# MAINLE: a Multi-Agent, Interactive, Natural Language Local Explainer of Classification Tasks

**Paulo Bruno Serafim**, Rômulo Férrer Filho, Stenio Freitas, Gizem Gezici, Fosca Giannotti, Franco Raimondi, and Alexandre Santos

17 September, 2025

GRAN SASSO SCIENCE INSTITUTE
SCHOOL OF ADVANCED STUDIES
Scuola Universitaria Superiore
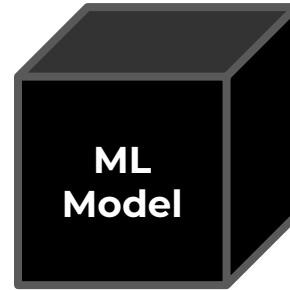
SCUOLA NORMALE SUPERIORE

UFC

# Motivation

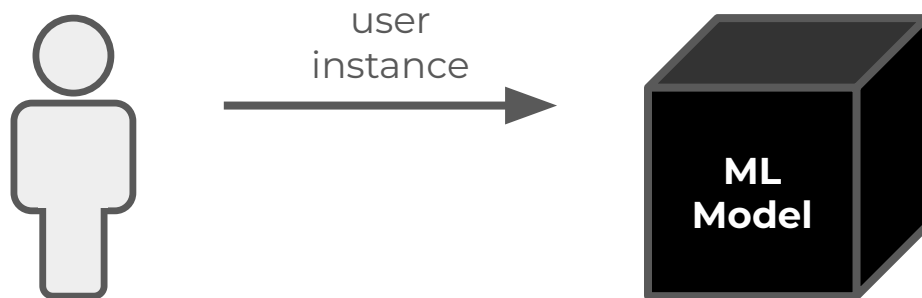- How to **explain** a Machine Learning **classification** to a **non-expert** user?

# Machine Learning is Everywhere

- Computers, phones, tvs, home appliances...
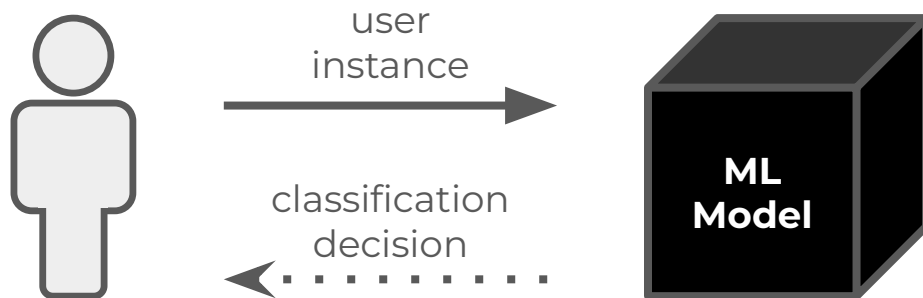- Customer support, recommendation system, classification, predictions...
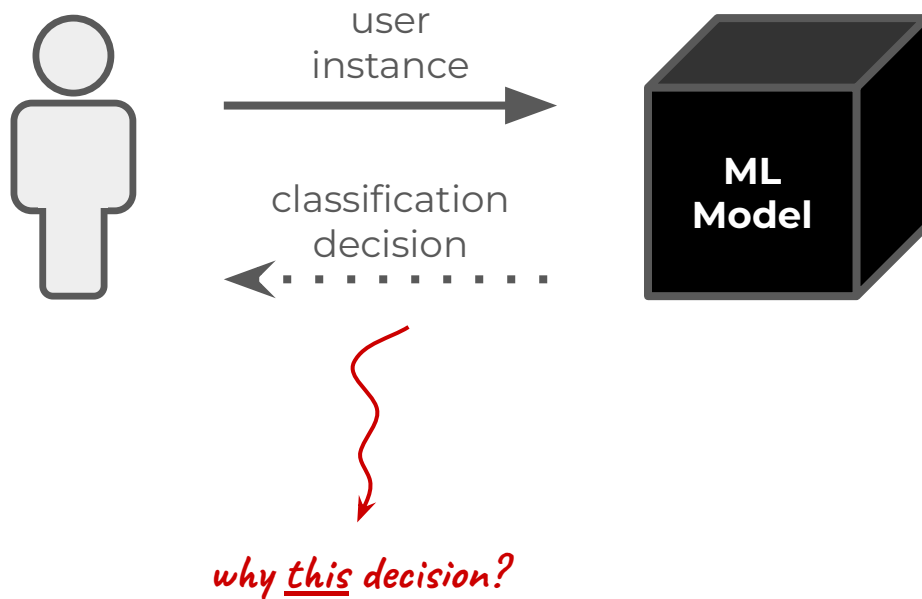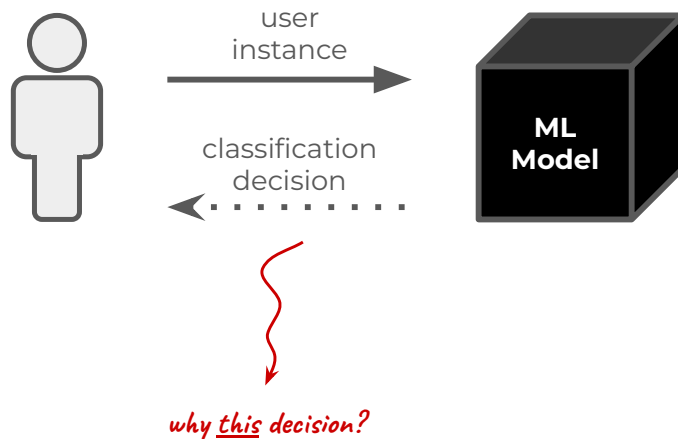
# Example Scenario

# Example Scenario



user
instance

**ML Model**

# Example Scenario

# Example Scenario

# Motivation

- How to explain a Machine Learning classification to a non-expert user?



user
instance

classification
decision

ML
Model

why *this* decision?

