

Causality and Causal Misperception in Dynamic Games

Sungmin Park

The Ohio State University

December 13, 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)

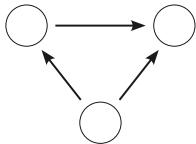
Is it Compatible with RE?

If agents have **perfect observation** of outcomes,

- **OE** \Leftrightarrow Self-confirming equilibrium
- **MOE** \Leftrightarrow Perfect Bayesian Equilibrium (PBE)

Literature

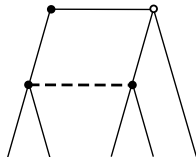
Bridging behavioral theory and standard game theory



Behavioral theory

(e.g. Spiegel, 2020, 2021)

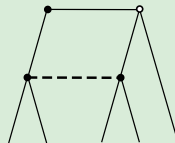
- Single-person decisions
- Directed Acyclic Graphs
- Subjective best responses



Standard game theory

(e.g. Kreps and Wilson, 1982)

- Multiple players
- Rational expectations
- Objective best responses



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

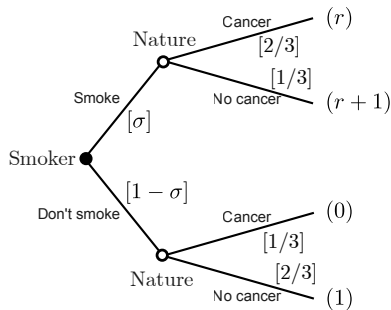
My paper (MOE)

- Multiple players
- Observational structure (C)
- Subjective best responses

Simplest Example

Simplest example

- Smoker 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.
- Smoker's **strategy** is the prob $\sigma \in [0, 1]$ of smoking.
- Smoker's **belief** is $\beta = (\beta_0, \beta_1)$ where β_s is the subjective probability of getting cancer given s .



Smoker's Problem

\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

Observational consistency

Assumption Smoker 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

$$\underbrace{\sigma\beta_1 + (1 - \sigma)\beta_0}_{\text{perceived marginal prob of cancer}} = \underbrace{\sigma \cdot \frac{2}{3} + (1 - \sigma) \cdot \frac{1}{3}}_{\text{actual marginal prob of cancer}}$$

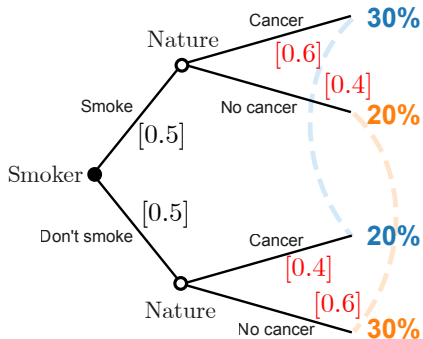
Interpretation Smoker sees a population of smokers choosing σ and sees the overall **rate of cancer** patients, but do not know the **conditional probabilities**.

Problem There are many observation-consistent beliefs.

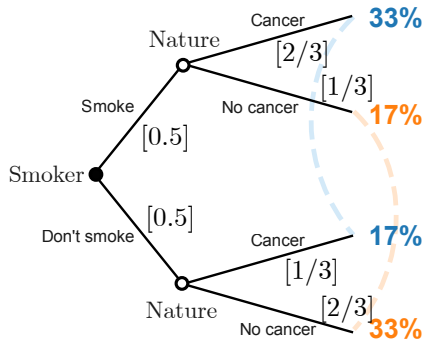
Illustration of an observational consistency

