# Causality and Causal Misperception in Dynamic Games

Sungmin Park

The Ohio State University

December 19, 2024

## Motivation

Limited observation of reality $\Rightarrow$ Varying perceptions of causality

- People have different perceptions about how actions affect outcomes



  **Smoking**          **Education**          **Police size**          **Social media**

- Subjects in lab experiments look at the same data and tell different causal narratives (Kendall and Charles, 2022)

- Yet, most applications of game theory continue to assume Rational Expectations (RE)

## What I do

| **Question** | What is a useful solution concept to incorporate people's misperceptions about causality in extensive-form games? |
|---|---|
| Answer | Let each player best respond to a belief about Nature and others' strategies consistent with observed outcomes |
| Even better | + let each player's belief be the simplest explanation consistent with observation |

"Maximum-entropy Observation-consistent Equilibrium" (MOE)

## What I do

| | |
|---|---|
| **Question** | What is a useful solution concept to incorporate people's misperceptions about causality in extensive-form games? |
| **Answer** | Let each player best respond to a belief about Nature and others' strategies consistent with observed outcomes |
| Even better | + let each player's belief be the simplest explanation consistent with observation |

"Maximum-entropy Observation-consistent Equilibrium" (MOE)

## What I do

**Question** What is a useful solution concept to incorporate people's misperceptions about causality in extensive-form games?

**Answer** Let each player best respond to a belief about Nature and others' strategies consistent with observed outcomes

**Even better** + let each player's belief be the simplest explanation consistent with observation

"Maximum-entropy Observation-consistent Equilibrium" (MOE)

## What I do

| | |
|---|---|
| **Question** | What is a useful solution concept to incorporate people's misperceptions about causality in extensive-form games? |
| **Answer** | Let each player best respond to a belief about Nature and others' strategies consistent with observed outcomes |
| **Even better** | + let each player's belief be the simplest explanation consistent with observation |

"Maximum-entropy Observation-consistent Equilibrium" (MOE)

## Main Results

**Does it Exist?**    Every finite extensive-form game with perfect recall and observational constraint has an MOE

**Is it Useful?**    MOE captures common causal misperceptions such as

- Correlation neglect
- Omitted-variable bias (selection neglect)
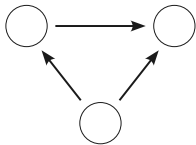- Simultaneity bias (reverse causality bias)

**Is it Compatible with RE?**    If agents have perfect observation of outcomes,

- OE ⇔ Self-confirming equilibrium
- MOE ⇔ Perfect Bayesian Equilibrium (PBE)

## Literature

Bridging behavioral theory and standard game theory



| **Behavioral theory** | **Standard game theory** | **My paper (MOE)** |
| --- | --- | --- |
| (e.g. Spiegler, 2020, 2021) | (e.g. Kreps and Wilson, 1982) | |

- Single-person decisions
- Directed Acyclic Graphs
- Subjective best responses

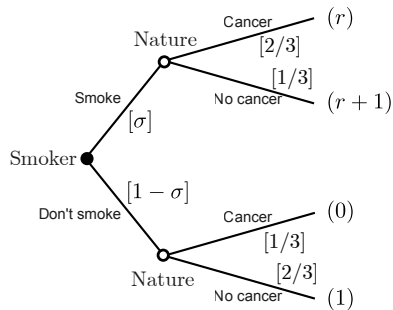- Multiple players
- Rational expectations
- Objective best responses

- Multiple players
- Observational structure ($C$) + maximum entropy
- Subjective best responses

Simplest Example

## Simplest example

- Player chooses to smoke ($s = 1$) or not ($s = 0$).
  - If he smokes, he gets cancer with prob $\pi_1 = 2/3$.
  - If not, he gets cancer with prob $\pi_0 = 1/3$.
  - He gets $r < \frac{1}{3}$ if he smokes and loses $1$ if he gets cancer.

- Player's **strategy** is the prob $\sigma \in [0,1]$ of smoking.

- Player's **belief** is $\beta = (\beta_0, \beta_1)$ where $\beta_s$ is the subjective probability of getting cancer given $s$.

