

Thesis Project Portfolio

MARL In Persuasion Games

(Technical Report)

Combating Online Misinformation Through Game-Theoretic Incentives

(STS Research Paper)

An Undergraduate Thesis

Presented to the Faculty of the School of Engineering and Applied Science

University of Virginia • Charlottesville, Virginia

In Fulfillment of the Requirements for the Degree

Bachelor of Science, School of Engineering

Gustavo Moreira

Spring, 2022

Department of Computer Science

Table of Contents

Sociotechnical Synthesis

MARL In Persuasion Games

Combating Online Misinformation Through Game-Theoretic Incentives

Prospectus

Sociotechnical Synthesis

In the study of complex multi-agent environments, understanding what incentive structures result in cooperation is crucial to effectively encouraging cooperation and preventing undue competition. Said understanding would be particularly useful in contexts in which large groups of agents act largely independently from each other and in such a way where their actions cannot be controlled, but rather indirectly influenced. One such real-world scenario is social media communication online, as users act largely outside of the control of the platform holders and various previous attempts to control their behavior have been ineffective. The study of this indirect influence in the language of mathematics and economics is Bayesian persuasion, and developing optimal strategies under a Bayesian persuasion setting is notoriously difficult. If one were able to develop a learning algorithm for the Bayesian persuasion setting, one could apply insights from it in the context of social media in order to incentivize certain user behaviors

My technical report covers the application of multi-agent reinforcement learning algorithms to the Bayesian persuasion setting. To be more formally concrete, in the Bayesian persuasion setting, one agent (the “Sender”) sends signals to another agent (the “Receiver”) and each agent’s utility is determined by the actions of the Receiver and some “state of nature” visible only to the Sender, usually some random value drawn from a set distribution. Thus, the Receiver is “warily trusting” of the Sender to send true information based on the state of nature that the Sender observes, and the Sender in turn can decide whether to send a true or false signal based on that observation. The conclusion of our findings is that we can construct a state of nature distribution such that agents using the EXP3-DH learning algorithm will fall into an equilibrium in which the Sender will always truthfully report to the Receiver, regardless of the utility values each agent receives for each action-state pair.

The STS research paper discusses the application of these insights to the context of social media, with preventing misinformation as the primary target. The paper begins with an overview of the existing social media landscape and the various attempts that have been previously made to prevent the spread of misinformation on platforms like Twitter, and why such measures have failed. I then go through various “intuitive” alternative proposals for ways to incentivize the spreading of truthful information. Finally, I construct a proposal for social media platforms under which users are incentivized to post truthful information (either through monetary means or on-platform content rewards). These incentives would be limited to a certain tier of users that are initially selected manually by the platform holder, and afterwards would select themselves in order to avoid potential abuse of this system to continue posting misinformation while still getting the incentives

While the application of the results from the technical report are here limited to the domain of social media misinformation, the potential applications are massive in scope, precisely because they apply to any situation in which rational agents must indirectly influence each others choices via aligned incentives. Despite this, social media misinformation is a notable real-world problem that has only been growing in severity over the past few years. If the changes suggested by the STS research paper were implemented into sites like Twitter or Facebook, the technical research suggests that the spread of misinformation would be significantly reduced, clearly demonstrating the research’s real-world value.\