

Effect of Adaptation Rate and Cost Display in a Human-Al Interaction Game

STUDENTS: Jason T. Isa, Bohan Wu, Qirui Wanq, Yilin Zhang

Introduction

- As interactions between humans and AI become more prevalent, it is critical to have better predictors of human behavior in these interactions to create more safe and socially beneficial outcomes.
- Previous work has shown that AI can unilaterally change their adaptive algorithm to drive the outcome of a simple Human-AI game to distinctly different game theoretic equilibria [1].

Problem

 Without changing the Al agent's adaptive algorithm, can we influence the Human-Al game outcome by changing the Human's feedback information and model these behavior changes within a game theoretic model?



Game Theoretic Equilibria

Nash Equilibrium

$h^{\mathbb{N}} = rg\min_{h\in\mathcal{H}} c_H(h, m^{\mathbb{N}}),$ $m^{\mathtt{N}} = rg\min_{m \in \mathcal{M}} c_M(h^{\mathtt{N}}, m)$

Stackelberg Equilibrium

 $h^{\mathsf{S}} = \arg\min_{h \in \mathcal{H}} \{ c_H(h, m^{\mathsf{S}}) \mid m^{\mathsf{S}} = \arg\min_{m} c_M(h, m) \},\$ $m^{\mathsf{S}} = \arg\min_{m \in \mathcal{M}} \{c_M(h^{\mathsf{S}}, m)\}.$

Experimental Algorithm

Algorithm 1 Experiments with human subjects

Require: initial $(h_{(0,:)}, m_{(0,:)})$, parameter $\gamma > 0$ for $t = 0, 1, \dots, T - 1$ do switch Experiment do case 1 display $c_H(h_{(t,:)}, m_{(t,:)})$ case 2

 $h_{(t+1,:)} = \texttt{manual}_i\texttt{nput}(t)$ $m_{(t+1,:)} = m_t - lpha rac{\partial c_M}{\partial m}(h_{(t,:)},m_{(t,:)})$ end for

Where both Human and AI cost functions were defined as quadratic functions, $c_{H}(h,m) = rac{1}{2}h^{ op}A_{H}h + h^{ op}B_{H}m + rac{1}{2}m^{ op}D_{H}m + h^{ op}a_{H} + m^{ op}b_{H}$ $c_M(h,m) = rac{1}{2}m^ op A_Mm + m^ op B_Mh + rac{1}{2}h^ op D_Mh + m^ op a_M + h^ op b_M$



ADVISERS: Samuel A. Burden, Lillian J. Ratliff **COLLABORATOR:** Benjamin J. Chasnov

Empirical Results

Questions?

email: jisa@uw.edu