$\Rightarrow$ Under RE, one shouldn't smoke because the causal effect of smoking on cancer ($\frac{2}{3} - \frac{1}{3} = \frac{1}{3}$) is larger than the reward $r$



**Smoker's Problem**

## Observational consistency

**Assumption** Player observes only the marginal prob of cancer.

### Definition

Given strategy $\sigma \in [0,1]$, a belief $\beta \in [0,1]^2$ is **observation-consistent** if

$$\underbrace{\sigma\beta_1 + (1-\sigma)\beta_0}_{\text{perceived marginal prob of cancer}} = \underbrace{\sigma \cdot \frac{2}{3} + (1-\sigma) \cdot \frac{1}{3}}_{\text{actual marginal prob of cancer}}$$
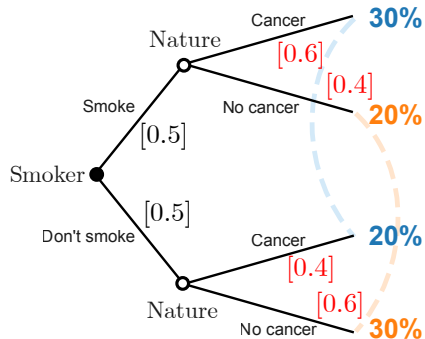
**Interpretation** Player sees a population of players choosing $\sigma$ and sees the overall rate of cancer patients, but do not know the conditional probabilities.

**Problem** There are many observation-consistent beliefs.
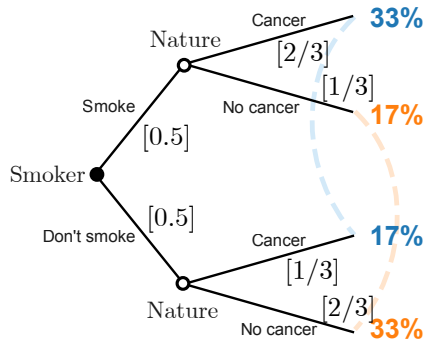
# Illustration of an observational consistency

Suppose I smoke half of the time ($\sigma = 0.5$).



What I **think** Nature does

What Nature **really** does

## Principle of Maximum Entropy

**Notation**

- $\mathbf{p}(\sigma, \beta)$: vector of probabilities over the 4 terminal nodes.
- $G(\cdot)$: Shannon entropy function, i.e. $G(\mathbf{q}) = \sum -q \log q$

### Definition

Given strategy $\sigma \in (0, 1)$, an observation-consistent belief $\beta^* \in [0, 1]^2$ **maximizes the entropy** if
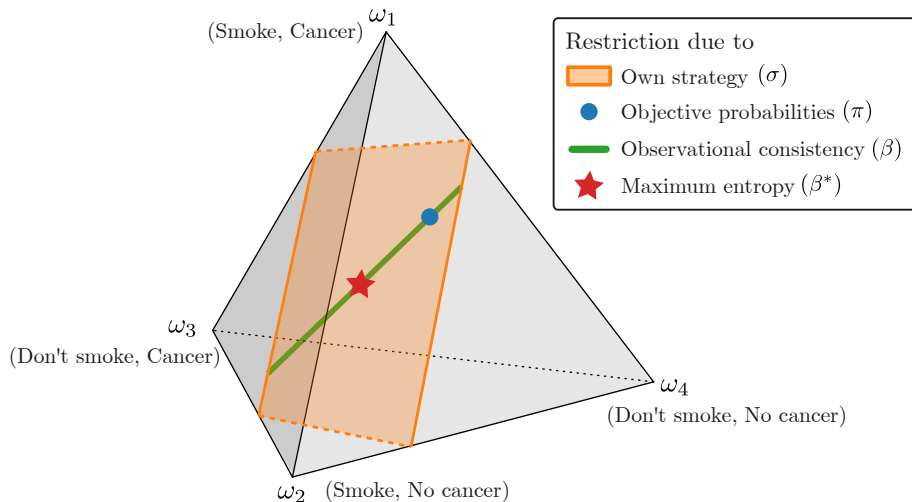
$$\beta^* \in \underset{\beta \text{ is observation-consistent}}{\operatorname{argmax}} G(\mathbf{p}(\sigma, \beta)).$$

**Interpretation**

- Among many worldviews consistent with observation, the agent believes in the the one that assumes the least information

