

# A Novel Approach to Evaluate Robustness of Incentive Mechanism Against Bounded Rationality (Doctoral Consortium)

Zehong Hu  
Rolls-Royce@NTU Corporate Lab, School of Computer Engineering  
Nanyang Technological University  
Singapore  
HUZE0004@e.ntu.edu.sg

## ABSTRACT

In this abstract, we first propose a general robustness definition as the upper-bound of the stable region of an equilibrium strategy by generalizing existing bounded rationality models. Then, we develop a robustness evaluation framework, of which a key component is the stability test given a certain level of bounded rationality.

## Keywords

Mechanism Design; Bounded Rationality; Robustness

## 1. INTRODUCTION

Incentive mechanisms are designed to produce desired outcomes through incentivizing agents to perform expected strategies. Incentive is created by inducing a Bayesian game and having the desired outputs as its Nash equilibrium. However, in real-life situations, agents often cannot behave fully rationally, causing mechanisms to fail. To ensure the practical usability of incentive mechanisms, it is crucial to evaluate their robustness. Quite a few qualitative studies have been conducted to judge whether a mechanism is robust. For example, Cabrales [1] and Tumennasan [5] considered robustness from testing the existence of convergent trajectories in all *ex-post* games under the better response learning process and the limited logit quantal response process, respectively.

However, one common limitation of these studies is that they cannot quantitatively evaluate robustness as to what extent an incentive mechanism can be kept under the effects of bounded rationality [3]. Furthermore, most of these studies do not consider precise bounded rationality models, and thus, are not applicable to real-world scenarios where a certain kind of bounded rationality dominates [2]. Besides, the evolutionary processes involved in these studies can only handle the *ex-post* Nash equilibrium with fixed agent's types, while mechanism design is based on Bayesian Nash equilibrium with all agents' types unknown, i.e. the *ex-ante* equilibrium. Thus, we aim to providing a general framework to quantitatively evaluate the robustness of incentive mechanisms against different kinds of bounded rationality.

**Appears in:** *Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2016)*, J. Thangarajah, K. Tuyls, C. Jonker, S. Marsella (eds.), May 9–13, 2016, Singapore.  
Copyright © 2016, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

## 2. ROBUSTNESS DEFINITION

To quantitatively evaluate robustness, we need to define robustness at first. However, the literature closely linked to this topic is limited. Thus, we extend the targeted domain of our definition to common Nash equilibrium in arbitrary kinds of games, and enforce it to be consistent with existing ones for correctness. Bounded rationality causes agents to take sub-optimal strategies of the optimal strategy. To depict the strategy variations, the uncertainty set  $\mathcal{US}$  is employed to generalize existing bounded rationality models as

$$\mathcal{US}(s_i^*, s_{-i}, \alpha, G) = \{s_i^1, \dots, s_i^{I_i}\}$$

where  $G$  denotes the targeted game and  $s_{-i}$  represents the strategy of other agents.  $s_i^j$  stands for agent's one strategy which is a probability distribution over all possible actions.  $I_i$  represents the total number of possible strategies. Besides,  $\alpha$  is a parameter to describe agents' bounded rationality level. To keep consistent among all models, we require  $\alpha \in [0, 1]$ , where  $\alpha = 0$  and  $\alpha = 1$  denote the cases where agents are fully rational and irrational, respectively. Then, for a game with  $N$  players, the uncertainty set involves more than one agent and can be formulated as:

$$\mathcal{US}_G(\mathbf{s}^*, \alpha_G, G) = \{\mathbf{s} = (s_1, \dots, s_N) \mid s_i \in \mathcal{US}_i(Pr(s_i^*, s_{-i}^*, \alpha_i, G), \alpha_i \in \Pi(\alpha_G))\} \quad (1)$$

where  $\alpha = (\alpha_1, \dots, \alpha_N)$  represents the bounded rationality level profile. The parameter  $\alpha_G \in [0, 1]$ , termed as system bounded rationality level, is used to measure agents' bounded rationality from system perspective, and it is required that  $\mathcal{US}_G(\mathbf{s}^*, 0, G) = \{\mathbf{s}^*\}$ .  $\Pi(\alpha_G)$  denotes all possible bounded rationality level profiles.