# Motivation

- How to explain a Machine Learning classification to a non-expert user?
    - **Explainable AI (XAI) can help**

# Problems

It is hard to **use** XAI methods

# Problems

It is hard to **use** XAI methods

- Command Line
- Git
- Programming
- Python
- PyTorch / Tensor Flow
- Machine Learning
- …

# Problems

It is hard to **present** Explanations

# Problems

It is hard to **present** Explanations
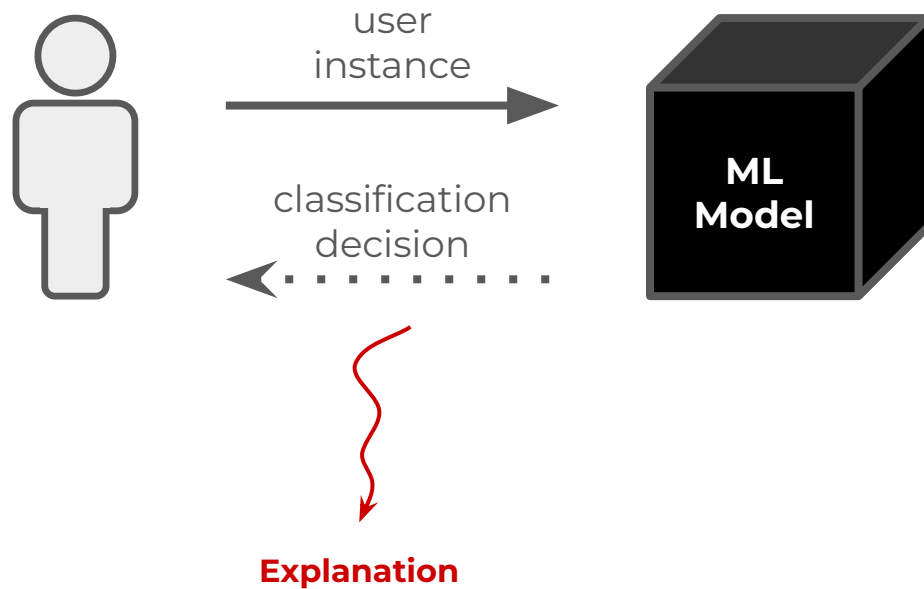
- List of numbers
- Tables
- Spreadsheets
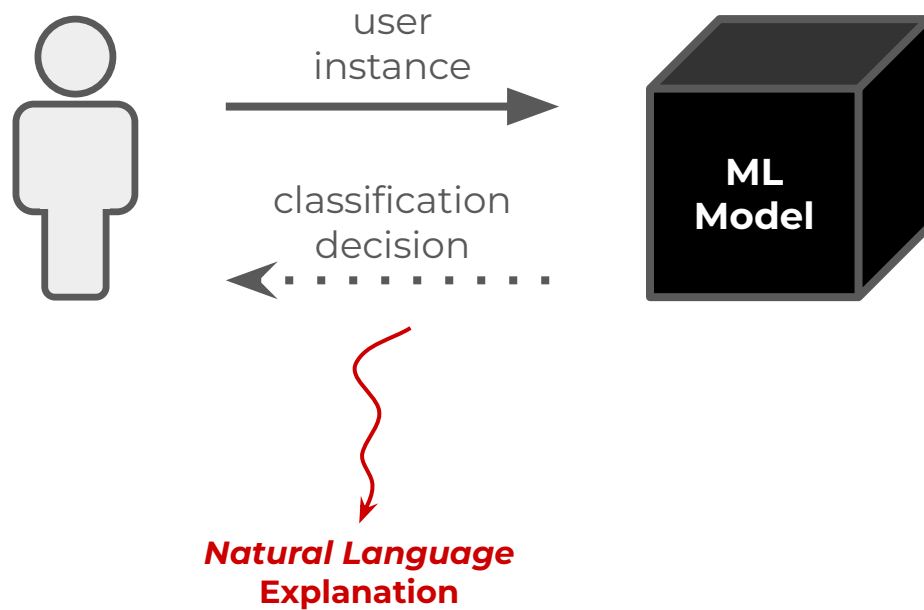- Scrollers
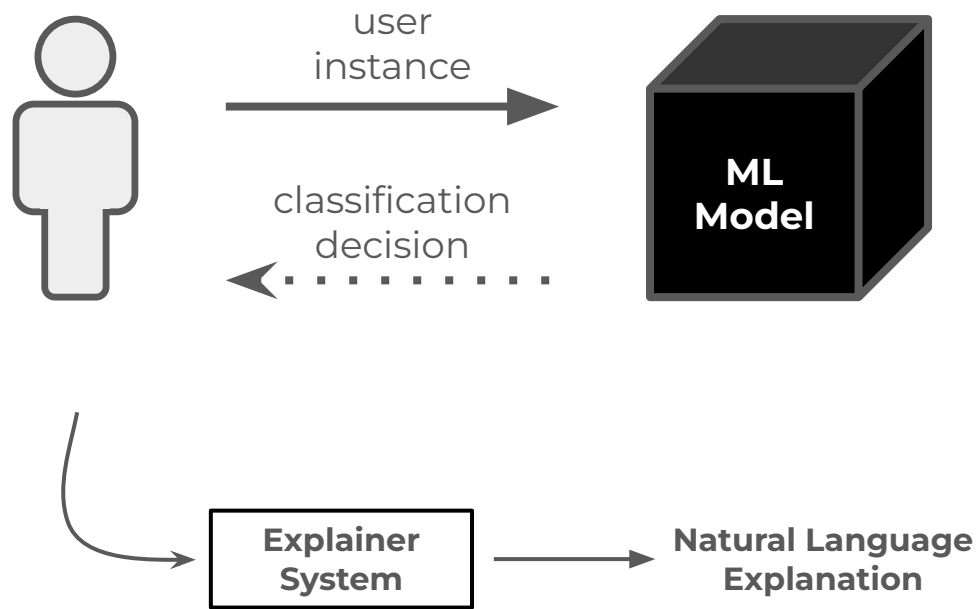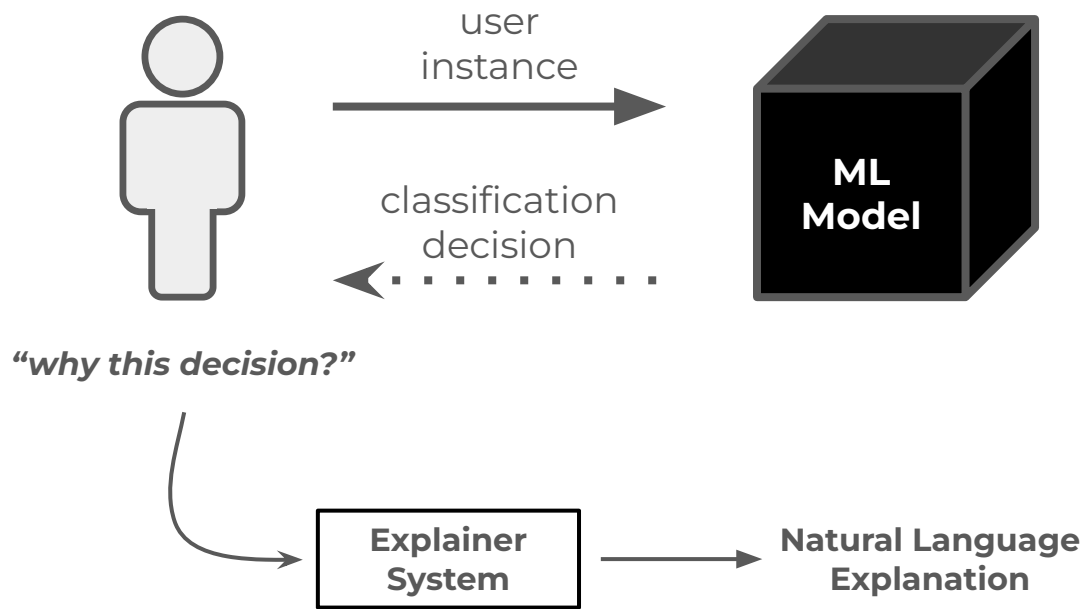- …

