

異なる価値観をもつエージェント間での合意形成戦略の提案及びその評価

Proposal of a Strategy for an Agreement between Agents with Their Inherent Values and Its Evaluation

森野尊行^{1*} 高橋和子¹
Takayuki Morino¹ Kazuko Takahashi¹

¹ 関西学院大学大学院理工学研究科

¹ Graduate School of Science and Technology, Kwansai Gakuin University

Abstract: We model a dialogue based on an argumentation framework between agents with their own inherent values and consider strategies and conditions leading to solutions that increase the satisfaction of both agents. As an evaluation criterion, we define satisfaction in agreement of the acceptability of subject argument at the end of dialogue and the number of acceptable arguments and propose a strategy to increase satisfaction. Implementing this dialogue model, we evaluate the experiment and show the effectiveness of the strategy.

1 はじめに

対話モデルの研究は人工知能の一分野として古くから行われている。近年の Dung の提案した議論フレームワーク [1] は論理プログラミングや非単調推論との関連から大きな注目をあび、これに基づく対話モデルについて多くの研究がなされてきている。横浜は、その中の一つ Amgoud らの対話モデル [2] に予測知識を導入した説得対話モデルを提案した [3]。そこでは、ある議題について対立するエージェントが互いに意見を述べ、一方が他方を説得するための戦略が提案されている。しかし、横浜の提案した対話モデルでは、価値観の異なるエージェント同士の対話は対象としていない。実際の対話では、価値観が異なるためにエージェントの意見衝突が解消しなかったり、合意がなされたとしても互いが十分納得できずに終わる場合がある。

価値観を導入した対話モデルも研究されているが [4]、戦略や結果の満足度にまでは言及されていない。対話の結果たとえ合意が得られたとしてもエージェントがどのくらい自分の意見を出せたかによって満足度は異なる、これについてはゲーム理論に基づいて対話の評価をする方法が提案されている [5]。

本発表では、これらの研究をもとに、価値観の違いによって立場が分かれているエージェント同士での対話をモデル化し、両者が合意をし満足度が高くなる結

果が得られるような戦略について考察する。対話終了時の議題の受理可能性の一致とエージェントが持つ議論フレームワーク内の受理可能な論証の数を対話の満足度として設定する。まず、基本戦略を作成し、この対話モデルを実装してシミュレーション実験を行う。次に、実験結果を解析することで、新たな戦略を発見し、それらを反映して再度シミュレーションを行う。

本発表は以下のように構成される。第2節は対話モデルの基礎となる議論フレームワークについて説明する。第3節は本研究で提案した対話モデルについて説明する。第4節は評価実験を行い対話を解析することで新たな戦略について考察し、その有効性を調べる。最後に第5節でまとめる。

2 議論フレームワーク

Dung が提案した議論フレームワークは論証の集合とその集合上の関係である攻撃関係の二項組 $AF = \langle AR, AT \rangle$ で定義される。論証はエージェントの発言であり、攻撃関係はある論証がどの論証に反論をしているかを表したものである。また議論フレームワークは論証をノード、攻撃関係をエッジとする有向グラフで表すことができる。本研究で扱う議論フレームは全て連結グラフで表されるもの限定して考える。

議論フレームワーク $AF_1 = \langle AR_1, AT_1 \rangle$, $AF_2 = \langle AR_2, AT_2 \rangle$ を考える。 $AR_1 \subseteq AR_2$, $AT_1 \subseteq AT_2$ を満たすとき、 AF_1 を AF_2 の部分議論フレームワークと

*連絡先：関西学院大学大学院理工学研究科
〒669-1337 兵庫県三田市学園2丁目1番地 関西学院大学理工学部
E-mail: fvo65684@kwansai.ac.jp

よび, $AF_1 \subseteq AF_2$ と表記する.

議論フレームワーク $AF_1 = \langle AR_1, AT_1 \rangle$, $AF_2 = \langle AR_2, AT_2 \rangle$ を考える. $AF_1 \subseteq AF_2$ かつ $|AT_1| = \lfloor \frac{|AT_2|}{2} \rfloor$ を満たすとき, AF_1 は AF_2 の half であるといい, $AF_1 = \text{half}(AF_2)$ と表記する.

