

Formal Explanations in Artificial Intelligence

ESSAI 2025

Julien Girard-Satabin
Zakaria Chihani
Dorin Doncenco

CEA LIST

2025-07-02

This work was supported by the French Agence Nationale de la Recherche (ANR) through SAIF (ANR-23-PEIA-0006) and DeepGreen (ANR-23-DEGR-0001) as part of the France 2030 programme.

About lecturer



Figure 1: Questions may be asked during the lecture!

- BSc in Computer Science & Engineering
- MSc in Artificial Intelligence (AI)
- Interested in applied eXplainable AI (XAI)
- PhD Student on Formal and Actionable XAI
- Likes applying the handshake protocol to human conversations



Session goal

- Understand importance of eXplainability in AI systems
- Discuss what makes an explanation
- Learn limitations of common XAI methods
- Address these limitations with Formal XAI (FXAI)
- Discover new FXAI concepts
- Practical results with FXAI

Motivation

Omnipresence of automated systems



Omnipresence of automated systems



Why is explainability important in XAI? Please give me an answer in only 2 sentences, nobody in my classroom is going to read your very long texts...

Explainability in XAI (Explainable AI) is important because it helps humans understand how and why AI systems make decisions, which builds trust and accountability. It also allows for easier detection of errors, biases, or unethical outcomes in AI behavior.



Figure 2: LLMs are large models that have seen quick adoption within society.

Omnipresence of automated systems



Why is explainability important in XAI? Please give me an answer in only 2 sentences, nobody in my classroom is going to read your very long texts...

Explainability in XAI (Explainable AI) is important because it helps humans understand how and why AI systems make decisions, which builds trust and accountability. It also allows for easier detection of errors, biases, or unethical outcomes in AI behavior.



Figure 2: LLMs are large models that have seen quick adoption within society.

“AI Is a Better Teacher Than Humans”
- Duolingo CEO



Figure 3: AI might be adopted in education.

Omnipresence of automated systems



Why is explainability important in XAI? Please give me an answer in only 2 sentences, nobody in my classroom is going to read your very long texts...

Explainability in XAI (Explainable AI) is important because it helps humans understand how and why AI systems make decisions, which builds trust and accountability. It also allows for easier detection of errors, biases, or unethical outcomes in AI behavior.



Figure 2: LLMs are large models that have seen quick adoption within society.

“AI Is a Better Teacher Than Humans”
- Duolingo CEO



Figure 3: AI might be adopted in education.



Figure 4: High stakes for avoiding harm from autonomous vehicles.

Omnipresence of automated systems



Why is explainability important in XAI? Please give me an answer in only 2 sentences, nobody in my classroom is going to read your very long texts...

Explainability in XAI (Explainable AI) is important because it helps humans understand how and why AI systems make decisions, which builds trust and accountability. It also allows for easier detection of errors, biases, or unethical outcomes in AI behavior.



Figure 2: LLMs are large models that have seen quick adoption within society.

“AI Is a Better Teacher Than Humans”
- Duolingo CEO



Figure 3: AI might be adopted in education.



Figure 4: High stakes for avoiding harm from autonomous vehicles.

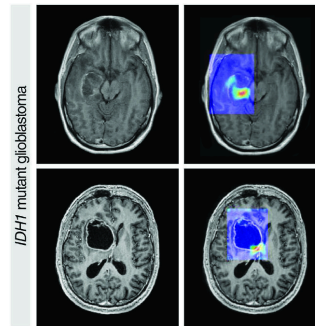


Figure 5: Automatic tumor detection has the potential to save lives.

Omnipresence of automated systems



Why is explainability important in XAI? Please give me an answer in only 2 sentences, nobody in my classroom is going to read your very long texts...

Explainability in XAI (Explainable AI) is important because it helps humans understand how and why AI systems make decisions, which builds trust and accountability. It also allows for easier detection of errors, biases, or unethical outcomes in AI behavior.



Figure 2: LLMs are large models that have seen quick adoption within society.

"AI Is a Better Teacher Than Humans"
- Duolingo CEO



Figure 3: AI might be adopted in education.



Figure 4: High stakes for avoiding harm from autonomous vehicles.

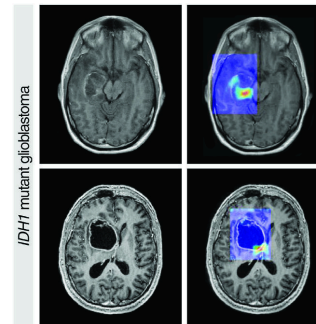


Figure 5: Automatic tumor detection has the potential to save lives.

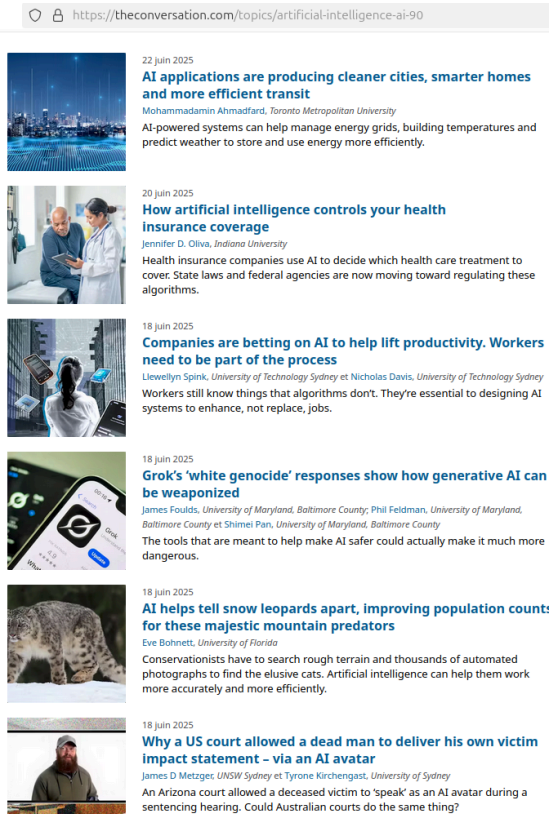


Figure 6: AI is a discussion point on many important topics.