# Problems
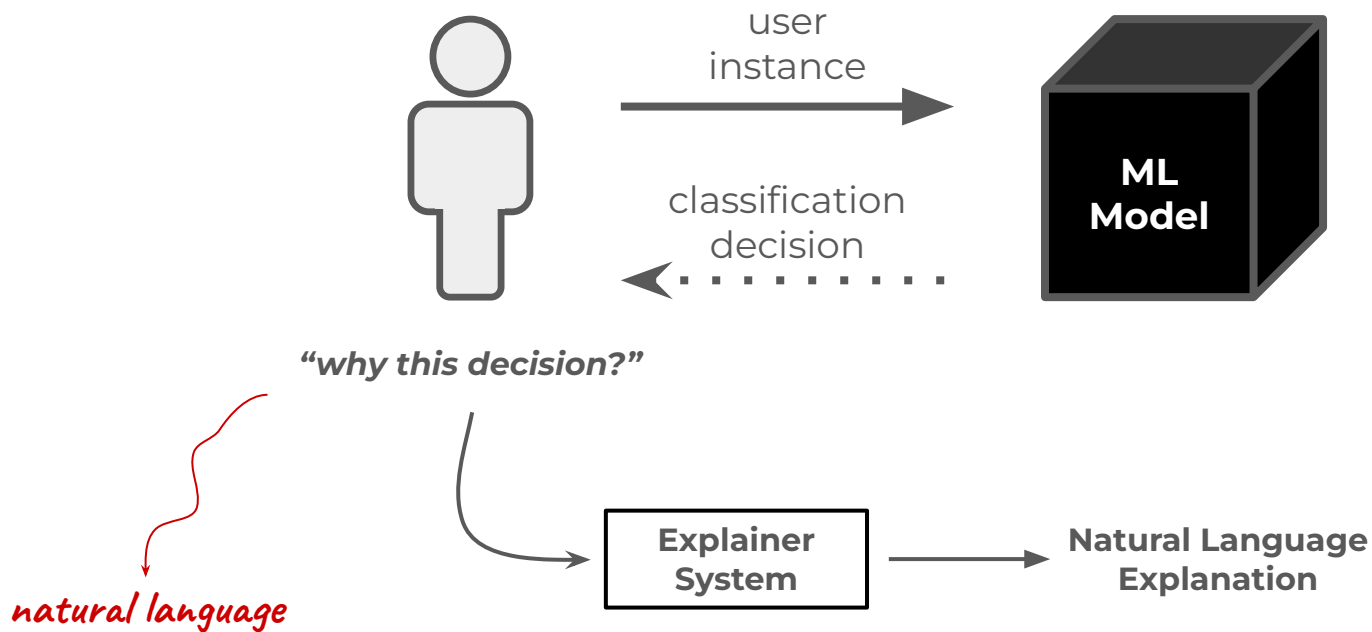
Our objective is to remove these barriers

*How?*

# Example Scenario

# Example Scenario

# Example Scenario

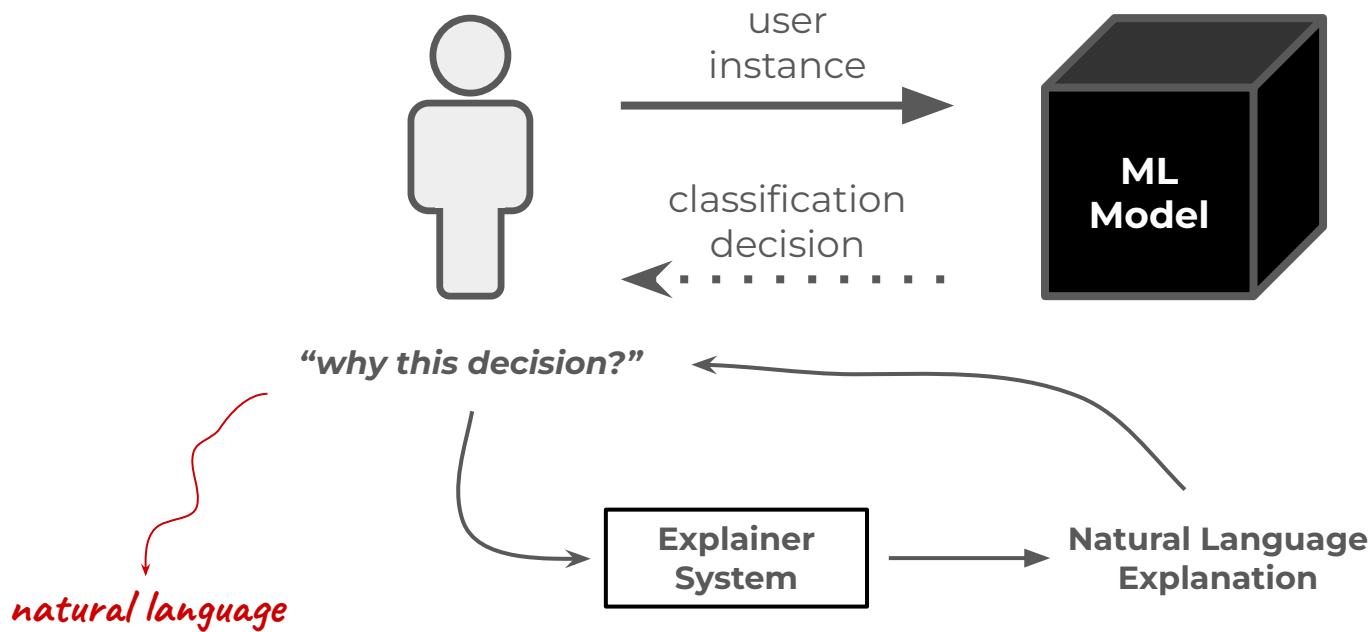# Example Scenario