The most classical way to define robustness is to measure to what extent parameter variation can make systems unstable. Following this idea, we give a general definition of robustness for Nash equilibrium:

*Definition 1.* Given a game  $G$  with a desired equilibrium  $\mathbf{s}^E$  and the uncertainty set  $\mathcal{US}_G(\mathbf{s}^*, \alpha_G, G)$  denoting agents' bounded rationality, the robustness of  $\mathbf{s}^E$  is  $R$  such that

$$R = \frac{Dist(\mathcal{US}_G(\mathbf{s}^E, 0, G), \mathcal{US}_G(\mathbf{s}^E, \alpha_G^M, G))}{Dist(\mathcal{US}_G(\mathbf{s}^E, 0, G), \mathcal{US}_G(\mathbf{s}^E, 1, G))} \quad (2)$$

where  $\alpha_G^M$  is determined by the stability function  $\mathcal{S}$  as

$$\alpha_G^M = \max\{\alpha_G \in [0, 1] \mid \forall \alpha \in [0, \alpha_G], \mathcal{S}(\alpha, G, \mathbf{s}^E, \mathcal{US}_G) \leq 0\} \quad (3)$$

and the distance function  $Dist(\cdot)$  is defined as

$$Dist(\mathcal{US}_G(\mathbf{s}^E, p, G), \mathcal{US}_G(\mathbf{s}^E, q, G)) = \max_{\mathbf{s}^p \in \mathcal{US}_G(p)} \max_{\mathbf{s}^q \in \mathcal{US}_G(q)} \max_{i \in N} \frac{\|s_i^p - s_i^q\|_1}{2} \quad (4)$$

wherein  $\|\cdot\|_1$  denotes the 1-norm function.

In other words, the robustness of a desired Nash equilibrium corresponds to the maximum system bounded rationality level, at which the employed stability function  $\mathcal{S}$  can keep non-positive and the game remains stable. This definition provides general guidelines to instantiate more concrete robustness formulations for specific applications.

### 3. ROBUSTNESS COMPUTATION

An incentive mechanism  $\mathcal{M}$  is defined by  $\{\mathbf{x}(\mathbf{a}), \mathbf{t}(\mathbf{a})\}$ , where  $\mathbf{x}$  and  $\mathbf{t}$  are the allocation vector and the payoff function, respectively. The Bayesian game induced by  $\mathcal{M}$  can be defined as  $G = \langle N, \{u_i, \Theta_i, A_i, f_i\}_{i \in N} \rangle$ . Here,  $N$  denotes the set of  $n$  agents. Each agent draws its type  $\theta_i \in \Theta_i$  independently from a commonly known distribution over  $\Theta_i$  with density distribution  $f_i$ . The utility of agent  $i$  can be calculated as  $u_i = v_i(\mathbf{x}(\mathbf{a}), \theta_i) - t_i(\mathbf{a})$ , where  $v_i$  represents of the value obtained from the allocation  $\mathbf{x}$ . According to Definition 1, given a certain kind of bounded rationality depicted in the form of Equation (1), the remaining two tasks for computing robustness are to: 1) test the stability of the desired equilibrium with a concrete stability function  $\mathcal{S}$ ; 2) solve the maximum value of  $\alpha_G$  according to Equation (3). Focusing on the Bayesian game induced by incentive mechanisms, we develop a general robustness evaluation framework:

1) **Robustness Solver** iteratively searches for  $\alpha_M^G$  in  $[0, 1]$ . In order to achieve an efficient search, Newton's Dichotomy method is used to skip some useless regions. Specifically, the inputs of our robustness solver include Bayesian game  $G$ , the desired action distribution profile  $\mathbf{s}^E$  and the uncertainty set  $\mathcal{US}_G$ . At first, the left boundary  $\alpha_L$  and right boundary  $\alpha_R$  are set to be 0 and 1. The fully irrational case with  $\alpha = 1$  is tested. If the test result is stable, we can conclude that the equilibrium is stable for all bounded rationality levels; otherwise, the real boundary should be some value between  $\alpha_L$  and  $\alpha_R$ . Then, we repeatedly conduct dichotomy and test the middle value  $(\alpha_L + \alpha_R)/2$  until the distance between two boundaries is smaller than the acceptable threshold.