The impact of automated decision systems

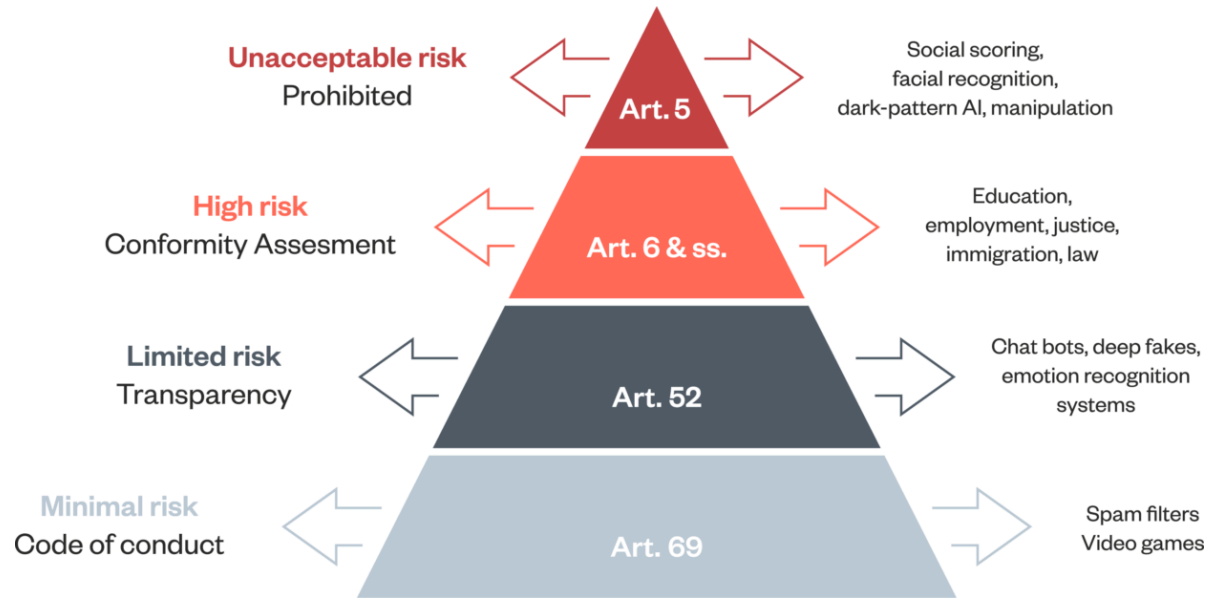
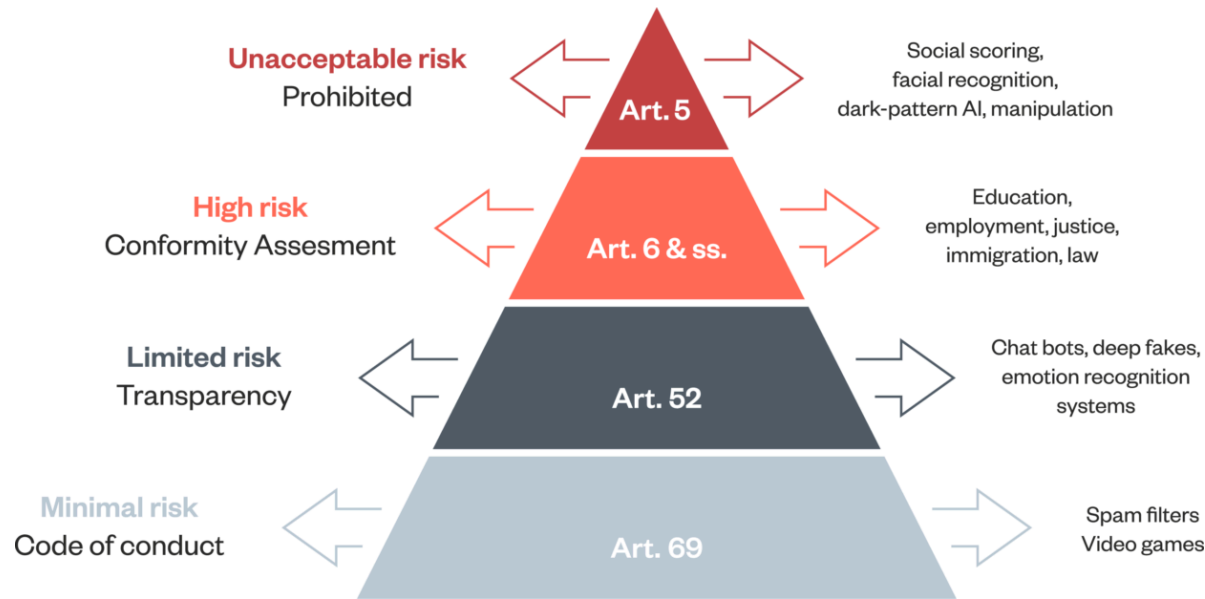


Figure 7: Risk hierarchy of AI use-cases¹. Various institutions have differing definitions.

¹<https://www.adalovelaceinstitute.org/resource/eu-ai-act-explainer/>

²https://en.wikipedia.org/wiki/Right_to_explanation

The impact of automated decision systems



Regulating high risks involves:

- transparency and information to users
- right to explanation²
- ensuring the safety of these systems

Figure 7: Risk hierarchy of AI use-cases¹. Various institutions have differing definitions.

¹<https://www.adalovelaceinstitute.org/resource/eu-ai-act-explainer/>

²https://en.wikipedia.org/wiki/Right_to_explanation

Explanations | Introduction



Proof vs Explanations

- The concept of explanation is illusive
- Not easy to explain something



Proof vs Explanations

- The concept of explanation is illusive
- Not easy to explain something

Let's demonstrate with an example:



Proof vs Explanations

- The concept of explanation is illusive
- Not easy to explain something

Let's demonstrate with an example:

$$12345679 \times 36 = 444444444$$



Proof vs Explanations

- The concept of explanation is illusive
- Not easy to explain something

Let's demonstrate with an example:

$$12345679 \times 36 = 444444444$$

Why is this result particular such that it consists only of the digit 4?

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline \end{array}$$

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \end{array}$$

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \end{array}$$

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \\ \hline \end{array}$$

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \\ \hline 444444444 \end{array}$$

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \\ \hline 444444444 \end{array}$$

- What is your take on such an explanation?

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \\ \hline 444444444 \end{array}$$

- What is your take on such an explanation?
- What is lacking about the explanation?

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \\ \hline 444444444 \end{array}$$

- What is your take on such an explanation?
- What is lacking about the explanation?
- We have proven the truthfulness (good)
- It doesn't help understand the result (bad)

First attempt – Demonstration

$$\begin{array}{r} 12345679 \\ \times \quad 36 \\ \hline 74074074 \\ +37037037 \\ \hline 444444444 \end{array}$$

- What is your take on such an explanation?
- What is lacking about the explanation?
- We have proven the truthfulness (good)
- It doesn't help understand the result (bad)
- Additional questions:
- Why does this result work for 444,444,444?
- Can we make it work for 555,555,555?



Second attempt – Explanation #1

$$444444444 = 12345679 \times 36$$

$$444444444 = 111111111 \times 4$$

$$111111111 = 12345679 \times 9$$

- This explanation allows us to generalize the example;



Second attempt – Explanation #1

$$444444444 = 12345679 \times 36$$

$$444444444 = 111111111 \times 4$$

$$111111111 = 12345679 \times 9$$

- This explanation allows us to generalize the example;
- What other questions are left out?



Second attempt – Explanation #1

$$444444444 = 12345679 \times 36$$

$$444444444 = 111111111 \times 4$$

$$111111111 = 12345679 \times 9$$

- This explanation allows us to generalize the example;
- What other questions are left out?
- We still don't know why 12345679 is interesting here
- Why does 111,111,111 have 9 digits?
- Why is the number 8 missing?

Third attempt – Explanation #2

11111111|9

Third attempt – Explanation #2

$$\begin{array}{r} 11111111|9 \\ \hline 1 \end{array}$$

Third attempt – Explanation #2

11111111|9

21

12

Third attempt – Explanation #2

11111111|9

21 _____

31 123

Third attempt – Explanation #2

11111111|9

21	_____
31	1234
41	

Third attempt – Explanation #2

11111111|9

21	_____
31	1234
41	



Third attempt – Explanation #2

1111 11111|9

21	_____
31	1234
41	

$$10 \times n + 1 = (9 \times n) + (n + 1)$$



Third attempt – Explanation #2

1111 1111|9

21	_____
31	1234
41	

$$10 \times n + 1 = (9 \times n) + (n + 1)$$

- left side: tens and ones
- right side: quotient + remainder



Third attempt – Explanation #2

$$\begin{array}{r} 11111111|9 \\ 21 \quad \text{-----} \\ 31 \quad 1234 \\ 41 \end{array}$$

$$10 \times n + 1 = (9 \times n) + (n + 1)$$

- left side: tens and ones
- right side: quotient + remainder
- At $n = 8$, we get the remainder 0

Third attempt – Explanation #2

11111111|9

21	_____
31	1234567
41	
51	
61	
71	
81	

$$10 \times n + 1 = (9 \times n) + (n + 1)$$

- left side: tens and ones
- right side: quotient + remainder
- At $n = 8$, we get the remainder 0

Third attempt – Explanation #2

11111111|9

21	_____
31	12345679
41	
51	
61	
71	
81	
0	

$$10 \times n + 1 = (9 \times n) + (n + 1)$$

- left side: tens and ones
- right side: quotient + remainder
- At $n = 8$, we get the remainder 0

Takeaways

- There are multiple ways of explaining something
- Proving a statement does not necessarily lead to an insight into the problem

<https://www.youtube.com/watch?v=6j8Vwbss038> (by Gilles Dowek, in French)



Takeaways

- There are multiple ways of explaining something
- Proving a statement does not necessarily lead to an insight into the problem
- Exercise source¹

¹<https://www.youtube.com/watch?v=6j8Vwbss038> (by Gilles Dowek, in French)

Interactive Time

- What makes an explanation?
- [Skip](#)

Interactive Time

- What makes an explanation?
- Case study
- *I have gotten funds for a new scholarship program, and I am looking for **good students** who would benefit from it. Conveniently, I have **decided** that I am a **good student** who should benefit from this scholarship.*

Interactive Time

- What makes an explanation?
- Case study
- *I have gotten funds for a new scholarship program, and I am looking for **good students** who would benefit from it. Conveniently, I have **decided** that I am a **good student** who should benefit from this scholarship.*
- Should this decision be explained?
- What kind of explanations would you like?

