Explainable Machine Learning: Mind the Users and their Knowledge

Freddy Lécué Inria, France CortAlx@Thales, Canada @freddylecue



Intelligent Systems Conference (IntelliSys) 2020

September 4th, 2020 Amsterdam, The Netherlands



Freddy Lecue. On the Role of Knowledge Graphs in Explainable AI. Semantic Web Journal (2020) http://www.semantic-web-journal.net/content/role-knowledge-graphs-explainable-ai



Al Adoption: Requirements



Explanation in Al

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.

Motivation

Business to Customer





Gary Chavez added a photo you might ... be in. about a minute ago · 🔐



Critical Systems





Markets We Serve (Critical Systems)



Trusted Partner For A Safer World

But not Only Critical Systems

COMPAS recidivism black bias

Opinion

By Rebecca Wexle

When a Computer Program Keeps You in Jail

DYLAN FUGETT

Prior Offense 1 attempted burglary

Subsequent Offenses 3 drug possessions

BERNARD PARKER

Prior Offense 1 resisting arrest without violence

Subsequent Offenses None

LOW RISK

HIGH RISK

10

Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.

3

XAI in a Nutshell

Source: https://www.cc.gatech.edu/~alanwags/DLAI2016/(Gunning)%20IJCAI-16%20DLAI%20WS.pdf

An Example of an end-to-end XAI System

Green regions argue for FISH, while RED pushes towards DOG. There's more green.

H: What happens if the 4 background anemones are removed? E.g.,

C: I still predict FISH. because of these green superpixels:

- Humans may have follow-up questions
- Explanations cannot answer all users' concerns

Weld, D., and Gagan Bansal. "The challenge of crafting intelligible intelligence." Communications of ACM (2018).

How to Explain? Accuracy vs. Explanability

- Challenges:
 - Supervised
 - Unsupervised learning

Learning

- Approach:
 - Representation Learning
 - Stochastic selection
- Output:
 - Correlation
 - No causation

XAI Objective

Supporting Industrialization of Al at Scale

Explainability by Design for AI Products

KDD 2019 Tutorial on Explainable AI in Industry - 5https://sites.google.com/view/kdd19-explainable-ai-tutorial

On Role of Data

In XAI

Interpretable Data for Interpretable Models

Table of baby-name data (baby-2010.csv)

	_			Field
name	rank	gender	year -	names
Jacob	1	boy	2010	One row
Isabella	1	girl	2010	(4 fields)
Ethan	2	boy	2010	
Sophia	2	girl	2010	
Michael	3	boy	2010	
2000 all	rows told			-

Text

Tabular

Images

What about the

Evaluation?

XAI: One Objective, Many Metrics

XAI in AI

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

On the Role of Knowledge Graphs in Explainable AI A Machine Learning Perspective

On the Role of Knowledge Graph in Explainable AI - Semantic Web Journal - <u>http://www.semantic-web-journal.net/content/role-knowledge-graphs-explainable-ai</u>

Knowledge Graph (1)

- Set of (*subject, predicate, object SPO*) **triples** *subject* and *object* are **entities**, and *predicate* is the **relationship** holding between them.
- Each SPO **triple** denotes a **fact**, i.e. the existence of an actual relationship between two entities.

Knowledge Graph in Machine Learning (1)

Augmenting (input) features with more semantics such as knowledge graph embeddings / entities

https://stats.stackexchange.com/questions/230581/decision -tree-too-large-to-interpret

Knowledge Graph in Machine Learning (2)

Knowledge Graph in Machine Learning (3)

Knowledge Graph in Machine Learning (4)

Knowledge Graph in Machine Learning (5)

