

LOGIC-BASED EXPLAINABILITY IN MACHINE LEARNING

Glimpse into ANITI's DeepLever Chair

Joao Marques-Silva

IRIT/CNRS & ANITI DeepLever Chair, Toulouse, France

September 2022

Recent & ongoing ML successes



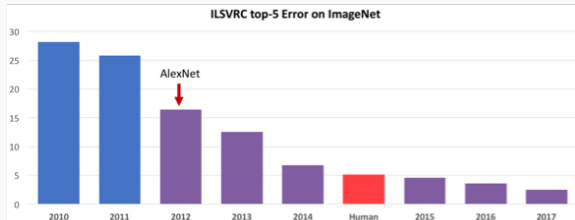
<https://en.wikipedia.org/wiki/Waymo>



AlphaGo Zero & Alpha Zero



Image & Speech Recognition

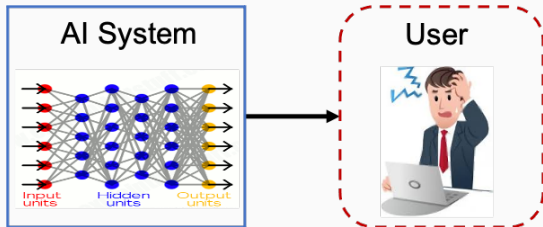


http://gradientscience.org/intro_adversarial/

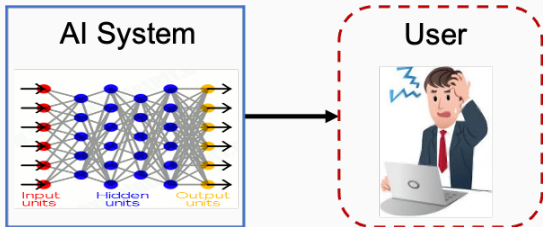


[https://fr.wikipedia.org/wiki/Pepper_\(robot\)](https://fr.wikipedia.org/wiki/Pepper_(robot))

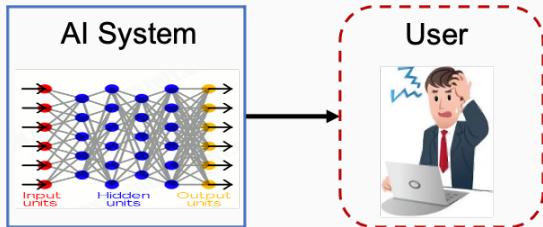
eXplainable AI (XAI)



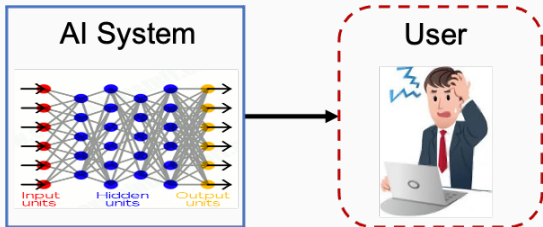
- ML models are most often **opaque**
- Goal of XAI: **to help humans understand ML models**
- Many questions to address:



- ML models are most often **opaque**
- Goal of XAI: **to help humans understand ML models**
- Many questions to address:
 - Properties of explanations
 - How to be human understandable?
 - How to answer **Why?** questions? I.e. Why the prediction?
 - How to answer **Why Not?** questions? I.e. Why not some other prediction?
 - Which guarantees of rigor?



- ML models are most often **opaque**
- Goal of XAI: **to help humans understand ML models**
- Many questions to address:
 - Properties of explanations
 - How to be human understandable?
 - How to answer **Why?** questions? I.e. Why the prediction?
 - How to answer **Why Not?** questions? I.e. Why not some other prediction?
 - Which guarantees of rigor?
 - Other queries: enumeration, membership, preferences, etc.



- ML models are most often **opaque**
- Goal of XAI: **to help humans understand ML models**
- Many questions to address:
 - Properties of explanations
 - How to be human understandable?
 - How to answer **Why?** questions? I.e. Why the prediction?
 - How to answer **Why Not?** questions? I.e. Why not some other prediction?
 - Which guarantees of rigor?
 - Other queries: enumeration, membership, preferences, etc.
 - Links with robustness, fairness, learning

REGULATION (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

of 27 April 2016

on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)

(Text with EEA relevance)

European Union regulations on algorithmic decision-making and a “right to explanation”

Bryce Goodman,^{1*} Seth Flaxman,²

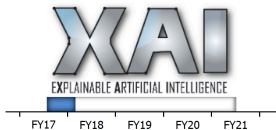
Proposal for a

REGULATION OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

LAYING DOWN HARMONISED RULES ON ARTIFICIAL INTELLIGENCE (ARTIFICIAL INTELLIGENCE ACT) AND AMENDING CERTAIN UNION LEGISLATIVE ACTS

■ We summarize the potential impact that the European Union's new General Data Protection Regulation will have on the routine use of machine-learning algorithms. Slated to take effect as law across the European Union in 2018, it will place restrictions on automated individual decision making (that is, algorithms that make decisions based on user-level predictors) that “significantly affect” users. When put into practice, the law may also effectively create a right to explanation, whereby a user can ask for an explanation of an algorithmic decision that significantly affects them. We argue that while this law may pose large challenges for industry, it highlights opportunities for computer scientists to take the lead in designing algorithms and evaluation frameworks that avoid discrimination and enable explanation.

Explainable Artificial Intelligence (XAI)



David Gunning
DARPA/I2O
Program Update November 2017



©DARPA

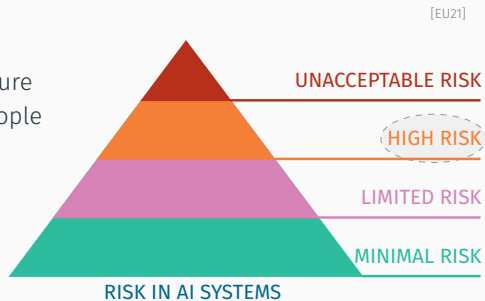
XAI for high-risk & safety-critical applications

- **High-risk** (EU regulations):

- Law enforcement
- Management and operation of critical infrastructure
- Biometric identification and categorization of people
- ...

- **Safety-critical:**

- Self-driving cars
- Autonomous vehicles
- Autonomous aerial devices
- ...



XAI for high-risk & safety-critical applications

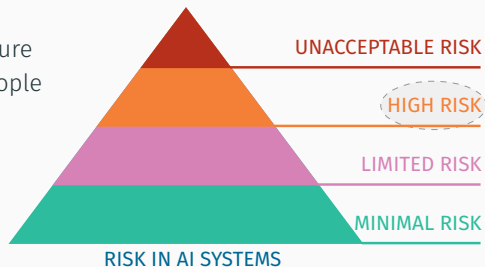
- **High-risk** (EU regulations):

- Law enforcement
- Management and operation of critical infrastructure
- Biometric identification and categorization of people
- ...

- **Safety-critical:**

- Self-driving cars
- Autonomous vehicles
- Autonomous aerial devices
- ...

[EU21]



PERSPECTIVE

<https://doi.org/10.1038/s42256-019-0048-x>

nature
machine intelligence

Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead

Cynthia Rudin

May 2019

XAI for high-risk & safety-critical applications

- **High-risk** (EU regulations):

- Law enforcement
- Management and operation of critical infrastructure
- Biometric identification and categorization of people
- ...

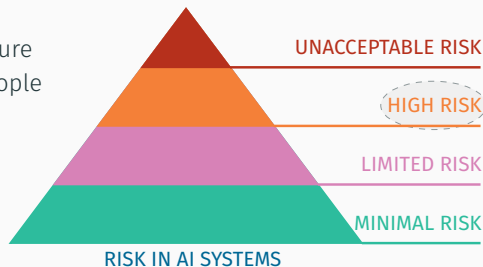
- **Safety-critical:**

- Self-driving cars
- Autonomous vehicles
- Autonomous aerial devices
- ...

- **Correctness of explanations is paramount!**

- To build trust
- To help debug AI systems
- To prevent accidents
- ...

[EU21]



PERSPECTIVE

<https://doi.org/10.1038/s42256-019-0048-x>

nature
machine intelligence

Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead

Cynthia Rudin

May 2019

XAI for high-risk & safety-critical applications

- **High-risk** (EU regulations):

- Law enforcement
- Management and operation of critical infrastructure
- Biometric identification and categorization of people
- ...

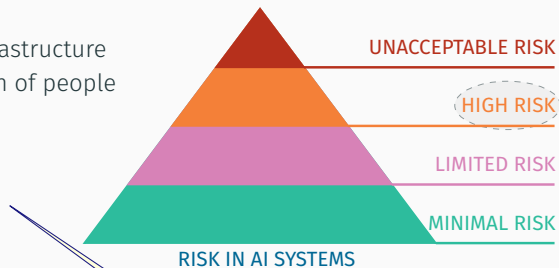
- **Safety-critical:**

- Self-driving cars
- Autonomous vehicles
- Autonomous aerial devices
- ...

- **Correctness of explanations is paramount!**

- To build trust
- To help debug AI systems
- To prevent accidents
- ...

[EU21]



Main motivation
for our work!

Outline

Basic Definitions

Limitations of Non-Formal XAI

Formal Explainability in AI

Progress in Formal Explainability

Beyond Computing Explanations

Beyond ML Explanations

Definitions – classification problems

- Set of features $\mathcal{F} = \{1, 2, \dots, m\}$, each feature i taking values from domain D_i
 - Features can be categorical, discrete or real-valued
 - Feature space: $\mathbb{F} = \prod_{i=1}^m D_i$
- Set of classes $\mathcal{K} = \{c_1, \dots, c_K\}$

Definitions – classification problems

- Set of features $\mathcal{F} = \{1, 2, \dots, m\}$, each feature i taking values from domain D_i
 - Features can be categorical, discrete or real-valued
 - Feature space: $\mathbb{F} = \prod_{i=1}^m D_i$
- Set of classes $\mathcal{K} = \{c_1, \dots, c_K\}$
- ML model \mathbb{M} computes a (non-constant) classification function $\kappa : \mathbb{F} \rightarrow \mathcal{K}$

Definitions – classification problems

- Set of features $\mathcal{F} = \{1, 2, \dots, m\}$, each feature i taking values from domain D_i
 - Features can be categorical, discrete or real-valued
 - Feature space: $\mathbb{F} = \prod_{i=1}^m D_i$
- Set of classes $\mathcal{K} = \{c_1, \dots, c_K\}$
- ML model \mathbb{M} computes a (non-constant) classification function $\kappa : \mathbb{F} \rightarrow \mathcal{K}$
- Instance (\mathbf{v}, c) for point $\mathbf{v} = (v_1, \dots, v_m) \in \mathbb{F}$, with prediction $c = \kappa(\mathbf{v})$, $c \in \mathcal{K}$
 - **Goal:** to compute explanations for (\mathbf{v}, c)

Outline

Basic Definitions

Limitations of Non-Formal XAI

Formal Explainability in AI

Progress in Formal Explainability

Beyond Computing Explanations

Beyond ML Explanations

Non-formal XAI approaches – among best known

- **Model-agnostic (MA) explainers:**

- **LIME & SHAP**

[RSG16, LL17]

- **Goal:** learn a simple interpretable ML model, e.g. linear classifier, decision tree, etc.
 - Approach: train classifier, sample-based vs. game theory

- **Anchor:**

[RSG18]

- **Goal:** learn features deemed **more** relevant for prediction
 - Anchor is sample-based

- **No formal guarantees of rigor in computed explanations**

- **Intrinsic interpretability (II)**

[Rud19, Mol20]

- (Interpretable) model is the explanation
 - E.g., DTs, DLs, DSs, etc.