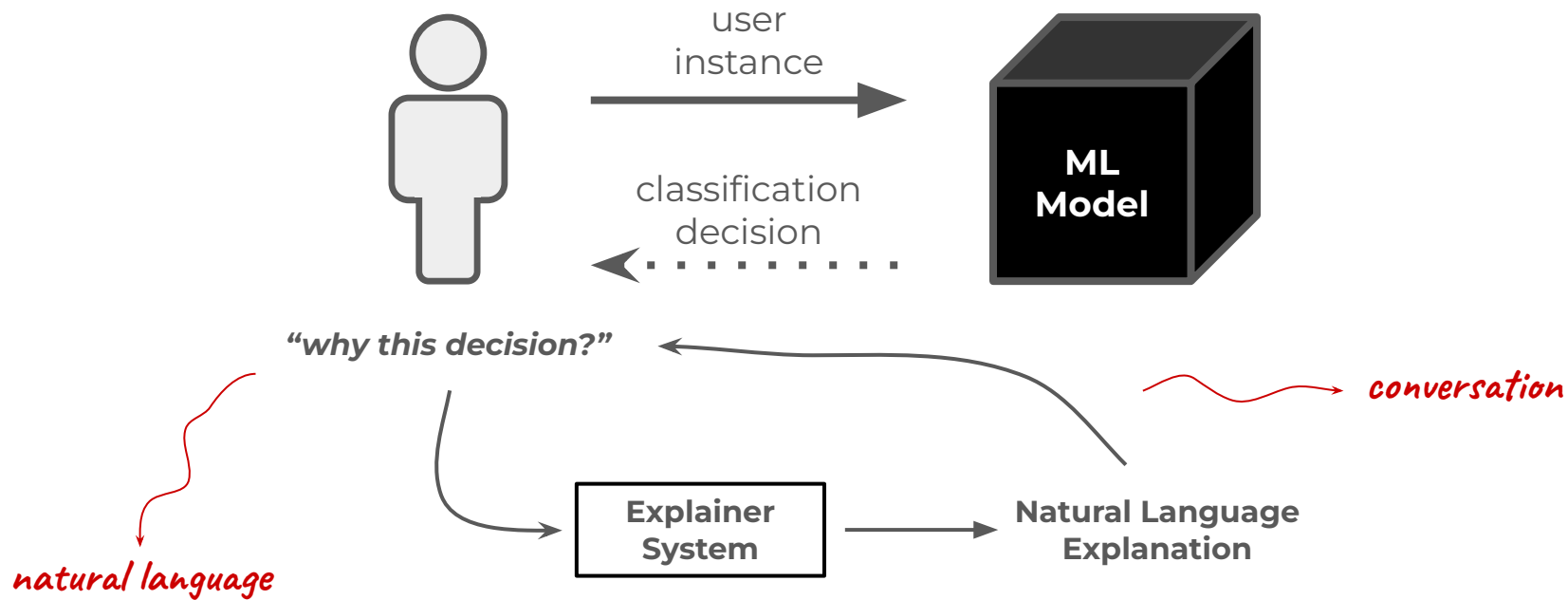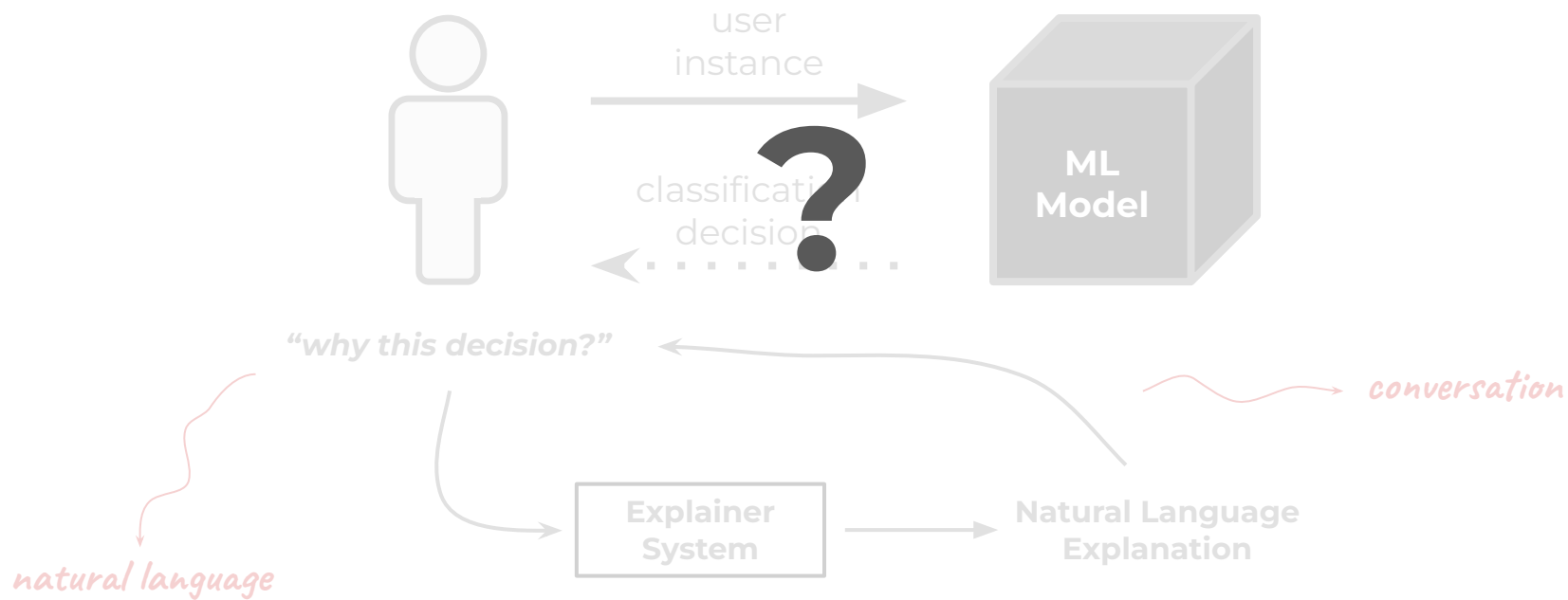
# Example Scenario

# Example Scenario

# Example Scenario

# Example Scenario



user instance

classification decision

**ML Model**

?