General consensus

- An explanation is the product of interaction between a stakeholder and a system, with the goal of the stakeholder extracting knowledge about the system, and its decision.¹²

¹"Metrics for explainable AI: Challenges and prospects", Hoffman et al.

²"Algorithms and evaluation metrics for improving trust in machine learning : application to visual object recognition", Romain Xu-Darme

³[Interpretable Machine Learning](#) by Christoph Molnar

⁴"Interpretable to Whom? A Role-based Model for Analyzing Interpretable Machine Learning Systems", Tomsett et al.

General consensus

- An explanation is the product of interaction between a stakeholder and a system, with the goal of the stakeholder extracting knowledge about the system, and its decision.¹²
- Explanations are not the goal, but a way of achieving other goals³⁴:

¹"Metrics for explainable AI: Challenges and prospects", Hoffman et al.

²"Algorithms and evaluation metrics for improving trust in machine learning : application to visual object recognition", Romain Xu-Darme

³[Interpretable Machine Learning](#) by Christoph Molnar

⁴"Interpretable to Whom? A Role-based Model for Analyzing Interpretable Machine Learning Systems", Tomsett et al.

General consensus

- An explanation is the product of interaction between a stakeholder and a system, with the goal of the stakeholder extracting knowledge about the system, and its decision.¹²
- Explanations are not the goal, but a way of achieving other goals³⁴:
 - Discover insights about the problem studied
 - Improve/debug models
 - Justify predictions (and models) to stakeholders

¹"Metrics for explainable AI: Challenges and prospects", Hoffman et al.

²"Algorithms and evaluation metrics for improving trust in machine learning : application to visual object recognition", Romain Xu-Darme

³[Interpretable Machine Learning](#) by Christoph Molnar

⁴"Interpretable to Whom? A Role-based Model for Analyzing Interpretable Machine Learning Systems", Tomsett et al.



Stakeholders for explanations

- developers
- operators
- executors
- auditors
- data subjects
- decision subjects

Interesting properties of explanations

Co-12 Property		Description
Content	Correctness	Describes how faithful the explanation is w.r.t. the black box. Key idea: Nothing but the truth
	Completeness	Describes how much of the black box behavior is described in the explanation. Key idea: The whole truth
	Consistency	Describes how deterministic and implementation-invariant the explanation method is. Key idea: Identical inputs should have identical explanations
	Continuity	Describes how continuous and generalizable the explanation function is. Key idea: Similar inputs should have similar explanations
	Contrastivity	Describes how discriminative the explanation is w.r.t. other events or targets. Key idea: Answers “why not?” or “what if?” questions
	Covariate complexity	Describes how complex the (interactions of) features in the explanation are. Key idea: Human-understandable concepts in the explanation
	Compactness	Describes the size of the explanation. Key idea: Less is more
Presentation	Composition	Describes the presentation format and organization of the explanation. Key idea: <i>How</i> something is explained
	Confidence	Describes the presence and accuracy of probability information in the explanation. Key idea: Confidence measure of the explanation or model output
User	Context	Describes how relevant the explanation is to the user and their needs. Key idea: How much does the explanation matter in practice?
	Coherence	Describes how accordant the explanation is with prior knowledge and beliefs. Key idea: Plausibility or reasonableness to users
	Controllability	Describes how interactive or controllable an explanation is for a user. Key idea: Can the user influence the explanation?

Figure 34: Properties of explanations, as identified in¹

¹From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI” by Nauta et al.

Interesting properties of explanations

	Co-12 Property	Description
Content	Correctness	Describes how faithful the explanation is w.r.t. the black box. Key idea: Nothing but the truth
	Completeness	Describes how much of the black box behavior is described in the explanation. Key idea: The whole truth
	Consistency	Describes how deterministic and implementation-invariant the explanation method is. Key idea: Identical inputs should have identical explanations
	Continuity	Describes how continuous and generalizable the explanation function is. Key idea: Similar inputs should have similar explanations
	Contrastivity	Describes how discriminative the explanation is w.r.t. other events or targets. Key idea: Answers “why not?” or “what if?” questions
	Covariate complexity	Describes how complex the (interactions of) features in the explanation are. Key idea: Human-understandable concepts in the explanation
	Compactness	Describes the size of the explanation. Key idea: Less is more
Presentation	Composition	Describes the presentation format and organization of the explanation. Key idea: <i>How</i> something is explained
	Confidence	Describes the presence and accuracy of probability information in the explanation. Key idea: Confidence measure of the explanation or model output
User	Context	Describes how relevant the explanation is to the user and their needs. Key idea: How much does the explanation matter in practice?
	Coherence	Describes how accordant the explanation is with prior knowledge and beliefs. Key idea: Plausibility or reasonableness to users
	Controllability	Describes how interactive or controllable an explanation is for a user. Key idea: Can the user influence the explanation?

Figure 34: Properties of explanations, as identified in¹



Figure 35: Poll: Which properties do you view as most important? Link: <https://whale5.noiraudes.net/polls/337317df-8895-402b-9ddb-0dea7564a1e9>

¹From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI" by Nauta et al.

Classical vs Formal XAI

Classical XAI

- Problem: system opacity¹

¹Further reading: "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead" by Cynthia Rudin

²<https://christophm.github.io/interpretable-ml-book/overview.html>

³<https://dept.utc2.edu.vn/bomoncntt/doi-tac/decision-trees-explained-with-a-practical-example-57.html>

Classical XAI

- Problem: system opacity¹

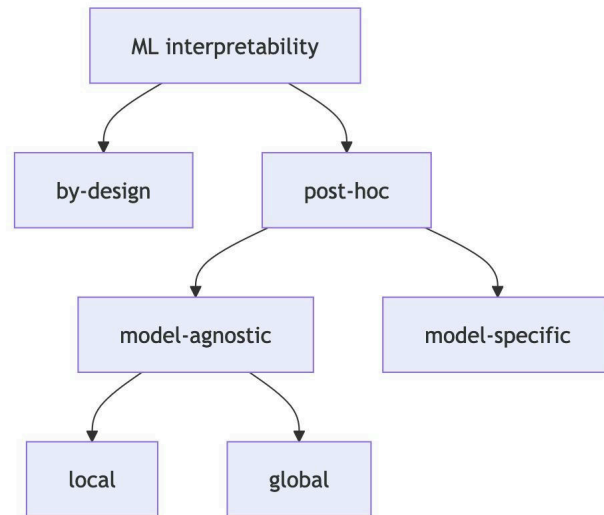


Figure 36: Taxonomy of interpretability methods²

¹Further reading: "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead" by Cynthia Rudin

²<https://christophm.github.io/interpretable-ml-book/overview.html>

³<https://dept.utc2.edu.vn/bomoncntt/doi-tac/decision-trees-explained-with-a-practical-example-57.html>