Non-formal XAI approaches – among best known

- **Model-agnostic (MA) explainers:**

- **LIME & SHAP**

[RSG16, LL17]

- **Goal:** learn a simple interpretable ML model, e.g. linear classifier, decision tree, etc.
 - Approach: train classifier, sample-based vs. game theory

- **Anchor:**

[RSG18]

- **Goal:** learn features deemed **more** relevant for prediction
 - Anchor is sample-based

- **No formal guarantees of rigor in computed explanations**

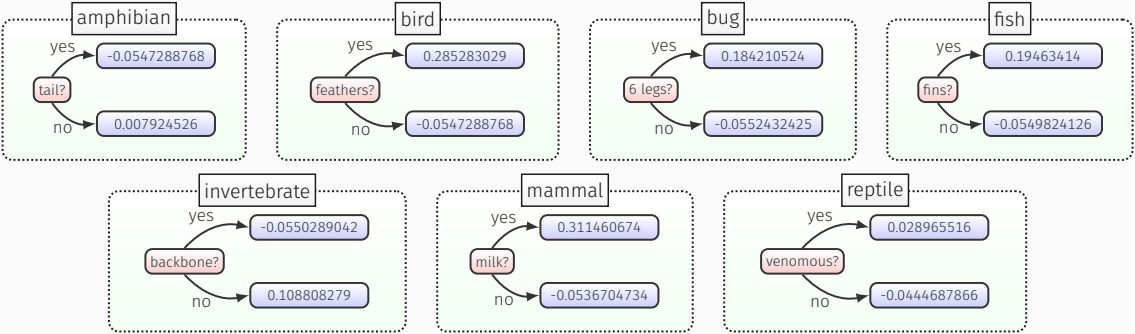
- **Intrinsic interpretability (II)**

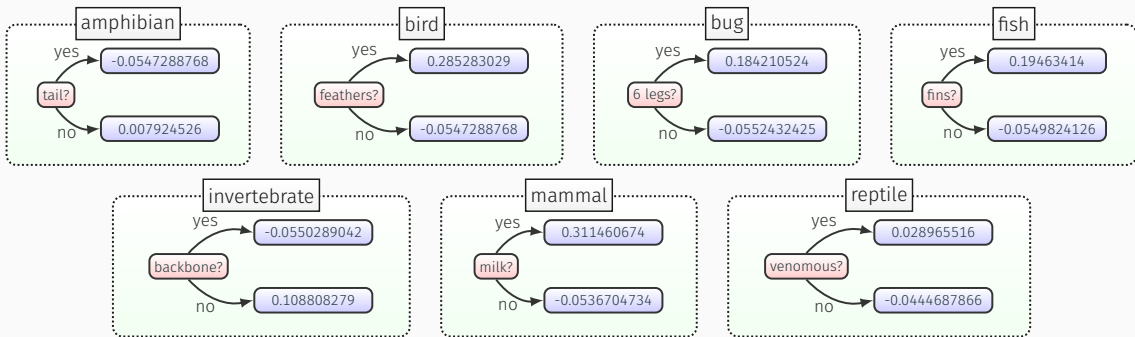
[Rud19, Mol20]

- (Interpretable) model is the explanation
 - E.g., DTs?, DLs?, DSs?, etc.?

Can MA explainers be trusted?
Is II indeed *interpretable*?

An example – BT for zoo dataset





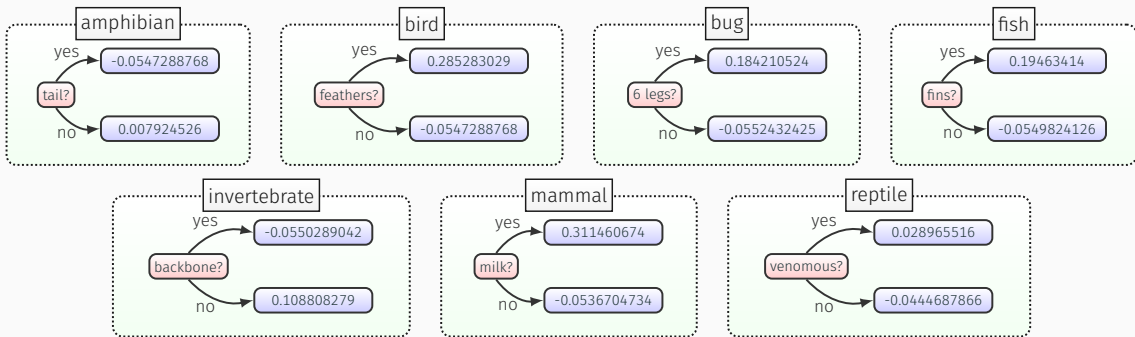
- Example instance:

IF $(\text{animal_name} = \text{pitviper}) \wedge \neg \text{hair} \wedge \neg \text{feathers} \wedge \text{eggs} \wedge \neg \text{milk} \wedge$
 $\neg \text{airborne} \wedge \neg \text{aquatic} \wedge \text{predator} \wedge \neg \text{toothed} \wedge \text{backbone} \wedge \text{breathes} \wedge$
 $\text{venomous} \wedge \neg \text{fins} \wedge (\text{legs} = 0) \wedge \text{tail} \wedge \neg \text{domestic} \wedge \neg \text{catsize}$

THEN $(\text{class} = \text{reptile})$

An example – BT for zoo dataset & Anchor

[INM19c, Ign20]



- Example instance (& Anchor picks):

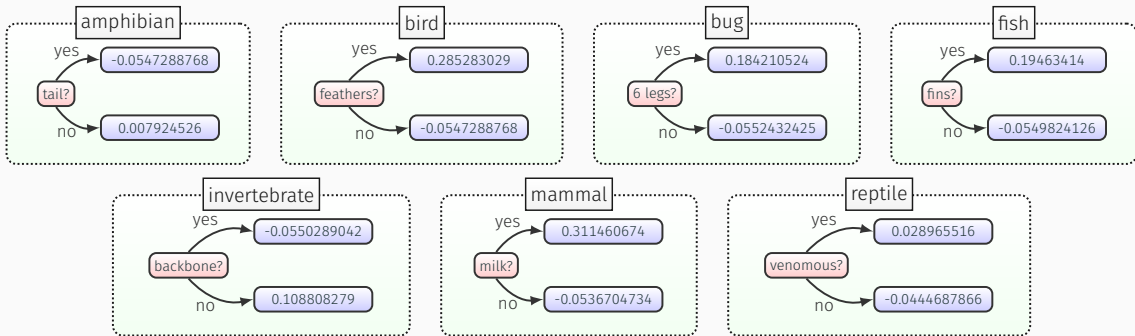
[RSG18]

IF (animal_name = pitviper) \wedge \neg hair \wedge \neg feathers \wedge eggs \wedge \neg milk \wedge \neg airborne \wedge \neg aquatic \wedge predator \wedge \neg toothed \wedge backbone \wedge breathes \wedge venomous \wedge \neg fins \wedge (legs = 0) \wedge tail \wedge \neg domestic \wedge \neg catsize

THEN (class = reptile)

An example – BT for zoo dataset & Anchor

[INM19c, Ign20]



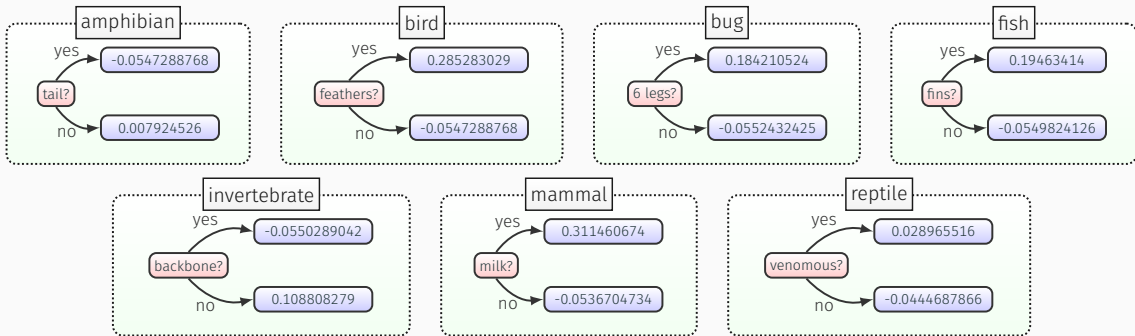
- Explanation obtained with Anchor:

[RSG18]

IF $\neg hair \wedge \neg milk \wedge \neg toothed \wedge \neg fins$
THEN (class = reptile)

An example – BT for zoo dataset & Anchor

[INM19c, Ign20]



- But, explanation **incorrectly “explains”** another instance (from **training data!**)

IF (animal_name = toad) \wedge \neg hair \wedge \neg feathers \wedge eggs \wedge \neg milk \wedge
 \neg airborne \wedge \neg aquatic \wedge \neg predator \wedge \neg toothed \wedge backbone \wedge breathes \wedge
 \neg venomous \wedge \neg fins \wedge (legs = 4) \wedge \neg tail \wedge \neg domestic \wedge \neg catsize

THEN (class = amphibian)

Incorrect explanations:

Classifier for deciding bank loans

Incorrect explanations:

Classifier for deciding bank loans

Two samples: Bessie := (v_1, \mathbf{Y}) and Clive := (v_2, \mathbf{N})

Incorrect explanations:

Classifier for deciding bank loans

Two samples: Bessie := (v_1, \mathbf{Y}) and Clive := (v_2, \mathbf{N})

Explanation X: age = 45, salary = 50K

Incorrect explanations:

Classifier for deciding bank loans

Two samples: Bessie := $(\mathbf{v}_1, \mathbf{Y})$ and Clive := $(\mathbf{v}_2, \mathbf{N})$

Explanation X: age = 45, salary = 50K

And,

X is consistent with Bessie := $(\mathbf{v}_1, \mathbf{Y})$

X is consistent with Clive := $(\mathbf{v}_2, \mathbf{N})$

Incorrect explanations:

Classifier for deciding bank loans

Two samples: Bessie := (v_1, \mathbf{Y}) and Clive := (v_2, \mathbf{N})

Explanation X: age = 45, salary = 50K

And,

X is consistent with Bessie := (v_1, \mathbf{Y})

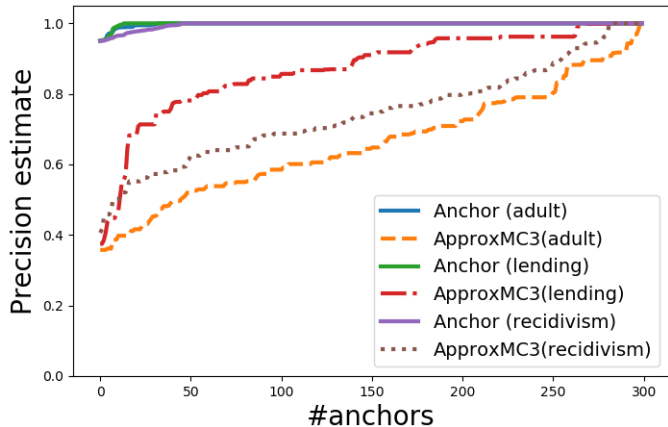
X is consistent with Clive := (v_2, \mathbf{N})

∴ different outcomes & same explanation !?