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

Augmenting models with semantics to support personalized explanation

Description 3: This is a public transportation accident <------

On One Industrial Application in **Thales**

State of the Art Machine Learning **Applied to Critical Systems**

Object (Obstacle) Detection Task

Object (Obstacle) Detection Task Stateof-the-art ML Result

Lumbermill - .59

Object (Obstacle) Detection Task Stateof-the-art ML Result

Lumbermill - .59

Boulder - .09

Railway - .11

State of the Art XAI **Applied to Critical**

Systems

Object (Obstacle) Detection Task State-of-the-art XAI Result

Lumbermill - .59

Unfortunately, this is of NO use for a human behind the system

Let's stay back

Why this Explanation? (meta explanation)

After Human Reasoning...

lum	hormil	1_ 50	
Lann			

💐 DBpedia		Formats -	C Faceted Browser	C Sparql Endpoint
dbo:wikiPageID		 352327 (xsd:integer) 		
dbo:wikiPageRevisior	nID	 734430894 (xsd:integer) 		
dct:subject		 dbc:Sawmills dbc:Saws dbc:Ancient_Roman_technology dbc:Timber_preparation dbc:Timber_industry 		
http://purl.org/linguis	tics/gold/hypernym	dbr:Facility		
rdf: type		owl:Thingdbo:ArchitecturalStructure		
rdfs:COmment		 A sawmill or lumber mill is a facility where logs are cut into lumber. Prior to the invention planed, or more often sawn by two men with a whipsaw, one above and another in a saw mill is the Hierapolis sawmill, a Roman water-powered stone mill at Hierapolis, Asia Mino water-powered mills followed and by the 11th century they were widespread in Spain an Asia, and in the next few centuries, spread across Europe. The circular motion of the what the saw blade. Generally, only the saw was powered, and the logs had to be loaded a was the developm (en) 	of the sawmill, boards we v pit below. The earliest k or dating back to the 3rd d North Africa, the Middl eel was converted to a re nd moved by hand. An ea	ere rived (split) and nown mechanical century AD. Other e East and Central ciprocating motion arly improvement
rdfs:label		 Sawmill (en) 		
owl:sameAs		 wikidata:Sawmill dbpedia-cs:Sawmill dbpedia-de:Sawmill dbpedia-es:Sawmill 		

What is missing?

Lumbermill - .59

Context

matters

Boulder - .09

Railway - .11

Srowse using - Formats -

C Faceted Browser C Sparql Endpoint

About: Boulder

An Entity of Type : place, from Named Graph : http://dbpedia.org, within Data Space : dbpedia.org

In geology, a boulder is a rock fragment with size greater than 25.6 centimetres (10.1 in) in diameter. Smaller pieces are called cobbles and pebbles, depending on their "grain size". While a boulder may be small enough to move or roll manually, others are extremely massive. In common usage, a boulder is too large for a person to move. Smaller boulders are usually just called rocks or stones. The word boulder is short for boulder stone, from Middle English bulderston or Swedish bullersten. Boulder sized clasts are found in some sedimentary rocks, such as coarse conglomerate and boulder clay.

Property	Value
dbo:abstract	In geology, a boulder is a rock fragment with size greater than 25.6 centimetres (10.1 in) in diameter. Smaller pieces are called cobbles and pebbles, depending on their "grain size". While a boulder may be small enough to move or roll manually, others are extremely massive. In common usage, a boulder is too large for a person to move. Smaller boulders are usually just called rocks or stones. The word boulder is short for boulder stone, from Middle English bulderston or Swedish bullersten. In places covered by ice sheets during lce Ages, such as Scandinavia, northern North America, and Russia, glacial erratics are common. Erratics are boulders picked up by the ice sheet during its advance, and deposited during its retreat. They are called "erratic" because they typically are of a different rock type than the bedrock on which they are deposited. One of them is used as the pedestal of the Bronze Horseman in Saint Petersburg, Russia. Some noted rock formations involve giant boulders exposed by erosion, such as the Devil's Marbles in Australia's Northern Territory, the Horeke basalts in New Zealand, where an entire valley contains only boulders, and The Baths on the island of Virgin Gorda in the British Virgin Islands. Boulder sized clasts are found in some sedimentary rocks, such as coarse conglomerate and boulder clay. The climbing of large boulders is called bouldering. (err)
dbo:thumbnail	wiki-commons:Special:FilePath/Balanced_Rock.jpg?width=300
dbo:wikiPageID	60784 (xsd:integer)
dbo:wikiPageRevisionID	• 743049914 (xsd:integer)
det:subject	dbc:Rock_formations dbc:Rocks

