

# Empirical-evidence Equilibria in Stochastic Games

Nicolas Dubebout

# Outline

- Stochastic games
- Empirical-evidence equilibria (EEEs)
- Open questions in EEEs

# Stochastic Games

- Game theory
- Markov decision processes

# Game Theory

## Decision making

$$u : \mathcal{A} \rightarrow \mathbb{R} \implies a^* \in \arg \max_{a \in \mathcal{A}} u(a)$$

## Game theory

$$u_1 : \mathcal{A}_1 \times \mathcal{A}_2 \rightarrow \mathbb{R}$$

$$u_2 : \mathcal{A}_1 \times \mathcal{A}_2 \rightarrow \mathbb{R}$$

## Nash Equilibrium

$$\begin{cases} a_1^* \in \arg \max_{a_1 \in \mathcal{A}_1} u_1(a_1, a_2^*) \\ a_2^* \in \arg \max_{a_2 \in \mathcal{A}_2} u_2(a_1^*, a_2) \end{cases}$$

# Example: Battle of the Sexes

	F	O
F	2, 2	0, 1
O	0, 0	1, 3

## Nash equilibria

- $(F, F)$
- $(O, O)$
- $(\frac{3}{4}F \frac{1}{4}O, \frac{1}{3}F \frac{2}{3}O)$

# Markov Decision Process (MDP)

Dynamic  $x^+ \sim f(x, a) \iff x^{t+1} \sim f(x^t, a^t)$

Stage cost  $u(x, a)$

History  $h^t = (x^0, x^1, \dots, x^t, a^0, a^1, \dots, a^t)$

Strategy  $\sigma : \mathcal{H} \rightarrow \mathcal{A}$

$$\text{Utility } U(\sigma) = \mathbb{E}_{f, \sigma} \left[ \sum_{t=0}^{\infty} \delta^t u(x^t, a^t) \right]$$

Bellman's equation

$$U^*(x) = \max_{a \in \mathcal{A}} \left\{ u(x, a) + \delta \mathbb{E}_f [U^*(x^+) \mid x, a] \right\}$$

Dynamic programming use knowledge of  $f$

Reinforcement learning learn  $f$  from repeated interaction

# Markov Decision Process (MDP)

Dynamic  $x^+ \sim f(x, a) \iff x^{t+1} \sim f(x^t, a^t)$

Stage cost  $u(x, a)$

History  $h^t = (x^0, x^1, \dots, x^t, a^0, a^1, \dots, a^t)$

Strategy  $\sigma : \mathcal{H} \rightarrow \mathcal{A}$

$$\text{Utility } U(\sigma) = \mathbb{E}_{f, \sigma} \left[ \sum_{t=0}^{\infty} \delta^t u(x^t, a^t) \right]$$

Bellman's equation

$$U^*(x) = \max_{a \in \mathcal{A}} \left\{ u(x, a) + \delta \mathbb{E}_f [U^*(x^+) \mid x, a] \right\}$$

Dynamic programming use knowledge of  $f$

Reinforcement learning learn  $f$  from repeated interaction

# Markov Decision Process (MDP)

Dynamic  $x^+ \sim f(x, a) \iff x^{t+1} \sim f(x^t, a^t)$

Stage cost  $u(x, a)$

History  $h^t = (x^0, x^1, \dots, x^t, a^0, a^1, \dots, a^t)$

Strategy  $\sigma : \mathcal{X} \rightarrow \mathcal{A}$

$$\text{Utility } U(\sigma) = \mathbb{E}_{f, \sigma} \left[ \sum_{t=0}^{\infty} \delta^t u(x^t, a^t) \right]$$

Bellman's equation

$$U^*(x) = \max_{a \in \mathcal{A}} \left\{ u(x, a) + \delta \mathbb{E}_f [U^*(x^+) \mid x, a] \right\}$$