Incorrect explanations are ubiquitous...

[INM19c, Ign20]

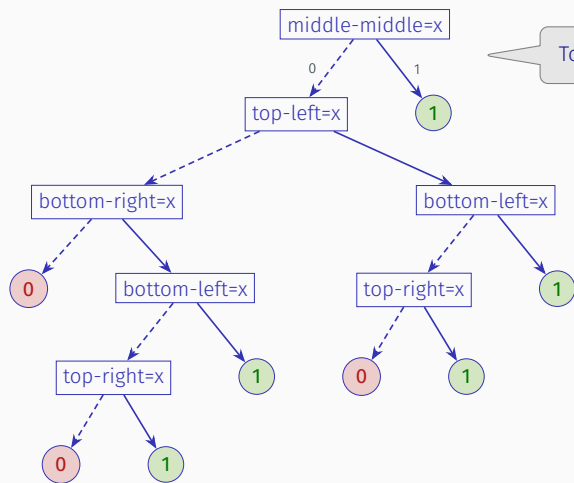
Dataset	# unique	Explanations								
		incorrect			redundant			correct		
		LIME	Anchor	SHAP	LIME	Anchor	SHAP	LIME	Anchor	SHAP
adult	(5579)	61.3%	80.5%	70.7%	7.9%	1.6%	10.2%	30.8%	17.9%	19.1%
lending	(4414)	24.0%	3.0%	17.0%	0.4%	0.0%	2.5%	75.6%	97.0%	80.5%
rcdv	(3696)	94.1%	99.4%	85.9%	4.6%	0.4%	7.9%	1.3%	0.2%	6.2%
compas	(778)	71.9%	84.4%	60.4%	20.6%	1.7%	27.8%	7.5%	13.9%	11.8%
german	(1000)	85.3%	99.7%	63.0%	14.6%	0.2%	37.0%	0.1%	0.1%	0.0%



“On interpretability, trees rate an A+.” [Bre01]

General agreement on interpretability of
decision trees [Rud19, Mol20, ANS20, Int21]

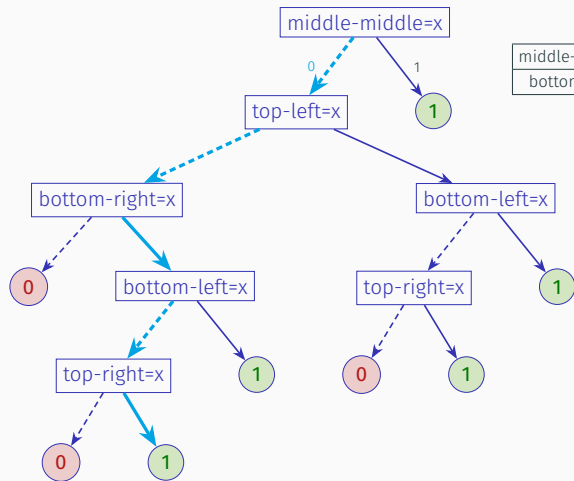
Decision tree explanations can be redundant



Tool: **OSDT**; train accuracy: 82.881%

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal Sparse Decision Trees**. **NeurIPS 2019**: 7265-7273

Decision tree explanations can be redundant

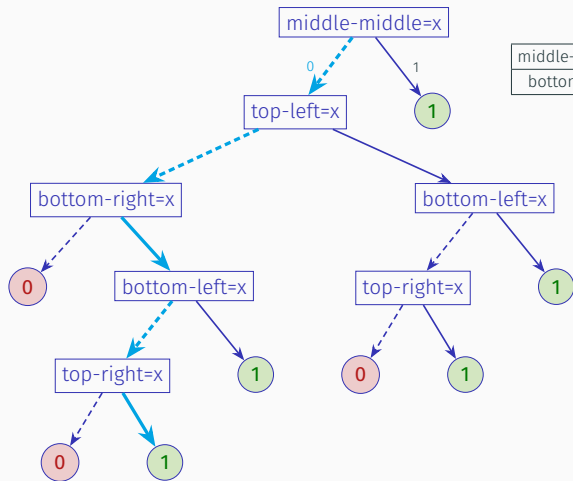


middle-middle = x iff MM = 1	top-left = x iff TL = 1	
bottom-right = x iff BR = 1	bottom-left = x iff BL = 1	top-right = x iff TR = 1

Q: What is explanation for
 $(MM = 0) \wedge (TL = 0) \wedge (BR = 1) \wedge$
 $(BL = 0) \wedge (TR = 1)$?

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal Sparse Decision Trees**. **NeurIPS 2019**: 7265-7273

Decision tree explanations can be redundant



middle-middle = x iff MM = 1	top-left = x iff TL = 1	
bottom-right = x iff BR = 1	bottom-left = x iff BL = 1	top-right = x iff TR = 1

Q: What is explanation for

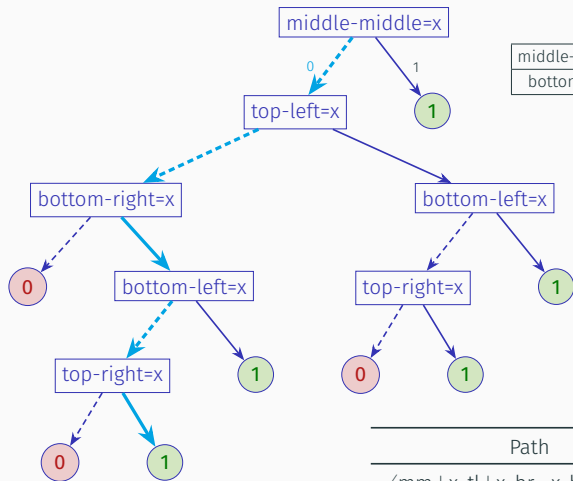
$(MM = 0) \wedge (TL = 0) \wedge (BR = 1) \wedge$

$(BL = 0) \wedge (TR = 1)?$

IF $(\neg MM \wedge \neg TL \wedge BR \wedge \neg BL \wedge TR)$ THEN $(\kappa(\cdot) = 1)?$

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal Sparse Decision Trees**. NeurIPS 2019: 7265-7273

Decision tree explanations can be redundant



middle-middle = x iff MM = 1	top-left = x iff TL = 1	
bottom-right = x iff BR = 1	bottom-left = x iff BL = 1	top-right = x iff TR = 1

Q: What is explanation for

$(MM = 0) \wedge (TL = 0) \wedge (BR = 1) \wedge$

$(BL = 0) \wedge (TR = 1)?$

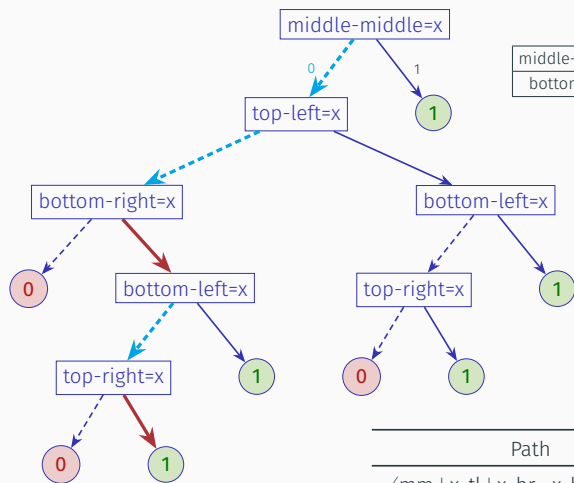
IF $(\neg MM \wedge \neg TL \wedge BR \wedge \neg BL \wedge TR)$ THEN $(\kappa(\cdot) = 1)?$

Path	Reduced XP	Dropped
$\langle mm \neq x, tl \neq x, br = x, bl \neq x, tr = x \rangle$		

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal**

Sparse Decision Trees. NeurIPS 2019: 7265-7273

Decision tree explanations can be arbitrarily redundant



middle-middle = x iff MM = 1	top-left = x iff TL = 1	
bottom-right = x iff BR = 1	bottom-left = x iff BL = 1	top-right = x iff TR = 1

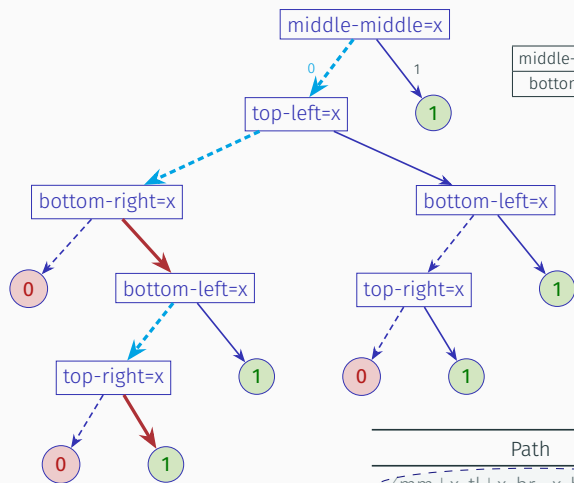
BR = 1, TR = 1			$\kappa(\cdot)$
MM	BL	TL	
0	0	0	1
0	0	1	1
0	1	0	1
0	1	1	1
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

Path	Reduced XP	Dropped
$\langle mm \neq x, tl \neq x, br = x, bl \neq x, tr = x \rangle$	$\langle br = x, tr = x \rangle$	$\{mm \neq x, tl \neq x, bl \neq x\}$

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal**

Sparse Decision Trees. NeurIPS 2019: 7265-7273

Decision tree explanations can be arbitrarily redundant



middle-middle = x iff MM = 1	top-left = x iff TL = 1	
bottom-right = x iff BR = 1	bottom-left = x iff BL = 1	top-right = x iff TR = 1