"why this decision?"

natural language

**Explainer System**

**Natural Language Explanation**

conversation

# Related Work

## Generación Automática de Explicaciones en Lenguaje Natural para Árboles de Decisión de Clasificación

B. López-Trigo, Jose M. Alonso, A. Bugarín
Centro Singular de Investigación en Tecnoloxías da Información (CiTIUS),
Universidade de Santiago de Compostela,
Campus Vida, E-15782, Santiago de Compostela, Spain
Email: bruno.lopez.trigo@rai.usc.es, {josemaria.alonso.moral, alberto.bugarin.diz}@usc.es

*Resumen—*En este trabajo describimos un modelo de explicaciones en lenguaje natural para árboles de decisión para clasificación. Las explicaciones incluyen aspectos globales del clasificador y aspectos locales de la clasificación de una instancia concreta. La propuesta está implementada en el servicio Web de código abierto ExpliClas [1], que en su versión actual opera sobre árboles construidos con Weka y conjuntos de datos con atributos numéricos. Ilustramos la viabilidad de la propuesta con dos casos de ejemplo, donde mostramos paso a paso cómo el modelo explica los respectivos árboles de clasificación.

*Index Terms—*Explicabilidad, Soft Computing, Árboles de decisión para Clasificación, Generación de Lenguaje Natural

## I. INTRODUCCIÓN

La generalización del uso de las nuevas tecnologías ha hecho que hoy trabajemos y vivamos rodeados de sistemas inteligentes [2]. Términos como ciudad inteligente, fábrica, casa, coche o teléfono inteligentes, son cada vez más populares. En realidad, existen multitud de dispositivos dotados de cierta inteligencia que nos asisten en el día a día, muchas veces sin que seamos totalmente conscientes de ello. Mención especial merece el teléfono móvil, que nos ofrece multitud de aplicaciones casi para cualquier cosa que podamos imaginar y va con nosotros a todas partes. Se puede afirmar que, si bien en el pasado vivimos una revolución industrial, ahora estamos viviendo una revolución social impulsada por la Inteligencia Artificial (IA).

Cuando un sistema inteligente toma decisiones que nos afectan (ej. filtrar llamadas, diagnóstico médico, concesión de un préstamo, etc.), surgen multitud de preguntas que deberíamos hacernos [3]: ¿quién es el responsable de las consecuencias colaterales que pudieran derivarse de las decisiones tomadas? ¿cuáles son las consecuencias éticas? ¿puede haber consecuencias legales?

Desde el punto de vista legal, el Parlamento Europeo aprobó una nueva Regulación General de Protección de Datos [4] que entró en vigor el 25 de mayo de 2018. La nueva regulación enfatiza el derecho de los ciudadanos a pedir explicaciones, independientemente de que las decisiones que les afectan sean tomadas por una persona o un programa informático. Esto significa que los ciudadanos pueden pedir a las empresas que den explicaciones asociadas a las decisiones tomadas por los sistemas inteligentes que utilizan.

Desde un punto de vista técnico: ¿puede explicarnos la aplicación que tomó una decisión por qué tomó esa decisión y no otra? Para eso, hay básicamente dos opciones [5]: (1) el sistema inteligente está construido siguiendo un modelo interpretable (también llamado de caja blanca) que un operario experto puede analizar y entender a fin de elaborar una explicación; o (2) el sistema está construido siguiendo un modelo explicable que genera explicaciones por sí mismo. La DARPA planteó en 2016 las siguientes cuestiones técnicas [5]: ¿puede una máquina inteligente aprender de forma autónoma a explicar su comportamiento? ¿está preparada la generación actual de sistemas inteligentes para dar explicaciones de forma clara, sin ambigüedades, tanto a públicos especializados como no especializados? Y lanzó el reto de crear una nueva generación de sistemas inteligentes explicables entre 2017 y 2021. El reto fue lanzado inicialmente a universidades y centros de investigación americanos, con énfasis en la creación de equipos multidisciplinares que abordasen no sólo aspectos algorítmicos sino también de implementación y evaluación con personas. Los equipos seleccionados empezaron a trabajar en mayo de 2017 pero a día de hoy sólo hemos encontrado resultados muy preliminares (ej. [6], [7]).

Hasta donde nosotros sabemos, en la práctica, la responsabilidad de generar explicaciones recae directamente en el operario asociado al sistema inteligente, si está disponible para ello [8]. Aunque hay sistemas basados en conocimiento que son interpretables, en los últimos años son cada vez más populares las técnicas de IA para aprendizaje automático y minería de datos, supervisadas y no supervisadas (es decir, con o sin intervención humana). Estos sistemas se están demostrando ciertamente útiles y versátiles, pero la mayoría no suelen tener ninguna capacidad explicativa ni tampoco pueden ser interpretados fácilmente por personas (en cuyo caso se dice que son sistemas de caja negra).

Por tanto, el nuevo marco legal demanda que los expertos en IA desarrollen nuevos algoritmos que proporcionen explicaciones de forma automática.

En este trabajo, presentamos un modelo para la interpre-

**López-Trigo, Alonso and Bugarín (2018)**
Generación Automática de Explicaciones en Lenguaje Natural para Árboles de Decisión de Clasificación

# Related Work



López-Trigo, Alonso and Bugarín (2018)
Generación Automática de Explicaciones en Lenguaje Natural para Árboles de Decisión de Clasificación



Ziems, Liu, Flanagan and Jian (2023)
Explaining Tree Model Decisions in Natural Language for Network Intrusion Detection