Classical XAI

- Problem: system opacity¹

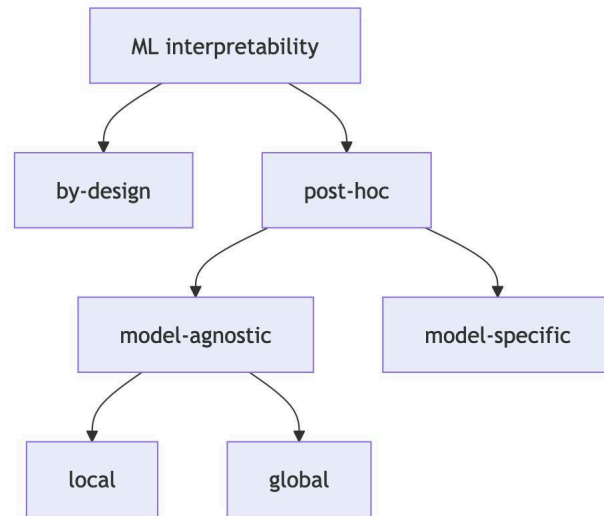


Figure 36: Taxonomy of interpretability methods²

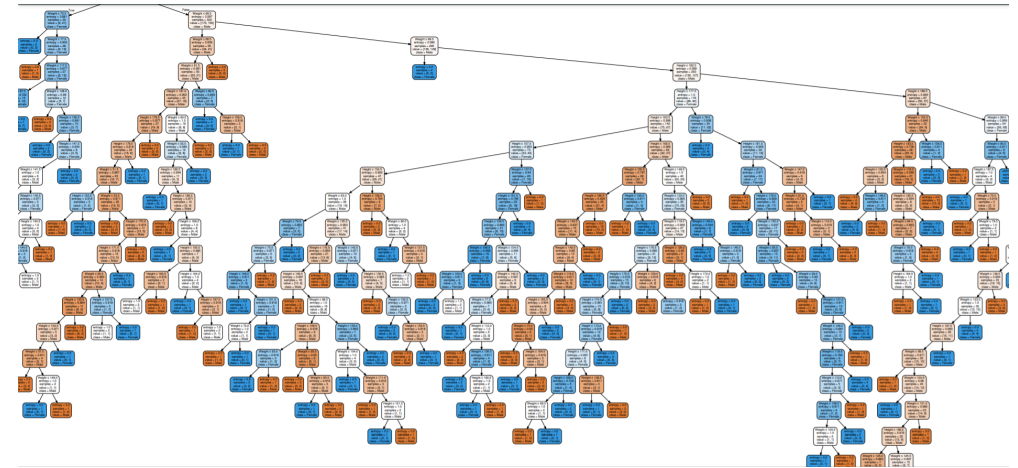


Figure 37: Deep decision tree.³

¹Further reading: "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead" by Cynthia Rudin

²<https://christophm.github.io/interpretable-ml-book/overview.html>

³<https://dept.utc2.edu.vn/bomoncntt/doi-tac/decision-trees-explained-with-a-practical-example-57.html>



What classical XAI had promised to us:

Classical XAI aims to¹:

- help us improve our model
- justify answers to various stakeholders
- discover insights about the problem

¹<https://christophm.github.io/interpretable-ml-book/goals.html>

AI risk scenarios (again)

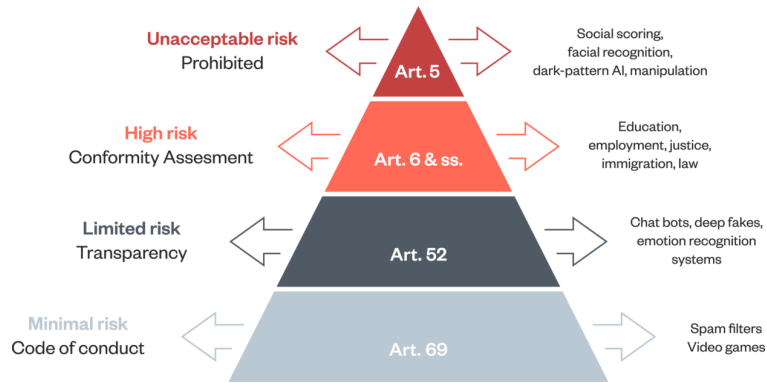


Figure 38: Reminder of AI risk hierarchy.

It is intolerable to:

- have systems that cannot be understood
- have mistakes in explanations



Limitations of traditional XAI

The explanations should be:

- compact (no redundant features)
- faithful to the model



Limitations of traditional XAI

The explanations should be:

- compact (no redundant features)
- faithful to the model

Limitations of XAI

- Not respecting these requirements



Limitations of traditional XAI

The explanations should be:

- compact (no redundant features)
- faithful to the model

Limitations of XAI

- Not respecting these requirements
- many others (but we will focus on the mentioned ones)

Limitations of traditional XAI

The explanations should be:

- compact (no redundant features)
- faithful to the model

Limitations of XAI

- Not respecting these requirements
- many others (but we will focus on the mentioned ones)

Case	Instance	Relevant	Irrelevant	Shapley values	Justification
14	$((0, 0, 1, 1), 0)$	1, 2, 4	3	$Sv(1) = -0.13$ $Sv(2) = 0.33$ $Sv(3) = 0.08$ $Sv(4) = 0.00$	$Irrelevant(3) \wedge Sv(3) \neq 0 \wedge$ $Relevant(4) \wedge Sv(4) = 0$
15	$((1, 1, 1, 1), 0)$	1, 2, 3	4	$Sv(1) = -0.12$ $Sv(2) = -0.12$ $Sv(3) = -0.12$ $Sv(4) = 0.17$	$Irrelevant(4) \wedge$ $\forall (j \in \{1, 2, 3\}). Sv(j) < Sv(4) $

TABLE IV: Examples of issues with Shapley for explainability for boolean classifiers of Figure 3

Figure 39: Shapley Values can return non-zero scores for irrelevant features, and zero scores for relevant ones.¹

¹"Disproving XAI Myths with Formal Methods – Initial Results" by Joao Marques-Silva



Advantages of Formal XAI

Formal XAI (FXAI) is able to:

- compute compact explanations
- compute model-faithful explanations
- validate the robustness of a model wrt. requirements

Formal XAI | Building Safe AI Systems

Toy problem

- Let's formalize a system that detects mammals!



Figure 40: Platypodes are weird mammals.

Toy problem

- Let's formalize a system that detects mammals!
- (disclaimer: not a biologist)



Figure 42: Platypodes are weird mammals.

Classification problem

Consider a classification instance,

- X : Feature space
- C : classification space
- x, c : an instance of a vector of values in X and corresponding class in C
- f : classification function, eg.
 $f(x) = c$

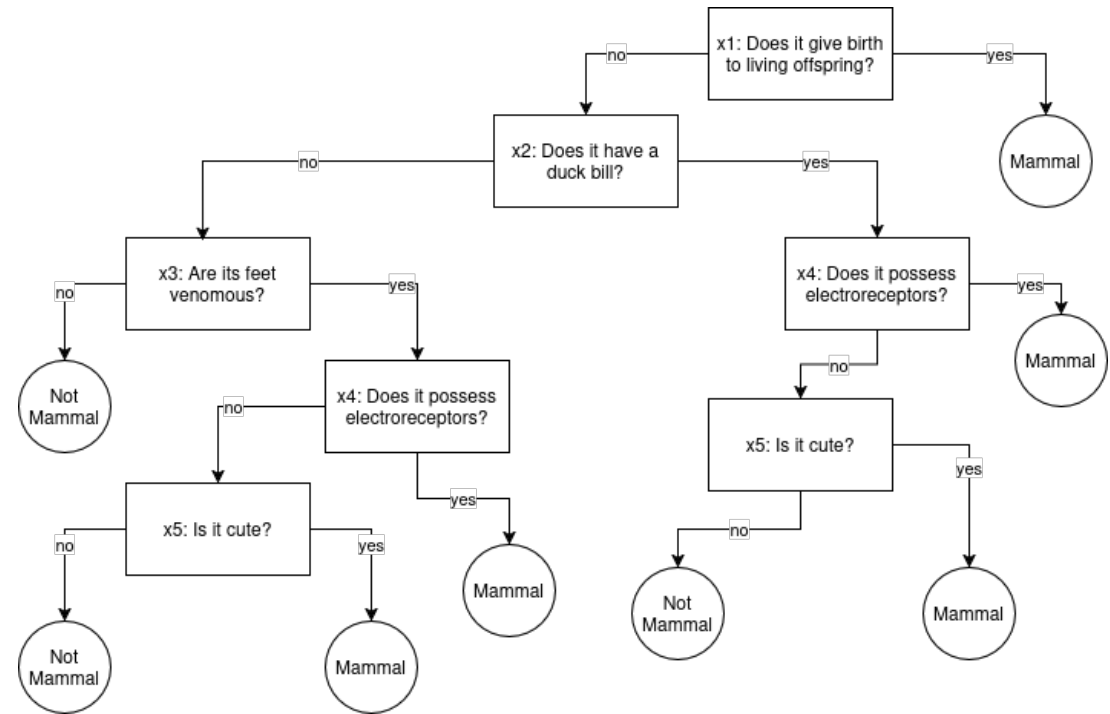


Figure 43: Adapted decision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Classification problem