$AF = \langle AR, AT \rangle$ において

1. $A, B, C \in AR$, $(A, B), (B, C), (C, A) \in AT$ が成り立つとき, AF 内に正論証三角形が存在するといひ, $\langle A, B, C \rangle$ は正論証三角形の辺とよぶ.
2. $A, B, C \in AR$, $(A, B), (B, C), (A, C) \in AT$ が成り立つとき, AF 内に逆論証三角形が存在するといひ, $\langle A, B, C \rangle$ は逆論証三角形の辺とよぶ.
3. $A, B, C \in AR$, $(A, B), (B, C), (A, C), (C, A) \in AT$ が成り立つとき, AF 内に双方向論証三角形が存在するという.

議論フレームワーク $AF = \langle AR, AT \rangle$, 関数 $L^{AF} : AR \rightarrow \{in, out, undec\}$ を考える. AR に属する全ての論証 A が以下の条件を満たすとき, L^{AF} を AF に対する完全ラベリングという [6].

1. $L^{AF}(A) = in \leftrightarrow A$ を攻撃する論証のラベルが全て "out" である.
2. $L^{AF}(A) = out \leftrightarrow A$ を攻撃する論証のラベルが少なくとも一つは "in" である.
3. $L^{AF}(A) = undec \leftrightarrow A$ の論証のラベルが "in", "out" とラベル付けできない.

また AF に対する完全ラベリング L^{AF} において $L^{AF}(A) = in$ のとき A は AF 内で受理可能とよび, 受理可能である論証の集合を $in(L^{AF})$ と表記する. $in(L^{AF})$ はそのエージェントの信じている論証の集合に相当する.

議論フレームワーク $AF = \langle AR, AT \rangle$ の完全ラベリング L^{AF} において, $in(L^{AF})$ が集合の包含関係において極小である完全ラベリングを基礎ラベリングとよぶ. なお本研究の議論フレームワークで用いられる意味論には基礎ラベリングが用いられるとする.

3 対話モデル

本研究の説得対話は議題 ρ に関して二人のエージェント (提案者 P , 対立者 Q) 間で行われる. 各エージェントは結果として ρ に関しての賛否が一致し, かつ, 相手を納得させる情報を多く述べることを目標とする.

各エージェントは知識として自分自身の議論フレームワークと相手の予測議論フレームワークの二つを持つ. 自分自身が持つ議論フレームワークは, エージェントがプロトコルに従って, 発言可能な論証の中から,

予測議論フレームワークを用いて説得のために最適な論証を発言する. エージェントの持つ二つの議論フレームワークはエージェントが発言することにより, 新たに論証が加わることで更新される. 対話終了時に二人のエージェントの議論フレームワークにおける議題のラベルが一致すれば説得は成功したと見なされる. また, そのときに信じている論証の数が多ほど, エージェントがより納得のいく対話だと考えられる.

3.1 エージェントの価値観

二人のエージェントは各論証について独自の価値観をもつ. そのため, 論証 A, B に対して, 片方のエージェントは論証 A から論証 B へ攻撃できると考えており, もう片方のエージェントは B から A への攻撃できると考えている場合がある.

エージェント X , その対話相手である Y 間で 対象とする全ての情報から構成される議論フレームワークを全体議論フレームワークを $UAF = \langle AR, AT \rangle$ とする. エージェント X の価値観である全体議論フレームワーク $UAF_X = \langle UAR_X, UAT_X \rangle$ とエージェント Y の価値観である全体議論フレームワーク $UAF_Y = \langle UAR_Y, UAT_Y \rangle$ はいずれでも UAF の部分議論フレームワークであり, UAF に属する論証 A, B 間がお互いに攻撃しあっている場合は, エージェント X は相手のエージェント Y とは 異なった攻撃関係を UAF から取り出し, それ以外の場合は共通の攻撃関係をとる.

3.2 知識の前提条件

任意のエージェントを X, X の全体議論フレームワーク $UAF_X = \langle UAR_X, UAT_X \rangle$ とする. X が持つ初期議論フレームワーク $AF_X = \langle AR_X, AT_X \rangle$ は, $AR_X \subseteq UAR_X$, $AT_X = (AR_X \times AR_X) \cap UAT_X$ という条件を満たす. つまり, X が持つ初期議論フレームワークの論証は, 自分の全体議論フレームワーク内の論証の範囲内で構成される.

