

# Adverse and Advantageous Selection in the Lab

with Max Mihm, Lucas Siga, and Chloe Tergiman



## motivating questions

*How do people behave when they know they are asymmetrically informed?*

*Do people account for selection better in some settings than in others?*

### Theory

- Simple class of games that allows for adverse and advantageous selection.
- Will use predictions of Bayes-Nash, cursed, and self-confirming equilibrium.

### Experiment

- Compare asymmetric info with otherwise identical symmetric information games.
- Study effects of removing strategic uncertainty & the role of payoff feedback.

theory: games of **adverse** / **advantageous** selection

## the game: uncertainty and information

Alice and Bob collectively choose between a safe ( $S$ ) and risky ( $R$ ) option.

- **Safe:** Each player obtains  $s$ .
- **Risky:** Each player obtains  $\ell$  with prob  $\frac{1}{2}$  and  $h$  with prob  $\frac{1}{2}$ , where  $0 < \ell < s < h$ .

We vary whether the **risky option** is **positively** or **negatively** correlated.

- Payoffs are in  $(\pi_A, \pi_B)$ , where  $\pi_i$  specifies  $i$ 's payoffs.
- **Positive correlation:** Risky option pays  $(\ell, \ell)$  if **Heads**,  $(h, h)$  if **Tails**.
- **Negative correlation:** Risky option pays  $(\ell, h)$  if **Heads**,  $(h, \ell)$  if **Tails**.

Alice and Bob both know the correlation of the risky option.

**Asymmetric Info:** Alice learns outcome of coin toss & Bob doesn't. This is CK.

## the game: choices and payoffs

Players vote simultaneously for the safe or risky option.

If **both** vote for risky option  $\implies$  Risky option is selected.

If **at least** one votes for safe option  $\implies$  Safe option is selected.

Each player has selfish preferences, with a utility that is increasing in own wealth.

What would “standard game theory” predict?

## solving the game via iterated weak dominance

Recall that Alice observes  $(\pi_A, \pi_B)$  and Bob does not.

Alice's weakly dominant strategy: vote for  $R$  if  $\pi_A = h$ , and for  $S$  if  $\pi_A = \ell$ .

Bob's vote for  $R$  matters only if Alice is also doing so.

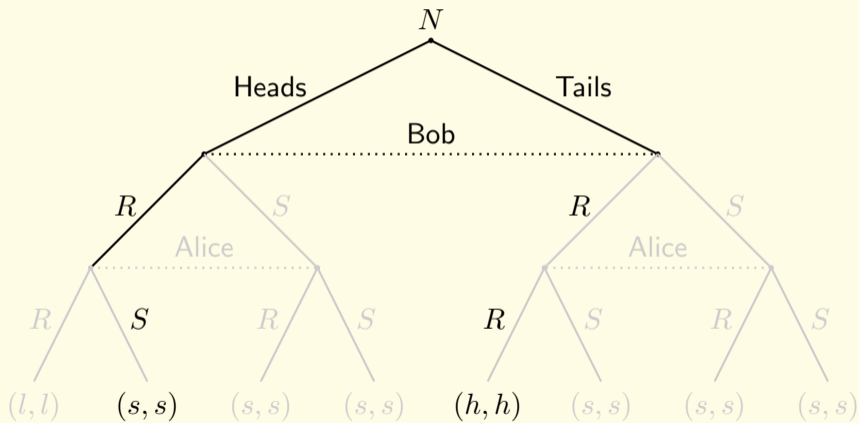
- **Positive correlation:**  $\pi_A = h \implies \pi_B = h \implies$  Bob should vote for  $R$ .
- **Negative correlation:**  $\pi_A = h \implies \pi_B = \ell \implies$  Bob should vote for  $S$ .

