# Explaining Classifiers in Ontology-Based Data Access

**Federico Croce** and Maurizio Lenzerini

DIPARTIMENTO DI INGEGNERIA INFORMATICA
AUTOMATICA E GESTIONALE ANTONIO RUBERTI

SAPIENZA
UNIVERSITÀ DI ROMA

Explainable Logic-Based Knowledge Representation (XLoKR 2020)

Workshop at KR2020 (17th International Conference on Principles of Knowledge Representation and Reasoning)
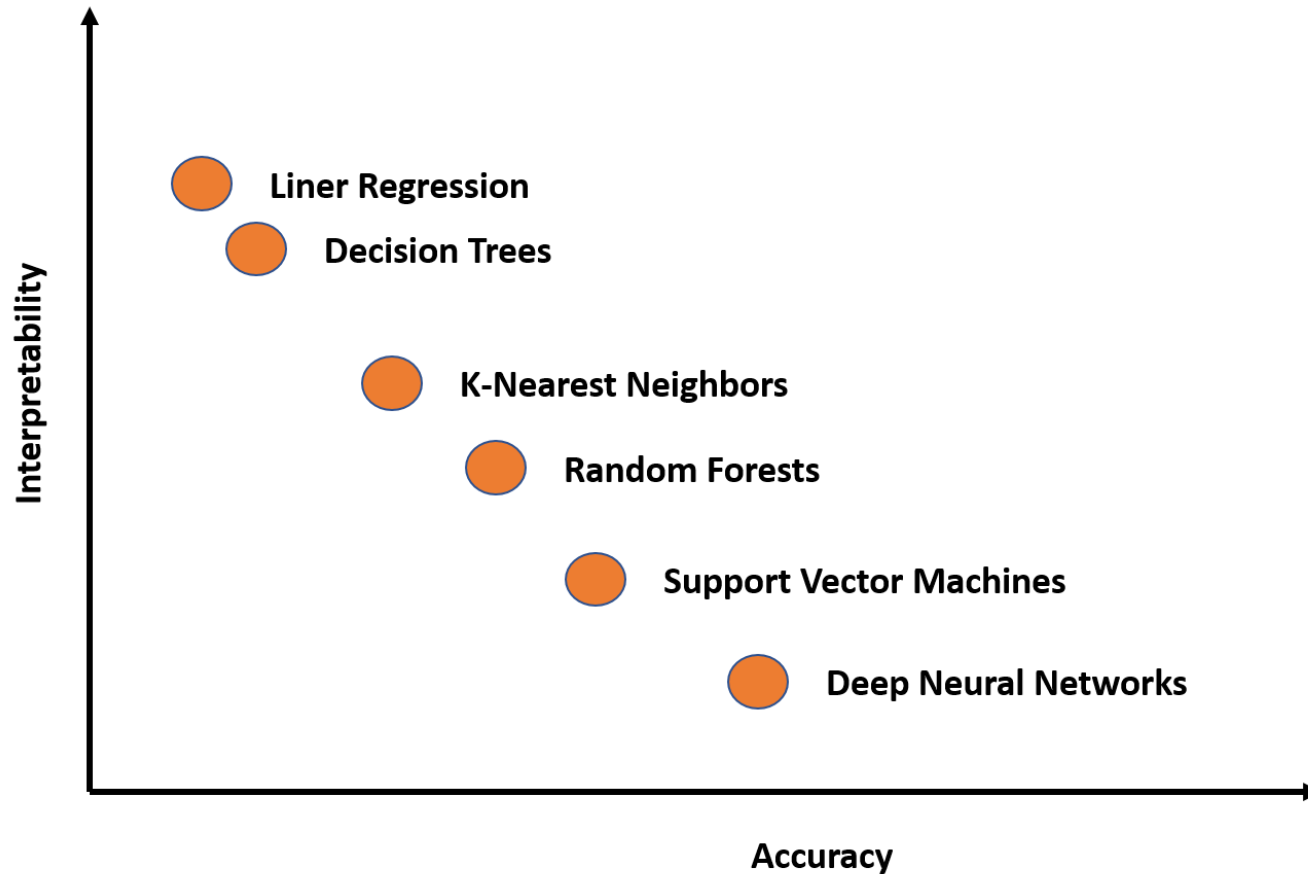
September 13th, 2020

