

Multi-Agent Systems

Albert-Ludwigs-Universität Freiburg



Bernhard Nebel, Felix Lindner, and Thorsten Engesser

Winter Term 2018/19

- Modeling agents exchanging arguments
 - Argumentation frameworks
 - Semantics
 - Algorithms

- A: My government cannot negotiate with your government because your government does not even recognize my government.
- B: Your government does not recognize my government either.
- A: But your government is a terrorist government.
- Which arguments should be accepted?

- A: Ralph goes fishing, because it is sunday.
- B: Ralph does not go fishing, because it is Mother's day, so he visits his parents.
- C: Ralph cannot visit his parents, because it is a leap year, so they are on vacation.
- Which arguments should be accepted?

- A statement is accepted if it can be successfully defended against attacking arguments.

Definition (Argument)

An **argument** is a pair (S, φ) , such that S is a set of formulae and φ can be derived from S . S is also called the **support** for the **claim** φ .

Definition (Attack)

Two definitions of attack:

Undercut Argument $A_1 = (S_1, \varphi_1)$ **undercuts** argument $A_2 = (S_2, \varphi_2)$ iff $\neg\varphi_2$ can be derived from S_1 .

Rebuttal Argument $A_1 = (S_1, \varphi_1)$ **rebutts** argument $A_2 = (S_2, \varphi_2)$ iff $\varphi_1 \equiv \neg\varphi_2$.

We can decide what to believe while looking at arguments at the abstract level (Dung, 1995):

- Disregarding internal structures of arguments
- Focus on the attack relation between arguments
 (a, b, c, d, \dots) : a **attacks** b or $a \rightsquigarrow b$
- Not concerned with the origin of arguments or the attack relation

Abstract argumentation framework

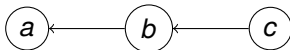
An **argumentation framework** is a pair $\mathcal{AF} = (Arg, \rightsquigarrow)$ where Arg is a set of arguments and $\rightsquigarrow \subseteq Arg \times Arg$. We say that $a \in Arg$ attacks $b \in Arg$ iff $(a, b) \in \rightsquigarrow$.

■ Remember:

- A: Ralph goes fishing, because it is sunday.
- B: Ralph does not go fishing, because it is Mother's day, so he visits his parents.
- C: Ralph cannot visit his parents, because it is a leap year, so they are on vacation.

■ Representation as an argumentation framework:

$$\mathcal{AF} = \langle \{a, b, c\}, \{(b, a), (c, b)\} \rangle,$$



$$\mathcal{AF} = \langle \{a, b, c\}, \{(b, a), (c, b)\} \rangle,$$

\mathcal{AF} can also be understood as a logic program $\Pi_{\mathcal{AF}}$:

$a :- \text{not } -a.$

$b :- \text{not } -b.$

$c :- \text{not } -c.$

$-a :- b.$

$-b :- c.$

The output of $\Pi_{\mathcal{AF}}$ is just its models. Those atoms that are true in the model correspond to the accepted arguments. Which ones?

- **Argument-based semantics** get as input an argumentation framework and output zero or more sets of acceptable arguments.

Definition: Labelling

Let $\mathcal{AF} = (Arg, \rightsquigarrow)$ be an argumentation framework. A **labelling** of \mathcal{AF} is a total function $\mathcal{Lab} : Arg \rightarrow \{in, out, undec\}$. The set of all labellings will be denoted by $\mathcal{L}(\mathcal{AF})$.

- $in(\mathcal{Lab}) = \{a \mid \mathcal{Lab}(a) = \mathbf{in}\}$
- $out(\mathcal{Lab}) = \{a \mid \mathcal{Lab}(a) = \mathbf{out}\}$
- $undec(\mathcal{Lab}) = \{a \mid \mathcal{Lab}(a) = \mathbf{undec}\}$
- To refer to a labelling \mathcal{Lab} we will also write $\langle in(\mathcal{Lab}), out(\mathcal{Lab}), undec(\mathcal{Lab}) \rangle$

$$\mathcal{AF} = \langle \{a, b, c\}, \{(b, a), (c, b)\} \rangle,$$



$$\mathcal{L}(\mathcal{AF}) = \{ \langle \emptyset, \emptyset, \{a, b, c\} \rangle, \langle \emptyset, \{a\}, \{b, c\} \rangle \dots \}$$

- How to identify the appropriate labellings?
- E.g., we do not want to accept both a and b , thus if $\mathcal{Lab}(a) = \mathbf{in}$ then $\mathcal{Lab}(b) \neq \mathbf{in}$.