X が持つ Y の初期予測議論フレームワーク $PAF_Y = \langle PAR_Y, PAT_Y \rangle$ は, 条件 $PAR_Y \subseteq AR_X \wedge PAR_Y \subseteq AR_Y$, $PAT_Y = (PAR_Y \times PAR_Y) \cap UAT_Y$ を満たすものとする. つまり, 予測議論フレームワークは自分が知っていて, なおかつ, 相手が知っている論証の範囲で構成される.

3.3 対話の進行

提案者 P , 対立者 Q の間の説得対話は, 最初に P が議題 ρ を述べた後 Q, P が交互に手を出すことで進行する. 手は論証を述べるか, 何も言わずパスをするかで

ある。また、対話が終了する条件は双方が続けてパスをすることである。

エージェントを X , act を T とする。手とは X と T の組 (X, T) である。act は論証 A を発言する $assert(A)$, 何も述べずパスする $pass$ のいずれかである。

提案者 P と対立者 Q 間の対話 $d_k (k \geq 1)$ とは以下の条件を満たす手の有限列 $[m_0, \dots, m_{k-2}, m_{k-1}]$ である。手 m_i は $(X_i, T_i) (0 \leq i \leq k-1)$ の形をなす。

1. 対話の始まりに P が議題 ρ を述べる。
2. 対話は P, Q が交互に手を出して進行する。
3. 手は、可能手の条件を満たさないと発言することはできない。可能手については後で述べる。

対話 $d_k = [m_0, \dots, m_{k-2}, m_{k-1}] (k \geq 1)$ における X のコミットメントストア $CS_X^{d_k}$ は対話中にエージェント X が発言した論証の集合であり、

$\bigcup_{i=0, \dots, k-1, X_i=X, T_i=assert(A)} \{A\}$ と定義する。

エージェント X の全体議論フレームワークを $UAF_X = (UAR_X, UAT_X)$ とする。議論フレームワーク $AF_X = (AR_X, AT_X)$ が論証 A を受け取ったときの UAF_X に関する更新を $AF \circ A$ と表記し、以下のよう

- $AF_X \circ A = (AR_X \circ A, AT_X \circ A)$
 1. $AR_X \circ A = AR \cup \{A\}$
 2. $AT_X \circ A = AT \cup \{(A, B), (C, A) \mid B, C \in AR, (A, B), (C, A) \in UAT_X\}$

X がもつ AF_X の更新は UAF_X に関するものになり、 X がもつ PAF_Y の更新は UAF_Y に関するものになる。これらの議論フレームワークは対話の進行によって更新される。

エージェントを X , その相手を Y とする。対話 $d_k = [m_0, \dots, m_{k-1}]$ における X が持つ議論フレームワーク $AF_X^{d_k}$, X が持つ Y の予測議論フレームワーク $PAF_Y^{d_k}$ は以下になる。

1. $AF_X^{d_k}$
 - (a) $m_{k-1} = (Y, assert(A))$ の場合
 $AF_X^{d_k} = AF_X^{d_{k-1}} \circ A$
 - (b) それ以外の場合は $AF_X^{d_k} = AF_X^{d_{k-1}}$
2. $PAF_Y^{d_k}$
 - (a) $m_{k-1} = (X, assert(A))$, または,
 $m_{k-1} = (Y, assert(A))$ の場合
 $PAF_Y^{d_k} = PAF_Y^{d_{k-1}} \circ A$
 - (b) それ以外の場合は $PAF_Y^{d_k} = PAF_Y^{d_{k-1}}$

P, Q の ρ に関する対話 $d_k = [m_0, \dots, m_{k-1}]$ における可能手 m_k の条件は以下の通りである。

1. $m_0 = (P, assert(\rho))$ (対話の初手は必ず P が議題 ρ を述べる.)
2. $k \geq 1$ のとき
 - $assert(A)$
 - (a) $B \in CS_Y^{d_k} \wedge (A, B) \in AT_X$ (相手が発言した論証 B に対して、反論できる論証 A を述べる.)
 - (b) $m_k \neq m_i (0 \leq i \leq k-1)$ (同じ論証 A を発言することはできない.)

