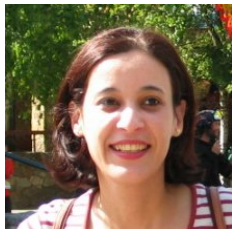
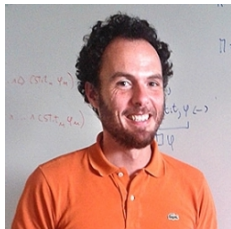


# Empowering Data-driven AI by Argumentation



Leila Amgoud (CNRS-IRIT)



Emiliano Lorini (CNRS-IRIT)



Philippe Muller (UPS-IRIT)

- **Argumentation** aims at increasing acceptability of claims by supporting them with **arguments**
- **Argument** is a set of *premises* intended to establish a *claim*

Premise 1

⋮

Premise n

Generally, birds fly

Tweety is a bird

Claim

Therefore, Tweety flies

# Several Types of Claims

- **Categorical** arguments (X is a Y)
- **Definitional** arguments (X is a Y; the definition of Y is contested)
- **Cause/Consequence** arguments (X causes Y; Y is a consequence of X)
- **Resemblance** arguments (X is like Y)
- **Evaluation** arguments (X is good or bad; X is true or false)
- **Proposal** arguments (One should do X)
- ...

⇒ Several Arguments Schemes

## ■ Practical applications

- Diagnosis in medical domain
- Online dispute resolution (e.g., CyberSettle)
- Online debate (e.g., DebateGraph, debate.org)
- Committees
- ...

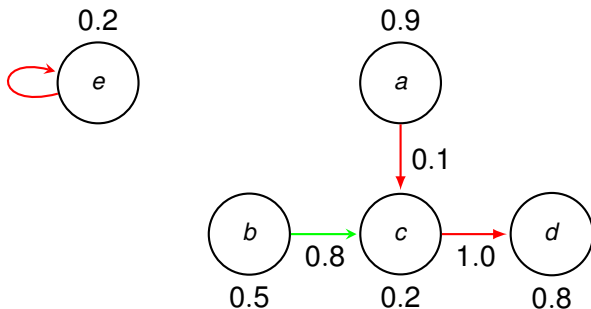
## ■ Theoretical applications

- Reasoning with (inconsistent, defeasible) information
- Decision making
- Negotiation
- Classification
- ...

# Argumentation Process

Given a **problem** (making a decision, classifying an object, ...)

- Construct **arguments**
- Identify their **basic strengths** + their **interactions**



# Example of Arguments

Let  $\Sigma$  be a finite propositional knowledge base.

## Definition

An **argument** is a pair  $\langle \Psi, \psi \rangle$  such that

- $\Psi \subseteq \Sigma$
- $\Psi$  is consistent
- $\Psi \vdash \psi$
- $\nexists \Psi'$  s.t.  $\Psi' \subset \Psi$  and  $\Psi' \vdash \psi$

$$\Sigma = \{p \wedge q, \neg p \wedge t\}$$

- $A = \langle \{p \wedge q\}, p \vee \neg t \rangle$
- $B = \langle \{p \wedge q\}, q \rangle$
- $C = \langle \{\neg p \wedge t\}, t \vee p \rangle$
- ...

# Example of Attacks

## Definition

$(\Psi, \psi)$  **attacks**  $(\Psi', \psi')$  iff  $\exists \phi \in \Psi'$  such that  $\psi \vdash \neg \phi$ .

$$\Sigma = \{p \wedge q, \neg p \wedge t\}$$

$$A = \langle \{p \wedge q\}, p \vee \neg t \rangle \text{ attacks } C = \langle \{\neg p \wedge t\}, t \vee p \rangle$$

Given a **problem** (making a decision, classifying an object, ...)

- Construct **arguments**
- Identify their **basic strengths** + their **interactions**  $\Rightarrow$  Graph
- **Analyse** the arguments  $\Rightarrow$  Semantics
- Conclude (the chosen option, the class of the object, ...)



## Individual arguments

- **Strength** concerns the **quality** of argument's components (premises, link, conclusion)

**Characteristics:** Uniqueness, Precise vs Vague

- **Acceptability** states whether an argument can be **accepted** so that its claim can safely be used for drawing conclusions, . . .

**Characteristics:** Uniqueness, Binary (Accepted, Rejected)

## Collections of arguments

- **Coalitions** Prevailing **viewpoints** expressed in an arg. graph

**Characteristics:** Multiple sets

# Three Families of Semantics

- A **semantics** is a function  $\pi$  that assigns to every  $\mathbf{G} = \langle \mathcal{A}, w, \mathcal{R}, \pi \rangle$ ,
  - a **set**  $\text{Ext}_{\mathbf{G}}^{\pi} \in 2^{2^{\mathcal{A}}}$  *(Extension Semantics)*
  - a **weighting**  $\text{Deg}_{\mathbf{G}}^{\pi} : \mathcal{A} \rightarrow \mathcal{D}$  *(Weighting Semantics)*
  - a **preorder**  $\succeq_{\mathbf{G}}^{\pi} \subseteq \mathcal{A} \times \mathcal{A}$  *(Ranking Semantics)*

$\mathcal{D}$  is a totally ordered scale.