# Introduction

# Problem statement

- Machine Learning (ML) has many elegant and efficient solutions to very difficult problems: Machine Translation, Vision, Autonomous Driving, and more

- An empiric rule shows that the more a ML algorithm is accurate, the less we understand its "magic"

- Deep learning is an extreme example of a high accuracy, black-box model

# ML interpretability (empiric)

# Why should we care?

- Caring only about performances is not the right choice in many fields: finance, justice, healthcare, privacy

- One famous example is COMPAS algorithm [1], used across the US to predict future criminals, and proved to be biased against black people

[1] https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing
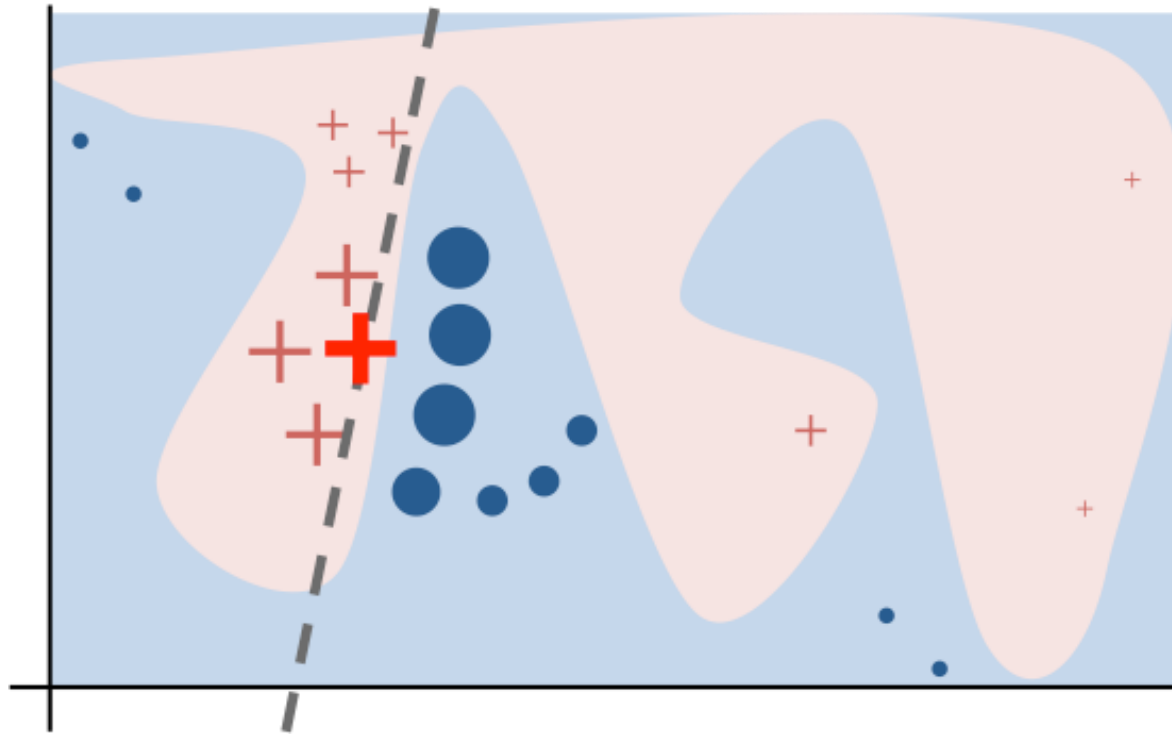
# Why should we care? (cont.)

*AARON HOLMES SEP 11, 2020*

- A sheriff launched an algorithm to predict who might commit a crime. Dozens of people said they were harassed by deputies for no reason [2].

- But according to a six-month investigation published this week by the Tampa Bay Times, the high-tech tool deployed by the Pasco Sheriff's Office didn't lead to a reduction in violent crimes. Instead, 21 families singled out by the algorithm said they were routinely harassed by deputies, even when there was no evidence of a specific crime.