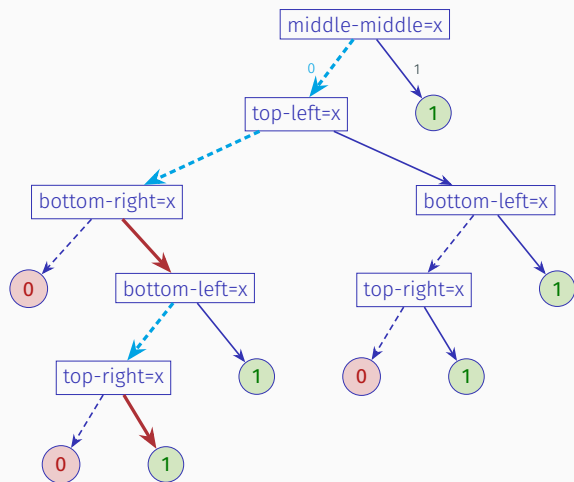
BR = 1, TR = 1			$\kappa(\cdot)$
MM	BL	TL	
0	0	0	1
0	0	1	1
0	1	0	1
0	1	1	1
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

Path	Reduced XP	Dropped
$\langle \langle mm \neq x, tl \neq x, br = x, bl \neq x, tr = x \rangle \rangle$	$\langle \langle br = x, tr = x \rangle \rangle$	$\{mm \neq x, tl \neq x, bl \neq x\}$

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal**

Sparse Decision Trees. NeurIPS 2019: 7265-7273

Decision tree explanations can be arbitrarily redundant



# tree paths	8
# red. paths	5
% red. paths	62.5%

Source: Xiyang Hu, Cynthia Rudin, Margo I. Seltzer: **Optimal Sparse Decision Trees**. *NeurIPS 2019*: 7265-7273

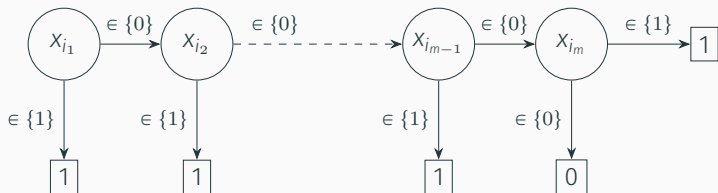
- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigwedge_{i=1}^m x_i$$

- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigvee_{i=1}^m x_i$$

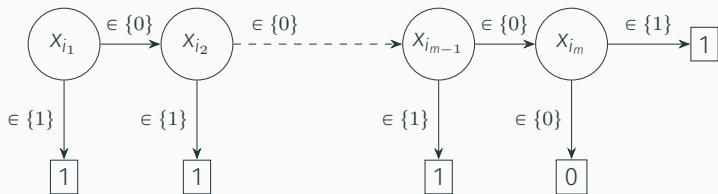
- Decision tree (DT), by picking variables in order $\langle i_1, i_2, \dots, i_m \rangle$, permutation of $\langle 1, 2, \dots, m \rangle$:



- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigvee_{i=1}^m x_i$$

- Decision tree (DT), by picking variables in order $\langle i_1, i_2, \dots, i_m \rangle$, permutation of $\langle 1, 2, \dots, m \rangle$:

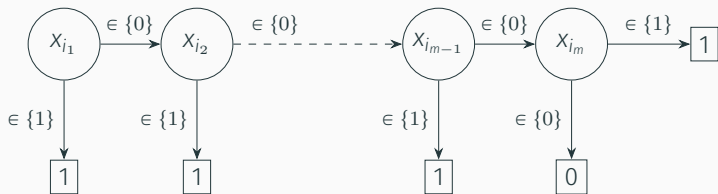


- Point: $(x_{i_1}, x_{i_2}, \dots, x_{i_{m-1}}, x_{i_m}) = (0, 0, \dots, 0, 1)$, and prediction 1

- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigvee_{i=1}^m x_i$$

- Decision tree (DT), by picking variables in order $\langle i_1, i_2, \dots, i_m \rangle$, permutation of $\langle 1, 2, \dots, m \rangle$:



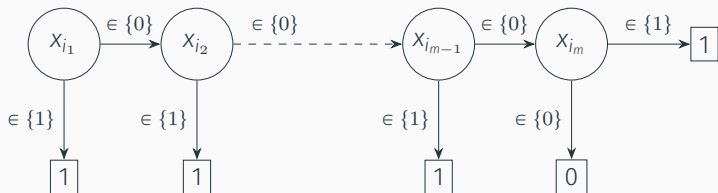
- Point: $(x_{i_1}, x_{i_2}, \dots, x_{i_{m-1}}, x_{i_m}) = (0, 0, \dots, 0, 1)$, and prediction 1
- Explanation using path in DT: $\{i_1, i_2, \dots, i_m\}$, i.e.

$$(x_{i_1} = 0) \wedge (x_{i_2} = 0) \wedge \dots \wedge (x_{i_{m-1}} = 0) \wedge (x_{i_m} = 1) \rightarrow \kappa(x_1, \dots, x_m)$$

- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigvee_{i=1}^m x_i$$

- Decision tree (DT), by picking variables in order $\langle i_1, i_2, \dots, i_m \rangle$, permutation of $\langle 1, 2, \dots, m \rangle$:



- Point: $(x_{i_1}, x_{i_2}, \dots, x_{i_{m-1}}, x_{i_m}) = (0, 0, \dots, 0, 1)$, and prediction 1
- Explanation using path in DT: $\{i_1, i_2, \dots, i_m\}$, i.e.

$$(x_{i_1} = 0) \wedge (x_{i_2} = 0) \wedge \dots \wedge (x_{i_{m-1}} = 0) \wedge (x_{i_m} = 1) \rightarrow \kappa(x_1, \dots, x_m)$$

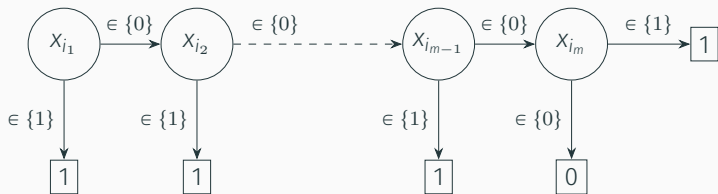
- But $\{i_m\}$ suffices for prediction, i.e. $\forall (\mathbf{x} \in \{0, 1\}^m). (x_{i_m}) \rightarrow \kappa(\mathbf{x})$

Minimal decision trees are **not** interpretable!

- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigvee_{i=1}^m x_i$$

- Decision tree (DT), by picking variables in order $\langle i_1, i_2, \dots, i_m \rangle$, permutation of $\langle 1, 2, \dots, m \rangle$:



- Point: $(x_{i_1}, x_{i_2}, \dots, x_{i_{m-1}}, x_{i_m}) = (0, 0, \dots, 0, 1)$, and prediction 1
- Explanation using path in DT: $\{i_1, i_2, \dots, i_m\}$, i.e.

$$(x_{i_1} = 0) \wedge (x_{i_2} = 0) \wedge \dots \wedge (x_{i_{m-1}} = 0) \wedge (x_{i_m} = 1) \rightarrow \kappa(x_1, \dots, x_m)$$

- But $\{i_m\}$ suffices for prediction, i.e. $\forall (\mathbf{x} \in \{0, 1\}^m). (x_{i_m}) \rightarrow \kappa(\mathbf{x})$**

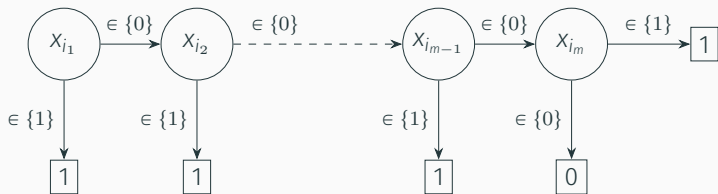
True explanations arbitrarily smaller than paths in DT; \therefore DTs are **not** interpretable*

Minimal decision trees are **not** interpretable!

- Classifier, with $x_1, \dots, x_m \in \{0, 1\}$:

$$\kappa(x_1, x_2, \dots, x_{m-1}, x_m) = \bigvee_{i=1}^m x_i$$

- Decision tree (DT), by picking variables in order $\langle i_1, i_2, \dots, i_m \rangle$, permutation of $\langle 1, 2, \dots, m \rangle$:



- Point: $(x_{i_1}, x_{i_2}, \dots, x_{i_{m-1}}, x_{i_m}) = (0, 0, \dots, 0, 1)$, and prediction 1
- Explanation using path in DT: $\{i_1, i_2, \dots, i_m\}$, i.e.

$$(x_{i_1} = 0) \wedge (x_{i_2} = 0) \wedge \dots \wedge (x_{i_{m-1}} = 0) \wedge (x_{i_m} = 1) \rightarrow \kappa(x_1, \dots, x_m)$$

- But $\{i_m\}$ suffices for prediction, i.e. $\forall (\mathbf{x} \in \{0, 1\}^m). (x_{i_m}) \rightarrow \kappa(\mathbf{x})$**

True explanations arbitrarily smaller than paths in DT; \therefore DTs are **not** interpretable*

*Pick any definition of interpretability that correlates with succinctness ...

Outline

Basic Definitions

Limitations of Non-Formal XAI

Formal Explainability in AI

Progress in Formal Explainability

Beyond Computing Explanations

Beyond ML Explanations

Formal explanations

- Instance (\mathbf{v}, c) , i.e. $c = \kappa(\mathbf{v})$

Formal explanations

- Instance (\mathbf{v}, c) , i.e. $c = \kappa(\mathbf{v})$
- **Abductive explanation** (AXp, PI-explanation) – answer to **Why?** question:
 - **Subset-minimal set of features** $\mathcal{X} \subseteq \mathcal{F}$ sufficient for **ensuring prediction**

[SCD18, INM19a]

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$$

Formal explanations

- Instance (\mathbf{v}, c) , i.e. $c = \kappa(\mathbf{v})$
- **Abductive explanation** (AXp, PI-explanation) – answer to **Why?** question:
 - **Subset-minimal set of features** $\mathcal{X} \subseteq \mathcal{F}$ sufficient for **ensuring prediction**

[SCD18, INM19a]

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$$

- **Contrastive explanation** (CXp) – answer to **Why Not?** question:
 - **Subset-minimal set of features** $\mathcal{Y} \subseteq \mathcal{F}$ sufficient for **changing prediction**

[Mil19, INAM20]

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\kappa(\mathbf{x}) \neq c)$$

Formal explanations

- Instance (\mathbf{v}, c) , i.e. $c = \kappa(\mathbf{v})$
- **Abductive explanation** (AXp, PI-explanation) – answer to **Why?** question:
 - **Subset-minimal set of features** $\mathcal{X} \subseteq \mathcal{F}$ sufficient for **ensuring prediction**

[SCD18, INM19a]

$$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$$

- **Contrastive explanation** (CXp) – answer to **Why Not?** question:
 - **Subset-minimal set of features** $\mathcal{Y} \subseteq \mathcal{F}$ sufficient for **changing prediction**

[Mil19, INAM20]

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\kappa(\mathbf{x}) \neq c)$$

- **Subset-minimal**, given predicate \mathbb{P} : $\mathbb{P}(\mathcal{Z}) \wedge \forall(\mathcal{Z}' \subsetneq \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z}')$

Formal explanations

- Instance (\mathbf{v}, c) , i.e. $c = \kappa(\mathbf{v})$
- **Abductive explanation** (AXp, PI-explanation) – answer to **Why?** question:
 - **Subset-minimal set of features** $\mathcal{X} \subseteq \mathcal{F}$ sufficient for **ensuring prediction**