# Illustration of maximum entropy

## A point prediction on belief

# Maximum entropy $\Rightarrow$ correlation neglect

### Claim

For every $\sigma \in (0, 1)$, the maximum-entropy belief $\beta^*$ satisfies

$$\beta_0^* = \beta_1^* = (1 - \sigma) \cdot \frac{1}{3} + \sigma \cdot \frac{2}{3}.$$

**Meaning**   Player doesn't think smoking causes cancer

**Intuition**   Player observes no evidence of dependence between smoking and cancer, so he believes in none.

**General result (Shore and Johnson, 1980; Csiszar, 1991)**

Correlation neglect $\Leftrightarrow$ maximum entropy, whenever agents observe only the marginal prob. distribution between two variables

## Equilibrium

Defined just for the Smoker's Problem

> ### Definition
>
> A strategy-belief pair $(\sigma, \beta)$ is an **observation-consistent equilibrium (OE)** if
>
> ❶ Given the belief $\beta$, the strategy $\sigma$ is a best response, and
>
> ❷ Given the strategy $\sigma$, the belief $\beta$ is observation-consistent.

**Interpretation**

- OE is a prediction of how the smoker behaves, given his possibly wrong but observationally consistent belief

# OE is too permissive

Every strategy is rationalizable by some observation-consistent belief

### Claim

Every strategy $\sigma$ has a belief $\beta$ such that $(\sigma, \beta)$ is an OE.

**Note**: Specifically, the OEs are

1. $\sigma = 0$, $\beta_0 = \frac{1}{3}$, and $\beta_1 - \beta_0 \geq r$,
2. $\sigma = 1$, $\beta_1 = \frac{2}{3}$, and $\beta_1 - \beta_0 \leq r$, and
3. $\sigma \in (0, 1)$, $\beta_0 = \sigma \cdot (\frac{2}{3} - r) + (1 - \sigma) \cdot \frac{1}{3}$, and $\beta_1 = \sigma \cdot \frac{2}{3} + (1 - \sigma)(\frac{1}{3} + r)$.

**Idea** Because there are many observation-consistent beliefs, there are many OEs.

## Definition of MOE

> ### Definition
>
> An OE $(\sigma, \beta)$ is a **maximum-entropy observation-consistent equilibrium (MOE)** if $\beta$ maximizes the entropy given $\sigma \in (0, 1)$.

* For $\sigma \notin (0, 1)$, an OE is an MOE if some $\{(\sigma^k, \beta^k)\}_{k=1}^{\infty} \to (\sigma, \beta)$ and each $\beta^k$ maximizes the entropy given $\sigma^k$

**Interpretation**

- MOE is an OE with the extra requirement that the smoker believes in the simplest explanation consistent with observation

## MOE provides a sharper prediction

### Claim

A strategy-belief pair $(\sigma, \beta)$ is an MOE if and only if

$$\sigma = 1 \quad \text{and} \quad \beta_0 = \beta_1 = \frac{2}{3}.$$

**Meaning**

- Player keeps smoking while thinking that smoking doesn't cause cancer

**Intuition**

- Maximum-entropy belief features correlation neglect, so no other strategy is a best response.

General Framework

## General framework

| | |
|---|---|
| **Model** | $(\Gamma, C)$ where |
| | • $\Gamma$: a finite extensive-form game with perfect recall, and |
| | • $C$: **observational structure**, a linear map from outcomes $(\Delta(\Omega))$ to observable outcomes $(\mathbb{R}^{\ell})$ |
| | |
| **Observational consistency** | Given a strategy $\sigma_i$, a belief $\beta_i$ is observation-consistent if |

$$C\mathbf{p}(\sigma_i, \beta_i) = C\mathbf{p}(\sigma_i, (\sigma_{-i}, \pi)).$$

| | |
|---|---|
| **Equilibrium (MOE)** | A profile of strategies, beliefs, and posterior functions such that |
| | • each strategy is (subjectively) sequentially rational, |
| | • each belief maximizes the entropy s.t. obs consistency, and |
| | • each posterior function satisfies Bayes rule |

▸ Definiton

# Existence of MOE

## Theorem

Every finite extensive-form game with perfect recall and observational constraint has an MOE.