[2] https://www.businessinsider.com/predictive-policing-algorithm-monitors-harasses-families-report-2020-9

# Possible Solutions

State-of-the-art: LIME, SHAP, Scoped Rules, Counterfactual and Adversarial Examples, Feature Visualization



Tulio Ribeiro, M., Singh, S., & Guestrin, C. (2016). " Why Should I Trust You?": Explaining the Predictions of Any Classifier. arXiv, arXiv-1602.
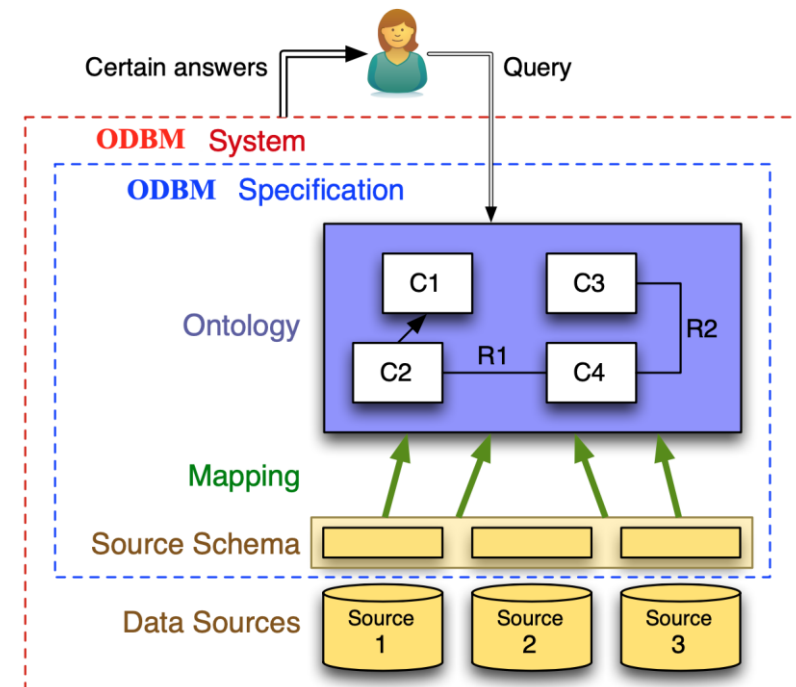
# Possible Solutions (cont.)

- **Our solution** is a form of *reverse engineering* of an Ontology-Based Data Management (OBDM) system*:* finding a query over the ontology that semantically describes the tagged individuals in the dataset

# Preliminaries

# Ontology-Based Data Management

It is a **three-layered architecture**:

- The ontology is a declarative and explicit representation of the domain of interest
- The data layer is constituted by the existing dataset
- The mapping layer is a set of declarative assertions specifying how the sources in the data layer relate to the ontology

# The notion of *certain answers*

- Let $\mathcal{O}$ be an ontology, $\mathcal{S}$ a dataset, and $\mathcal{M}$ a set of mappings, we call $\mathcal{J} = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ an OBDM *specification*

- Let $q_{\mathcal{O}}$ be a query over $\mathcal{O}$, we define the *certain answers* of $q_{\mathcal{O}}$ w.r.t. $\mathcal{J}$ and a database $D$, denoted by $cert_{q_{\mathcal{O}}, \mathcal{J}}^{D}$ as the set of tuples $\vec{t}$ of $\mathcal{S}$-constants, such that

$\vec{t} \in q_{\mathcal{O}}^{B}$ **for every possible interpretation $B$** that satisfies $\mathcal{J}$ for an $\mathcal{S}$-database $D$ (called a *model* of $\mathcal{J}$ w.r.t. *D)*

# The Classifier

Given a dataset *D*, we consider a binary classifier:

$$\lambda : dom(D)^n \rightarrow \{+1, -1\}$$

Also, we will denote the set of tuples that have been classified positively (resp. negatively) as:

$$\lambda^+ = \left\{\vec{t} \in dom(D)^n \mid \lambda\big(\vec{t}\big) = +1\right\}$$

$$(\text{resp. } \lambda^- = \left\{\vec{t} \in dom(D)^n \mid \lambda\big(\vec{t}\big) = -1\right\})$$

# The Framework

# The Notion of Border

- For each tuple $\vec{t} \in D$ and natural number $r$, we define $\mathcal{B}_{\vec{t},r}(D)$ as the ***Border*** of radius $r$ for $t$ in $D$, representing all the atoms in $D$ that are *reachable* from $\vec{t}$ in at most $r$ joins

*Example:* Let a database be $D = \{R(a,b), S(a,c), Z(c,d), W(d,e), W(e,h), R(f,g)\}$ and let **t** = $\langle a \rangle$. By denoting with $\mathcal{W}_{t,n}(D)$ the atoms in D that are reachable from **t** in at most *n* joins, we have that:

- $\mathcal{W}_{t,0}(D) = \{R(a,b), S(a,c)\}$
- $\mathcal{W}_{t,1}(D) = \{Z(c,d)\}$
- $\mathcal{W}_{t,2}(D) = \{W(d,e)\}$

Therefore, the border of radius 2 of **t** in D is:

$$\mathcal{B}_{t,2}(D) = \{R(a,b), S(a,c), Z(c,d), W(d,e)\}$$

# The $\mathcal{J}$-match

- A query $q_{\mathcal{O}}$ $\mathcal{J}$-*matches* a Border $\mathcal{B}_{\vec{t},r}(D)$ of radius *r* of a tuple $\vec{t}$ in a source database *D,* if $\vec{t}$ is in the *certain answers* of $q_{\mathcal{O}}$ w.r.t to $\mathcal{J}$ and *D,* i.e. if

$$ t \in cert_{q_{\mathcal{O}}, \mathcal{J}}^{\mathcal{B}_{t,r}(D)} $$

# The goal of the framework

- The goal of our framework, is to find a semantic description of $\lambda$ that is as close as possible to a set of user-defined criteria.

- Each criterion has a function associated to it, that returns a quantitative measure of how much a given query meets the criteria

- The user also defines an expression to compute, for a given query, a unique value out of all the measures returned by the functions of each criterion

# The criteria, the functions and the expression

- $\delta_1$ = "Maximize the number of tuples $\boldsymbol{t} \in \lambda^+$ such that $q_{\mathcal{O}}$ $\mathcal{J}$-matches $\mathcal{B}_{\boldsymbol{t},r}(D)$"
- $\delta_2$ = "Minimize the number of tuples $\boldsymbol{t} \in \lambda^-$ such that $q_{\mathcal{O}}$ $\mathcal{J}$-matches $\mathcal{B}_{\boldsymbol{t},r}(D)$"
- $\delta_3$ = "Minimize the number of disjuncts of the query $q_{\mathcal{O}}$"

- $f_{\delta_1}(q_{\mathcal{O}}) = \dfrac{|\{\boldsymbol{t} \in \lambda^+ \ s.t. \ q_{\mathcal{O}} \ \mathcal{J}-matches \ \mathcal{B}_{\boldsymbol{t},r}(D)\}|}{|\lambda^+|}$

- $f_{\delta_2}(q_{\mathcal{O}}) = 1 - \dfrac{\left|\{\boldsymbol{t} \in \lambda^- \ s.t. \ q_{\mathcal{O}} \ \mathcal{J}-matches \ \mathcal{B}_{\boldsymbol{t},r}(D)\}\right|}{|\lambda^-|}$

- $f_{\delta_3}(q_{\mathcal{O}}) = \dfrac{1}{|CQs \ in \ q_{\mathcal{O}}|}$

- $\mathcal{Z}_{\mathcal{F}}(q_{\mathcal{O}}) = \dfrac{\alpha f_{\delta_1}(q_{\mathcal{O}}) + \beta f_{\delta_2}(q_{\mathcal{O}}) + \gamma f_{\delta_3}(q_{\mathcal{O}})}{\alpha + \beta + \gamma}$ (we call this the $\mathcal{Z}$ score of $q_{\mathcal{O}}$ under $\mathcal{F}$)

  where $\alpha, \beta, \gamma$ represents the importance of criterion $\delta_1, \delta_2, \delta_3$ respectively

# The Algorithm

# Example (1/7)

Consider the following database *D*

| STUD | | $\lambda$ |
|---|---|---|
| | A10 | +1 |
| | B80 | +1 |
| $\lambda^+$ | C12 | +1 |
| | D50 | +1 |
| $\lambda^-$ | E25 | -1 |

| LOC | |
|---|---|
| Sap | Rome |
| TV | Rome |
| Pol | Milan |

| ENR | | |
|---|---|---|
| A10 | Math | TV |
| B80 | Math | Sap |
| C12 | Science | Norm |
| D50 | Science | TV |
| E25 | Arts | Pol |

# Example (2/7)

Let the ontology be:

$$\mathcal{O} = \{\text{MathStudent} \sqsubseteq \text{ScientificStudent},$$
$$\text{ScienceStudent} \sqsubseteq \text{ScientificStudent}\}$$