[SCD18, INM19a]

$$\text{WeakAXp}(\mathcal{X}) := \forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$$

- **Contrastive explanation** (CXp) – answer to **Why Not?** question:
 - **Subset-minimal set of features** $\mathcal{Y} \subseteq \mathcal{F}$ sufficient for **changing prediction**

[Mil19, INAM20]

$$\text{WeakCXp}(\mathcal{Y}) := \exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\kappa(\mathbf{x}) \neq c)$$

- **Subset-minimal**, given predicate P : $P(\mathcal{Z}) \wedge \forall(\mathcal{Z}' \subsetneq \mathcal{Z}). \neg P(\mathcal{Z}')$

$$\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall(\mathcal{X}' \subsetneq \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X}')$$

$$\text{CXp}(\mathcal{Y}) := \text{WeakCXp}(\mathcal{Y}) \wedge \forall(\mathcal{Y}' \subsetneq \mathcal{Y}). \neg \text{WeakCXp}(\mathcal{Y}')$$

Formal explanations

- Instance (\mathbf{v}, c) , i.e. $c = \kappa(\mathbf{v})$
- **Abductive explanation (AXp, PI-explanation)** – answer to **Why?** question:
 - **Subset-minimal set of features** $\mathcal{X} \subseteq \mathcal{F}$ sufficient for **ensuring prediction**

[SCD18, INM19a]

$$\text{WeakAXp}(\mathcal{X}) := \forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$$

- **Contrastive explanation (CXp)** – answer to **Why Not?** question:
 - **Subset-minimal set of features** $\mathcal{Y} \subseteq \mathcal{F}$ sufficient for **changing prediction**

[Mil19, INAM20]

$$\text{WeakCXp}(\mathcal{Y}) := \exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\kappa(\mathbf{x}) \neq c)$$

- **Subset-minimal**, given predicate \mathbb{P} : $\mathbb{P}(\mathcal{Z}) \wedge \forall(\mathcal{Z}' \subsetneq \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z}')$

$$\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall(\mathcal{X}' \subsetneq \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X}')$$

$$\text{CXp}(\mathcal{Y}) := \text{WeakCXp}(\mathcal{Y}) \wedge \forall(\mathcal{Y}' \subsetneq \mathcal{Y}). \neg \text{WeakCXp}(\mathcal{Y}')$$

- **AXp's are minimal hitting sets (MHSEs) of CXp's and vice-versa** (& more)

[INAM20, INM19b]

- Builds on work of R. Reiter

[Rei87]

General approach

- Weak AXp/CXp conditions are **monotone** wrt set inclusion

General approach

- Weak AXp/CXp conditions are **monotone** wrt set inclusion \implies suffices to check: $\mathbb{P}(\mathcal{Z}) \wedge \forall (t \in \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z} \setminus \{t\})$
- E.g. AXp is defined by: $\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall (t \in \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X} \setminus \{t\})$

General approach

- Weak AXp/CXp conditions are **monotone** wrt set inclusion \implies suffices to check: $\mathbb{P}(\mathcal{Z}) \wedge \forall(t \in \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z} \setminus \{t\})$
 - E.g. AXp is defined by: $\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall(t \in \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X} \setminus \{t\})$
- Approach:
 - Encode classifier into suitable logic representation \mathcal{T}
 - For **AXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while AXp condition holds
 - For **CXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while CXp condition holds (*)

General approach

- Weak AXp/CXp conditions are **monotone** wrt set inclusion \implies suffices to check: $\mathbb{P}(\mathcal{Z}) \wedge \forall (t \in \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z} \setminus \{t\})$
 - E.g. AXp is defined by: $\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall (t \in \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X} \setminus \{t\})$
- Approach:
 - Encode classifier into suitable logic representation \mathcal{T}
 - For **AXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while AXp condition holds
 - For **CXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while CXp condition holds (*)
 - **Q**: which **reasoner** for \mathcal{T} to use?
- **Monotone** predicates for AXp & CXp:

$$\mathbb{P}_{\text{axp}}(\mathcal{S}) \triangleq \neg \text{CO} \left(\left[\left(\bigwedge_{i \in \mathcal{S}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \right] \right) \quad \mathbb{P}_{\text{cxp}}(\mathcal{S}) \triangleq \text{CO} \left(\left[\left[\left(\bigwedge_{i \in \mathcal{F} \setminus \mathcal{S}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \right] \right] \right)$$

General approach

- Weak AXp/CXp conditions are **monotone** wrt set inclusion \implies suffices to check: $\mathbb{P}(\mathcal{Z}) \wedge \forall(t \in \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z} \setminus \{t\})$
 - E.g. AXp is defined by: $\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall(t \in \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X} \setminus \{t\})$
- Approach:
 - Encode classifier into suitable logic representation \mathcal{T}
 - For **AXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while AXp condition holds
 - For **CXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while CXp condition holds (*)
 - **Q**: which **reasoner** for \mathcal{T} to use?
- **Monotone** predicates for AXp & CXp:

$$\mathbb{P}_{\text{axp}}(\mathcal{S}) \triangleq \neg \text{CO} \left(\left[\left(\bigwedge_{i \in \mathcal{S}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \right] \right)$$

$$\mathbb{P}_{\text{cxp}}(\mathcal{S}) \triangleq \text{CO} \left(\left[\left(\bigwedge_{i \in \mathcal{F} \setminus \mathcal{S}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \right] \right)$$

Input: Predicate \mathbb{P} , parameterized by $\mathcal{T}, \mathcal{F}, \kappa, \mathbf{v}$

Output: One XP \mathcal{S}

1: **procedure** oneXP(\mathbb{P})

2: $\mathcal{S} \leftarrow \mathcal{F}$

▷ Initialization: $\mathbb{P}(\mathcal{S})$ holds

3: **for** $i \in \mathcal{F}$ **do**

▷ Loop invariant: $\mathbb{P}(\mathcal{S})$ holds

4: **if** $\mathbb{P}(\mathcal{S} \setminus \{i\})$ **then**

5: $\mathcal{S} \leftarrow \mathcal{S} \setminus \{i\}$

▷ Update \mathcal{S} only if $\mathbb{P}(\mathcal{S} \setminus \{i\})$ holds

6: **return** \mathcal{S}

▷ Returned set \mathcal{S} : $\mathbb{P}(\mathcal{S})$ holds

General approach

- Weak AXp/CXp conditions are **monotone** wrt set inclusion \implies suffices to check: $\mathbb{P}(\mathcal{Z}) \wedge \forall(t \in \mathcal{Z}). \neg \mathbb{P}(\mathcal{Z} \setminus \{t\})$
 - E.g. AXp is defined by: $\text{AXp}(\mathcal{X}) := \text{WeakAXp}(\mathcal{X}) \wedge \forall(t \in \mathcal{X}). \neg \text{WeakAXp}(\mathcal{X} \setminus \{t\})$
- Approach:
 - Encode classifier into suitable logic representation \mathcal{T}
 - For **AXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while AXp condition holds
 - For **CXp**: start from $\mathcal{S} = \mathcal{F}$ and drop features from \mathcal{S} while CXp condition holds (*)
 - **Q**: which **reasoner** for \mathcal{T} to use, and how to **encode** classifiers?
- **Monotone** predicates for AXp & CXp:

$$\mathbb{P}_{\text{axp}}(\mathcal{S}) \triangleq \neg \text{CO} \left(\left[\left(\bigwedge_{i \in \mathcal{S}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \right] \right)$$

$$\mathbb{P}_{\text{cxp}}(\mathcal{S}) \triangleq \text{CO} \left(\left[\left(\bigwedge_{i \in \mathcal{F} \setminus \mathcal{S}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \right] \right)$$

Input: Predicate \mathbb{P} , parameterized by $\mathcal{T}, \mathcal{F}, \kappa, \mathbf{v}$

Output: One XP \mathcal{S}

```
1: procedure oneXP( $\mathbb{P}$ )
2:    $\mathcal{S} \leftarrow \mathcal{F}$ 
3:   for  $i \in \mathcal{F}$  do
4:     if  $\mathbb{P}(\mathcal{S} \setminus \{i\})$  then
5:        $\mathcal{S} \leftarrow \mathcal{S} \setminus \{i\}$ 
6:   return  $\mathcal{S}$ 
```

How to encode
NNs, RFs, BTs, etc. ?

▷ Initialization: $\mathbb{P}(\mathcal{S})$ holds
▷ Loop invariant: $\mathbb{P}(\mathcal{S})$ holds

▷ Update \mathcal{S} only if $\mathbb{P}(\mathcal{S} \setminus \{i\})$ holds
▷ Returned set \mathcal{S} : $\mathbb{P}(\mathcal{S})$ holds

A simple example – AXp's

- Classifier:

$$\kappa(X_1, X_2, X_3, X_4) = \bigvee_{i=1}^4 X_i$$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$?

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$?

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$?

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 4 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \top \rightarrow \kappa(x_1, x_2, x_3, x_4)$?

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 4 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \top \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **No**

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 4 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \top \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **No**
- AXp $\mathcal{X} = \{4\}$

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 4 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \top \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **No**
- AXp $\mathcal{X} = \{4\}$
- Validity/consistency checked with SAT/SMT/MILP/CP reasoners**

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 4 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \top \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **No**
- AXp $\mathcal{X} = \{4\}$
- **Validity/consistency checked with SAT/SMT/MILP/CP reasoners**
 - **Obs:** for some classes of classifiers, poly-time algorithms exist

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

A simple example – AXp's

- Classifier:

$$\kappa(x_1, x_2, x_3, x_4) = \bigvee_{i=1}^4 x_i$$

- Point $\mathbf{v} = (0, 0, 0, 1)$ with prediction $\kappa(\mathbf{v}) = 1$. **AXp?**
- Define $\mathcal{X} = \{1, 2, 3, 4\} = \mathcal{F}$
- Can feature 1 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_2 \wedge \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 2 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \neg x_3 \wedge x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 3 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). x_4 \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **Yes**
- Can feature 4 be removed, i.e. $\forall(\mathbf{x} \in \{0, 1\}^4). \top \rightarrow \kappa(x_1, x_2, x_3, x_4)$? **No**
- AXp $\mathcal{X} = \{4\}$
- Validity/consistency checked with SAT/SMT/MILP/CP reasoners**
 - Obs:** for some classes of classifiers, poly-time algorithms exist
- Similar approach for CXp's

Recap weak AXp: $\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\kappa(\mathbf{x}) = c)$

Relating formal explainability with MBD – a first attempt

- Weak abductive explanation (with $\mathcal{X} = \mathcal{F}$):

- Basic statement:

$$\forall(\mathbf{x} \in \mathbb{F}). \left(\bigwedge_{i \in \mathcal{F}} (x_i = v_i) \right) \rightarrow (\kappa(\mathbf{x}) = c)$$

- As entailment:

$$\left(\bigwedge_{i \in \mathcal{F}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \models \perp$$