また、 $pass$ は初手を除いていつでも出すことができる。

対話 $d_k = [m_0, \dots, m_{k-2}, m_{k-1}]$ が $m_{k-2} = (X_{k-2}, pass), m_{k-1} = (X_{k-1}, pass)$ を満たす場合、その対話を終了対話とよぶ。対話はエージェントがお互いに何も言わず $pass$ を続けると終了する。

終了対話 d_k における、エージェント P, Q の議論フレームワークの基礎ラベリングをそれぞれ $L^{AF_P^{d_k}}, L^{AF_Q^{d_k}}$ とする。 $L^{AF_P^{d_k}}(\rho) = L^{AF_Q^{d_k}}(\rho)$ のとき終了対話 d_k は議題解決、 $L^{AF_P^{d_k}}(\rho) \neq L^{AF_Q^{d_k}}(\rho)$ のとき終了対話 d_k は議題不解決という。

終了対話 d_k における満足度 $I(d_k)$ を終了対話 d_k が議題解決の場合、

$I(d_k) = |in(L^{AF_P^{d_k}})| + |in(L^{AF_Q^{d_k}})|$, 終了対話 d_k が議題不解決の場合、 $I(d_k) = 0$ と定義する。

すべての可能な終了対話を D とする。終了対話 $d_k \in D$ において、満足度 $I(d_k)$ が最大となる対話 d_k を提案者 P と対立者 Q の議題 ρ に関する説得対話の最適解とよぶ。

4 戦略と実験

前節で提案した対話モデルはエージェントの価値観が異なるために対話の進行に伴う各議論フレームワークの更新も複雑で、議題解決するための戦略や議論フレームワークの条件を理論的に見出すのは大変困難である。本研究では対話モデルをシミュレートすることで、必要な戦略を見つける方法をとる。

まず、全体議論フレームワークの形や各議論フレームワークの関係を単純なものに限定してシミュレーションを行い、それにしたがうとどんな対話においても議題解決をし最適解を得られる戦略について考察する。

4.1 基本戦略

横浜の説得対話モデルでは、対話の目標は提案者が対立者に議題を受理させることであり、提案者 P は対立者 Q の議論フレームワーク内で議題が受理可能になるような手の出し方を戦略としてとっていた。一方、本研究では必ずしも議題が受理可能でなくてもよく、議題に対するラベルが一致していることを目的とする。また、 P と Q がいずれも同一な戦略を用いて議題解決し、最適解を求めることも目標とする。

そのためにまず、横浜の戦略をもとにした以下のような基本戦略を考える。ある時点において複数の論証 A が発言可能な場合、エージェントの手の選び方の優先順位を以下の順番とする。

基本戦略

1. 論証 A を発言した結果、自分の議論フレームワークと自分が持つ相手の予測議論フレームワークにおける議題のラベルが一致する。
2. 現在、自分の議論フレームワークと自分が持つ相手の予測議論フレームワークにおける議題のラベルが一致していれば何も述べず pass する。
3. 論証 A を発言する。
4. パスする。

4.2 テストデータの準備

UAF を全体議論フレームワーク、 UAF_P, UAF_Q をそれぞれ P, Q の全体議論フレームワークとする。 ρ を議題、 $AF_P = \langle AR_P, AT_P \rangle$ 、 $AF_Q = \langle AR_Q, AT_Q \rangle$ をそれぞれ初期状態で P, Q がもつ議論フレームワーク、 PAF_Q, PAF_P をそれぞれ P がもつ Q の予測議論フレームワーク、 Q がもつ P の予測議論フレームワークとする。これらの中で以下の条件を満たすものを予備テストデータとして考える。

1. $\rho \in AR_P \wedge \rho \in AR_Q$
2. $AF_P = UAF_P, AF_Q = half(UAF_Q)$
3. $L^{UAF}(\rho) = undec, L^{AF_P}(\rho) = in, L^{AF_Q}(\rho) = out$
4. UAF 内に双方向論証三角形が一つだけ存在する。

また、計算時間を考慮し、予備テストデータとして用いた UAF は論証の数が 10 個、攻撃関係の数が 11 個、かつ相互攻撃する論証のペアは双方向論証三角形の辺のみとする。

この条件を満たす議論フレームワークの組を予備テストデータとして対話を実行し、1000 個の予備テストデータに対してそれぞれ可能な全対話と満足度を求める。この結果から戦略の有効性を調べるために妥当な条件を満たすものを 300 個のテストデータとして取り出した。予備テストデータを妥当なテストデータとして使用するための条件は以下のものである。