Definition: Admissible labelling

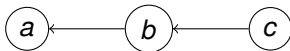
Definition

Let \mathcal{Lab} be a labelling of argumentation framework \mathcal{AF} . An **in**-labelled argument is said to be **legally in** iff all its attackers are labelled **out**. An **out**-labelled argument is said to be **legally out** iff it has at least one attacker that is labelled **in**.

Definition

Let \mathcal{AF} be an argumentation framework. An **admissible labelling** is a labelling where each **in**-labelled argument is legally **in** and each **out**-labelled argument is legally **out**.

$$\mathcal{AF} = \langle \{a, b, c\}, \{(b, a), (c, b)\} \rangle,$$



Admissible labellings

- $\langle \emptyset, \emptyset, \{a, b, c\} \rangle$
- $\langle \{a, c\}, \{b\}, \emptyset \rangle$

Definition

Given an argumentation framework $\mathcal{AF} = (Arg, \rightsquigarrow)$, a **labelling semantics** S associates with \mathcal{AF} a subset of $\mathcal{L}(\mathcal{AF})$, denoted as $\mathcal{L}_S(\mathcal{AF})$.

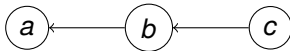
Definition

Let $\mathcal{AF} = (Arg, \rightsquigarrow)$ be an argumentation framework and $\mathcal{Lab} : Arg \rightarrow \{in, out, undec\}$ be a total function. We say that \mathcal{Lab} is a **complete labelling** iff it satisfies the following:

$\forall a \in Arg : (\mathcal{Lab}(a) = \mathbf{out} \leftrightarrow \exists b \in Arg : (b \rightsquigarrow a \wedge \mathcal{Lab}(b) = \mathbf{in}))$

$\forall a \in Arg : (\mathcal{Lab}(a) = \mathbf{in} \leftrightarrow \forall b \in Arg : (b \rightsquigarrow a \rightarrow \mathcal{Lab}(b) = \mathbf{out}))$

$$\mathcal{AF} = \langle \{a, b, c\}, \{(b, a), (c, b)\} \rangle,$$



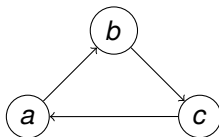
Complete labellings:

1 $\langle \{a, c\}, \{b\}, \emptyset \rangle$

Why not $\langle \emptyset, \emptyset, \{a, b, c\} \rangle$?

- **A:** Bert says that Ernie is unreliable, therefore everything that Ernie says cannot be relied on.
- **B:** Ernie says that Elmo is unreliable, therefore everything that Elmo says cannot be relied on.
- **C:** Elmo says that Bert is unreliable, therefore everything that Bert says cannot be relied on.

$$\mathcal{AF} = \langle \{a, b, c\}, \{(a, b), (b, c), (c, a)\} \rangle,$$

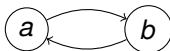


Complete labellings:

1 $\mathcal{Lab}_1 : \langle \emptyset, \emptyset, \{a, b, c\} \rangle$

- A: Nixon is a pacifist, because he is a quaker.
- B: Nixon is not a pacifist, because he is a republican.

$$\mathcal{AF} = \langle \{a, b\}, \{(a, b), (b, a)\} \rangle,$$



Complete labellings:

- 1 $\mathcal{Lab}_1 : \langle \emptyset, \emptyset, \{a, b\} \rangle$
- 2 $\mathcal{Lab}_2 : \langle \{a\}, \{b\}, \emptyset \rangle$
- 3 $\mathcal{Lab}_3 : \langle \{b\}, \{a\}, \emptyset \rangle$

⇒ Three resonable positions a rational agent can take.

Definition

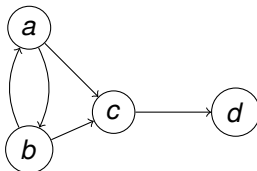
Let \mathcal{AF} be an argumentation framework. The **grounded labelling** of \mathcal{AF} is a complete labelling \mathcal{Lab} where $in(\mathcal{Lab})$ is minimal w.r.t. set inclusion.

- Grounded semantics picks the complete labelling with minimal **in**, minimal **out**, and maximal **undec**.
- Intuitively, the arguments in **in** are those that must be accepted by every rational agent.
- These arguments are in the **in** set of every complete labelling.
- The grounded labelling is unique.