- MBD mapping:

- Components, $C_i, i \in \mathcal{F}$: $(x_i = v_i)$
- System description, SD : $\bigwedge_i (C_i \vee Ab_i)$
- Observation, Obs : $(\kappa(\mathbf{x}) \neq c)$
- Hence,

$$SD \wedge Obs \wedge \bigwedge_{i \in \mathcal{F}} (\neg Ab_i) \models \perp$$

Relating formal explainability with MBD – a first attempt

- Weak abductive explanation (with $\mathcal{X} = \mathcal{F}$):

- Basic statement:

$$\forall(\mathbf{x} \in \mathbb{F}). \left(\bigwedge_{i \in \mathcal{F}} (x_i = v_i) \right) \rightarrow (\kappa(\mathbf{x}) = c)$$

- As entailment:

$$\left(\bigwedge_{i \in \mathcal{F}} (x_i = v_i) \right) \wedge (\kappa(\mathbf{x}) \neq c) \models \perp$$

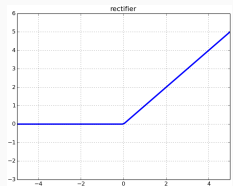
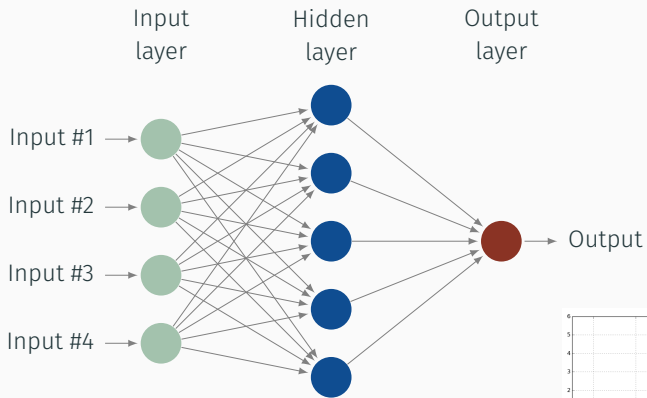
- MBD mapping:

- Components, $C_i, i \in \mathcal{F}$: $(x_i = v_i)$
- System description, SD : $\bigwedge_i (C_i \vee Ab_i)$
- Observation, Obs : $(\kappa(\mathbf{x}) \neq c)$
- Hence,

$$SD \wedge Obs \wedge \bigwedge_{i \in \mathcal{F}} (\neg Ab_i) \models \perp$$

- **Q:** Could MBD tools be efficient for XAI??

Encoding NNs



- Each layer (except first) viewed as a **block**, and
 - Compute \mathbf{x}' given input \mathbf{x} , weights matrix \mathbf{A} , and bias vector \mathbf{b}
 - Compute output \mathbf{y} given \mathbf{x}' and activation function
- Each unit uses a **ReLU** activation function

Encoding NNs (using MILP)

Computation for a NN ReLU **block**, in two steps:

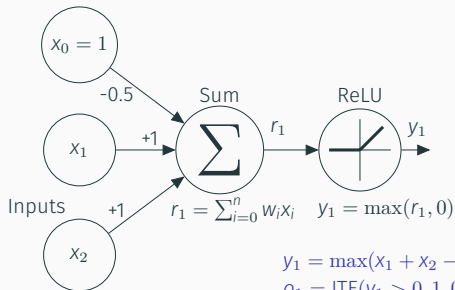
$$\begin{aligned} \mathbf{x}' &= \mathbf{A} \cdot \mathbf{x} + \mathbf{b} \\ \mathbf{y} &= \max(\mathbf{x}', \mathbf{0}) \end{aligned}$$

Encoding each **block**:

[F18]

$$\begin{aligned} \sum_{j=1}^n a_{i,j}x_j + b_i &= y_i - s_i \\ z_i = 1 &\rightarrow y_i \leq 0 \\ z_i = 0 &\rightarrow s_i \leq 0 \\ y_i \geq 0, s_i \geq 0, z_i &\in \{0, 1\} \end{aligned}$$

Encoding a simple NN in MILP



$$y_1 = \max(x_1 + x_2 - 0.5, 0)$$

$$o_1 = \text{ITE}(y_1 > 0, 1, 0)$$

x_1	x_2	r_1	y_1	o_1
0	0	-0.5	0	0
1	0	0.5	0.5	1
0	1	0.5	0.5	1
1	1	1.5	1.5	1

MILP encoding:

$$x_1 + x_2 - 0.5 = y_1 - s_1$$

$$z_1 = 1 \rightarrow y_1 \leq 0$$

$$z_1 = 0 \rightarrow s_1 \leq 0$$

$$o_1 = (y_1 > 0)$$

$$x_1, x_2, z_1, o_1 \in \{0, 1\}$$

$$y_1, s_1 \geq 0$$

Instance: $(\mathbf{x}, c) = ((1, 0), 1)$

$$1 + 0 - 0.5 = 0.5 - 0$$

$$1 \vee 0.5 \leq 0$$

$$0 \vee 0 \leq 0$$

$$1 = (0.5 > 0)$$

$$x_1 = 1, x_2 = 0, z_1 = 0, o_1 = 1$$

$$y_1 = 0.5, s_1 = 0$$

Checking: $\mathbf{x} = (0, 0)$

$$0 + 0 - 0.5 = 0 - 0.5$$

$$0 \vee 0 \leq 0$$

$$1 \vee 0.5 \leq 0$$

$$0 = (0 > 0)$$

$$x_1 = 0, x_2 = 0, z_1 = 1, o_1 = 0$$

$$y_1 = 0, s_1 = 0.5$$

Dataset			Minimal explanation			Minimum explanation		
			size	SMT (s)	MILP (s)	size	SMT (s)	MILP (s)
australian	(14)	m	1	0.03	0.05	—	—	—
		a	8.79	1.38	0.33	—	—	—
		M	14	17.00	1.43	—	—	—
backache	(32)	m	13	0.13	0.14	—	—	—
		a	19.28	5.08	0.85	—	—	—
		M	26	22.21	2.75	—	—	—
breast-cancer	(9)	m	3	0.02	0.04	3	0.02	0.03
		a	5.15	0.65	0.20	4.86	2.18	0.41
		M	9	6.11	0.41	9	24.80	1.81
cleve	(13)	m	4	0.05	0.07	4	—	0.07
		a	8.62	3.32	0.32	7.89	—	5.14
		M	13	60.74	0.60	13	—	39.06
hepatitis	(19)	m	6	0.02	0.04	4	0.01	0.04
		a	11.42	0.07	0.06	9.39	4.07	2.89
		M	19	0.26	0.20	19	27.05	22.23
voting	(16)	m	3	0.01	0.02	3	0.01	0.02
		a	4.56	0.04	0.13	3.46	0.3	0.25
		M	11	0.10	0.37	11	1.25	1.77
spect	(22)	m	3	0.02	0.02	3	0.02	0.04
		a	7.31	0.13	0.07	6.44	1.61	0.67
		M	20	0.88	0.29	20	8.97	10.73

First rigorous approach
for explaining NNs !

			Minimal explanation			Minimum explanation		
			size	SMT (s)	MILP (s)	size	SMT (s)	MILP (s)
australian	(14)	m	1	0.03	0.05	—	—	—
		a	8.79	1.38	0.33	—	—	—
		M	14	17.00	1.43	—	—	—
backache	(32)	m	13	0.13	0.14	—	—	—
		a	19.28	5.08	0.85	—	—	—
		M	26	22.21	2.75	—	—	—
breast-cancer	(9)	m	3	0.02	0.04	3	0.02	0.03
		a	5.15	0.65	0.20	4.86	2.18	0.41
		M	9	6.11	0.41	9	24.80	1.81
cleve	(13)	m	4	0.05	0.07	4	—	0.07
		a	8.62	3.32	0.32	7.89	—	5.14
		M	13	60.74	0.60	13	—	39.06
hepatitis	(19)	m	6	0.02	0.04	4	0.01	0.04
		a	11.42	0.07	0.06	9.39	4.07	2.89
		M	19	0.26	0.20	19	27.05	22.23
voting	(16)	m	3	0.01	0.02	3	0.01	0.02
		a	4.56	0.04	0.13	3.46	0.3	0.25
		M	11	0.10	0.37	11	1.25	1.77
spect	(22)	m	3	0.02	0.02	3	0.02	0.04
		a	7.31	0.13	0.07	6.44	1.61	0.67
		M	20	0.88	0.29	20	8.97	10.73

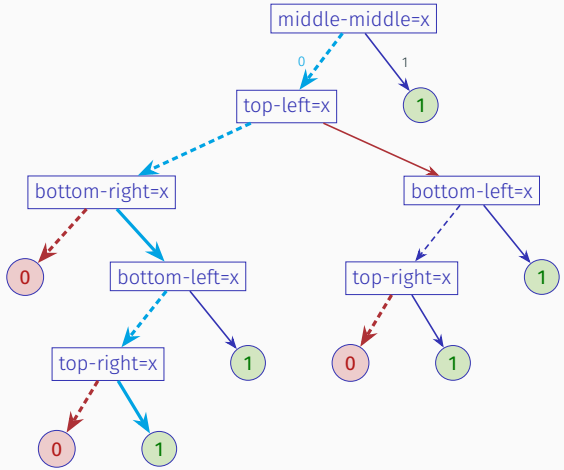
First rigorous approach
for explaining NNs !

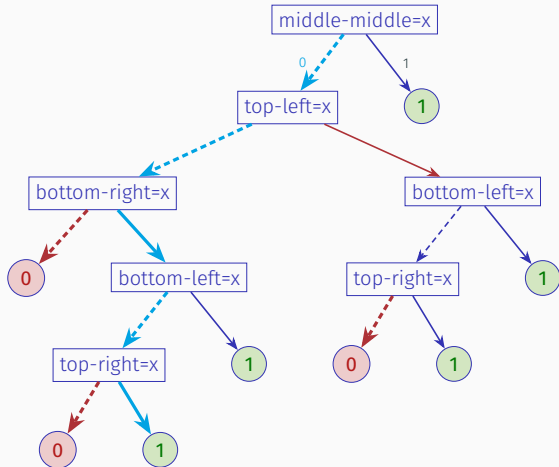
			Minimal explanation			Minimum explanation		
			size	SMT (s)	MILP (s)	size	SMT (s)	MILP (s)
australian	(14)	m	1	0.03	0.05	—	—	—
		a	8.79	1.38	0.33	—	—	—
		M	14	17.00	1.43	—	—	—
backache	(32)	m	13	0.13	0.14	—	—	—
		a	19.28	5.08	0.85	—	—	—
		M	26	22.21	2.75	—	—	—
breast-cancer	(9)	m	3	0.02	0.04	3	0.02	0.03
		a	5.15	0.65	0.20	4.86	2.18	0.41
		M	9	6.11	0.41	9	24.80	1.81
cleve	(13)	m	4	0.05	0.07	4	—	0.07
		a	8.62	3.32	0.32	7.89	—	5.14
		M	13	60.74	0.60	13	—	39.06
hepatitis	(19)	m	6	0.02	0.04	4	0.01	0.04
		a	11.42	0.07	0.06	9.39	4.07	2.89
		M	19	0.26	0.20	19	27.05	22.23
voting	(16)	m	3	0.01	0.02	3	0.01	0.02
		a	4.56	0.04	0.13	3.46	0.3	0.25
		M	11	0.10	0.37	11	1.25	1.77
spect	(22)	m	3	0.02	0.02	3	0.02	0.04
		a	7.31	0.13	0.07	6.44	1.61	0.67
		M	20	0.88	0.29	20	8.97	10.78

Scales to (a few)
tens of neurons...