1. 可能な全対話の中に議題解決するものと議題不解決なもの必ず一つ含まれる。
2. 可能な全対話のうち議題不解決する数が議題解決する数をこえる。

また、各テストデータに対して、満足度の最大値を求めてそれぞれのデータの最適解の満足度とする。

4.3 実験による戦略の考察

一組のテストデータに対し以下三つの条件についてそれぞれ対話を 200 回実行する。これを 300 組のテストデータに対して行う。

戦略あり P の予測知識あり エージェント P, Q の持つ予測知識は $PAF_Q = AF_Q, PAF_P = \emptyset$ で P, Q はともに戦略を用いる。

戦略あり P の予測知識なし エージェント P, Q が持つ予測知識は $PAF_Q = \emptyset, PAF_P = \emptyset$ で P, Q はともに戦略を用いる。

戦略なし エージェント P, Q は戦略を用いず、可能手を一様分布に従って出す。

その結果、複数の論証 A が発言可能なとき、手の選び方について以下の五つの戦略が新たに得られた。

A を発言した結果、自分が持つ相手の予測議論フレームワーク内において以下の条件を満たすときに論証 A を発言しない。

- (a) これまでに自分が発言した論証のラベルが "in" から "out" に変わる。
- (b) これまでに自分が発言した論証のラベルが "in" から "undec" に変わる。
- (c) A がこれまでに自分が発言した論証に対して攻撃をする。
- (d) 正論証三角形または逆論証三角形の辺が存在する。
- (e) A が攻撃する論証が存在しない。

これらは全ての価値観の相違に起因するものだが、(c),(e) は論証のラベルの変化がないため、この条件を見つけるのは困難だった。

以下で (a),(e) が得られた理由について例を用いて説明する。P の全体議論フレームワークを UAF_P 、Q の全体議論フレームワークを UAF_Q 、P のもつ自分の議論フレームワークを AF_P 、Q のもつ自分の議論フレームワークを AF_Q 、P のもつ Q の予測議論フレームワークを PAF_Q 、Q のもつ P の予測議論フレームワークを PAF_P とする。

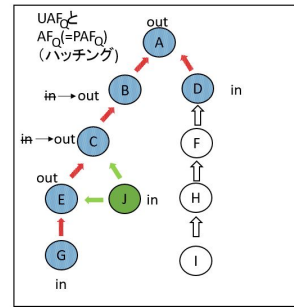


図 2: 条件 (a): 議題不解決に至る手

4.3.1 条件 (a)

図 1 のような P と Q の初期議論フレームワークが与えられるとする。図 1(1) は初期 UAF_P で、ハッチングされたノードとエッジはそれぞれ初期 AF_P に含まれるとする。図 1(2) は初期 UAF_Q で、ハッチングされたノードとエッジはそれぞれ初期 AF_Q に含まれるとする。初期 PAF_Q は、初期 AF_Q と一致しており、初期 PAF_P は \emptyset である。また、論証 A を議題とする。このとき、P,Q の初期 AF_P 、 AF_Q において、議題のラベルはそれぞれ "in", "out" となっている。

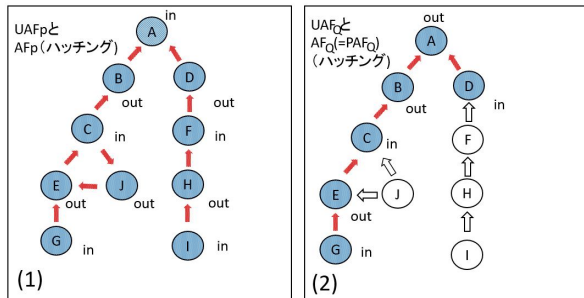


図 1: 条件 (a): 初期状態

対話 $d_4 = [(P, assert(A)), (Q, assert(B)), (P, assert(C)), (Q, assert(E))]$ において、P が手 $(P, assert(J))$ を出すと、その時点での PAF_Q では、P の発言した論証 C のラベルが "in" から "out" になり (図 2)、それ以降のように対話が進行しても議題不解決になる。したがってこのような手を出すのは避ける必要がある。

4.3.2 条件 (e)