To put it differently, from Bob's perspective,

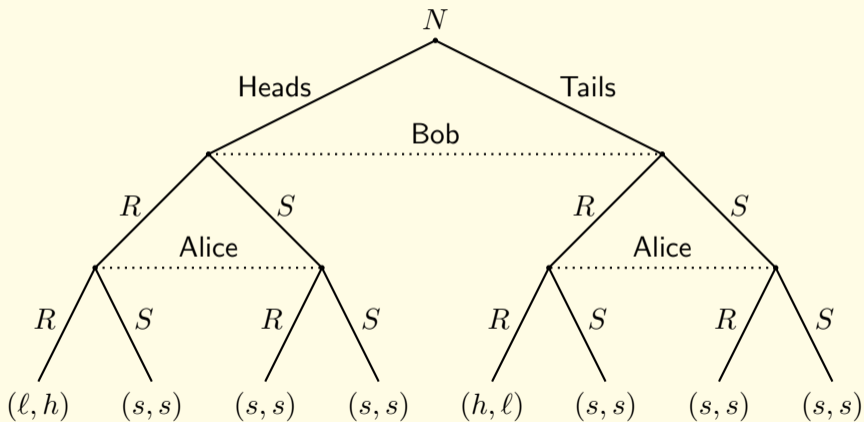
Positive correlation: Risky option is selected **advantageously**.

Negative correlation: Risky option is selected **adversely**.

# extensive-form for advantageous selection



## extensive-form for adverse selection





**Proposition 1.**  $\exists$  a unique strategy profile that survives two rounds of elimination of weakly dominated strategies. In this strategy profile:

- Alice (informed) votes for  $R$  iff  $\pi_A = h$
- Bob (uninformed) votes for  $R$  ( $S$ ) if payoffs are **positively** (**negatively**) correlated.

Also unique weakly undominated Bayes-Nash & trembling-hand perfect equilibrium.

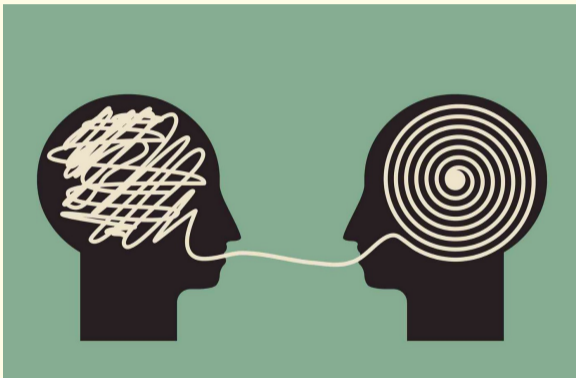
*Qualitative prediction:* Uninformed player should vote for risky far more often when payoffs are positively correlated than when payoffs are negatively correlated.

**Potential confound:** Perhaps that player is averse to **ex post inequality**.

So we compare to a game where both players are symmetrically uninformed.

- Both uninformed  $\implies$  Neither adverse nor advantageous selection.
- Features risk and inequity. (We also have a risk-elicitation stage).
- Compare with main game to see if “Bob” adjusts for “Alice” having private info.

experiment: the [human-human](#) treatment



The main part of our experiment are **asymmetric information** games in which:

- Risky option payoffs:  $h = 20, \ell = 10$ .
- Safe option payoff:  $s \in \{12, 16\}$ .
- Subjects play in the role of both Alice (informed) and Bob (uninformed).

We compare behavior of subjects in the role of Bob in these AI games with those of:

- **Symmetric information**: both subjects are uninformed.
- **Dictator**: each subject chooses between options with and without uncertainty.

Each session lasted  $\approx 50$  minutes, and average earnings were \$15.

86 subjects, Penn State's LEMA, Spring 2019.

For each part, subjects also answered incentivized comprehension questions.

## aggregate results: fraction choosing risky

| Round | Safe Option  | Risky Option                 | AI Game<br>(uninformed "Bob") | SI Game | Dictator Game |
|-------|--------------|------------------------------|-------------------------------|---------|---------------|
| 12N   | (\$12; \$12) | (\$10, \$20) or (\$20, \$10) | 47.7%                         | 77.9%   | 72.1%         |
| 12P   | (\$12; \$12) | (\$10, \$10) or (\$20, \$20) | 86.0%                         | 88.4%   | 82.6%         |
| 16N   | (\$16; \$16) | (\$10, \$20) or (\$20, \$10) | 1.2%                          | 3.5%    | 0%            |
| 16P   | (\$16; \$16) | (\$10, \$10) or (\$20, \$20) | 32.6%                         | 8.1%    | 7.0%          |

In AI, diff between 12N and 12P and between 16N and 16P are significant at  $p < 0.001$ .  
 In 12N and 16P, diff between AI and SI is significant at  $p < 0.001$ .