Serafim, Crescenzi, Gezici, Cappuccio, Rinzivillo, and Giannotti (2024)
Exploring Large Language Models Capabilities to Explain Decision Trees
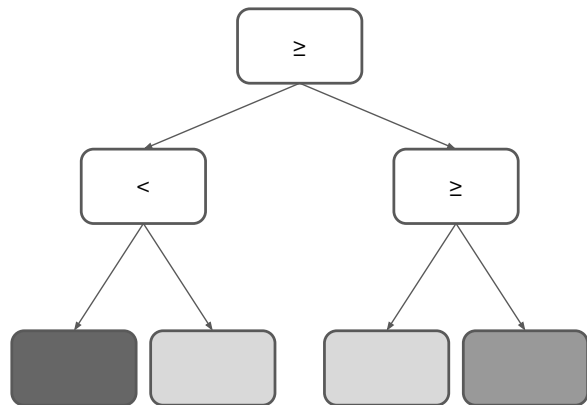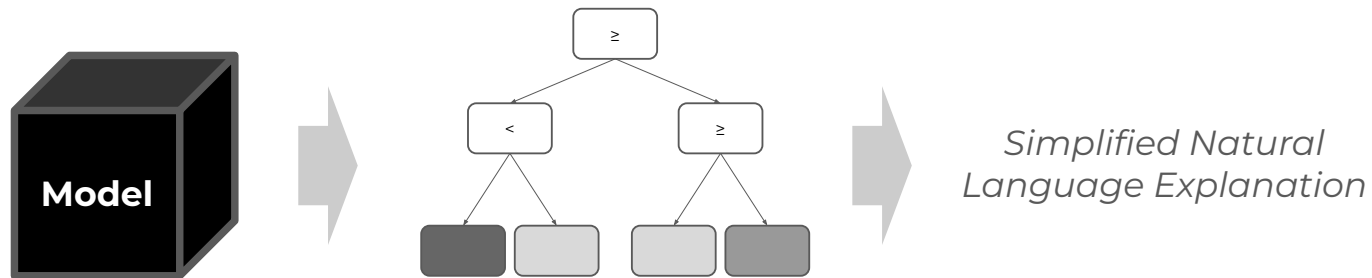
# Decision Tree to NLE

# Decision Tree to NLE



*Natural Language Explanation*

# Decision Tree to NLE



Simplified *Natural Language Explanation*

# Decision Tree to NLE



**Model** → [Decision Tree] → *Simplified Natural Language Explanation*

# Decision Tree to NLE



Simplified Natural Language Explanation

# LORE

## Factual and Counterfactual Explanations
## for Black-Box Decision Making

Riccardo Guidotti, Anna Monreale, Fosca Giannotti, Dino Pedreschi, Salvatore Ruggieri, and Franco Turini

*Guidotti, Monreale, Giannotti, Pedreschi, Ruggieri, and Turini (2019)*
Factual and Counterfactual Explanations for Black-Box Decision Making

---

## Stable and actionable explanations of black-box models through factual and counterfactual rules

Riccardo Guidotti[1] · Anna Monreale[1] · Salvatore Ruggieri[1] ·
Francesca Naretto[2] · Franco Turini[1] · Dino Pedreschi[1] · Fosca Giannotti[2]

**Abstract**
Recent years have witnessed the rise of accurate but obscure classification models that hide the logic of their internal decision processes. Explaining the decision taken by a black-box classifier on a specific input instance is therefore of striking interest. We propose a local rule-based model-agnostic explanation method providing stable and actionable explanations. An explanation consists of a factual logic rule, stating the reasons for the black-box decision, and a set of actionable counterfactual logic rules, proactively suggesting the changes in the instance that lead to a different outcome. Explanations are computed from a decision tree that mimics the behavior of the black-box locally to the instance to explain. The decision tree is obtained through a bagging-like approach that favors stability and fidelity: first, an ensemble of decision trees is learned from neighborhoods of the instance under investigation; then, the ensemble is merged into a single decision tree. Neighbor instances are synthetically generated through a genetic algorithm whose fitness function is driven by the black-box behavior. Experiments show that the proposed method advances the state-of-the-art towards a comprehensive approach that successfully covers stability and actionability of factual and counterfactual explanations.

**Keywords** Explainable AI · Local explanations · Model-agnostic explanations · Rule-based explanations · Counterfactuals

*Guidotti, Monreale, Ruggieri, Naretto, Turini, Pedreschi, and Giannotti (2022)*
Stable and actionable explanations of black-box models through factual and counterfactual rules

---

**Algorithm 1** Explainer agent process

**Input:** $x \leftarrow$ instance, $\bar{y} \leftarrow$ target classification,
  $b \leftarrow$ black-box classifier, $D \leftarrow$ dataset
**Output:** $i \leftarrow$ explanation data

1: $y \leftarrow b(x)$;
2: **assert** $y = \bar{y}$;
3: $LORE_{sa} = LORE_{exp}(b, D)$;
4: $e = \langle c, r, \Phi \rangle \leftarrow LORE_{sa}(x, b)$;
5: $c_t \leftarrow text(c)$;
6: $i = \langle$ feature names, target values, confidence, $c_t, x, r, \Phi, \bar{y} \rangle$
7: **return** $i$

# Decision Tree to NLE

# Decision Tree to NLE

# Decision Tree to NLE

# Goal

*Develop a **conversational interface** to provide **simplified explanations** for **non-expert** users*

# Goal

*Develop a conversational interface to provide simplified explanations for non-expert users*

- Internally, an explanation is generated from a Decision Tree built from the classifier using LORE

# Goal

*Develop a conversational interface to provide simplified explanations for non-expert users*

- Internally, an explanation is generated from a Decision Tree built from the classifier using LORE

- The generation process is distributed among **multiple specialized agents**

# MAINLE's Simplified Architecture

**User**

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

asks for an explanation

Instance,
Classification

**User**

**Parser
Agent**

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

# MAINLE's Simplified Architecture

# MAINLE

a Multi-Agent architecture, interactive, that provides Natural Language Local Explanations of classification outcomes

# Example

**user**
Hello, please explain the classification of a flower that has sepal length of 7.3, sepal width of 2.9, petal length of 6.3, and petal width of 1.8

# Example

**user**
Hello, please explain the classification of a flower that has sepal length of 7.3, sepal width of 2.9, petal length of 6.3, and petal width of 1.8

**assistant**
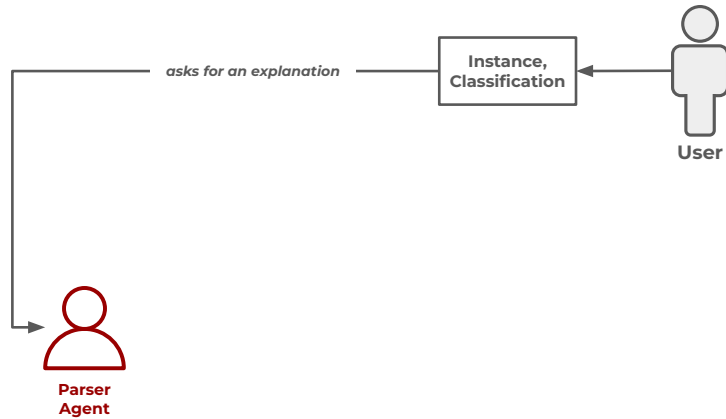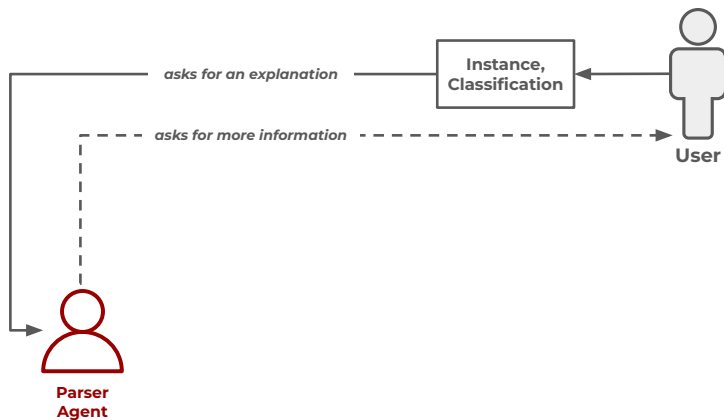gpt-4o-2024-08-06

I see that you have provided the following information:

- Sepal length: 7.3
- Sepal width: 2.9
- Petal length: 6.3
- Petal width: 1.8

However, I still need the classification of the flower, specifically the class (setosa, versicolor, or virginica). Could you please provide the class for the flower? Thank you!

# Example

**user**
Hello, please explain the classification of a flower that has sepal length of 7.3, sepal width of 2.9, petal length of 6.3, and petal width of 1.8

**assistant**
gpt-4o-2024-08-06

I see that you have provided the following information:

- Sepal length: 7.3
- Sepal width: 2.9
- Petal length: 6.3
- Petal width: 1.8

However, I still need the classification of the flower, specifically the class (setosa, versicolor, or virginica). Could you please provide the class for the flower? Thank you!

**user**
The class is virginica

# Example

**user**

Hello, please explain the classification of a flower that has sepal length of 7.3, sepal width of 2.9, petal length of 6.3, and petal width of 1.8

**assistant**
gpt-4o-2024-08-06

I see that you have provided the following information:

- Sepal length: 7.3
- Sepal width: 2.9
- Petal length: 6.3
- Petal width: 1.8

However, I still need the classification of the flower, specifically the class (setosa, versicolor, or virginica). Could you please provide the class for the flower? Thank you!

**user**

The class is virginica

**assistant**
gpt-4o-2024-08-06

The classifier identified the example as belonging to the 'virginica' class with a confidence level of approximately 81.78%. This conclusion was based on evaluating the features of your instance.

Here's a simplified explanation of why it was classified this way:

1. The sepal width is relatively high, indicating that it falls above a certain threshold.
2. The petal length is also high, reinforcing the previous point.
3. The sepal length is greater than a specific value, further supporting the classification.
4. Lastly, the petal width is above the necessary level, which aligns with the characteristics typical of the 'virginica' class.

Overall, these feature values closely matched the conditions required for an instance to be classified as 'virginica', leading the model to make this prediction with a strong level of certainty.

# Example

**assistant**
gpt-4o-2024-08-06

The classifier identified the example as belonging to the 'virginica' class with a confidence level of approximately 81.78%. This conclusion was based on evaluating the features of your instance.

Here's a simplified explanation of why it was classified this way:

1. The sepal width is relatively high, indicating that it falls above a certain threshold.
2. The petal length is also high, reinforcing the previous point.
3. The sepal length is greater than a specific value, further supporting the classification.
4. Lastly, the petal width is above the necessary level, which aligns with the characteristics typical of the 'virginica' class.

Overall, these feature values closely matched the conditions required for an instance to be classified as 'virginica', leading the model to make this prediction with a strong level of certainty.

# Example

The classifier identified the example as belonging to the 'virginica' class with a confidence level of
approximately 81.78%. This conclusion was based on evaluating the features of your instance.

Here's a simplified explanation of why it was classified this way:

1. The sepal width is relatively high, indicating that it falls above a certain threshold.
2. The petal length is also high, reinforcing the previous point.
3. The sepal length is greater than a specific value, further supporting the classification.
4. Lastly, the petal width is above the necessary level, which aligns with the characteristics typical of the
'virginica' class.

Overall, these feature values closely matched the conditions required for an instance to be classified as
'virginica', leading the model to make this prediction with a strong level of certainty.

**user**

If the sepal width was low, what would be the classification?

# Example

# Example

**assistant**
gpt-4o-2024-08-06

The classifier identified the example as belonging to the 'virginica' class with a confidence level of approximately 81.78%. This conclusion was based on evaluating the features of your instance.

Here's a simplified explanation of why it was classified this way:

1. The sepal width is relatively high, indicating that it falls above a certain threshold.
2. The petal length is also high, reinforcing the previous point.
3. The sepal length is greater than a specific value, further supporting the classification.
4. Lastly, the petal width is above the necessary level, which aligns with the characteristics typical of the 'virginica' class.

Overall, these feature values closely matched the conditions required for an instance to be classified as 'virginica', leading the model to make this prediction with a strong level of certainty.

**user**

If the sepal width was low, what would be the classification?

**assistant**
gpt-4o-2024-08-06

If the sepal width was low, that would likely lead to a different classification. More specifically, if the sepal width were below or equal to a certain threshold, it could change the outcome altogether.

In this case, the premises that would support a different classification would include:

- A low sepal width (below that specific threshold) would not support the conditions for 'virginica'.

Since the current classification of your instance relies on the relatively high sepal width, altering this feature to be low could lead to a classification of either 'versicolor' or 'setosa' depending on the other values of the instance.

To sum it up, a low sepal width would disrupt the conditions that led to the classification of 'virginica' and instead shift it towards one of the other classes.

**user**

Thank you

# Validation

- User Study

    25 people

# Validation

- User Study

    25 people

- Form with raw and simplified explanations

# Validation

- User Study

    25 people

- Form with raw and simplified explanations

- Likert scale rating

    1: Strongly Disagree

    2: Disagree

    3: Neutral

    4: Agree

    5: Strongly Agree

# Evaluation Criteria

- Technical Jargon

- Simplicity

- Completeness

- Conciseness

# Results - User Study - Explanation Rating

| Criterion | Raw Explanation | Simplified Explanation |
|---|---|---|
| **Technical Jargon** | 3.1 ± 1.3 | **4.0 ± 1.0** |
| **Simplicity** | 2.8 ± 1.3 | **4.4 ± 0.8** |
| **Completeness** | **3.4 ± 1.3** | 3.1 ± 1.3 |
| **Conciseness** | 3.0 ± 1.2 | **3.9 ± 1.2** |