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

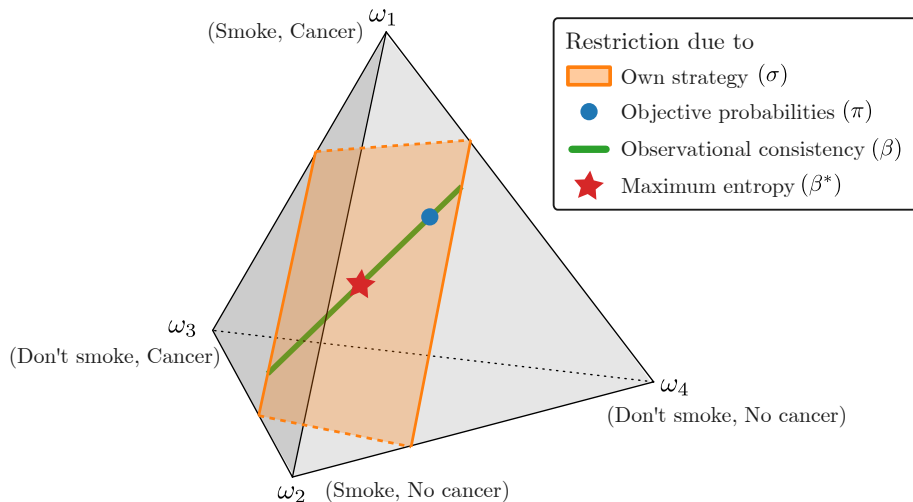
$$\beta^* \in \underset{\beta \text{ is observation-consistent}}{\operatorname{argmax}} G(\mathbf{p}(\sigma, \beta)).$$

Interpretation

- Among many worldviews consistent with observation, players believe 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 β^* satisfies

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

Meaning The smoker doesn't think smoking **causes** cancer

Intuition The smoker 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

Definition

A strategy-belief pair (σ, β) is an **observation-consistent equilibrium (OE)** if

- ① Given the belief β , the strategy σ is a **best response**, and
- ② Given the strategy σ , the belief β 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 σ has a belief β such that (σ, β) is an OE.

Note: Specifically, the OCE equilibria 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 (σ, β) is a **maximum-entropy observation-consistent equilibrium (MOE)** if β maximizes the entropy given $\sigma \in (0, 1)$.

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

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

Meaning

- Smoker **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.

I. General Framework

General framework

Model

(Γ, C) where

- Γ : a finite extensive-form game with perfect recall, and
- C : **observational structure**, a linear map from **outcomes** $(\Delta(\Omega))$ to **observable outcomes** (\mathbb{R}^ℓ)

Observational consistency

Given a strategy σ_i , a belief β_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**

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.

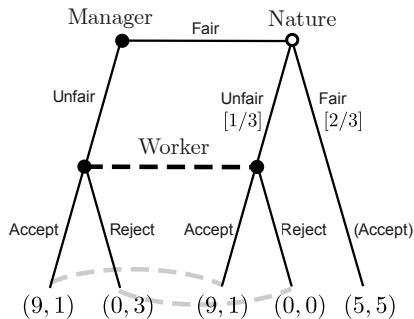
Key proof step

- With **ϵ -constrained strategies**, mappings from a strategy profile σ to a maximum-entropy beliefs β_i and posterior functions are well-behaved.

Example: An ultimatum-game-like scenario

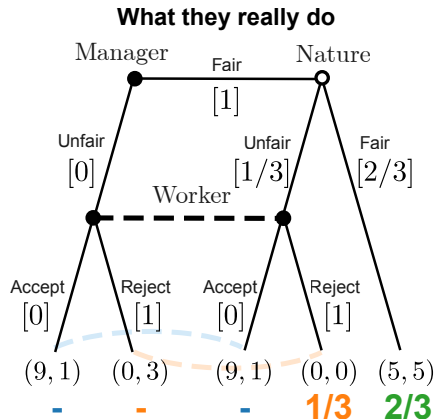
Manager-Worker game

- 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)



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

Manager always tries to be fair



20 / 33

II. 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

$$N = \{1, 2, \dots, n\}$$

Stages

1. Players choose **actions** $x = (x_i)_{i \in N}$.
2. Nature chooses a **consequence** $y = (y_1, y_2)$
with conditional probability $\pi(y|x) > 0$ for all (x, y) .

Payoffs

$$u_i(x, y)$$

Obs. structure

Marginal probabilities of pairs (x, y_1) and (x, y_2)

Correlation neglect

Proposition

An OE (σ, β, μ) is a MOE if and only if for every player i ,

$$\beta_i(x_{-i}) = \sigma_{-i}(x_{-i}) \quad \text{for all } x_{-i}, \text{ and}$$