## individual-level results

### What fraction of subjects conform to theory?

- **Adverse Selection:** 52.3% of subjects choose S in both 12N and 16N.
- **Advantageous Selection:** 30.2% of subjects choose R in both 12P and 16P.
- **All Selection:** 20.9% of subjects do both.

Higher fraction of subjects account for adverse selection ( $p = 0.001$ ).

### Role of Social Preferences:

- **Aversion to Ex-Post Inequality:** In both SI and dictator games, vast majority of subjects make identical choices across correlation.
- **Preferences for Efficiency:** We find no subject that behaves consistent with a strong preference for efficiency across all the games.

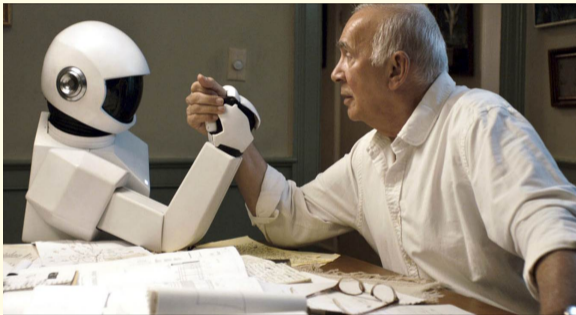
We see evidence that (some) subjects do account for selection, and some fail to do so, and more subjects account for adverse selection.

What factors lead subjects to not account for selection?

- **Strategic uncertainty:** Bob is uncertain about what informed Alice is doing.
- **Failures to reason about contingencies:** Bob does not think about “pivotality,” or draw inferences from it.
- **Cursedness:** Bob fails to realize how Alice’s choices reflect her information.

To study these possibilities, we turn to a “Human-Robot” treatment that pairs subjects with computer players whose strategies are revealed ahead of time.

experiment: the human-robot treatment



Recruited 82 subjects; each is paired with a **robot player** whose strategy is revealed ahead of time.

Idea: all human subjects are in the role of Bob; robot in the role of Alice.

Subjects receive close to identical instructions as in HH treatment, except told that robot players earned “virtual dollars” and were programmed to:

- In **symmetric information** games: vote for the risky option
- In **asymmetric information** games: vote for option that maximizes virtual \$.



## SI games in human-robot treatment

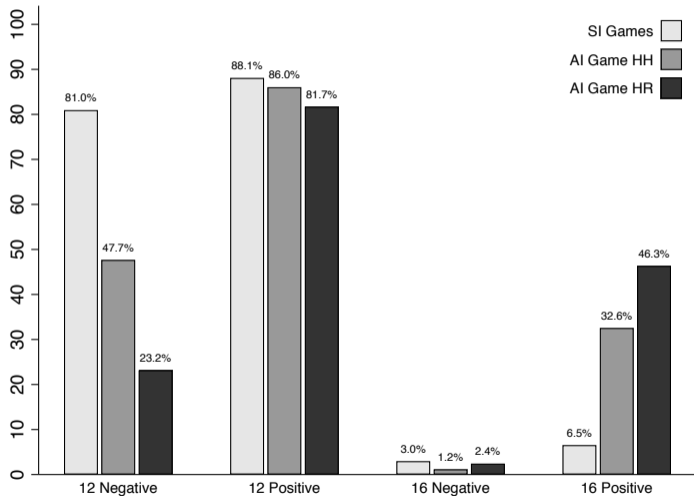
We see no significant differences in behavior in SI games in the HR treatment from those in the HH treatment.