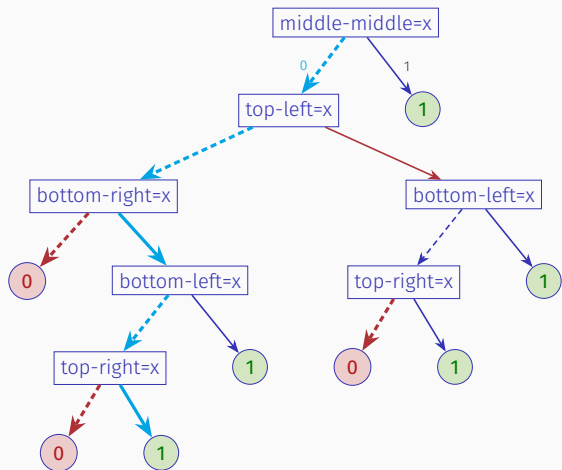
DT explanations

[11M20]

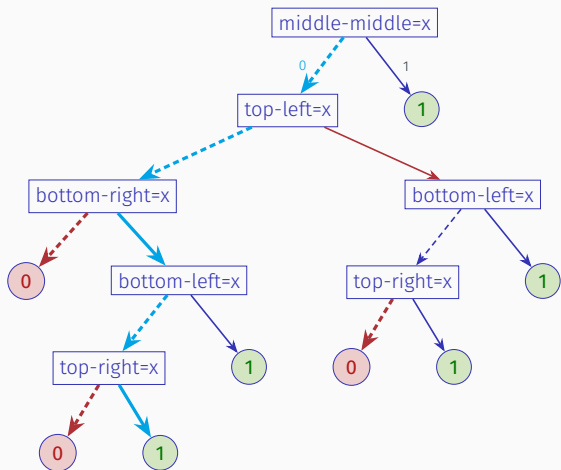




- Run PI-explanation algorithm based on NP-oracles
 - Worst-case exponential time



- Run PI-explanation algorithm based on NP-oracles
 - Worst-case exponential time
- For prediction **1**, it suffices to ensure **all** paths with prediction **0** remain inconsistent



- Run PI-explanation algorithm based on NP-oracles
 - Worst-case exponential time
- For prediction **1**, it suffices to ensure **all** paths with prediction **0** remain inconsistent
 - I.e. find a **subset-minimal hitting set** of **all 0** paths; these are the features to **keep**
 - E.g. BR and TR suffice for prediction
 - Well-known to be solvable in **polynomial time**

Preliminary results for DTs

[IIM20, IIM21, IIM22]

Dataset	(#F	#S)	IAI									ITI								
			D	#N	%A	#P	%R	%C	%m	%M	%avg	D	#N	%A	#P	%R	%C	%m	%M	%avg
adult	(12	6061)	6	83	78	42	33	25	20	40	25	17	509	73	255	75	91	10	66	22
anneal	(38	886)	6	29	99	15	26	16	16	33	21	9	31	100	16	25	4	12	20	16
backache	(32	180)	4	17	72	9	33	39	25	33	30	3	9	91	5	80	87	50	66	54
bank	(19	36293)	6	113	88	57	5	12	16	20	18	19	1467	86	734	69	64	7	63	27
biodegradation	(41	1052)	5	19	65	10	30	1	25	50	33	8	71	76	36	50	8	14	40	21
cancer	(9	449)	6	37	87	19	36	9	20	25	21	5	21	84	11	54	10	25	50	37
car	(6	1728)	6	43	96	22	86	89	20	80	45	11	57	98	29	65	41	16	50	30
colic	(22	357)	6	55	81	28	46	6	16	33	20	4	17	80	9	33	27	25	25	25
compas	(11	1155)	6	77	34	39	17	8	16	20	17	15	183	37	92	66	43	12	60	27
contraceptive	(9	1425)	6	99	49	50	8	2	20	60	37	17	385	48	193	27	32	12	66	21
dermatology	(34	366)	6	33	90	17	23	3	16	33	21	7	17	95	9	22	0	14	20	17
divorce	(54	150)	5	15	90	8	50	19	20	33	24	2	5	96	3	33	16	50	50	50
german	(21	1000)	6	25	61	13	38	10	20	40	29	10	99	72	50	46	13	12	40	22
heart-c	(13	302)	6	43	65	22	36	18	20	33	22	4	15	75	8	87	81	25	50	34
heart-h	(13	293)	6	37	59	19	31	4	20	40	24	8	25	77	13	61	60	20	50	32
kr-vs-kp	(36	3196)	6	49	96	25	80	75	16	60	33	13	67	99	34	79	43	7	70	35
lending	(9	5082)	6	45	73	23	73	80	16	50	25	14	507	65	254	69	80	12	75	25
letter	(16	18668)	6	127	58	64	1	0	20	20	20	46	4857	68	2429	6	7	6	25	9
lymphography	(18	148)	6	61	76	31	35	25	16	33	21	6	21	86	11	9	0	16	16	16
mortality	(118	13442)	6	111	74	56	8	14	16	20	17	26	865	76	433	61	61	7	54	19
mushroom	(22	8124)	6	39	100	20	80	44	16	33	24	5	23	100	12	50	31	20	40	25
pendigits	(16	10992)	6	121	88	61	0	0	—	—	—	38	937	85	469	25	86	6	25	11
promoters	(58	106)	1	3	90	2	0	0	—	—	—	3	9	81	5	20	14	33	33	33
recidivism	(15	3998)	6	105	61	53	28	22	16	33	18	15	611	51	306	53	38	9	44	16
seismic_bumps	(18	2578)	6	37	89	19	42	19	20	33	24	8	39	93	20	60	79	20	60	42
shuttle	(9	58000)	6	63	99	32	28	7	20	33	23	23	159	99	80	33	9	14	50	30
soybean	(35	623)	6	63	88	32	9	5	25	25	25	16	71	89	36	22	1	9	12	10
spambase	(57	4210)	6	63	75	32	37	12	16	33	19	15	143	91	72	76	98	7	58	25
spect	(22	228)	6	45	82	23	60	51	20	50	35	6	15	86	8	87	98	50	83	65
splice	(2	3178)	3	7	50	4	0	0	—	—	—	88	177	55	89	0	0	—	—	—

Outline

Basic Definitions

Limitations of Non-Formal XAI

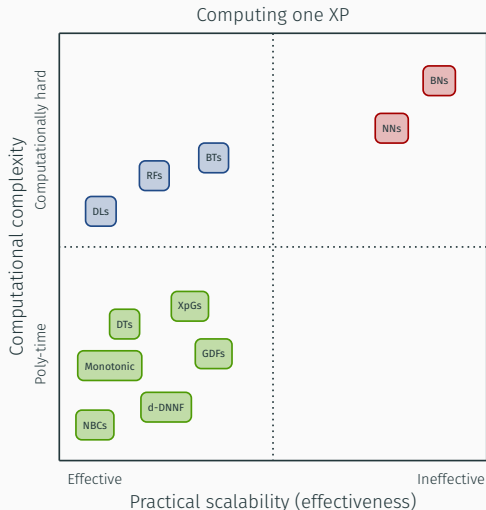
Formal Explainability in AI

Progress in Formal Explainability

Beyond Computing Explanations

Beyond ML Explanations

Efficacy map – current status



[INM19c, Ign20, IIM20, MGC⁺20, MGC⁺21, HIIM21, IMS21, IM21, CM21, HII⁺22, IISMS22]

• Formal explanations efficient for several families of classifiers

- Polynomial-time:
 - Naive-Bayes classifiers (NBCs) [MGC⁺20]
 - Decision trees (DTs) [IIM20, HIIM21]
 - XpG's: DTs, OBDDs, OMDDs, etc. [HIIM21]
 - Monotonic classifiers [MGC⁺21]
 - Propositional languages (e.g. d-DNNF, ...) [HII⁺22]
 - Additional results [CM21, HII⁺22]
- Comp. hard, but **effective** (efficient in practice):
 - Random forests (RFs) [IMS21]
 - Decision lists (DLs) [IM21]
 - Boosted trees (BTs) [INM19c, Ign20, IISMS22]
- Comp. hard, and **ineffective** (hard in practice):
 - Neural networks (NNs) [INM19a]
 - Bayesian networks (BNs) [SCD18]

Some recent results – RFs (with SAT)

[IMS21]

Dataset	(#F #C #I)	RF			CNF		SAT oracle				PI-expl (RFxp1)				Anchor	
		D	#N	%A	#var	#cl	MxS	MxU	#S	#U	Mx	m	avg	%w	avg	%w
ann-thyroid	(21 3 718)	4	2192	98	17854	29230	0.12	0.15	2	18	0.36	0.05	0.13	96	0.32	4
appendicitis	(7 2 43)	6	1920	90	5181	10085	0.02	0.02	4	3	0.05	0.01	0.03	100	0.48	0
banknote	(4 2 138)	5	2772	97	8068	16776	0.01	0.01	2	2	0.03	0.02	0.02	100	0.19	0
biodegradation	(41 2 106)	5	4420	88	11007	23842	0.31	1.05	17	22	2.27	0.04	0.29	97	4.07	3
heart-c	(13 2 61)	5	3910	85	5594	11963	0.04	0.02	6	7	0.07	0.01	0.04	100	0.85	0
ionosphere	(34 2 71)	5	2096	87	7174	14406	0.02	0.02	22	11	0.11	0.02	0.03	100	12.43	0
karhunen	(64 10 200)	5	6198	91	36708	70224	1.06	1.41	35	29	14.64	0.65	2.78	100	28.15	0
letter	(16 26 398)	8	44304	82	28991	68148	1.97	3.31	8	8	6.91	0.24	1.61	70	2.48	30
magic	(10 2 381)	6	9840	84	29530	66776	0.51	1.84	6	4	2.13	0.07	0.14	99	0.91	1
new-thyroid	(5 3 43)	5	1766	100	17443	28134	0.03	0.01	3	2	0.08	0.03	0.05	100	0.36	0
pendigits	(16 10 220)	6	12004	95	30522	59922	2.40	1.32	10	6	4.11	0.14	0.94	96	3.68	4
ring	(20 2 740)	6	6188	89	19114	42362	0.27	0.44	11	9	1.25	0.05	0.25	92	7.25	8
segmentation	(19 7 42)	4	1966	90	21288	35381	0.11	0.17	8	10	0.53	0.11	0.31	100	4.13	0
shuttle	(9 7 1160)	3	1460	99	18669	29478	0.11	0.08	2	7	0.34	0.05	0.14	99	0.42	1
sonar	(60 2 42)	5	2614	88	9938	20537	0.04	0.06	36	24	0.43	0.04	0.09	100	23.02	0
spectf	(44 2 54)	5	2306	88	6707	13449	0.07	0.06	20	24	0.34	0.02	0.07	100	8.12	0
texture	(40 11 550)	5	5724	87	34293	64187	0.79	0.63	23	17	3.24	0.19	0.93	100	28.13	0
twonorm	(20 2 740)	5	6266	94	21198	46901	0.08	0.08	12	8	0.28	0.06	0.10	100	5.73	0
vowel	(13 11 198)	6	10176	90	44523	88696	1.66	2.11	8	5	4.52	0.15	1.15	66	1.67	34
waveform-40	(40 3 500)	5	6232	83	30438	58380	0.50	0.86	15	25	7.07	0.11	0.88	100	11.93	0
wdbc	(33 2 78)	5	2432	76	9078	18675	1.00	1.53	20	13	5.33	0.03	0.65	79	3.91	21