### Meaning

- There always exists a prediction where everyone best responds to what they think how others play, with a belief that assumes the least information beyond observation.
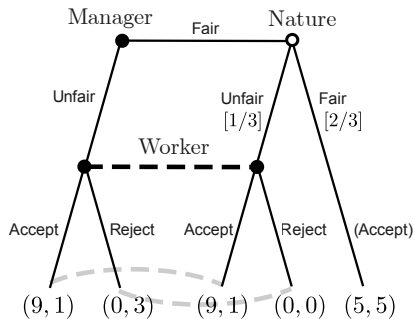
### Key proof step

- With $\epsilon$-constrained strategies, mappings from a strategy profile $\sigma$ to a maximum-entropy beliefs $\beta_i$ and posterior functions are well-behaved.

## Example: An ultimatum-game-like scenario

**Manager-Worker game**

- Manager decides a fair or unfair bonus to Worker

- Even if Manager chooses a fair bonus, Nature might change it to unfair or keep it fair

- If Worker receives fair bonus, he accepts. If not, he either accepts or rejects.
  - He gets a thrill for rejecting an unfair Manager

- Worker doesn't know how likely Manager treats him unfairly in the interim or ex post (in a population)



$$C = \begin{bmatrix} 1 & \cdot & 1 & \cdot & \cdot \\ \cdot & 1 & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{bmatrix}$$

19 / 33

## Standard prediction
Manager often treats Worker unfairly

### Claim
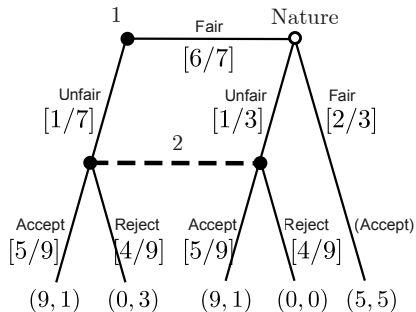
In the unique Perfect Bayesian Equilibrium (PBE),

- Manager offers an unfair bonus 1 out of 7 times

- Worker accepts an unfair bonus 5 out of 9 times
  - He infers (correctly) that any unfair offer is due to Manager 1 out of 3 times

**Perfect Bayesian Equilibrium**



**Intuition**

- There is no causal misperception, because there is no ex-ante uncertainty about others' strategies

## MOE prediction: Manager always tries to be fair

**What Worker thinks others do**



> ### Claim
>
> In the unique MOE,
>
> - Manager always offers the fair bonus
>   - She believes (correctly) that Worker will reject any unfair offer.
>
> - Worker always rejects an unfair offer.
>   - He believes (incorrectly) that Manager offers the unfair bonus 1 out of 6 times
>   - He infers (incorrectly) that any unfair offer is caused by Manager 1 out of 2 times

**What others really do**



**Intuition**   Worker has no clue about the causes of his unfair treatment

## Discussion: How to test MOE in the lab

**Ideal experiment**   Have lab subjects play a game with different observational structures (perfect and imperfect)

1. Randomly assign subjects into Control and Treated groups

2. Within each group, randomly match each subject with another and let them play 1 round of the game

3. Control players receive perfect feedback about all Control outcomes; Treated players receive imperfect feedback about all Treated outcomes

4. Repeat steps 2–3 for sufficiently many rounds

**Example**   A simplified poker game (work in progress)

## MOE and Common Causal Misperceptions

1. Correlation neglect

2. Omitted-variable bias (selection neglect)

3. Simultaneity bias (reverse causality bias)

# 1. A two-stage game of correlated consequences

**Players** $\qquad\qquad N = \{1, 2, \ldots, n\}$

**Stages** $\qquad\qquad$ **1.** Players choose actions $x = (x_i)_{i \in N}$.

$\qquad\qquad\qquad\qquad$ **2.** Nature chooses a consequence $y = (y_1, y_2)$

$\qquad\qquad\qquad\qquad\qquad$ with conditional probability $\pi(y|x) > 0$ for all $(x, y)$.

