

Introduction

Our work focuses on attempting to comprehend the complexity in human decisions, particularly in the context of strategic persuasion in personalized recommender systems. We have developed a mathematical framework to model strategic interaction between a personalized recommender system and a human user as a Stackelberg signaling game. We assume that the rationality of the user is based on expected utility maximization. Our goal is to compute the equilibrium strategies at both the recommender system and the user, and investigate conditions under which the recommender system reveals manipulated information, and user trust toward the recommender system deteriorates.

Problem Setup

We considered a human-AI interaction as shown in the following graphic, where Bob is presented with a set of choices $\mathcal{N} = \{1, \dots, N\}$, which is also known to Alice. As is the case with typical recommender systems, we assumed that Alice has access to extrinsic information, motivating Bob to rely on Alice's messages. Using Alice's belief and Bob's prior belief, Bob's posterior belief is constructed as

$$\phi(x) = \alpha \pi_p(x) + (1 - \alpha)q(x).$$

Letting $\psi = \{\psi_1, \dots, \psi_N\} \in S_N$ denote the probabilistic decision rule employed by Bob where S_N is the probability simplex on the choice set N and is the probability of picking the n^{th} choice based on Bob's ex-post belief, Alice's average reward utility is

$$U_A(\pi_p, \psi) = \sum_{n=1}^N \psi_n \cdot \mathbb{E}_p(X_n),$$

And Bob's ex-post utility is

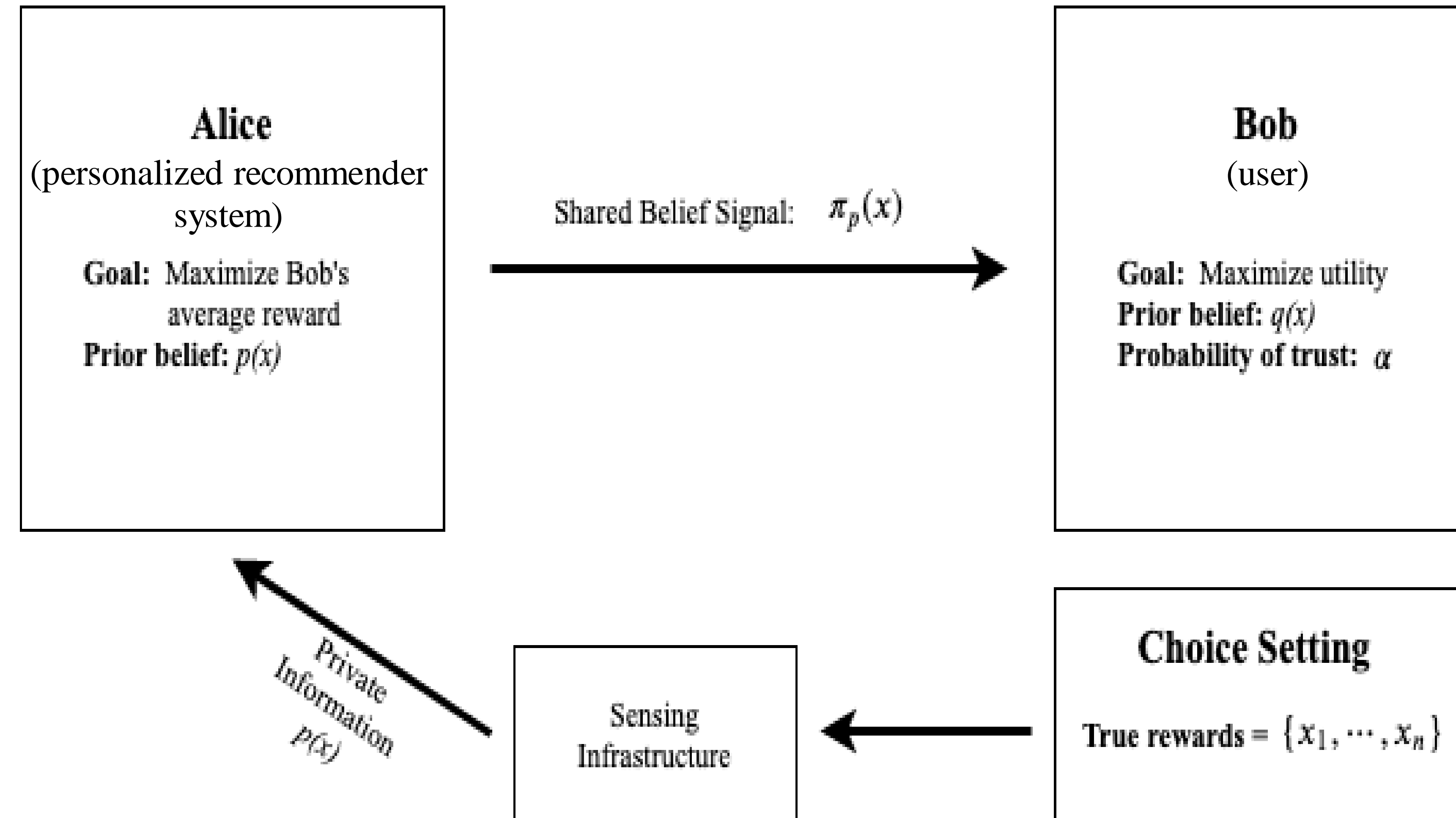
$$U_B(\pi_p, \psi) = \sum_{n=1}^N \psi_n \cdot \mathbb{E}_\phi(X_n),$$