Further suggestive evidence that social preferences play a limited role in our HH data.

| Round | Safe Option  | Risky Option                 | Asymmetric Information<br>Human-Human | Asymmetric Information<br>Human-Robot |
|-------|--------------|------------------------------|---------------------------------------|---------------------------------------|
| 12N   | (\$12; \$12) | (\$10, \$20) or (\$20, \$10) | 47.7%                                 | <b>23.2%</b>                          |
| 12P   | (\$12; \$12) | (\$10, \$10) or (\$20, \$20) | 86.0%                                 | <b>81.7%</b>                          |
| 16N   | (\$16; \$16) | (\$10, \$20) or (\$20, \$10) | 1.2%                                  | <b>2.4%</b>                           |
| 16P   | (\$16; \$16) | (\$10, \$10) or (\$20, \$20) | 32.6%                                 | <b>46.3%</b>                          |

- Difference between 12N HH and 12N HR is significant at  $p = 0.001$ .
- Difference between 16P HH and 16P HR is significant at  $p = 0.068$ .

## fraction choosing the risky option in SI and AI



## individual level

**Negative Correlation:** 74.4% vote for Safe in both rounds vs. 52.3% in HH ( $p = 0.003$ ).

**Positive Correlation:** 42.7% vote for Risky in both rounds vs. 30.2% in HH ( $p = 0.094$ ).

**40.2%** of subjects follow predictions in all rounds (vs. 20.9% in HH).

→ **25%** of deviations from theory in HH are due to strategic uncertainty.

Despite removing strategic uncertainty, gap between adverse and advantageous selection remains. ( $p < 0.001$ ).

## contingent reasoning

It appears that eliminating strategic uncertainty has a significant effect.

Another possibility is that subjects do not reason about contingencies.

To investigate this possibility, we have 4 high-stakes questions on contingent reasoning.

- We observe answers to these questions.
- Subjects then play the AI game again.

**Computer player is INFORMED and votes for option that gives it the HIGHEST number of “virtual dollars”**

| Safe Option  | Risky Option                 |
|--------------|------------------------------|
| (\$17; \$17) | (\$15, \$20) or (\$20, \$15) |

- “If the computer player votes for the option requiring 2 votes (the option on the right), what does that tell you about the outcome of the coin flip?”  
Multiple Choice: HEADS; TAILS; Nothing.
- “If the computer player votes for the option requiring 2 votes, and you vote for it too, how much will you earn?”  
Multiple Choice: \$15, \$17, \$20, \$15 or \$20 with equal chance of each.

## how well do subjects answer CR questions?

89% answer **all** of the **contingent reasoning** questions correctly.

Of the 11% who answered at least one CR question incorrectly, none played according to Proposition 1 in the preceding AI game.

How does it affect subsequent play?

| Round | Safe Option  | Risky Option                 | Asymmetric Information | Asymmetric Information(2) |
|-------|--------------|------------------------------|------------------------|---------------------------|
| 1     | (\$12; \$12) | (\$10, \$20) or (\$20, \$10) | 23.2%                  | <b>22.0%</b>              |
| 2     | (\$12; \$12) | (\$10, \$10) or (\$20, \$20) | 81.7%                  | <b>90.2%</b>              |
| 3     | (\$16; \$16) | (\$10, \$20) or (\$20, \$10) | 2.4%                   | <b>1.2%</b>               |
| 4     | (\$16; \$16) | (\$10, \$10) or (\$20, \$20) | 46.3%                  | <b>62.2%</b>              |

- **76.8%** vote for safe in both **adverse selection** rounds (vs. 74.4% in AI)
- **61.0%** vote for risky in both **advantageous selection** rounds (vs. 47.2% in AI).
- **57.3%** of subjects vote according to predictions in all rounds (vs. 40.2% in AI).
- **What can we say about the 42.7% of subjects who don't?**



## failure to understand or to apply?

Of the 27 subjects who deviate from theory, 76.5% answer every contingent reasoning question correctly.

[Sidenote: Barring 1 subject, all those whose behavior in AI(2) matches Prop 1 also answer every CR question correctly].

Evidence suggests that at high stakes, some subjects understand contingent reasoning but fail to apply it in their games.

## a puzzle

Throughout all of our treatments, we see that subjects are more likely to account for adverse selection than advantageous selection.

Our subsequent theory and experimental designs investigate why.

theory: limited strategic thinking + risk aversion

Can models of limited strategic thinking, combined with risk aversion, help us understand why people may be better at accounting for adverse selection?