図 3 のような P と Q の初期議論フレームワークが与えられるとする。図の説明は条件 (a) と同じであるので省略する。

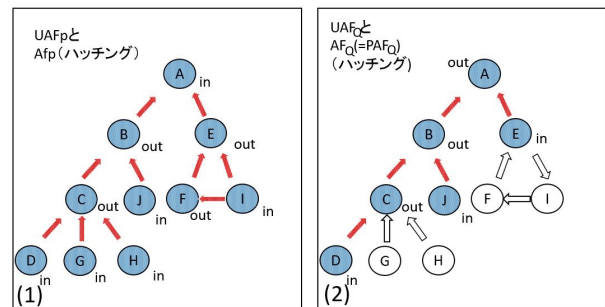


図 3: 条件 (e): 初期状態

対話 $d_2 = [(P, assert(A)), (Q, assert(C))]$ において、P が手 $(P, assert(F))$ を出すと、その時点での PAF_Q では、エージェント P やエージェント Q がどのような手を出そうと議題のラベルを "out" から "in" に変えることができない (図 4)。それ以降のように対話が進行しても議題不解決になる。したがってこのような手を出すのは避ける必要がある。

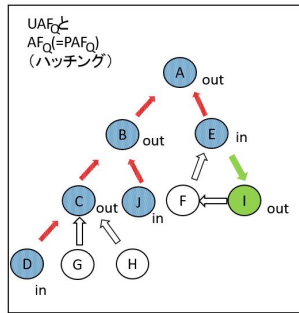


図 4: 条件 (e): 議題不解決に至る手

このようなことが起こるのはその時点での AF_P の中では論証 A は論証 E を攻撃しているが、その時点での AF_Q の中では論証 A は論証 E から攻撃を受けているためである。これは価値観の相違に起因するものである。

4.4 実験結果と評価

以上得られた戦略の有効性を評価するため、これらを基本戦略に加えて再度対話シミュレーションを行った。実行された 6 万回の対話において議題解決した回数、議題不解決した回数、対話の最適解が得られた回数を表 1 に示す。

表 1: 対話結果と回数 (回)

	議題解決	議題不解決	最適解が得られた対話
戦略あり (P の予測知識あり)	60000	0	51428
戦略あり (P の予測知識なし)	59975	25	51348
戦略なし	6171	53829	3744

実行された対話 6 万回のうち議題解決の割合は、戦略なしだと 12.2 %なのに対して、戦略ありだと、 P の予測知識がある場合は 100 %、 P の予測知識がない場合は 99.9 %である。したがって、得られた戦略は有効であるといえる。また、基本戦略として議題解決したあとでもできるだけ多くの論証を発言するようにすることで最適解を得る効果があると思われたが、「戦略あり (P の予測知識あり)」の条件のもとでもすべての場合で最適解を得ることはできなかった。

5 まとめ

本研究では、横浜の議論フレームワークに基づく対話モデルを、エージェントごとの価値観を導入したモデルに拡張し、このモデルに対する対話の評価基準を

設定した。そして、この基準に基づいて、議題が解決され、かつ、最適解が得られるような戦略を考察し、提案した戦略の有効性を示すために評価実験を行った。この結果、議論フレームワークのラベルの変化のような表層的なものでなく、理論的解析のみだと見つけにくい条件も、実験解析によって新たな戦略として得ることができた。特定の初期条件のもとで得られた戦略を用いると、全ての対話で議題解決をすることができたが、全ての対話での最適解を得ることはできなかった。

今後はこの特定の条件下で全ての対話で最適解を得る戦略について考察し、次に条件を変えた場合の戦略について考察する、さらにそれらの有効性を調べる評価実験を行う。

参考文献

- [1] Phan Minh Dung. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, Vol. 77, No. 2, pp. 321–357 (1995)
- [2] Leila Amgoud, Nicolas Maudet and Simon Parsons.: Modeling dialogues using argumentation. *ICMAS2000*, pp.31–38 (2000)
- [3] Shizuka Yokohama and Kazuko Takahashi.: What Should an Agent Know Not to Fail in Persuasion?. *EUMAS-AT2015*, pp. 219–233 (2015)
- [4] Trevor J.M. Bench-Capon, Sylvie Doutre and Paul E. Dunne.: Audiences in argumentation frameworks. *Artificial Intelligence*, Vol. 171, pp. 42–71 (2007)
- [5] Iyad Rahwan and Kate Larson.: Argumentation and game theory. *Argumentation in Artificial Intelligence*, pp. 321–340 (2009)
- [6] Pietoro Baroni, Martin Caminada and Massimiliano Giacomin.: An introduction to argumentation semantics. *The Knowledge Engineering Review*, Vol.26:4, pp.365–410 (2011)