Outline

Basic Definitions

Limitations of Non-Formal XAI

Formal Explainability in AI

Progress in Formal Explainability

Beyond Computing Explanations

Beyond ML Explanations

Additional computational problems – queries

- For a given instance (\mathbf{v}, c) :
 - **Enumeration:** list all/some/preferred explanations (AXp's/CXp's)
 - Exploit duality between AXp's & CXp's [INAM20]
 - For classifiers with explanations in P: **one SAT oracle call per computed explanation** [MGC⁺21, HIIM21]
 - For NBCs: enumeration with polynomial delay [MGC⁺20]
 - **Membership:** decide whether there exists explanation that includes target feature
 - Σ_2^P -hard for DNF classifiers [HIIM21]
 - In P for DTs [HIIM21]
 - In NP if finding one AXp/CXp in P [HM22]
 - **Probabilistic explanations:** compute set of features which, if fixed, the probability of predicting the target class is sufficiently large
 - NP^{PP}-hard for boolean circuit classifiers [WMHK21]
- Membership & enumeration for a prediction c , independently of point in feature space
- ...

Outline

Basic Definitions

Limitations of Non-Formal XAI

Formal Explainability in AI

Progress in Formal Explainability

Beyond Computing Explanations

Beyond ML Explanations

Explainability beyond ML explanations

- Model-based diagnosis & abduction
- AI planning
- Optimization
- Problem solving
- ...
- Any sort of algorithmic decision making

Conclusions

- Critical limitations of widely used XAI approaches:
 - Model-agnostic methods can compute **incorrect** explanations
 - Similar limitations for methods based on saliency maps for NNs
 - Intrinsic interpretability explanations can be (very) **redundant**

Conclusions

- Critical limitations of widely used XAI approaches:
 - Model-agnostic methods can compute **incorrect** explanations
 - Similar limitations for methods based on saliency maps for NNs
 - Intrinsic interpretability explanations can be (very) **redundant**
- **Formal explainability in AI (FXAI)**:
 - Logic-based, rigorous definitions of explanations
 - Initial theoretical insights, e.g. duality between AXp's and CXp's (and so between "Why?" and "Why not?" explanations)
 - Initial links with MBD
 - Other problems of crucial importance: **enumeration**, **membership**, etc.

Conclusions

- Critical limitations of widely used XAI approaches:
 - Model-agnostic methods can compute **incorrect** explanations
 - Similar limitations for methods based on saliency maps for NNs
 - Intrinsic interpretability explanations can be (very) **redundant**
- **Formal explainability in AI (FXAI)**:
 - Logic-based, rigorous definitions of explanations
 - Initial theoretical insights, e.g. duality between AXp's and CXp's (and so between "Why?" and "Why not?" explanations)
 - Initial links with MBD
 - Other problems of crucial importance: **enumeration**, **membership**, etc.
- Ongoing research related with:
 - Computation of AXp's & CXp's
 - Enumeration of explanations
 - Deciding membership of features in explanations
 - Probabilistic rigorous explanations
 - Also, complexity of explainability

Q & A

Acknowledgment: joint work with A. Ignatiev, Y. Izza, X. Huang, M. Cooper, N. Asher, N. Narodytska, E. Hebrard, M. Siala, et al.

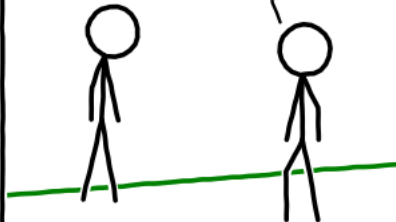
BLACK BOX MODELS

MY ML MODEL...

IS LIKE A
(BLACK) BOX OF
CHOCOLATES.

I NEVER KNOW WHAT
I'M GONNA GET.

BUT WHY?



<http://arxiv.org/abs/1901.01686> & <http://cmx.io/ed1/>

References i

- [ANS20] Gaël Aglin, Siegfried Nijssen, and Pierre Schaus.
PyDL8.5: a library for learning optimal decision trees.
In *IJCAI*, pages 5222–5224, 2020.
- [Bre01] Leo Breiman.
Statistical modeling: The two cultures.
Statistical science, 16(3):199–231, 2001.
- [CM21] Martin C. Cooper and Joao Marques-Silva.
On the tractability of explaining decisions of classifiers.
In *CP*, October 2021.
- [EG95] Thomas Eiter and Georg Gottlob.
Identifying the minimal transversals of a hypergraph and related problems.
SIAM J. Comput., 24(6):1278–1304, 1995.
- [EU21] EU.
European Artificial Intelligence Act – Proposal.
<https://eur-lex.europa.eu/legal-content/EN/TXT/?qid=1623335154975&uri=CELEX%3A52021PC0206>,
2021.

References ii

- [FJ18] Matteo Fischetti and Jason Jo.
Deep neural networks and mixed integer linear optimization.
Constraints, 23(3):296–309, 2018.
- [HII+22] Xuanxiang Huang, Yacine Izza, Alexey Ignatiev, Martin Cooper, Nicholas Asher, and Joao Marques-Silva.
Tractable explanations for d-DNNF classifiers.
In *AAAI*, February 2022.
- [HIIM21] Xuanxiang Huang, Yacine Izza, Alexey Ignatiev, and Joao Marques-Silva.
On efficiently explaining graph-based classifiers.
In *KR*, November 2021.
Preprint available from <https://arxiv.org/abs/2106.01350>.
- [HM22] Xuanxiang Huang and João Marques-Silva.
On deciding feature membership in explanations of SDD & related classifiers.
CoRR, abs/2202.07553, 2022.
- [Ign20] Alexey Ignatiev.
Towards trustable explainable AI.
In *IJCAI*, pages 5154–5158, 2020.

References iii

- [IIM20] Yacine Izza, Alexey Ignatiev, and Joao Marques-Silva.
On explaining decision trees.
CoRR, abs/2010.11034, 2020.
- [IIM22] Yacine Izza, Alexey Ignatiev, and João Marques-Silva.
On tackling explanation redundancy in decision trees.
J. Artif. Intell. Res., 2022.
In Press. Preprint available from <https://doi.org/10.48550/arXiv.2205.09971>.
- [IISMS22] Alexey Ignatiev, Yacine Izza, Peter J. Stuckey, and Joao Marques-Silva.
Using MaxSAT for efficient explanations of tree ensembles.
In *AAAI*, February 2022.
- [IM21] Alexey Ignatiev and Joao Marques-Silva.
SAT-based rigorous explanations for decision lists.
In *SAT*, pages 251–269, July 2021.
- [IMS21] Yacine Izza and Joao Marques-Silva.
On explaining random forests with SAT.
In *IJCAI*, pages 2584–2591, July 2021.

References iv

- [INAM20] Alexey Ignatiev, Nina Narodytska, Nicholas Asher, and João Marques-Silva.
From contrastive to abductive explanations and back again.
In *AlxIA*, pages 335–355, 2020.
- [INM19a] Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva.
Abduction-based explanations for machine learning models.
In *AAAI*, pages 1511–1519, 2019.
- [INM19b] Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva.
On relating explanations and adversarial examples.
In *NeurIPS*, pages 15857–15867, 2019.
- [INM19c] Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva.
On validating, repairing and refining heuristic ML explanations.
CoRR, abs/1907.02509, 2019.
- [Int21] Interpretable AI.
<https://www.interpretable.ai/>, 21.
- [Lip18] Zachary C. Lipton.
The mythos of model interpretability.
Commun. ACM, 61(10):36–43, 2018.

References v

- [LL17] Scott M. Lundberg and Su-In Lee.
A unified approach to interpreting model predictions.
In *NIPS*, pages 4765–4774, 2017.
- [MGC⁺20] Joao Marques-Silva, Thomas Gerspacher, Martin C. Cooper, Alexey Ignatiev, and Nina Narodytska.
Explaining naive bayes and other linear classifiers with polynomial time and delay.
In *NeurIPS*, 2020.
- [MGC⁺21] Joao Marques-Silva, Thomas Gerspacher, Martin C. Cooper, Alexey Ignatiev, and Nina Narodytska.
Explanations for monotonic classifiers.
In *ICML*, pages 7469–7479, July 2021.
- [Mil19] Tim Miller.
Explanation in artificial intelligence: Insights from the social sciences.
Artif. Intell., 267:1–38, 2019.
- [Mol20] Christoph Molnar.
Interpretable machine learning.
Lulu.com, 2020.
<https://christophm.github.io/interpretable-ml-book/>.

References vi

- [NH10] Vinod Nair and Geoffrey E. Hinton.
Rectified linear units improve restricted boltzmann machines.
In *ICML*, pages 807–814, 2010.
- [NSM⁺19] Nina Narodytska, Aditya A. Shrotri, Kuldeep S. Meel, Alexey Ignatiev, and Joao Marques-Silva.
Assessing heuristic machine learning explanations with model counting.
In *SAT*, pages 267–278, 2019.
- [Rei87] Raymond Reiter.
A theory of diagnosis from first principles.
Artif. Intell., 32(1):57–95, 1987.
- [RSG16] Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin.
"why should I trust you?": Explaining the predictions of any classifier.
In *KDD*, pages 1135–1144, 2016.
- [RSG18] Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin.
Anchors: High-precision model-agnostic explanations.
In *AAAI*, pages 1527–1535. AAAI Press, 2018.

References vii

- [Rud19] Cynthia Rudin.
Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead.
Nature Machine Intelligence, 1(5):206–215, 2019.
- [SCD18] Andy Shih, Arthur Choi, and Adnan Darwiche.
A symbolic approach to explaining bayesian network classifiers.
In *IJCAI*, pages 5103–5111, 2018.
- [WMHK21] Stephan Wäldchen, Jan MacDonald, Sascha Hauch, and Gitta Kutyniok.
The computational complexity of understanding binary classifier decisions.
J. Artif. Intell. Res., 70:351–387, 2021.