# Results - User Study - Best Explanation

| Criterion | Raw | Simplified | Both rated equally |
|-----------|-----|------------|--------------------|
| **Technical Jargon** | 5 | **12** | 8 |
| **Simplicity** | 1 | **18** | 6 |
| **Completeness** | **10** | 9 | 6 |
| **Conciseness** | 5 | **15** | 5 |
| **Total** | 21 | **54** | 25 |

# Critic Agent Evaluation

# Critic Agent Evaluation

# Critic Agent Evaluation

# Results - Critic Agents - Explanation Rating

| Criterion | GPT-4o | | Gemini 2.0 | | LLaMA 3.2 | | DeepSeek-R1 | |
|---|---|---|---|---|---|---|---|---|
| | Raw | Sim. | Raw | Sim. | Raw | Sim. | Raw | Sim. |
| **Technical Jargon** | 2.1 | **4.2** | 3.7 | **4.7** | 2.1 | **3.9** | 2.9 | **3.6** |
| **Simplicity** | 2.0 | **4.8** | 3.7 | **4.7** | 3.9 | **4.8** | 2.7 | **4.6** |
| **Completeness** | 3.4 | **4.1** | **3.9** | 3.4 | 3.5 | **4.6** | **3.8** | **3.8** |
| **Conciseness** | 3.3 | **4.6** | 3.4 | **4.6** | 4.0 | **4.7** | 2.7 | **4.4** |

# Results - Critic Agents - Best Explanation

| Criterion | GPT-4o | | Gemini 2.0 | | LLaMA 3.2 | | DeepSeek-R1 | |
|---|---|---|---|---|---|---|---|---|
| | Raw | Sim. | Raw | Sim. | Raw | Sim. | Raw | Sim. |
| **Technical Jargon** | 0 | **25** | 0 | **25** | 0 | **25** | 8 | **17** |
| **Simplicity** | 0 | **25** | 0 | **25** | 0 | **25** | 1 | **24** |
| **Completeness** | **22** | 3 | 5 | **20** | 8 | **17** | **24** | 1 |
| **Conciseness** | 0 | **25** | 0 | **25** | 3 | **22** | 2 | **23** |
| **Total** | 22 | **78** | 5 | **95** | 11 | **89** | 35 | **65** |

# MAINLE

# Quality Attributes

**Interactivity**

# Quality Attributes

# Quality Attributes

# Quality Attributes

# Potential Risks and Concerns

1. LLM-related problems

# Potential Risks and Concerns

1. LLM-related problems

2. Usage of a surrogate model

# MAINLE: a Multi-Agent, Interactive, Natural Language Local Explainer of Classification Tasks

**Paulo Bruno Serafim**, Rômulo Férrer Filho, Stenio Freitas, Gizem Gezici, Fosca Giannotti, Franco Raimondi, and Alexandre Santos

**Poster Section
(Board 36)**
Thursday, 18/Sep
19:00 - 21:00

# Thank you!

*pre-print version*

github.com/paulobruno/ecml-pkdd-2025

paulobruno.github.io

paulo.desousa@gssi.it

GSSI | GRAN SASSO SCIENCE INSTITUTE | SCHOOL OF ADVANCED STUDIES | Scuola Universitaria Superiore

SCUOLA NORMALE SUPERIORE

UFC

# Extra slides

# Contributions

1. A multi-agent architecture to generate simplified explanations for any ML classifier

2. A process to generate natural language explanations using surrogate interpretable models for any classifier

3. Interactive explanations through conversation

4. Validation of the architecture with user studies

5. Public availability of the code

# LORE Algorithm

**Algorithm 1** Explainer agent process

**Input:** $x \leftarrow$ instance, $\bar{y} \leftarrow$ target classification,
$\quad b \leftarrow$ black-box classifier, $D \leftarrow$ dataset

**Output:** $i \leftarrow$ explanation data

1: $y \leftarrow b(x)$;                                  // get instance classification
2: **assert** $y = \bar{y}$;                             // validate target classification
3: $LORE_{sa} = LORE_{exp}(b, D)$;                       // generate explainer
4: $e = \langle c, r, \Phi \rangle \leftarrow LORE_{sa}(x, b)$;   // generate explanation
5: $c_t \leftarrow text(c)$;                             // convert tree to text
6: $i = \langle$feature names, target values, confidence, $c_t, x, r, \Phi, \bar{y}\rangle$
7: **return** $i$

# Quality Attributes

- Interactivity

- Modularity

- Flexibility

- Transparency

  - Decision transparency

  - Information transparency

# Information Transparency

# Information Transparency

**assistant**
gpt-4o-mini-2024-07-18

The model used for this classification is a decision tree classifier. Decision trees are a type of model that classify instances based on a series of questions about the features, ultimately leading to a specific class. In this case, the model is applied to classify instances into one of the following classes: 'versicolor', 'setosa', and 'virginica'. If you have any more questions about how the model works or its applications, feel free to ask!

**user**

please provide the decision tree used

**assistant**
gpt-4o-mini-2024-07-18

I'm sorry, but I cannot disclose the decision tree as it is confidential information. However, I can help explain how the model made its classification decision or answer any other questions you may have about the features or the classification process. Please let me know how else I can assist you!

# Datasets, Classifiers and LLMs

- Iris Flower
- Breast Cancer
- Wine
- Adult Income
- Credit Approval

- Random Forest
- Gradient Boosting
- MLP

- Gpt-4o
- Gemini 2.0
- Llama 3.2
- DeepSeek-R1

# Non-LLM Translator

## Raw Explanation

The instance is classified as virginica because sepal length (cm) is greater than 6.96 and sepal width (cm) is 2.09 and petal length (cm) is greater than 2.26 and petal width (cm) is greater than 1.08.

## Simplified Explanation

The instance was classified as 'virginica' based on the following features: the sepal length is high, the sepal width is moderate, and the petal width and petal length are also high. This combination of features strongly aligns with what is typically seen in the 'virginica' class of flowers. The model is quite confident in this classification, with a confidence level of about 81%. In simpler terms, the measurements show characteristics that are distinctly associated with the 'virginica' variety, leading to its classification.

# Questions from the audience

*Here are the two questions asked by the audience during the presentation and their summarized answers.*

1. **Did you ask the users if they could understand the explanation and not only to rate them?**
   a. No, because that would require users to have domain knowledge about the dataset, which we did not want to require, since we were evaluating MAINLE from the perspective of non-expert users.
2. **Which LLM did you use for the simplifier and how did you implement it to simplify an explanation?**
   a. (I showed the extra slide with the LLMs used) We created a structured prompt and passed it as system instructions for the LLMs that accept system prompts or in the question for the LLMs that do not accept them. All prompts are publicly available in the repository.