 9. I. Mollas, N. Bassiliades, G. Tsoumakas. 2019. LioNets: Local Interpretation of Neural Networks through Penultimate Layer Decoding.
- 10. Čyras, Kristijonas and Satoh, Ken and Toni, Francesca. 2016. Abstract argumentation for case-based reasoning.
- 11. Čyras, Kristijonas and Satoh, Ken and Toni, Francesca. 2016. Explanation for case-based reasoning via abstract argumentation.
- 12. Martin Možina, Jure Žarbkar, Ivan Bratko. 2007. Argument based machine learning.



Some say: Future Is Model Agnostic Models

Others say: Future is New Models, Interpretable by Design

Before we go...

We Say Model Specific Implementations

Explainable Artificial Intelligence

Interrogating the AI systems



"Magic model on the core, Explain yourself in front of all"

Nick Bassiliades nbassili@csd.auth.gr The End



Intelligent Systems

TΒ

OOL OF INFORMATICS