- $X = \{X_i \in \{0, 1\}, \forall i \in \{1, 2, 3, 4, 5\}\}$
- $C = \{0, 1\}$
- $x = \{1, 0, 1, 0, 1\}$
- $c = 1$

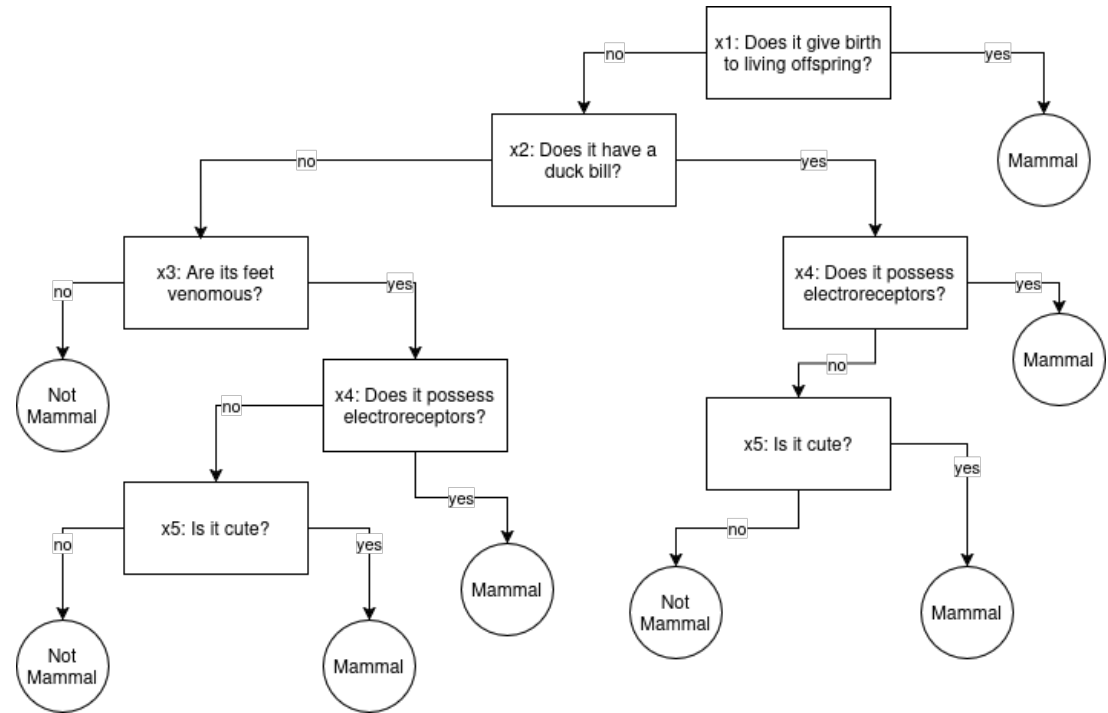


Figure 47: Adapted decision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Classification problem

- $X = \{X_i \in \{0, 1\}, \forall i \in \{1, 2, 3, 4, 5\}\}$
- $C = \{0, 1\}$
- $x = \{1, 0, 1, 0, 1\}$
- $c = 1$

Why did the model classify $f(x) = 1$?

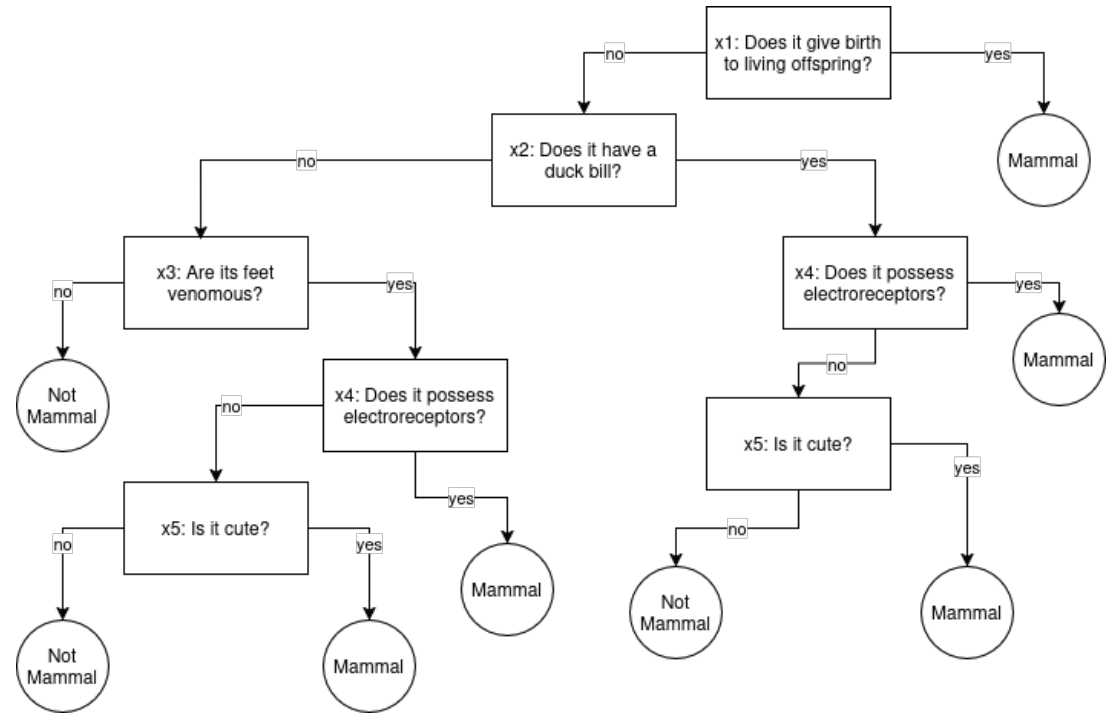


Figure 47: Adapted decision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Abductive explanations

- $x = \{1, 0, 1, 0, 1\}$
- Why $f(x) = 1$ (mammal)?

