

Q-learn Argumentation Schemes for Car Sales Dialogues

Adrian Groza

Department of Computer Science
Baritiu 28, RO-3400 Cluj-Napoca, Romania
adrian@cs-gw.utcluj.ro

Abstract

Agents need to argue with other agents many times, developing persuasion strategies that are effective over repeated situations. Applying reinforcement learning (RL) to the design of argumentation policies is appealing to dialogues where the counterpart can be modelled as a probability distribution. The idea of this research is to apply RL to speech acts in order to learn which discourse pattern is best to be conveyed during an argumentation game. Empowered by this learning mechanism, the persuasive agents gradually become more skillful through repeated argumentation.

1 Introduction

Over the past years, the research on argumentation theory has focused on identifying the most adequate technical instrumentation for modelling argumentation - defeasible logic seems to be the answer to these attempts. Recently, in the context of large scale argumentation of World Wide Argument Web (WWAW) [5], the interest has been shifted towards frameworks where all different inference mechanisms are able to co-exist under one umbrella, presently given by the Argument Interchange Format (AIF) ontology.

Learning and argumentation are related, sharing the same motivations of increasing efficiency or the range of solvable problems [4]. Agents have to adapt their persuasion policies depending on the failure or success of past dialogues, using a learning mechanism. The RL theory is quite advanced, but applications have been limited to problems of control, operations research, game playing. Two advantages of applying RL to argumentation sustain this research: i) the persuasive agent (lawyer, seller, teacher), can learn from its own experience; and ii) he can learn from simulation games [2].

2 Argumentation Framework

Extending the AIF ontology. We introduce a new node type, namely *context node* (*CO-node*) in the AIF ontology [5]. We argue the necessity of this node, due to the fact that context exists independent of any object in the system. Thus, one context may be used to evaluate different arguments, whilst the same argument can be evaluated in different contexts. Because of the separation of the argument structure from contexts, more power to re-use arguments and decide when to apply a specific argument in a given state is provided.

The extended-AIF ontology has five types of nodes: An *information node* (*I-node*) represents passive information of an argument, such as claim, premise, data, locution. A *scheme node* (*S-node*) captures active information or domain-independent patterns of reasoning. The schemes are split in three disjoint sets, whose elements are: rule of inference schemes (*RA-node*), conflict application node (*CA-node*), preference application node (*PA-node*). *Forms of arguments* (*F-nodes*) model argumentation schemes, by defining premises descriptor, conclusion descriptor, presumptions and exceptions. *Protocol interaction nodes* (*PIA-node*) are used to constrain the dialogue moves within an argumentation process. *CO-nodes* encapsulate all that is considered relevant to model the state of the argument, among the following contextual dimensions: *dialectical*, *intentional*, and *social*.

Argumentation schemes as protocol. Formally, an argumentation scheme (AS), encapsulated as a *F-node*, is composed of a set of premises A_i , a conclusion C , and a set of critical questions CQ_i , which enumerate specific ways to defeat the derivation of the consequent. One desiderata of the ASs is to simplify the argumentation process by hiding secondary premises and encapsulating them as *CQs* (figure 1). Based on the main premises A_1 , A_2 , and A_3 the consequent C is defeasibly inferred. During the process of gradually re-

<i>Argument from Analogy</i> \doteq <i>AS_AA</i>
$A_1 : X$ has properties P, Q, R, \dots
$A_2 : Y$ has properties P, Q, R, \dots
$A_3 : X$ has property T .
$C : Y$ has property T as well.
CQ_1^p : Does X have properties P, Q, R, \dots ? (evidence)
CQ_2^p : Are P, Q, R, \dots relevant to T ? (relevance)
CQ_3^e : Does Y have another property V which make it atypical, so that typical conclusion T does not follow? (distinguish)
CQ_4^e : Does X have missing features compared to Y that suggests that X is not typical? (exception)
CQ_5^p : Is a more typical case known? (increase confidence)
CQ_6^e : Is there another case W similar to X but in which T is false? (counterexample)

Figure 1. Simplifying argumentation.

vealing information in a debate, the conclusion might be retracted in the light of new information. Each *AS* sustaining a claim provides the correspondent *CQs* that the opponent may use to challenge the pleading. When a *CQ* is conveyed, the conclusion is suspended, until the subject of the dispute is clarified. Whoever is responsible for this clarification, in other words who has the burden of proof, depends on the type of the *CQ* [1].