We show that it does not in two steps.

### Theory:

- Level- $k$  requires a very high degree of risk aversion, and works for *only* level-1.
- Cursed equilibrium requires significant risk aversion.
- Formal arguments without parametric restrictions on the utility function.

### Data:

- Most of our subjects are not sufficiently risk averse for either to apply.
- Pattern remains even when excluding those subjects.

Theory abstracts from social preferences, and we use data from HR treatment.

Since our solution uniquely emerges from two rounds of elimination of weakly dominated strategies, it can apply *only* if Bob is a level-1 thinker.

Suppose Bob envisions an L0 player who chooses S with prob  $p$ , regardless of info.

For him to vote for safe in 12N / 12P:

$$u(12) \geq pu(12) + (1 - p) \left( \frac{1}{2}u(10) + \frac{1}{2}u(20) \right)$$

Re-arranging terms,  $u(12) \geq \frac{u(10)+u(20)}{2}$ , which is true for only 2.4% of our subjects.

## fully and partially cursed equilibrium

Suppose Bob's cursedness is  $\chi$  in  $[0, 1]$ .

In equilibrium, the **marginal probability** with which Alice votes for S is  $1/2$ .

So if Bob is  $\chi$ -cursed, then he attributes:

- Prob  $\chi$  to Alice voting S with prob  $1/2$  in every payoff state (conveying no info)
- Prob  $(1 - \chi)$  to Alice voting S if  $\pi_A = \ell$  and R if  $\pi_A = h$ .

For Bob to vote for Safe in 12N:

$$u(12) \geq (1 - \chi) \left( \frac{1}{2}u(12) + \frac{1}{2}u(10) \right) + \chi \left( \frac{1}{2}u(12) + \frac{1}{4}u(10) + \frac{1}{4}u(20) \right).$$

For Bob to vote for Safe in 16P:

$$u(16) \geq (1 - \chi) \left( \frac{1}{2}u(16) + \frac{1}{2}u(20) \right) + \chi \left( \frac{1}{2}u(16) + \frac{1}{4}u(10) + \frac{1}{4}u(20) \right).$$

**Proposition 2.** If  $u(14) \geq \frac{u(10)+u(20)}{2}$ , then there is no value of  $\chi$  for which both inequalities above are satisfied.

But in risk-elicitation task, less than 10% of our subjects switch from lottery to sure thing below \$14, and less than 30% do so below \$15.

theory: self-confirming equilibrium



Individuals likely enter AI games with various biases:

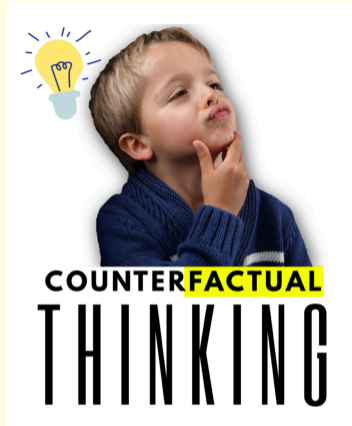
- Zero-sum thinking
- Cursedness / limited strategic thinking

But these biases are influenced by experiences.

A challenge we face is that we observe consequences of actions chosen and rarely observe counterfactuals:

*What would have happened had I chosen differently?*

This can contribute to gap between adverse & advantageous selection.



Suppose that Bob obtains payoff feedback *only* from on-path actions.

If payoffs are *negatively correlated*, then if he mistakenly chooses the risky option, he observes that R is selected only when  $\pi_b = \ell$ .

⇒ he sees that this is a mistake.

If payoffs are *positively correlated*, then if he mistakenly chooses the safe option, then the safe option is always selected.

⇒ he doesn't observe what Alice would have done.

⇒ he *does not* see that this is a mistake.

### Proposition 3.

- ① If payoffs are **negatively correlated**, then in every weakly undominated self-confirming equilibrium, the uninformed voter votes for the safe option.
- ② If payoffs are **positively correlated**, then there exists a weakly undominated self-confirming equilibrium in which the uninformed voter votes for the safe option.

Thus, incorrect beliefs off the path of play can rationalize “mistakes” if payoffs are positively correlated but not when payoffs are negatively correlated.