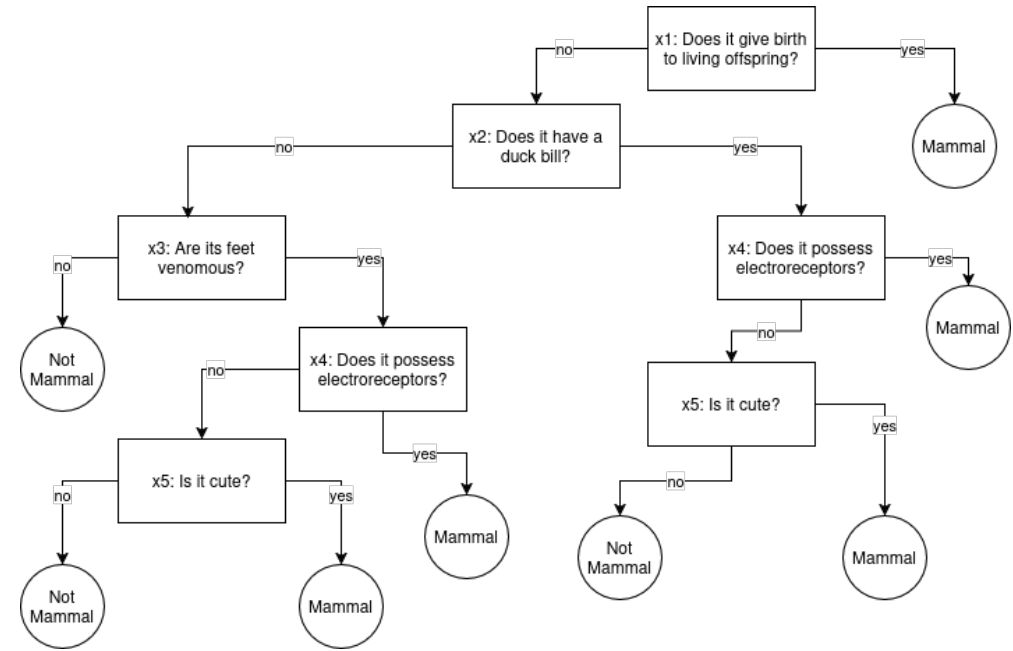


Figure 48: Adapted ecision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Abductive explanations

- $x = \{1, 0, 1, 0, 1\}$
- Why $f(x) = 1$ (mammal)?
- An abductive explanation (AXP) is a subset-minimal of features $i \in$ AXP, such that:

$$\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in \text{AXP}} (\bar{x}_i = x_i) \right) \rightarrow f(\bar{x}) = c \right]$$

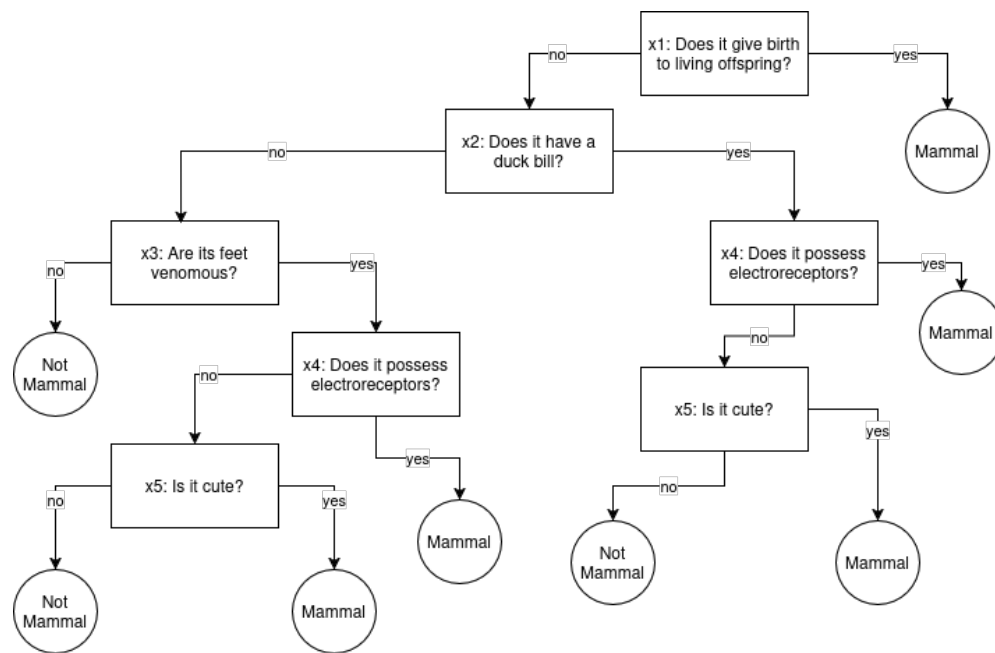


Figure 51: Adapted ecision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Abductive explanations

- $x = \{1, 0, 1, 0, 1\}$
- Why $f(x) = 1$ (mammal)?
- An abductive explanation (AXP) is a subset-minimal of features $i \in$ AXP, such that:

$$\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in \text{AXP}} (\bar{x}_i = x_i) \right) \rightarrow f(\bar{x}) = c \right]$$

- A weak abductive explanation (wAXP) is an AXP that is not subset minimal:

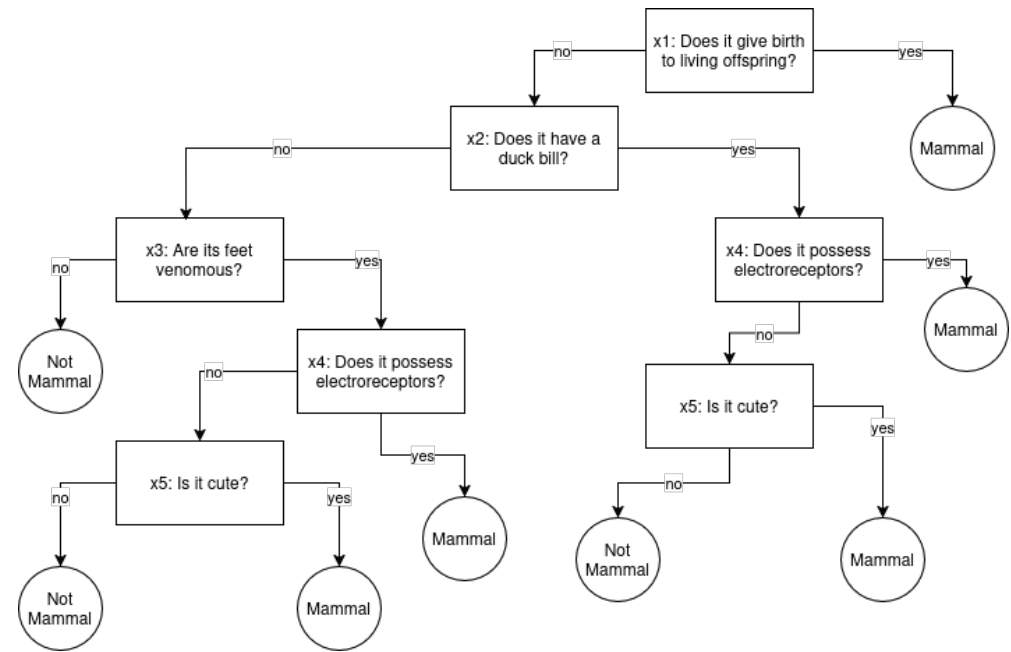


Figure 51: Adapted ecision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Abductive explanations

- $x = \{1, 0, 1, 0, 1\}$
- Why $f(x) = 1$ (mammal)?
- An abductive explanation (AXP) is a subset-minimal of features $i \in$ AXP, such that:

$$\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in \text{AXP}} (\bar{x}_i = x_i) \right) \rightarrow f(\bar{x}) = c \right]$$

- A weak abductive explanation (wAXP) is an AXP that is not subset minimal:
- Trivial wAXP: $\{x_1, x_2, x_3, x_4, x_5\}$