Definition

Let \mathcal{AF} be an argumentation framework. The **preferred labelling** of \mathcal{AF} is a complete labelling \mathcal{Lab} where $in(\mathcal{Lab})$ is maximal w.r.t. set inclusion.

- Preferred semantics picks the complete labelling with maximal **in**, maximal **out**, and minimal **undec**.
- For every argumentation framework one or more preferred labellings exists.



- Ground labelling: $\langle \emptyset, \emptyset, \{a, b, c, d\} \rangle$
- Preferred labellings: $\langle \{a, d\}, \{b, c\}, \emptyset \rangle, \langle b, d \rangle, \{a, c\}, \emptyset$

Observe: Ground labelling is not among the preferred labellings and none of the preferred labellings is the ground labelling. Also, it is not the case that the ground labelling coincides with the intersection of all preferred labellings.

Definition

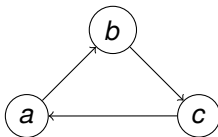
Let \mathcal{Lab} be a labelling of an argumentation framework \mathcal{AF} . \mathcal{Lab} is a **stable labelling** of \mathcal{AF} iff it is a complete labelling with $\mathbf{undec}(\mathcal{Lab}) = \emptyset$.

- Stable semantics decides for every argument if it is **in** or **out**, no **undec**.
- As it minimizes **undec** it maximizes **in** and **out**. Thus, every stable labelling is a preferred labelling.
- But not vice versa: Whereas a preferred labelling always exists, the existence of a stable labelling is not guaranteed.



Complete labellings:

$$1 \quad \mathcal{Lab}_1 : \langle \emptyset, \emptyset, \{a\} \rangle$$



Complete labellings:

$$1 \quad \mathcal{Lab}_2 : \langle \emptyset, \emptyset, \{a, b, c\} \rangle$$

$\Rightarrow \mathcal{Lab}_1, \mathcal{Lab}_2$ are complete, ground, preferred, but not stable.

restriction on complete labelling	resulting semantics
no restrictions	complete semantics
empty undec	stable semantics
maximal in	preferred semantics
maximal out	preferred semantics
maximal undec	grounded semantics
minimal in	grounded semantics
minimal out	grounded semantics

- Every complete labelling is admissible.
- Every ground labelling is complete.
- Every preferred labelling is complete.
- Every stable labelling is preferred.

Definition (Credulous Acceptance)

Given $\mathcal{AF} = (Arg, \rightsquigarrow)$ and $a \in Arg$: is a labelled **in** in at least one ground/preferred/stable/... labelling?

- NP-complete for stable, admissible, complete, preferred
- P-complete for ground

Definition (Skeptical Acceptance)

Given $\mathcal{AF} = (Arg, \rightsquigarrow)$ and $a \in Arg$: is a labelled **in** in every ground/preferred/stable/... labelling?

- stable: co-NP-complete
- admissible: trivially false (empty **in** is admissible)
- complete, ground: P-complete
- preferred: Π_2^p -complete

- Other interesting decision problems:
 - Given some labelling, is it ground/preferred/stable/...?
 - Does there some ground/preferred/stable/... labelling exist?
 - Does there some nonempty ground/preferred/stable/... labelling exist?

- Given an argument A and an argumentation framework \mathcal{AF} , is A in the **in** set of \mathcal{AF} 's ground labelling?
- Given an argument A and an argumentation framework \mathcal{AF} , is A in the **in** set of some of \mathcal{AF} 's preferred labellings?

Definition

A **partial labelling** is a partial function $\mathcal{Lab} : \text{Args} \rightarrow \{\mathbf{in}, \mathbf{out}\}$ such that

- if $\mathcal{Lab}(A) = \mathbf{in}$ then for each attacker B $\mathcal{Lab}(B) = \mathbf{out}$
- if $\mathcal{Lab}(A) = \mathbf{out}$ then for some attacker B $\mathcal{Lab}(B) = \mathbf{in}$
- Partial labellings are admissible labellings
- A partial labelling \mathcal{Lab} can be extended to a complete labelling $\mathcal{Lab}' \supseteq \mathcal{Lab}$
- For each complete labelling \mathcal{Lab}' there exists a partial labelling $\mathcal{Lab} \subseteq \mathcal{Lab}'$ (just remove the **undec** labels)

Definition

$extendin(\mathcal{L}ab) = \mathcal{L}ab \cup \{(A, \mathbf{in}) \mid \forall B[B \rightsquigarrow A \rightarrow \mathcal{L}ab(B) = \mathbf{out}]\}$