2) **Stability Tester** tests the stability of the desired equilibrium at the given level of bounded rationality. To achieve this function, we need to develop a convincing evolutionary process. Those utilized in the literature are not applicable to the *ex-ante* equilibrium. A good alternative is the fictitious play (FP) process which was initially proposed to compute equilibrium in normal-form games. Recently, through introducing new methods to compute the best response, Rabinovich *et al.* [4] used the FP process to compute the pure-strategy Bayesian Nash equilibrium. Although convergence of the FP process cannot be assured for all Bayesian games, it can still be used for our stability test because the divergent case can be used to identify the unstable region. Specifically, our algorithm computes the uncertainty set at first, which acts as the initial belief set of the FP process. Then, the FP process learns the Bayesian Nash equilibrium through gradually updating elements in the belief set with the the best response strategy  $s_i^*(a_i)$ . The stability function  $\mathcal{S}$  is  $Dist(\mathcal{US}_t, \{\mathbf{s}^E\}) - \epsilon$ , where  $\epsilon$  is the threshold of distance.

3) **Best Response Computing** is the most challenging step in stability tester. Existing methods depend on a specific target domain. To achieve an efficient and general evaluation of different mechanisms, we propose a sampling-based algorithm. Since  $\sum_{a_i} \mathbb{E}_{s_{-i}}(a_i|\theta_i) Pr_i(a_i|\theta_i) \leq \mathbb{E}_{s_{-i}}(a_i^*|\theta_i)$ , the best response for a specific  $\theta_i$  value should be an action  $a_i^*$  with the highest expected utility. Thus, the best response computation depends on the comparison between the expected utilities of different actions as  $\mathbb{E}_{s_{-i}}[u_i(a_i^l|\theta_i) - u_i(a_i^k|\theta_i)] \geq \tau \epsilon(l, k)$ , where  $l, k = 1, \dots, |A_i|$ , and  $|A_i|$  denotes the number of elements in  $A_i$ . Since the accurate comparison is impossible for sampling-based methods, the confidence level  $\tau$  is introduced, and  $a \geq_\tau b$  represents that  $a$  is bigger than or equal to  $b$  with a probability of  $\tau$ . Thus, if  $a_i^*$  is the desired best response with probability  $\tau$ , then  $\epsilon(l^*, k) \geq 0$  for  $k = 1, \dots, |A_i|$ . Furthermore, to compute  $\epsilon(l, k)$ , bootstrap statistics is employed because it does not rely on explicit assumptions about the shape of a distribution. This feature is useful because nothing is known about the distribution of agent's utility. Specifically, in our algorithm,  $K$  samples of opponents' actions are generated at first and the utility for every action  $a_i^k \in A_i$  is calculated. Then, the distribution of expectation is computed with bootstrap statistics and the lower bound  $\epsilon(l, k)$  is obtained. For those sub-optimal actions, there must exist at least one action that can achieve better utility, and  $\{k | \max_{j \in \{1, \dots, |A_i|\}} \epsilon(j, k) \leq 0\}$  represents the optimal set. If the size of  $I$  is 1, the optimal action is found; otherwise, we need to add another  $K$  samples.

### 4. CONCLUSION AND FUTURE WORK

In this abstract, we proposed a novel framework for evaluating robustness of incentive mechanisms. For future work, we will extend our framework to cover more kinds of bounded rationality and conduct more experiments on various incentive mechanisms to validate the framework.

### REFERENCES

- [1] A. Cabrales and R. Serrano. Implementation in adaptive better-response dynamics: Towards a general theory of bounded rationality in mechanisms. *Games and Economic Behavior*, 73(2):360–374, 2011.
- [2] L. Chen and P. Tang. Bounded rationality of restricted turing machines. In *Proceedings of International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 1673–1674, 2015.
- [3] Y. Liu, J. Zhang, B. An, and S. Sen. A simulation framework for measuring robustness of incentive mechanisms and its implementation in reputation systems. *Autonomous Agents and Multi-Agent Systems*, pages 1–20, 2015.
- [4] Z. Rabinovich, V. Naroditskiy, E. H. Gerding, and N. R. Jennings. Computing pure bayesian-nash equilibria in games with finite actions and continuous types. *Artificial Intelligence*, 195:106–139, 2013.
- [5] N. Tumennasan. To err is human: Implementation in quantal response equilibria. *Games and Economic Behavior*, 77(1):138 – 152, 2013.