Source of the second se

C Faceted Browser C Sparql Endpoint

About: Rail transport

Pror

An Entity of Type : software, from Named Graph : http://dbpedia.org, within Data Space : dbpedia.org

Rail transport is a means of conveyance of passengers and goods on wheeled vehicles running on rails, also known as tracks. It is also commonly referred to as train transport. In contrast to road transport, where vehicles run on a prepared flat surface, rail vehicles (rolling stock) are directionally guided by the tracks on which they run. Tracks usually consist of steel rails, installed on ties (sleepers) and ballast, on which the rolling stock, usually fitted with metal wheels, moves. Other variations are also possible, such as slab track, where the rails are fastened to a concrete foundation resting on a prepared subsurface.

erty	Value
bstract ·	• Rail transport is a means of conveyance of passengers and goods on wheeled vehicles running on rails, also known as tracks. It is also commonly referred to as train transport. In contrast to road transport, where vehicles running on rails, also known as tracks. It is also commonly referred to as train transport. In contrast to road transport, where vehicles run on a prepared flat surface, rail vehicles (rolling stock) are directionally guided by the tracks on which they run. Tracks usually consist of steel rails, installed on ties (sleepers) and ballast, on which ther onling stock, usually fitted with metal wheels, moves. Other variations are also possible, such as slab track, where the rails are fastened to a concrete foundation resting on a prepared subsurface. Rolling stock in a rail transport system generally encounters lower frictional resistance than road vehicles, so passenger and freight cars (carriages and wagons) can be coupled into longer trains. The operation is carried out by a railway company, providing transport between train stations or freight customer facilities. Power is provided by locomotives which either draw electric power from a railway electrification system or produce their own power, usually by diesel engines. Most tracks are accompanied by a signalling system. Railways are a safe land transport system when compared to other forms of transport. Railway transport is capable of high levels of passenger and cargo tuitization and energy efficiency. Jut is often less fits/like and more capital-intensive than mad transport hower target and cargo

considered. The oldest, man-hauled railways date back to the 6th century BC, with Periander, one of the Seven Sages of Greece

XAI Thales Platform

- Higher accuracy with no intensive fine-tuning
- Human interpretable explanation
- Running on the edge at inference time

 Hardware: High performance, scalable, generic (to different FGPA family) & portable CNN dedicated programmable processor implemented on an FPGA for real-time embedded inference

Software: Knowledge graph extension of object detection

X

This is an **Obstacle: Boulder** obstructing the train: XG142-R on **Rail_Track** from City: Cannes to City: Marseille at Location: Tunnel VIX due to **Landslide**

Conclusion

Conclusion

- Explainable AI is motivated by real-world applications in AI
- Not a new problem a reformulation of past research challenges in AI
- Knowledge graphs should be foundational for XAI
- But they are facing challenges related to their integration (data mapping)
- Many industrial applications already crucial for AI adoption in critical systems

Why do we Need Knowledge Graphs to Achieve XAI?

Because this is not an explanation from an intelligent system

This is even not interpretable, and then not actionable

Open Research Questions for the Semantic Web / Knowledge Graph Community

- [Data] Machine learning experts do not buy the data knowledge mapping
- [Explanation] There is *no agreement* on *what an explanation is*
- [Explanation] There is not a formalism for explanations (neither model nor output)
- [Model] There is very limited work in *machine learning modules* composability – and none from a semantics perspective
- [Model] There is no work on describing and representing models
- [Model] What are **disentangled representations** and how can its factors be quantified and detected?
- [Human-in-the-loop] There is *no work* that seriously addresses the problem of *quantifying* the grade of *comprehensibility* of an explanation for humans

