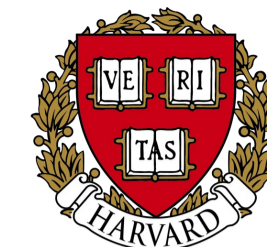




Persuading a Learning Agent

Generalized Principal-Agent Problem with a Learning Agent

Tao Lin,
Yiling Chen



Generalized Principal-Agent Problem

A general model for principal-agent interactions [1][2]. It includes:

- Contract design
- Bimatrix Stackelberg games
- Bayesian persuasion
- ...

Formally,

- The principal has a convex decision space \mathcal{X}
- The agent has a finite action set A
- Principal's utility $u: \mathcal{X} \times A \rightarrow \mathbb{R}$
- Agent's utility $v: \mathcal{X} \times A \rightarrow \mathbb{R}$
- u, v are assumed to be linear in $x \in \mathcal{X}$
- Additionally, there is a finite set of signals/messages S

Timeline:

- The principal commits to a strategy $\pi = (p_s, x_s)_{s \in S}$, which consists of:
 - $(p_s)_{s \in S}$: a distribution over S
 - x_s : a decision associated with signal s
- The agent chooses a strategy $\rho: S \rightarrow A$
 - Best response:

$$\rho(s) \in \operatorname{argmax}_{a \in A} v(x_s, a)$$

- The two players obtain (expected) utility

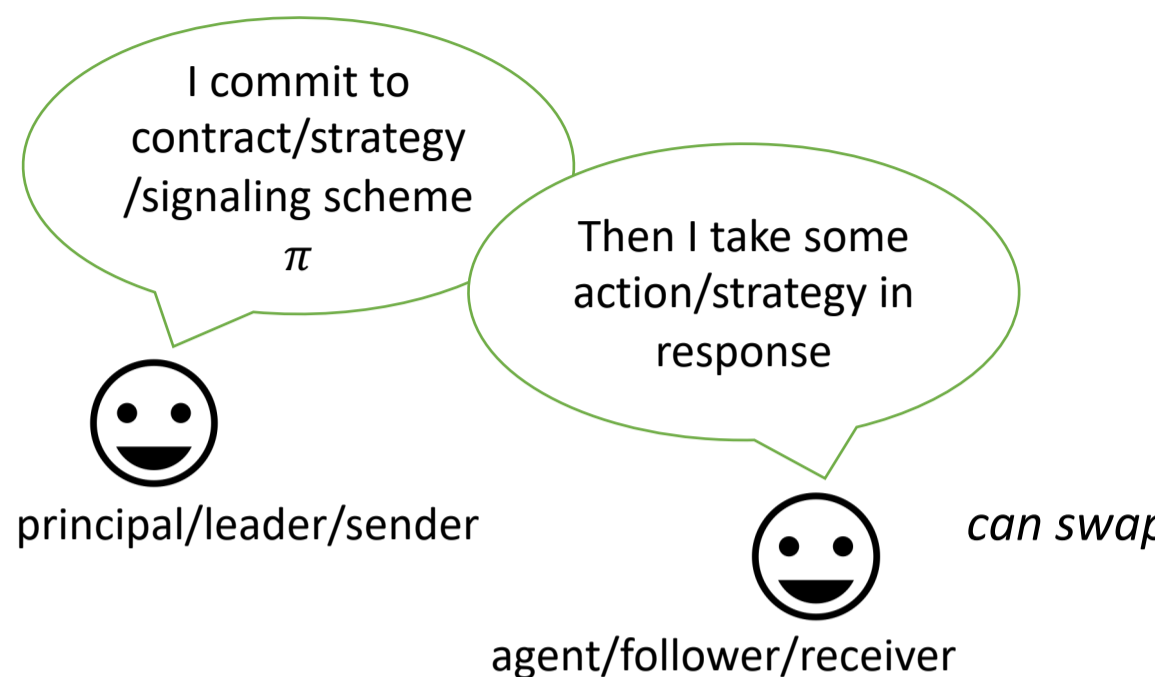
$$\mathbb{E}_{s \sim p} [u(x_s, \rho(s))]$$

$$\mathbb{E}_{s \sim p} [v(x_s, \rho(s))]$$

[1] Myerson (1982): Optimal Coordination Mechanism in Generalized Principal-Agent Problems.

[2] Gan, Han, Wu, Xu (2024): Generalized Principal-Agency: Contracts, Information, Games and Beyond.

[3] Deng, Schneider, Sivan (2019): Strategizing Against No-Regret Learners.



- In contract design, $x = (p_1, \dots, p_n)$ is a payment vector
- In bimatrix Stackelberg game, $x \in \mathcal{X} = \Delta([n])$ is a mixed strategy
- In Bayesian persuasion, $x \in \mathcal{X} = \Delta(\Omega)$ is a posterior belief

Additionally, π satisfies a constraint:

$$\sum_{s \in S} p_s x_s \in \text{convex set } \mathcal{C} \subseteq \mathcal{X}.$$

In Bayesian persuasion, $\mathcal{C} = \{\mu_0\}$ where μ_0 is the prior for the state of the world.

Learning Agent

Instead of best-responding, we consider an agent who *learns* which action to take for each signal.

T rounds of interactions. In each round t ,

- The principal chooses a strategy $\pi^t = (p_s^t, x_s^t)_{s \in S}$, can be *unknown* to the agent.
- Based on history, the agent chooses a (randomized) strategy $\rho^t: S \rightarrow \Delta(A)$.
- A signal $s^t \sim p^t$ is realized, the principal makes decision $x^t = x_{s^t}^t$, the agent samples action $a^t \sim \rho^t(s^t)$. The two players obtain (expected) utility $\mathbb{E}[u(x^t, a^t)]$ and $\mathbb{E}[v(x^t, a^t)]$

Definition: The agent's learning algorithm satisfies

- **No-contextual-regret** if

$$\forall d: S \rightarrow A, \quad \mathbb{E} \left[\sum_{t=1}^T (v(x^t, d(s^t)) - v(x^t, a^t)) \right] \leq \text{CReg}(T) = o(T).$$

- **No-contextual-swap-regret** if

$$\forall d: S \times A \rightarrow A, \quad \mathbb{E} \left[\sum_{t=1}^T (v(x^t, d(s^t, a^t)) - v(x^t, a^t)) \right] \leq \text{CSReg}(T) = o(T).$$

Main Results: Under some regularity conditions,

- Against a **no-contextual-regret** learning agent, the principal can obtain average utility at least $U^* - O\left(\sqrt{\frac{\text{CReg}(T)}{T}}\right)$; U^* is the principal's optimal utility against a best-responding agent.
- Against a **no-contextual-swap-regret** learning agent, the principal *cannot* obtain more than $U^* + O\left(\frac{\text{CSReg}(T)}{T}\right)$ (even knowing the agent's learning algorithm and using adaptive strategies).
- For *some* **no-contextual-regret** agent, the principal can obtain more than $U^* + O\left(\frac{\text{CReg}(T)}{T}\right)$.

Intuition: no contextual **swap**-regret learning \approx *approximately* best responding \approx best responding