**Payoffs** $\qquad\qquad u_i(x, y)$

**Obs. structure** $\qquad$ Marginal probabilities of pairs $(x, y_1)$ and $(x, y_2)$

## Correlation neglect

### Proposition

An OE $(\sigma, \beta, \mu)$ is a MOE if and only if for every player $i$,

$$\beta_i(x_{-i}) = \sigma_{-i}(x_{-i}) \qquad \text{for all } x_{-i}, \text{ and}$$
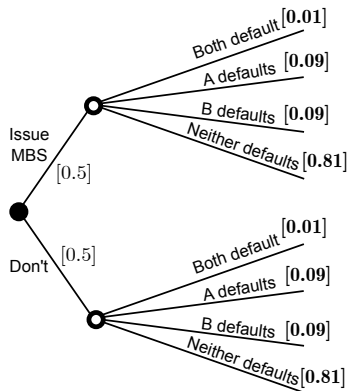
$$\beta_i(y_1, y_2|x) = \pi(y_1|x)\pi(y_2|x) \qquad \text{for all } x \text{ and } (y_1, y_2).$$

**Meaning**   In an MOE, players believe $y_1$ and $y_2$ remain (conditionally) independent regardless of their actions $x$.

**Example (Acharya and Richardson, 2009)**   Financial regulators neglect the correlation between bank failures under lenient regulation

## Stylized example of correlation neglect



**What I think how mortgage-backed securities (MBS) work**

**How they really work**

● Players
○ Nature

Issue MBS
[0.5]

Both default **[0.01]**
A defaults **[0.09]**
B defaults **[0.09]**
Neither defaults **[0.81]**

Don't [0.5]

Both default **[0.01]**
A defaults **[0.09]**
B defaults **[0.09]**
Neither defaults **[0.81]**

Issue MBS
[0.5]

Both default [0.09]
A defaults [0.01]
B defaults [0.01]
Neither defaults [0.89]

Don't [0.5]

Both default [0.01]
A defaults [0.09]
B defaults [0.09]
Neither defaults [0.81]

**Result**   Regulators neglect that issuing MBS causes correlated defaults

26 / 33

## 2. An omitted-variable game

| | |
|---|---|
| **Players** | $N = \{1, 2, \ldots, n\}$ |
| **Stages** | **1.** Nature assigns a state $t$ with probability $\pi(t)$. |
| | **2.** Players see the state $t$ and choose actions $x = (x_i)_{i \in N}$. |
| | **3.** Nature chooses a consequence $y$ with probability $\pi(y|t, x)$. |
| **Payoffs** | $u_i(t, x, y)$ |
| **Obs. structure** | Marginal probabilities of pairs $(t, x)$ and $(x, y)$ |

## Omitted-variable bias (selection neglect)

### Proposition

An OE $(\sigma, \beta, \mu)$ is an MOE if and only if every player's belief $\beta_i$ satisfies

$$\beta_i(t) = \pi(t),$$
$$\beta_i(x_{-i}|t) = \sigma_{-i}(x_{-i}|t), \text{ and}$$
$$\beta_i(y|t, x) = \sum_{t' \in \mathcal{T}} \pi(y|t', x)w(t', x) \qquad \text{for all } (t, x, y).$$

**Note**: $w(\cdot)$ is a weight function $w(t', x) = \lim_{k \to \infty} \frac{\sigma^k(x|t')\pi(t')}{\sum_{t'' \in \mathcal{T}} \sigma^k(x|t'')\pi(t'')}$,

**Meaning**   Players believe the effect of $x$ on $y$ is the same across states $t$

**Example**   High school graduates may overestimate or underestimate the value of college education

# Stylized example of omitted-variable bias



**Result**   High-ability students underestimate the value of college education.

Low-ability students overestimate it.

29 / 33

# 3. Game with simultaneous causality

**Players**          $N = \{1, 2, \ldots, n\}$

**Stages**          **(1)** Nature assigns a state $t \in \{\text{Forward}, \text{Reverse}\}$ with probability $\pi(t)$.

If $t = F$, **(2)** players learn $t$ and choose actions $x = (x_i)_{i \in N}$ and
**(3)** Nature chooses consequence $y$ with prob $\pi(y|F, x)$.

If $t = R$, **(2)** Nature chooses consequence $y$ with prob $\pi(y|R)$ and
**(3)** players learn $(t, y)$ and choose actions $x = (x_i)_{i \in N}$.