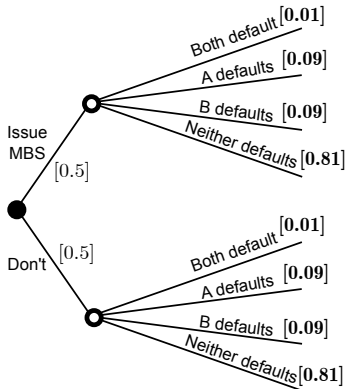
$$\beta_i(y_1, y_2 | x) = \pi(y_1 | x) \pi(y_2 | x) \quad \text{for all } x \text{ and } (y_1, y_2).$$

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

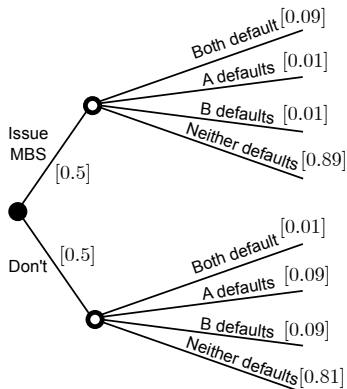
Example Let x be whether an investment bank issues **mortgage-backed securities (MBS)** or not. Let y be the **default outcomes** of two households.

Stylized example of correlation neglect

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



How they really work



Result Investment banks issue MBS, neglecting that those **cause financial crises**

2. An omitted-variable game

Players	$N = \{1, 2, \dots, n\}$
Stages	<ol style="list-style-type: none">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 (σ, β, μ) is an MOE if and only if every player's belief β_i satisfies

$$\beta_i(t) = \pi(t),$$

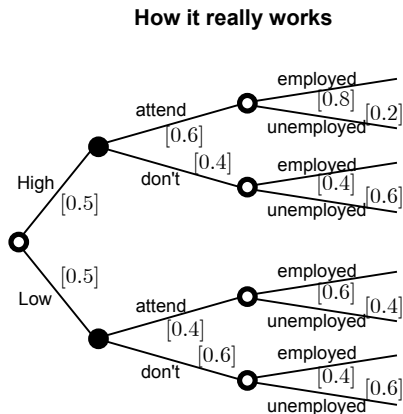
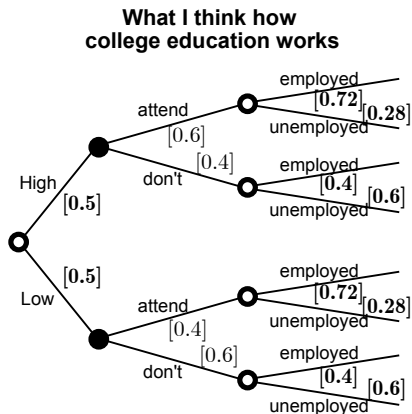
$$\beta_i(x_{-i}|t) = \sigma_{-i}(x_{-i}|t), \text{ and}$$

$$\beta_i(y|t, x) = \sum_{t' \in \mathcal{T}} \pi(y|t', x) w(t', x) \quad \text{for all } (t, x, y).$$

Note: $w(\cdot)$ is a weight function such that $w(t', x) = \lim_{k \rightarrow \infty} \frac{\sigma^k(x|t')\pi(t')}{\sum_{t'' \in \mathcal{T}} \sigma^k(x|t'')\pi(t'')}$, for some sequence $\{\sigma^k\}_{k=1}^{\infty}$ of totally mixed strategy profiles converging to σ .

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

Stylized example of omitted-variable bias



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

Low-ability students overestimate it.

3. Simultaneous causality game

Players $N = \{1, 2, \dots, n\}$

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** $x = (x_i)_{i \in N}$.

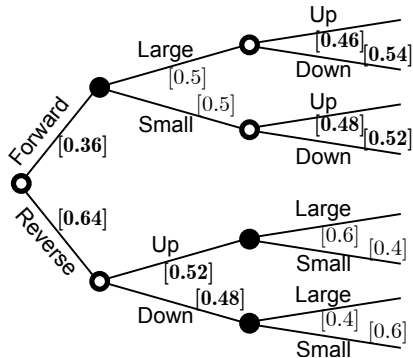
Payoffs $u_i(t, x, y)$

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

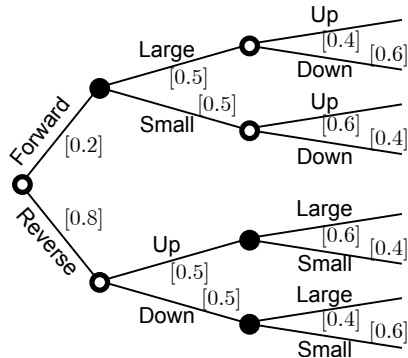
Stylized example of simultaneity (reverse causality) bias