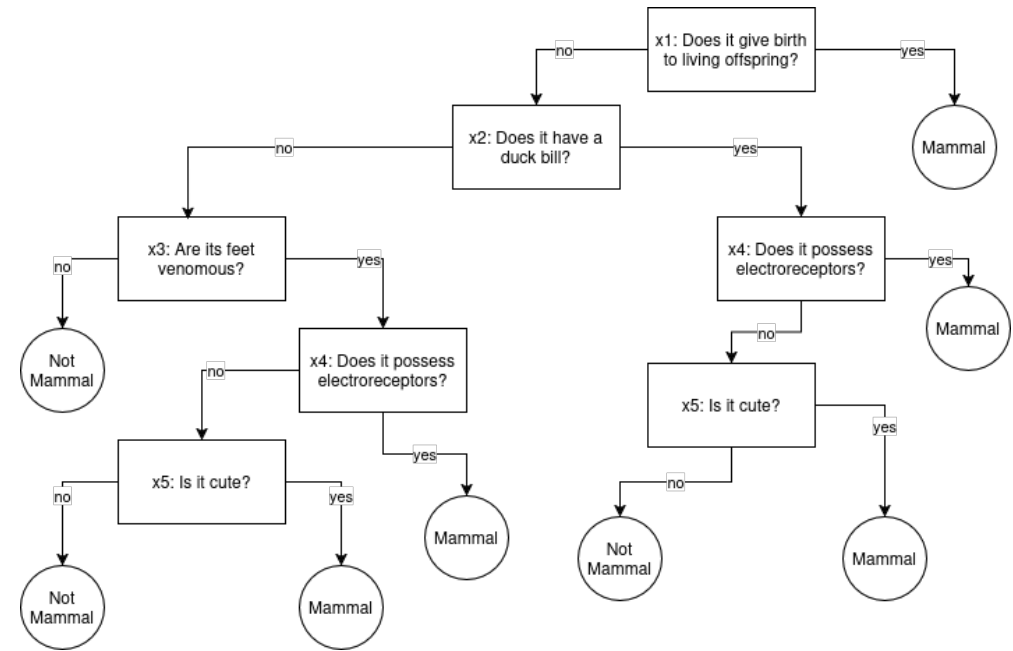


Figure 51: Adapted ecision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Abductive explanations

- $x = \{1, 0, 1, 0, 1\}$
- Why $f(x) = 1$?
- An abductive explanation (AXP) is a subset-minimal of features $i \in \text{AXP}$, such that:

$$\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in \text{AXP}} (\bar{x}_i = x_i) \right) \rightarrow f(\bar{x}) = c \right]$$

- Subset minimal AXP: $\{x_1\}$

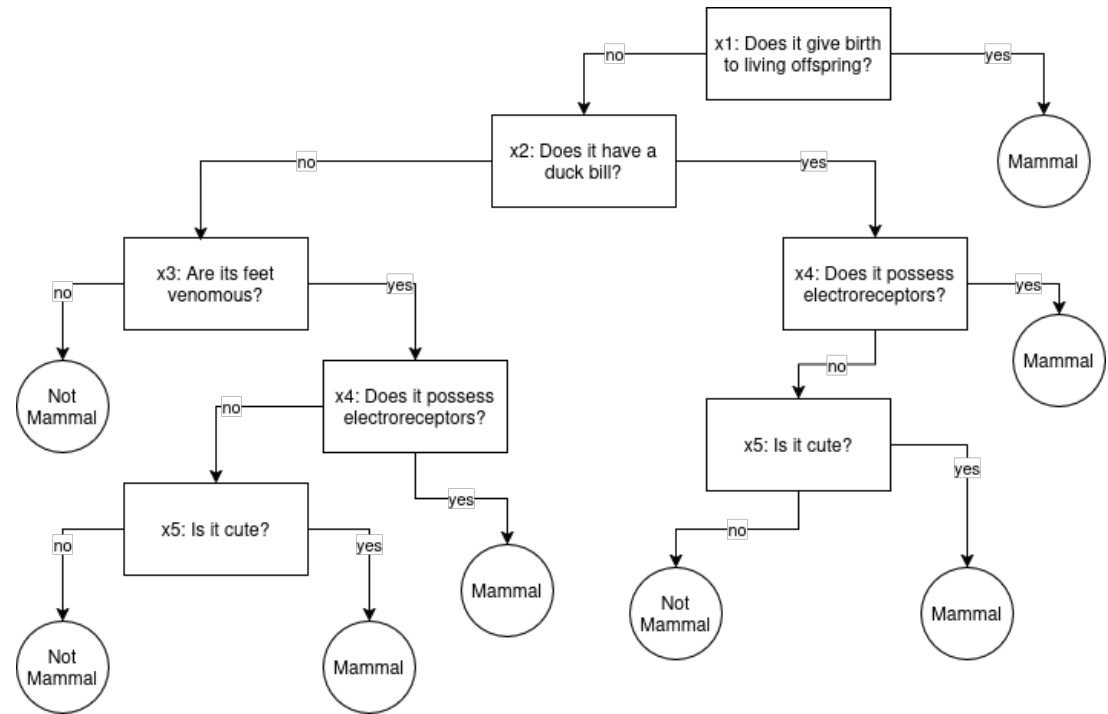


Figure 52: Adapted decision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva

Algorithm to a formal explanation

Input: function f , datapoint $x = \{x_1, x_2, \dots, x_n\}$

Output: Subset minimal abductive explanations

1. $\text{wAXP} \leftarrow \{1, 2, \dots, n\}$ (weak Abductive Explanation)
2. $c = f(x)$
3. for i in $\{1, 2, \dots, n\}$ do:
 4. $\text{wAXP}' \leftarrow \text{wAXP} - \{i\}$
 5. Check SAT of $\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in \text{wAXP}'} (\bar{x}_i = x_i) \right) \rightarrow f(\bar{x}) = c \right]$
 6. if SAT then $\text{wAXP} \leftarrow \text{wAXP}'$
7. $\text{AXP} \leftarrow \text{wAXP}$
8. Return AXP

Algorithm to a formal explanation

1. $wAXP \leftarrow \{1, 2, \dots, n\}$ (weak Abductive Explanation)
2. $c = f(x)$
3. for i in $\{1, 2, \dots, n\}$ do:
 4. $wAXP' \leftarrow wAXP - \{i\}$
 5. Check SAT of $\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in wAXP'} (\bar{x}_i = x_i) \right) \rightarrow f(\bar{x}) = c \right]$
 6. if SAT then $wAXP \leftarrow wAXP'$
7. $AXP \leftarrow wAXP$
8. Return AXP

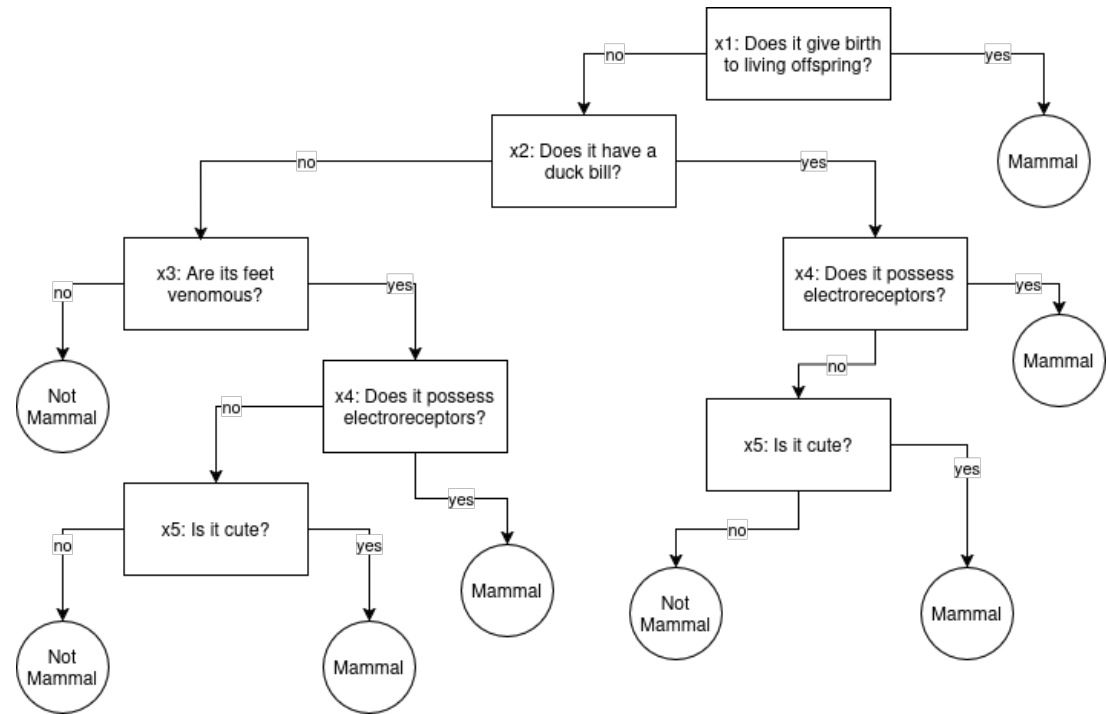


Figure 53: Adapted decision tree example¹

¹"Disproving XAI Myths with Formal Methods – Initial Results", Joao Marques-Silva



Contrastive explanation

- Similarly, we can ask “Why not “not mammal”?”