**Payoffs**          $u_i(t, x, y)$

**Obs. structure**   Marginal probabilities of the pair $(x, y)$

**Example**          City mayors may misperceive the effects of police on reducing crimes

# Stylized example of simultaneity (reverse causality) bias



**What I think how police and violent crimes work**

**How they really work**

**Result**   Mayor underestimates the effect of police on reducing crime

## Discussion: Implications for stuctural econometrics

**Rational expectations (RE) assumption**

- "Ubiquitous" even though it's a "very strong assumption"
  (Aguirregabiria and Mira, 2010)

- Relaxing it requires modeling and estimating beliefs
  (e.g., Aguirregabiria and Magesan, 2020)

**MOE assumption**

- A viable alternative to RE by providing a point-prediction on beliefs

- Only requires an existing model + observational structure $C$

- Example application: Models of education and occupational choice
  (e.g., Keane and Wolpin, 1997)

## Rest of the paper and takeaway

**Rest of the paper**

- Comparison with related concepts  ▸ Comparison
- Game-theoretic definition of causality  ▸ Causality
- Cooperation in Centipede games  ▸ Centipede game
- Games with infinite time horizons  ▸ Markov games

**Takeaway:**   MOE is useful if you want to

- allow causal misperception in a dynamic model,
- let misperception arise endogenously from the observational structure, and
- want narrow predictions.

## Rest of the paper and takeaway

**Rest of the paper**

- Comparison with related concepts  ▸ Comparison
- Game-theoretic definition of causality  ▸ Causality
- Cooperation in Centipede games  ▸ Centipede game
- Games with infinite time horizons  ▸ Markov games

**Takeaway:**   MOE is useful if you want to

- allow causal misperception in a dynamic model,
- let misperception arise endogenously from the observational structure, and
- want narrow predictions.

## Thank you!

Appendix

## Precise definitions in the general framework

**Strategy** $\sigma_i \in \mathcal{S}_i$        $\sigma_i(a|I_i)$ is player $i$'s objective prob of action $a$ by $i$ at info set $I_i$

**Belief** $\beta_i \in \mathcal{S}_{-i}$        $\beta_i(a|I_j)$ is player $i$'s subjective prob of action $a$ by Nature or an opponent at info set $I_j$.

**Posterior function** $\mu_i$        $\mu_i(h|I_i)$ is player $i$'s subjective prob of history $h \in I_i$ given $I_i$.

**"Assessment"**        $(\sigma, \beta, \mu) = \{(\sigma_i, \beta_i, \mu_i)\}_{i \in N}$

# Definition of OE

**Notation**   $\mathbf{p}(\sigma_i, \beta_i)$ is the subjective probability distribution over $\Omega$

---

### Definition

An assessment $(\sigma, \beta, \mu)$ is an **observation-consistent equilibrium (OE)** if for every player $i$,

❶ the strategy $\sigma_i$ is (subjectively) sequentially rational given $(\beta_i, \mu_i)$,

❷ the belief $\beta_i$ is observation-consistent given the strategy profile $\sigma$:

$$C\mathbf{p}(\sigma_i, \beta_i) = C\mathbf{p}(\sigma_i, (\sigma_{-i}, \pi)), \text{ and}$$

❸ the posterior function $\mu_i$ is Bayes-consistent given $(\sigma_i, \beta_i)$.

---

## Definition of MOE

Given a strategy profile $\sigma$, a player's observation-consistent belief $\beta_i$ **maximizes the entropy** if

$$\beta_i \in \underset{\beta_i' \text{ is obs-cons}}{\mathrm{argmax}} \; G(\mathbf{p}(\sigma_i, \beta_i')).$$

### Definition

An OE $(\sigma, \beta, \mu)$ is a **maximum-entropy observation-consistent equilibrium (MOE)** if there exists a sequence

$$\{\sigma^k, \beta^k\}_{k=1}^{\infty} \longrightarrow (\sigma, \beta)$$

where each $\sigma^k$ is a totally mixed strategy profile and each player's belief $\beta_i^k$ maximizes the entropy given $\sigma^k$.

# OE and MOE nest standard concepts as special cases

> **Proposition**
>
> Under perfect observation of outcomes ($C$ = identity),
>
> $$\text{OE} \iff \text{Self-confirming equilibrium}^*, \quad \text{and}$$
> $$\text{MOE} \iff \text{Perfect Bayesian equilibrium}.$$

**\*** Version with sequential rationality.

### Implication

- Varying the extent of misperception is straightforward: Take an existing model and vary the observational structure $C$.