Police size and violent crime rates

What I think how police and violent crimes work



How they really work



Result The police chief **underestimates** the effect of police on reducing crime

III. Relationship with existing solution concepts

OE and MOE nest standard concepts as special cases

Proposition

Under **perfect observation** of outcomes ($C = \text{identity}$),

$$\begin{array}{ll} \text{OE} & \iff \text{Self-confirming equilibrium}^*, \text{ and} \\ \text{MOE} & \iff \text{Perfect Bayesian equilibrium.} \end{array}$$

* 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**

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**.

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!



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

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** σ_i , and **belief** β_i .

Definition

Let (σ, β, μ) 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

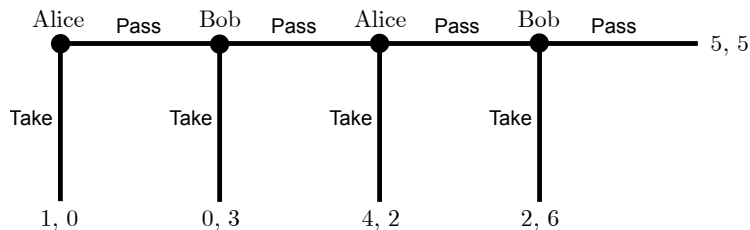


Figure: A four-node centipede game

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.

Each thinks the other mixes more than they really do



Extension: Stochastic (Markov) Games

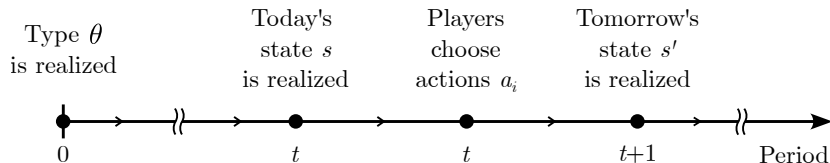


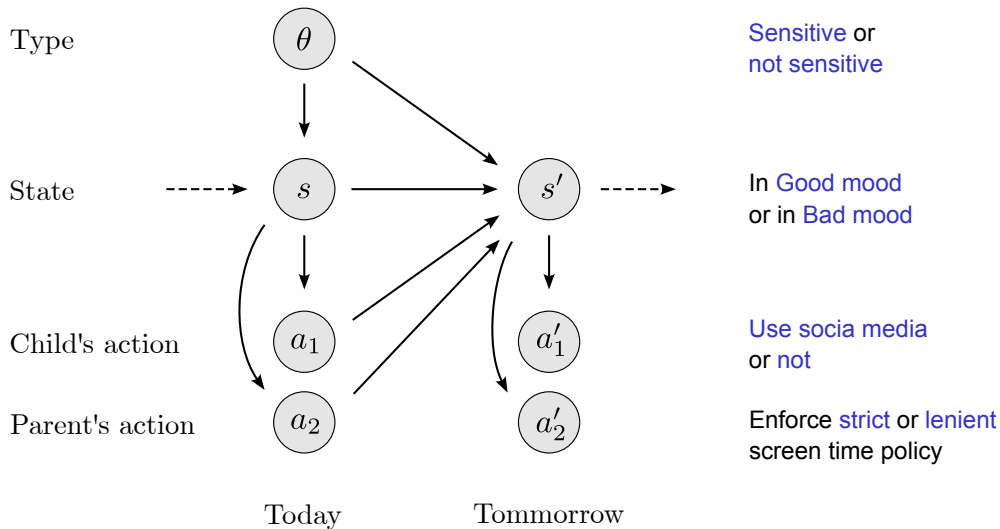
Figure: Stochastic game with permanent game types θ

Proposition

If players perfectly observe steady-state outcomes (θ, s, a, s') ,

MOE \iff Markov perfect equilibrium (MPE).

Illustration: Parent-Child game of social media use



Equilibrium in the Parent-Child game

Equilibrium	Type (θ)	Child's strategy (σ_1)		Parent's strategy (σ_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.

References I

- 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.
- 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.

References II

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.