$extendout(\mathcal{L}ab) = \mathcal{L}ab \cup \{(A, \mathbf{out}) \mid \exists B[B \rightsquigarrow A \wedge \mathcal{L}ab(B) = \mathbf{in}]\}$

$extendinout(\mathcal{L}ab) = extendin(\mathcal{L}ab) \circ extendout(\mathcal{L}ab)$

- If $\mathcal{L}ab$ is a partial labelling, then $extendin(\mathcal{L}ab)$, $extendout(\mathcal{L}ab)$, $extendinout(\mathcal{L}ab)$ are partial labellings.

function GROUNDLABELLING(\mathcal{AF})

$L \leftarrow \emptyset$

repeat

$L_{old} \leftarrow L$

$L \leftarrow extendin(L)$

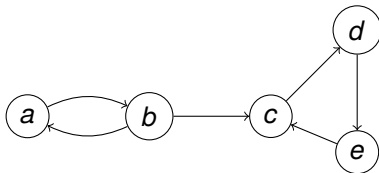
$L \leftarrow extendout(L)$

until $L = L_{old}$

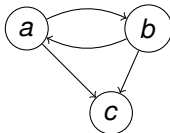
return $L \cup \{(A, \mathbf{undec}) \mid (A, \mathbf{in}) \notin L \text{ and } (A, \mathbf{out}) \notin L\}$

end function

- **Idea:** Take the other's opinion and then derive a contradiction:
 - Proponent (M) makes a statement (A)
 - Opponent (S) derives from A more statements M will be committed to
 - S aims at letting M commit himself to a contradiction
- **Dialog game**
 - M starts and claims the existence of a reasonable position (admissible labelling) in which a particular argument is accepted (labelled **in**).
 - S confronts M with the consequences of M's own position, and asks M to resolve these consequences.
 - S wins if she leads M to a contradiction.
- If M wins then his argument is in the **in** set of an admissible labelling, and thus in the **in** of a preferred labelling.



- M: in(D) *I have an admissible labelling in which D is in*
- S: out(C) *But then in your labelling C is out. Why?*
- M: in(B) *Because B is in*
- S: out(A) *But then A must be out. Why?*
- M: in(B) *Because B is in.*



- M: in(C) *I have an admissible labelling in which C is in*
- S: out(A) *But then in your labelling A is out. Why?*
- M: in(B) *Because B is in*
- S: out(B) *But B must be out!*

Definition

Let $\mathcal{AF} = (Arg, \rightsquigarrow)$ be an argumentation framework. An **admissible discussion** is a sequence of moves

$[\Delta_1, \dots, \Delta_n] (n \geq 0)$ such that:

- each move $\Delta_i (1 \leq i \leq n)$ where i is odd is called M-move and is of the form $in(A)$
- each move $\Delta_i (1 \leq i \leq n)$ where i is even is called S-move and is of the form $out(A)$
- for each S-move $\Delta_i = out(A) (2 \leq i \leq n)$ there exists an M-move $\Delta_j = in(B) (j < i)$ such that A attacks B
- for each M-move $\Delta_i = in(A) (3 \leq i \leq n)$ it holds that Δ_{i-1} is of the form $out(B)$, where A attacks B
- there exist no two S-moves $\Delta_i = \Delta_j$ with $i \neq j$

Definition

An admissible discussion $[\Delta_1, \dots, \Delta_n]$ is said to be **finished** iff

- 1 There exists no Δ_{n+1} such that $[\Delta_1, \dots, \Delta_n, \Delta_{n+1}]$ is an admissible discussion, or there exists a M-move and a S-move containing the same argument
- 2 No subsequence of the discussion is finished.

Definition

A finished admissible discussion is **won** by player S if there exist a M-move and a S-move containing the same argument. Otherwise, it is **won** by the player making the last move.

Theorem

Let g be an admissible discussion won by M and let $\mathcal{L}ab : Ar \rightarrow \{\mathbf{in}, \mathbf{out}, \mathbf{undec}\}$ be a function defined as follows. For every argument $B \in Ar$:

- $\mathcal{L}ab(B) = \mathbf{in}$ if B was labeled in during g
- $\mathcal{L}ab(B) = \mathbf{out}$ if B was labeled out during g
- $\mathcal{L}ab(B) = \mathbf{undec}$ otherwise

Then $\mathcal{L}ab$ is an admissible labelling.