## Other related concepts

### Analogy-based expectation equilibrium (ABEE)

Jehiel (2005); Jehiel and Koessler (2008); Jehiel (2022)

- Players believe others behave the same in "analogous" situations

### Cursed (sequential) equilibrium

Eyster and Rabin (2005, **CE**); Fong, Lin and Palfrey (2023, **CSE**); Cohen and Li (2022, **SCE**)

- Players believe others behave the same regardless of their types/info

### Berk-Nash equilibrium

Esponda and Pouzo (2016)

- Players' beliefs about the game are misspecified

## Wait... what do I even mean by causality?

**Notation**   $p(\sigma_i, \beta_i)(E|h)$ is the subjective probability of event $E \subset \Omega$ given history $h$, strategy $\sigma_i$ , and belief $\beta_i$.

### Definition

Let $(\sigma, \beta, \mu)$ be an OE. An action $a$ instead of $b$ is a **subjective cause** of an event $E \subset \Omega$ given history $h$ to player $i$ if

$$p(\sigma_i, \beta_i)(E|h, a) > p(\sigma_i, \beta_i)(E|h, b).$$

An action $a$ instead of $b$ is an **objective cause** of an event $E \subset \Omega$ given history $h$ to player $i$ if

$$p(\sigma_i, (\sigma_{-i}, \pi))(E|h, a) > p(\sigma_i, (\sigma_{-i}, \pi))(E|h, b).$$
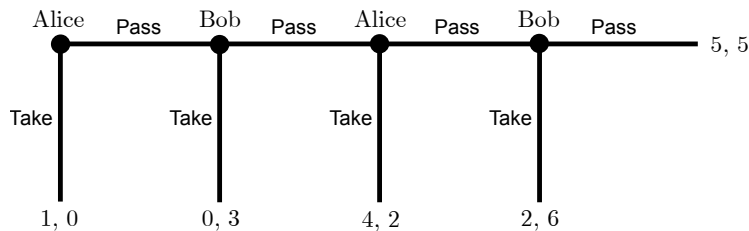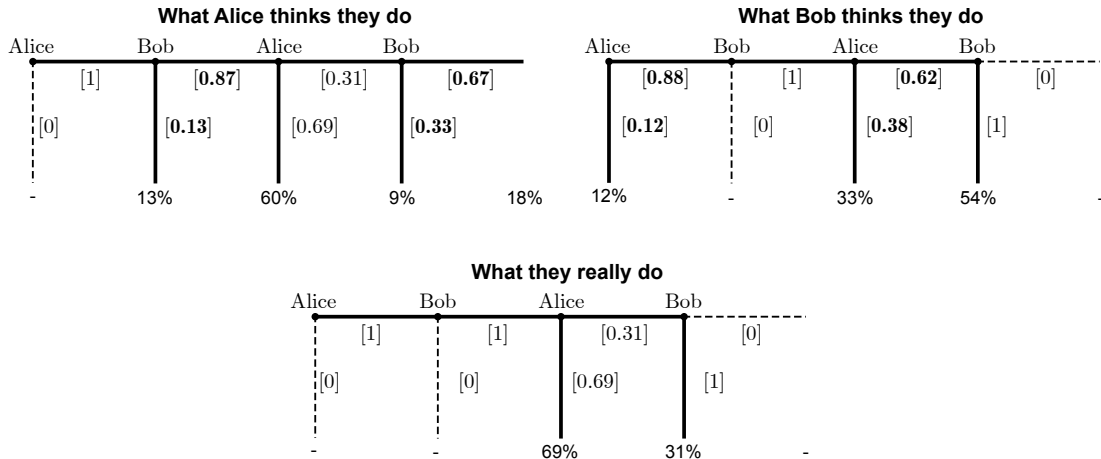
# Example: A centipede game



Figure: A four-node centipede game

### Claim

Suppose players observe only the average number of passes ($C = [0\ 1\ 2\ 3\ 4]$).

There exists no MOE in which Alice Takes immediately.