Contrastive explanation

- Similarly, we can ask “Why not “not mammal”?”
- A contrastive explanation (CXP) is a subset-minimal of features s.t.:

$$\exists \bar{x} \in X. \left[\left(\bigwedge_{i \in (X - \text{CXP})} (\bar{x}_i = x_i) \right) \wedge f(\bar{x}) \neq c \right]$$

Contrastive explanation

- Similarly, we can ask “Why not “not mammal”?”
- A contrastive explanation (CXP) is a subset-minimal of features s.t.:

$$\exists \bar{x} \in X. \left[\left(\bigwedge_{i \in (X - \text{CXP})} (\bar{x}_i = x_i) \right) \wedge f(\bar{x}) \neq c \right]$$

- Similarly, a weak CXP (wCXP) is a CXP that is not subset minimal.



Are CXPs and AXP unique? (No)

- we should've seen in the algorithm that we got a different AXP than before
- it depends on the order in which you traverse the features
- AXPs: $\{1\}$, $\{3, 5\}$

Are CXPs and AXP unique? (No)

- we should've seen in the algorithm that we got a different AXP than before
- it depends on the order in which you traverse the features
- AXPs: $\{1\}$, $\{3, 5\}$

Duality between CXPs and AXPs: Minimal hitting sets¹

- Interesting property: If you were to find all AXPs, it would allow you to compute a CXP aswell!
- CXPs: $\{1, 3\}$, $\{1, 5\}$

¹"Delivering Trustworthy AI through Formal XAI" by Joao Marques-Silva and Alexey Ignatiev



Are CXPs and AXP unique? (No)

- we should've seen in the algorithm that we got a different AXP than before
- it depends on the order in which you traverse the features
- AXP: {1}, {3, 5}

Duality between CXPs and AXP: Minimal hitting sets¹

- Interesting property: If you were to find all AXP, it would allow you to compute a CXP aswell!
- CXP: {1, 3}, {1, 5}

Open Problem: Multiple explanations for the same decision

- Given multiple formal explanations, which one to choose?

¹"Delivering Trustworthy AI through Formal XAI" by Joao Marques-Silva and Alexey Ignatiev

New paradigm: Local formal explanations

Reminder:

Local robustness (Katz et al. 2017)

Let a classifier $f : \mathcal{X} \mapsto \mathcal{Y}$. Given $x \in \mathcal{X}$ and $\varepsilon \in \mathbb{R} \ll 1$ the problem of *local robustness* is to prove that $\forall x^{\{\cdot\}}. \|x - x^{\{\cdot\}}\|_p < \varepsilon \rightarrow f(x) = f(x^{\{\cdot\}})$

Algorithm to a formal explanation (VeriX)

Input: function f , datapoint $x = \{x_1, x_2, \dots, x_n\}$, *perturbation level* ε

Output: Subset minimal local abductive explanations

1. $\text{wAXP} \leftarrow \{1, 2, \dots, n\}$ (weak Abductive Explanation)
2. $c = f(x)$
3. for i in $\{1, 2, \dots, n\}$ do:
 4. $\text{wAXP}' \leftarrow \text{wAXP} - \{i\}$
 5. Check SAT of:
$$\forall \bar{x} \in X. \left[\left(\bigwedge_{i \in \text{wAXP}'} (\bar{x}_i = x_i) \wedge \left(\bigwedge_{j \notin \text{wAXP}'} (\|\bar{x}_j - x_j\|_p \leq \varepsilon) \right) \rightarrow f(\bar{x}) = c \right) \right]$$
 6. if SAT then $\text{wAXP} \leftarrow \text{wAXP}'$
7. $\text{AXP} \leftarrow \text{wAXP}$
8. Return AXP

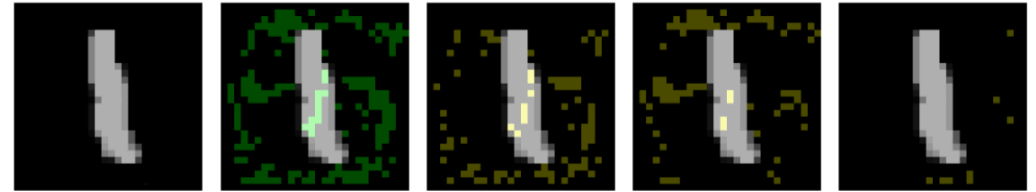


Applicable domains

- We've seen FXAI for decision trees and tabular data
- They can also work on more complex scenarios: neural networks, computer vision
- While some problems are quick to solve (instantaneous), complex problems take a long time:

Example explanations

- Explanations computed for the decision of a convolutional neural network



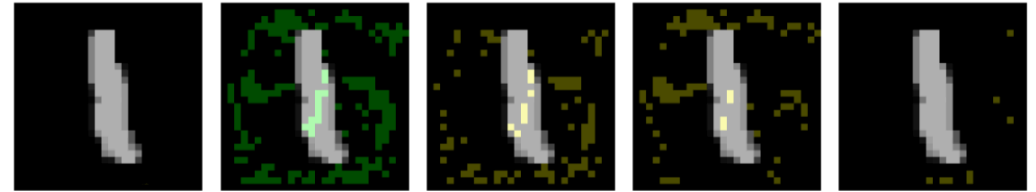
(b) Handwritten digit “1”, not “8”, not “5”, not “2”

Figure 57: Explanation from VeriX¹. Image resolution 1x28x28

¹“VERIX: Towards Verified Explainability of Deep Neural Networks”, by M. Wu, H. Wu and C. Barrett

Example explanations

- Explanations computed for the decision of a convolutional neural network
- Question: What is your take on this explanation?



(b) Handwritten digit “1”, not “8”, not “5”, not “2”

Figure 60: Explanation from VeriX¹. Image resolution 1x28x28

¹“VERIX: Towards Verified Explainability of Deep Neural Networks”, by M. Wu, H. Wu and C. Barrett

Example explanations

- Explanations computed for the decision of a convolutional neural network
- Question: What is your take on this explanation?



Figure 61: Example of original image (left) and of the explanation (right) for the current prediction (give way sign). Image resolution: 3x32x32

Example explanations

- Explanations can be computed for target datapoints, on larger networks
- Question: What is your take on these explanations?



Figure 62: Example of original image (top-left) and of three different explanations for the same decision (“Watch”) on Resnet18 (same weights). Image resolution: 3x64x64

Practicality



Dataset Method	MNIST	GTSRB	GTSRB, simplified	TinyImageNet, simplified
Satisfiability Modulo Theories	820s	2400s	32s	–
Abstract Interpretation	220s	1600s	8s	4050s
Interval Bound Propagation	40s	140s	1s	–



Limitations of formal XAI

From why to use it, now we go back to what are limitations for it:

- tabular data (or smaller in size) vs bigger data (eg visual)
- Scalability, memory usage, viability of explanations
- NLP tasks & transformers, approximating various activation functions
 - (eg Softmax is an unsolved problem that prevents us from verifying transformers)
 - (eg the discreteness of words is an unsolved problem that prevents us from verifying NLP models)
- + your opinions of the formal explanations you have seen



Open ended directions

- Disentangle the multiplicity of explanations
- Scalability
- Apply FXAI to multicriterion decision making (analysis of conflicting criteria)



Key takeaways from lecture

- Explainability is important in Artificial Intelligence
- Explanations take into account both the system explained and the user seeking the explanation
- No XAI method is flawless
- Formal methods address limitations in terms of model faithfulness and explanation compactness

Feedback

Thank you for listening to me!