This motivates varying payoff feedback and seeing how that affects behavior.

experiment: the **feedback** treatments



Study of self-confirming equilibria suggests that:

- **Negative correlation**: info about off-path **shouldn't** significantly affect behavior.
- **Positive correlation**: info about off-path **should** significantly affect behavior.
- Info about off-path histories should reduce gap in how well subjects account for adverse and advantageous selection.

We test these hypotheses in our feedback treatments:

- Subjects obtain partial or full feedback in a number of rounds.
- Partial Feedback Treatment = Feedback about only on-path choices.
- Full feedback Treatment = Feedback about on- and off-path choices.
- We then study choices in subsequent no-feedback rounds.

## partial feedback

After each partial feedback round, subject reminded of how he voted and told what payoffs would be if that round is selected for payment.

**You voted for the option on the LEFT.**

**Votes needed: 1**

Your earnings: \$12  
Computer player: \$12 virtual (imaginary) dollars

**Votes needed: 2**

**HEADS:**  
Your earnings: \$10  
Computer player: \$10 virtual (imaginary) dollars

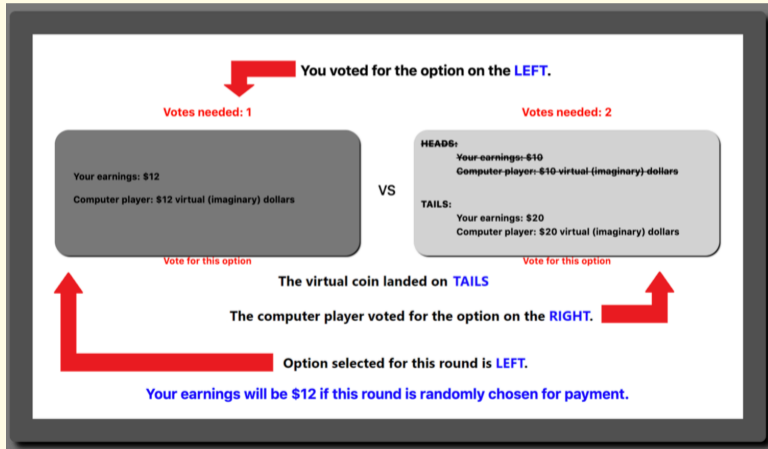
**TAILS:**  
Your earnings: \$20  
Computer player: \$20 virtual (imaginary) dollars

**VS**

**Your earnings will be \$12 if this round is randomly chosen for payment.**

# full feedback

In full feedback, subject also told robot player's choice & outcome of coin toss.



---

|                                    | Partial Feedback Treatment | Full Feedback Treatment |
|------------------------------------|----------------------------|-------------------------|
| % voting for S in both 12N and 16N | 77.9%                      | 81.9%                   |
| % voting for R in both 12P and 16P | 62.8%                      | 75.9%                   |
| % doing both                       | 55.8%                      | 71.1%                   |

---

- Full feedback: insignificant gap between negative & positive correlation ( $p = 0.166$ ).
- Partial feedback: significant gap between negative & positive correlation ( $p = 0.009$ )
- **Negative Correlation**: insignificant diff between partial and full feedback ( $p = 0.515$ ).
- **Positive Correlation**: significant diff between partial and full feedback. ( $p = 0.065$ ).



## what did we learn?

Theory + experiments on adverse and advantageous selection.

Significant fraction of subjects account for selection, and removing strategic uncertainty increases this fraction, particularly for adverse selection.

Systematic gap in how well subjects account for adverse versus advantageous selection.

We relate that gap to individuals' inability to learn about counterfactuals.

We believe that this has potentially important implications for political behavior, the perception of polarization, and the distrust of the elite and experts.

## related literature

We build on a vast literature on asymmetric information games, failures in contingent reasoning and selection-neglect, as well as learning. See paper for references.

Where we see our contributions as being:

- Unified treatment of adverse & advantageous selection, facilitating comparison.
- Comparison of AI games with otherwise identical SI games.
- Human-human vs. Human-Robot.
- Role of counterfactuals.

Thank you!