And the mappings:

$$\mathcal{M} = \quad ENR(x, Math, z) \rightsquigarrow MathStudent(x)$$
$$ENR(x, Science, z) \rightsquigarrow ScienceStudent(x)$$
$$ENR(x, y, z) \rightsquigarrow enrolledIn(x, z)$$
$$LOC(x, y) \rightsquigarrow locatedIn(x, y)$$

# Example (3/7)

The corresponding borders of radius 1, for each tuple are:

$$\mathcal{B}_{A10,1}(D) = \{STUD(A10),\ ENR(A10,\ Math,\ TV),\ LOC(TV,\ Rome)\}$$

$$\mathcal{B}_{B80,1}(D) = \{STUD(B80),\ ENR(B80,\ Math,\ Sap),\ LOC(Sap,\ Rome)\}$$

$$\mathcal{B}_{C12,1}(D) = \{STUD(C12),\ ENR(C12,\ Science,\ Norm)\}$$

$$\mathcal{B}_{D50,1}(D) = \{STUD(D50),\ ENR(D50,\ Science,\ TV),\ LOC(TV,\ Rome)\}$$

$$\mathcal{B}_{E25,1}(D) = \{STUD(E25),\ ENR(E25,\ Arts,\ Pol),\ LOC(Pol,\ Milan)\}$$

# Example (4/7)

Consider each border associated to the tuples in $\lambda^+$ as a CQ, and compute the complete s-to-o rewriting of each query, as described in [3]. In a nutshell, this means to apply all the mappings to the queries.

$$q_1(A10) \leftarrow \mathrm{MathStudent}(A10) \wedge \mathrm{enrolledIn}(A10, \mathrm{TV}) \wedge \mathrm{locatedIn}(\mathrm{TV}, \mathrm{Rome})$$

$$q_2(B80) \leftarrow \mathrm{MathStudent}(B80) \wedge \mathrm{enrolledIn}(B80, \mathrm{Sap}) \wedge \mathrm{locatedIn}(\mathrm{Sap}, \mathrm{Rome})$$

$$q_3(C12) \leftarrow \mathrm{ScienceStudent}(C12) \wedge \mathrm{enrolledIn}(C12, \mathrm{Norm})$$

$$q_4(D50) \leftarrow \mathrm{ScienceStudent}(D50) \wedge \mathrm{enrolledIn}(D50, \mathrm{TV}) \wedge \mathrm{locatedIn}(\mathrm{TV}, \mathrm{Rome})$$

[3] Cima, G., Lenzerini, M., & Poggi, A. (2019). Semantic Characterization of Data Services through Ontologies. In *IJCAI*.

# Example (5/7)

- To reduce the number of queries generated, we introduce the notion of **query patterns**

- We say that two CQs have the same *pattern*, if they are conjunctions of the same set of atoms

- Our intuition is that similar tuples of the database will be described by similar properties, and will form similar *query patterns* when processed by the previous steps of the algorithm

- For each *pattern*, we only keep the constants that are shared by all the queries of the pattern. All the other constants will be substituted by new variables.

# Example (6/7)

- The query patterns of the example are:

$$q_5(x) \leftarrow \text{MathStudent}(x) \wedge \text{enrolledIn}(x, y) \wedge \text{locatedIn}(y, \text{Rome})$$

$$q_6(C12) \leftarrow \text{ScienceStudent}(C12) \wedge \text{enrolledIn}(C12, \text{Norm})$$

$$q_7(D50) \leftarrow \text{ScienceStudent}(D50) \wedge \text{enrolledIn}(D50, \text{TV}) \wedge \text{locatedIn}(\text{TV}, \text{Rome})$$

# Example (7/7)

- Let $k$ be the highest number of atoms appearing in a query pattern. We enumerate and compute the $\mathcal{Z}$ score of all the possible UCQs such that:
  - i.    Each CQ only uses atoms that either belong to a query pattern, or are implied by one of such atoms and the ontology
  - ii.   Each CQ has at most $k$ atoms

One can verify that the query $q(x) \leftarrow$ **ScientificStudent**$(x)$ achieves the highest $\mathcal{Z}$ score of 1.0, and is therefore the best explanation of the classifier $\lambda$.

# Conclusions

- Our framework uses the Ontology-Based Data Management paradigm to provide an explanation to the behavior of a classifier

- The short-term goal is to explore possible optimizations of the algorithm drafted in this presentation

- The future work includes an evaluation of the framework to real world scenarios, as well as comparison with other similar works