Weighting Semantics  
Ranking Semantics

Strength

Extension Semantics

Acceptability

Coalitions

## Weighting semantics

Let  $\mathbf{G} = \langle \mathcal{A}, w, \mathcal{R}, \pi \rangle$ ,  $a \in \mathcal{A}$ ,  $b_1, \dots, b_n$  its attackers.

$$\text{Deg}(a) = f(w(a), g(h(\pi((b_1, a)), \text{Deg}(b_1)), \dots, h(\pi((b_n, a)), \text{Deg}(b_n))))$$

- $h : [0, 1] \times [0, 1] \rightarrow [0, 1]$
- $g : \bigcup_{n=0}^{+\infty} [0, 1]^n \rightarrow [0, +\infty)$  such that  $g$  is symmetric
- $f : [0, 1] \times \text{Range}(g) \rightarrow [0, 1]$

# Examples of Functions

$f_{comp}(x_1, x_2) = x_1(1 - x_2)$	$g_{sum}(x_1, \dots, x_n) = \sum_{i=1}^n x_i$	$h_{prod}(x_1, x_2) = x_1 x_2$
$f_{exp}(x_1, x_2) = x_1 e^{-x_2}$	$g_{sum, \alpha}(x_1, \dots, x_n) = \left(\sum_{i=1}^n (x_i)^\alpha\right)^{\frac{1}{\alpha}}$	$h_{prod, \alpha}(x_1, x_2) = x_1^\alpha x_2, \alpha > 0$
$f_{frac}(x_1, x_2) = \frac{x_1}{1+x_2}$	$g_{max}(x_1, \dots, x_n) = \max\{x_1, \dots, x_n\}$	$h_{min}(x_1, x_2) = \min\{x_1, x_2\}$
$f_{min}(x_1, x_2) = \min\{x_1, 1 - x_2\}$	$g_{psum}(x_1, \dots, x_n) = x_1 \oplus \dots \oplus x_n,$ where $x_1 \oplus x_2 = x_1 + x_2 - x_1 x_2$	$h_{Ham}(x_1, x_2) = \frac{x_1 x_2}{x_1 + x_2 - x_1 x_2};$ $h_{Ham}(x_1, x_2) = 0$ if $x_1 = x_2 = 0$

The choice of functions depends on **axioms** that need to be satisfied by a semantics

# Example of a Weighting Semantics

## *h*-Categorizer

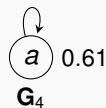
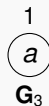
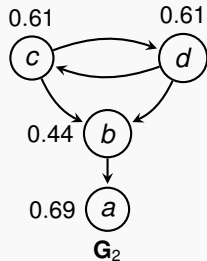
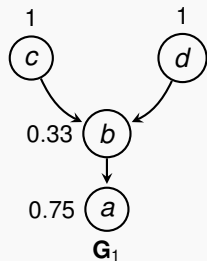
Let  $\mathbf{G} = \langle \mathcal{A}, \mathcal{R} \rangle$  and  $a \in \mathcal{A}$ .

$$\text{Deg}_{\mathbf{G}}^{\pi}(a) = \frac{1}{1 + \sum_{(b,a) \in \mathcal{R}} \text{Deg}_{\mathbf{G}}^{\pi}(b)}$$

$$\left\{ \begin{array}{l} g_{\text{sum}}(x_1, \dots, x_n) = \sum_{i=1}^n x_i \\ f_{\text{frac}}(x) = \frac{1}{1+x_2} \\ \mathcal{D} = [0, 1] \end{array} \right.$$

# Example of a Weighting Semantics

## Examples



- Argumentative counterparts of data-driven models
- Explanation theory and persuasion

# Argumentative Counterparts of Data-driven Models

- Define argumentative view of existing data-driven models, namely NNs
- Improve predictions by incorporating arguments given by experts
- Reduce the need for large amounts of data. The new arguments introduce crucial domain knowledge,
- Improve search performance. The new arguments will constrain the combinatorial search among possible hypotheses



- Links between explanation and argument
- Explanation schemes
- Evaluation of explanations
- Persuasive explanation
- Which explanation to present to users and under which format ?

- 2 PhD thesis (one for each part of the project)
- 1 (or 2) postdoc on the first part

- Joao Marques-Silva
- Louise Travé-Massouyès
- Jean-Michel Loubes