We then modeled the strategic interaction between Alice and Bob as a Stackelberg game with Alice as the leader and Bob as the follower

$$\psi^*(\pi_p) \triangleq \arg \max_{\psi} U_B(\pi_p, \psi), \text{ and}$$

$$\pi_p^* \triangleq \arg \max_{\pi_p} U_A(\pi_p, \psi^*(\pi_p)).$$

Problem Setup



Equilibrium Analysis

Bob's Best Response:

Given that Alice chooses a signalling strategy π_p , Bob's best response is to choose $\psi = \{\psi_1, \dots, \psi_N\}$ such that Bob's expected utility is

$$U_B(\pi_p, \psi) = \sum_{n=1}^N \psi_n y_n$$

where $y_n = \alpha \mathbb{E}_{\pi_p}(X_n) + (1 - \alpha) \mathbb{E}_q(X_n)$.

Theorem 1. For a given trust parameter α , signaling strategy $\pi_p(x)$ and prior belief $q(x)$, Bob's best response is given by

$$\psi_n^*(\pi_p) = \begin{cases} 1, & \text{if } n = \arg \max_{n \in N} y_n, \\ 0, & \text{otherwise,} \end{cases}$$

Optimal Signaling at Alice:

Alice's optimal signal strategy is to choose the maximum entry ψ_n in the vector $\psi = \{\psi_1, \dots, \psi_N\}$ such that

$$U_A(\pi_p, \psi) = \sum_{n=1}^N \psi_n \cdot \mathbb{E}_p(X_n)$$

is maximized.

Theorem 2. The optimal signaling strategy at Alice is to choose a distribution π_p such that $\alpha \mathbb{E}_{\pi_p}(X_{n^*}) + (1 - \alpha) \mathbb{E}_q(X_{n^*}) \geq \alpha \mathbb{E}_{\pi_p}(X_n) + (1 - \alpha) \mathbb{E}_q(X_n)$ holds true for all $n \in N$ where $n^* = \arg \max \mathbb{E}_{\pi_p}(X)$.

Trust Analysis

We defined **strategic manipulation** as follows: Alice employs strategic manipulation if she chooses $\mathbb{E}_{\pi_p}(X) \neq \mathbb{E}_p(X)$. We found the following to hold true:

Corollary 2. Alice adopts strategic manipulation if there exists at least one $n \in N$ such that $\alpha \mathbb{E}_p(X_{n^*}) + (1 - \alpha) \mathbb{E}_q(X_{n^*}) < \alpha \mathbb{E}_p(X_n) + (1 - \alpha) \mathbb{E}_q(X_n)$ where $n^* = \arg \max \mathbb{E}_{\pi_p}(X)$.

Corollary 3. When $\alpha = 1$, Alice has no incentive to share manipulated information to Bob.

We denoted Bob's trust as $\alpha^+ = \alpha + \epsilon \cdot R_B(\alpha)$ where $\epsilon > 0$ is the default step size and Bob's regret is

$R_B = \alpha \mathbb{E}_{\pi_p}(X_{n^*}) + (1 - \alpha) \mathbb{E}_q(X_{n^*}) - \mathbb{E}_q(X_{n^*})$ and $n^* = \arg \max_{n \in N} \alpha \mathbb{E}_{\pi_p}(X) + (1 - \alpha) \mathbb{E}_q(X)$.

Based on these results, we made the following finding:

Bob's trust deteriorates even though Alice reveals truthful signals, whenever

$$\alpha < \frac{\mathbb{E}_q(X_{n^*}) - \mathbb{E}_q(X_{n^*})}{\mathbb{E}_p(X_{n^*}) - \mathbb{E}_q(X_{n^*})}$$

Conclusion

The findings presented in this project regarding the evolution of trust between the recommender system and the user largely coincide with the behavior seen in the real-world interactions of humans with recommender systems. By analyzing the dynamics of Bob's trust, we found that Alice will employ strategic manipulation if her signaling strategy does not match her posterior belief, but Alice has no incentive to employ strategic manipulation if Bob does not distrust her. Further, we found that Bob's trust may deteriorate even when Alice reveals information truthfully, if Bob does not obtain a desired outcome after interacting with Alice.

We plan to extend this work by validating the results of this project, modeling this interaction as a repeated game, and considering the evolution of Bob's trust towards Alice when Bob's rationality is characterized by human decision models as opposed to expected utility maximization.