- Thus, if there is a winning game for M defending A then A is in the **in** set of some preferred labelling (add **undec** arguments to **in** as long as possible).

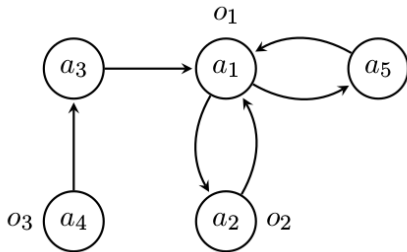
Proposition 6 *Let (A, R) be an argument system. A set $S \subseteq A$ is admissible iff S is a model of the formula*

$$\bigwedge_{a \in A} ((a \rightarrow \bigwedge_{b: (b,a) \in R} \neg b) \wedge (a \rightarrow \bigwedge_{b: (b,a) \in R} (\bigvee_{c: (c,b) \in R} c))).$$

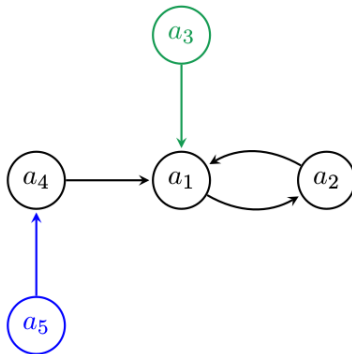
Example

- $\mathcal{AF} = (\{a, b\}, \{(a, b), (b, a)\})$
- $\varphi_{\mathcal{AF}} = (a \rightarrow \neg b) \wedge (a \rightarrow a) \wedge (b \rightarrow \neg a) \wedge (b \rightarrow b)$
- Is a credulously accepted: Is $a \wedge \varphi_{\mathcal{AF}}$ satisfiable?

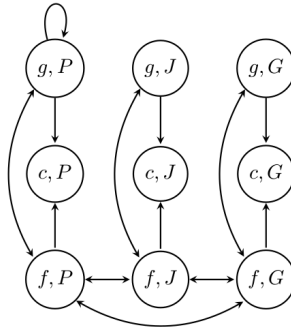
- Can be used to decide what to do next.
- Can be used to find perfect matchings [3]
 - Arg: The couples
 - $(m_1, w_1) \rightsquigarrow (m_2, w_2)$ iff
 - $m_1 = m_2$ and m_1 prefers w_1 to w_2 , or
 - $w_1 = w_2$ and w_1 prefers m_1 to m_2
- Ressource allocation
 - Arg: Pairs $(agent, task)$
 - $(agent_i, task_i) \rightsquigarrow (agent_j, task_j)$ iff one of:
 - $(agent_i, task_i)$ is preferred to $(agent_j, task_j)$
 - $(agent_i, task_i)$ excludes $(agent_j, task_j)$
 - Agent is unable to do $task_i$ (then self attack of $(agent_i, task_i)$)
- Can be used to compute the set of arguments an agent should utter / keep for itself (Persuasion).



Source: [4]



Source: [4]








Source: [4]

- In abstract argumentation systems all arguments are equally strong—relaxation
 ~> **Preference-based argumentation systems** (e.g., Amgoud et al. 1998f) model preference (weights) of arguments.
- Acceptability of arguments can depend on the target audience (e.g., newspaper vs. scientific article)
 ~> **Value-based argumentation systems** (Bench-Capon et al, 2003ff)
- Arguments in abstract argumentation systems do not have an internal (logical) structure
 ~> **Deductive argumentation systems**

Some projects available for BA/MA:

- SAT reductions of argumentation-framework semantics.
Complexity preserving?
- Dialogue-based algorithms for the other semantics
- More types of semantics: semi-stable, stage, ideal, eager, cf2, ...
- Building argumentation frameworks from utterances
- Planning utterances based on argumentation frameworks
- Verbalization of argumentation frameworks

-  M. Caminada, A gentle introduction to argumentation semantics. Technical report, University of Luxembourg, Summer 2008.
-  M. Caminada, W. Dvorak, S. Vesic, Preferred semantics as socratic discussion. Journal of Logic and Computation, 2014.
-  P.M. Dung. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. Artificial Intelligence 77, pp. 321-357, 1995.
-  J.-G. Maily, Dynamics of Argumentation Frameworks, PhD Thesis, 2015.
-  Philippe Besnard and Sylvie Doutre. Checking the acceptability of a set of arguments. In Proceedings of the 10th International Workshop on Non-Monotonic Reasoning (NMR'04), pages 59-64, 2004.