# Unique MOE of the centipede game

Each thinks the other mixes more than they really do

**What Alice thinks they do**

| Alice | Bob | Alice | Bob | |
|---|---|---|---|---|
| [1] | [**0.87**] | [0.31] | [**0.67**] | |
| [0] | [**0.13**] | [0.69] | [**0.33**] | |
| - | 13% | 60% | 9% | 18% |

**What Bob thinks they do**

| Alice | Bob | Alice | Bob | |
|---|---|---|---|---|
| [**0.88**] | [1] | [**0.62**] | [0] | |
| [**0.12**] | [0] | [**0.38**] | [1] | |
| 12% | - | 33% | 54% | - |

**What they really do**

| Alice | Bob | Alice | Bob | |
|---|---|---|---|---|
| [1] | [1] | [0.31] | [0] | |
| [0] | [0] | [0.69] | [1] | |
| - | - | 69% | 31% | - |

# Extension: Stochastic (Markov) Games



Figure: Stochastic game with permanent game types $\theta$

---

### Proposition

If players perfectly observe steady-state outcomes $(\theta, s, a, s')$,

$$\text{MOE} \iff \text{Markov perfect equilibrium (MPE)}.$$

# Illustration: Parent-Child game of social media use



| | | |
|---|---|---|
| Type | $\theta$ | Sensitive or not sensitive |
| State | $s \to s'$ | In Good mood or in Bad mood |
| Child's action | $a_1$, $a_1'$ | Use socia media or not |
| Parent's action | $a_2$, $a_2'$ | Enforce strict or lenient screen time policy |

Today          Tommorow

# Equilibrium in the Parent-Child game

| Equilibrium | Type ($\theta$) | Child's strategy ($\sigma_1$) | | Parent's strategy ($\sigma_2$) | |
|---|---|---|---|---|---|
| | | Bad mood | Good mood | Bad mood | Good mood |
| MPE | Not sensitive | Use | Use | Lenient | Lenient |
| | Sensitive | Don't | Use | Lenient | Lenient |
| MOE | Not sensitive | Use | Use | Strict | Lenient |
| | Sensitive | Use | Use | Strict | Lenient |

**Note**: MPE refers to Markov perfect equilibrium. MOE refers to maximum-entropy observation-consistent equilibrium.

Acharya, Viral V. and Matthew Richardson (2009) "Causes of the financial crisis," *Critical Review*, 21 (2-3), 195–210.

Aguirregabiria, Victor and Arvind Magesan (2020) "Identification and estimation of dynamic games when players' beliefs are not in equilibrium," *The Review of Economic Studies*, 87 (2), 582–625.

Aguirregabiria, Victor and Pedro Mira (2010) "Dynamic discrete choice structural models: A survey," *Journal of Econometrics*, 156 (1), 38–67.

Cohen, Shani and Shengwu Li (2022) "Sequential Cursed Equilibrium," *arXiv preprint arXiv:2212.06025*.

Csiszar, Imre (1991) "Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems," *The Annals of Statistics*, 2032–2066.

Esponda, Ignacio and Demian Pouzo (2016) "Berk–Nash equilibrium: A framework for modeling agents with misspecified models," *Econometrica*, 84 (3), 1093–1130.

Eyster, Erik and Matthew Rabin (2005) "Cursed equilibrium," *Econometrica*, 73 (5), 1623–1672.

## References II

Fong, Meng-Jhang, Po-Hsuan Lin, and Thomas R Palfrey (2023) "Cursed sequential equilibrium," *arXiv preprint arXiv:2301.11971*.

Jehiel, Philippe (2005) "Analogy-based expectation equilibrium," *Journal of Economic Theory*, 123 (2), 81–104.

———— (2022) "Analogy-based expectation equilibrium and related concepts: Theory, applications, and beyond."

Jehiel, Philippe and Frédéric Koessler (2008) "Revisiting games of incomplete information with analogy-based expectations," *Games and Economic Behavior*, 62 (2), 533–557.

Keane, Michael P and Kenneth I Wolpin (1997) "The career decisions of young men," *Journal of Political Economy*, 105 (3), 473–522.

Kendall, Chad W and Constantin Charles (2022) "Causal narratives,"Technical report.

Kreps, David M and Robert Wilson (1982) "Sequential equilibria," *Econometrica: Journal of the Econometric Society*, 863–894.

Shore, John and Rodney Johnson (1980) "Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy," *IEEE Transactions on Information Theory*, 26 (1), 26–37.

Spiegler, Ran (2020) "Behavioral implications of causal misperceptions," *Annual Review of Economics*, 12, 81–106.

——— (2021) "Modeling players with random "data access"," *Journal of Economic Theory*, 198, 105374.