Definition. A *presumption-type* CQ^p is a hidden premise encapsulated as a *CQ*, at the instantiation of which the burden of proof is shifted to the proponent of the argument. An *exception-type* CQ^e is a hidden premise encapsulated as a *CQ*, at the instantiation of which the burden of proof remains to the opponent.

Enacting the *Argument from Analogy* scheme as a *F – node* is depicted in figure 2. Based on the activation of a *CA – node*, the derivation of the conclusion is blocked until the subject is clarified. The activation of an exception leads to the instantiation of an *AS* that supports that exception. The state variables of the dialogue are encapsulated as *CO – nodes*.

The presumptions CQ_1^p, CQ_2^p, CQ_5^p request from the opponent more justification for supporting the claim, while the exceptions CQ_3^e, CQ_4^e, CQ_6^e challenge the argument by attacking the link between the premises and the conclusion. By instantiating the CQ_1^p , the opponent requests evidence that the case X is actually characterized by the mentioned properties P, Q, R . The burden of proof is shifted back to the proponent of argument who has to provide evidence supporting its claim. The CQ_2^p represents an undercutting defeater by attacking the link between the premises and the conclusion: the

premises are not relevant enough to infer the consequent. When the CQ_5^p is uttered, the opponent does not challenge the arguments provided by the proponent; it accepts them, but they are not considered sufficient for the consequent to be accepted. This *CQ* suggests to the proponent to engage in a persuasion dialogue, where the accrual of different arguments supporting the same conclusion might hold. In case of an exception-type *CQ*, the oponent identifies evidence that Y is an exception (CQ_3^e), X is an exception (CQ_4^e), or there have been similar situations where the same premises lead to the opposite conclusion (CQ_6^e), the opponent being the one responsible for providing data to distinguish between the case X and Y .

3 Learning from Experience

Learning how to argue is based on two ideas: i) if an argument in a given situation causes something bad to happen (at the end of the game the argument is not accepted by the opponent), the agent must learn not to convey that argument in a similar context; and ii) if all arguments in a given situation cause something bad to happen, the agent should avoid utterances that would lead the dialogue in that state. The *RL* model for argumentation consists of a set of environment states s_i , a set of actions a_i , and a set of rewards r_i .

The *states* of the dialogue game are defined as a collection of state variables, encapsulated as *CO – nodes*, which include all the information that is relevant in determining what *AS* to be conveyed next. The *actions* represent a finite number of speech acts. The agent learns which locution is best to be conveyed in a given context. The communication language L_c consists of a set of *ASs* used when arguing from experience: analogical argument, argument from example, argument from generalization, argument from correlation to cause, argument from common cause [6].

The *reward function* models the goal of the proponent, representing the cost of the argumentation game defined in terms of (i) the cost to obtain information O_I (the cost to populate *ASs* by querying the dataset); (ii) the cost of information disclosure R_I (cost of sharing its own experiences in order to alter the opponent’s cognitive state); and (iii) the solution cost (defined either in terms of time left and opportunity cost in case the persuasion fails, or in terms of reward or expected profit if the argument is accepted). According to the ϵ – *greedy learning policy*, the agent must decide between choosing an *AS* that is most likely to influence the dispute, according to past experience (*exploit step*), or to instantiate a scheme or a *CQ* that stimulates discussion and extracts information (*exploration step*).

so far: $R_1 = -0.1 - 1 = -1.1$. In the next step the seller conveys another argument $r_2^s(0.8, 0.3) : a, b \rightarrow acc$. The acceptance value being $A_2 = 0.5(0.8 - 0.6) + 0.5(0.3 - 0.2) = 0.15 > \tau$, the argument is not accepted. The resulting penalty is the same as in the previous step $R_2 = -1.1$. In the step 3, the seller keeps the same association rule as argument, but he decides to alter the cognitive state of the buyer by providing a past experience. If the buyer trusts the example, he incorporates it in his theory. The association rule r_2^b will be altered, resulting $r_2^b(0.69, 0.21) : a, b \rightarrow acc$. The accept value being $A_3 = 0.5(0.8 - 0.69) + 0.5(0.3 - 0.22) = 0.095 < \tau$, the buyer accepts the argument r_2^s . Thus, the received reward in this state is $R_3 = -0.3 + 10 = 9.7$. The reward represents the immediate value of an action in a given state. After a number of dialogue games, the long running value of each state will be estimated based on:

$$Q_n(s, a) = (1 - \alpha_n)Q_{n-1}(s, a) + \alpha_n[r + \gamma \max_{a'} Q_{n-1}(s', a')]$$

where $\alpha_n = 1/(1 + \text{visits}_n(s, a))$ and $\text{visits}(s, a)$ is the number of visits that the state-action pair has received so far. The learning rule considers a weighted average of the current Q value and the revised estimate in order to avoid the repeated alternations of the values of Q in the nondeterministic case². Thus, the revision of the Q value is made more smoothly, updates becoming smaller as training takes place [3].

CQs-based feedback. One advantage of *RL* is that it can accommodate different levels of feedback amount. For instance, the buyer can share its viewpoint in terms of its own certainty factor cf_o and coverage cov_o . If this is the case, it is much easier for the seller to decide which argumentation policy is better to follow. If the provided values are close to their own, he anticipates if sharing an example impacts on the cognitive state of the buyer at the required level. Otherwise, the seller may decide to change the argument type and to approach the buyer from a different perspective. The argumentation process is a trajectory in the chosen state space driven by the conveyed *CQs* and *ASs*. When planning for this trajectory, the first idea is to propose rules having the subject of debate as consequent of the rule. If this fails, the seller has to get away from the goal by starting to prove the premises of the above rules. The *RL* provides agents exactly with this mechanism: *learning when it is better to get away from the goal*. For instance, if the CQ_1^p is uttered, the seller can start to prove the requested properties or it can instantiate

²To identical pairs of state and arguments, buyers react in different ways.

another rule having the same initial consequent. The main drawback is given by the number of interactions needed for the agent to learn a persuasion strategy.

5 Related Work and Conclusion

Synergy of argumentation and machine learning is an open research issue. The work in [2] adopts *RL* to adjust negotiation strategies for a sales agent. The persuasion dialogue is based on Dung's abstract argumentation framework and the status of arguments according to the dialectical approach of Vreeswijk and Prakken. By using *ASs*, our work is more oriented towards the interaction with human agents. Combining argumentation and datamining has been investigated in [7]. Rising arguments is based on PADUA (Protocol for Argumentation Dialogue Using Association Rules), designed to empower agents with the ability to extract arguments from a dataset. In our approach, the argumentation protocol is directly encapsulated as *CQs*.

The current research³ can be integrated into the area of *action language perspective*, in the larger context of newly arising *Pragmatic Web* paradigm. If the *WWAW* is successful, we anticipate an explosion of public structured arguments, available on the Internet, that can be used in the learning process to better estimate the Q -value of each *AS*. The strategy will be shifted towards re-using arguments in different contexts and learning from others' experiences.

References

- [1] T. F. Gordon, H. Prakken, and D. Walton. The Carneades model of argument and burden of proof. *Artificial Intelligence*, 171(10-15):875–896, 2007.
- [2] S. Huang and F. Lin. The design and evaluation of an intelligent sales agent for online persuasion and negotiation. *Electron. Comm. Rec. Appl.*, 6(3):285–296, 2007.
- [3] T. M. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [4] S. Onta and E. Plaza. Learning and joint deliberation through argumentation in multiagent systems. In *AA-MAS '07*, pages 1–8, New York, NY, USA, 2007. ACM.
- [5] I. Rahwan, F. Zablith, and C. Reed. Laying the foundations for a World Wide Argument Web. *Artificial Intelligence*, 171(10-15):897–921, 2007.
- [6] C. Reed and G. Rowe. Araucaria: Software for argument analysis, diagramming and representation. *Intern. J. on Artif. Intell. Tools*, 13(4):961–979, 2004.
- [7] M. Wardeh, T. J. M. Bench-Capon, and F. Coenen. PADUA protocol: Strategies and tactics. In K. Mellouli, editor, *ECSQARU*, volume 4724 of *LNCS*, pages 465–476. Springer, 2007.

³This work was supported by the grant TD7 CNCSIS 534 from the National Research Council of the Romanian Ministry for Education and Research.