Dynamic programming use knowledge of  $f$

Reinforcement learning learn  $f$  from repeated interaction



# Imperfect Information (POMDP)

Dynamic  $w^+ \sim n(w, a)$

Signal  $s \sim v(w)$

History  $h^t = (s^0, s^1, \dots, s^t, a^0, a^1, \dots, a^t)$

Strategy  $\sigma : \mathcal{H} \rightarrow \mathcal{A}$

Belief  $\mathbb{P}_{n,v,\sigma}[w | h]$

# Imperfect Information (POMDP)

Dynamic  $w^+ \sim n(w, a)$

Signal  $s \sim v(w)$

History  $h^t = (s^0, s^1, \dots, s^t, a^0, a^1, \dots, a^t)$

Strategy  $\sigma : \mathcal{H} \rightarrow \mathcal{A}$

Belief  $\mathbb{P}_{n,v,\sigma}[w | h]$

# Imperfect Information (POMDP)

Dynamic  $w^+ \sim n(w, a)$

Signal  $s \sim v(w)$

History  $h^t = (s^0, s^1, \dots, s^t, a^0, a^1, \dots, a^t)$

Strategy  $\sigma : \Delta(\mathcal{W}) \rightarrow \mathcal{A}$

Belief  $\mathbb{P}_{n,v,\sigma}[w | h]$

# Stochastic Games

Dynamic  $w^+ \sim n(w, a_1, a_2)$

Signals  $\begin{cases} s_1 \sim v_1(w) \\ s_2 \sim v_2(w) \end{cases}$

Histories  $\begin{cases} h_1^t = (s_1^0, s_1^1, \dots, s_1^t, a_1^0, a_1^1, \dots, a_1^t) \\ h_2^t = (s_2^0, s_2^1, \dots, s_2^t, a_2^0, a_2^1, \dots, a_2^t) \end{cases}$

Strategies  $\begin{cases} \sigma_1 : \mathcal{H}_1 \rightarrow \mathcal{A}_1 \\ \sigma_2 : \mathcal{H}_2 \rightarrow \mathcal{A}_2 \end{cases}$

Beliefs  $\begin{cases} \mathbb{P}_{n, v_1, \sigma_1, v_2, \sigma_2} [w, h_2 \mid h_1] \\ \mathbb{P}_{n, v_1, \sigma_1, v_2, \sigma_2} [w, h_1 \mid h_2] \end{cases}$

# Stochastic Games

Dynamic  $w^+ \sim n(w, a_1, a_2)$

Signals  $\begin{cases} s_1 \sim v_1(w) \\ s_2 \sim v_2(w) \end{cases}$

Histories  $\begin{cases} h_1^t = (s_1^0, s_1^1, \dots, s_1^t, a_1^0, a_1^1, \dots, a_1^t) \\ h_2^t = (s_2^0, s_2^1, \dots, s_2^t, a_2^0, a_2^1, \dots, a_2^t) \end{cases}$

Strategies  $\begin{cases} \sigma_1 : \mathcal{H}_1 \rightarrow \mathcal{A}_1 \\ \sigma_2 : \mathcal{H}_2 \rightarrow \mathcal{A}_2 \end{cases}$

Beliefs  $\begin{cases} \mathbb{P}_{n, v_1, \sigma_1, v_2, \sigma_2} [w, h_2 \mid h_1] \\ \mathbb{P}_{n, v_1, \sigma_1, v_2, \sigma_2} [w, h_1 \mid h_2] \end{cases}$

# Existing Approaches

- (Weakly) belief-free equilibrium
- Mean-field equilibrium
- Incomplete theories

# Empirical-evidence Equilibria

# Motivation



0. Pick arbitrary strategies
1. Formulate simple but *consistent* models
2. Design strategies optimal w.r.t. models, then, back to 1.

Empirical-evidence equilibrium is a fixed point:

- Strategies **optimal** w.r.t. models
- Models **consistent** with strategies



# Example: Asset Management

Trading one asset on the stock market

Model based on

- information published by the company
- observed trading activity

Model very different for each agent

# Multiple to Single Agent



# Multiple to Single Agent



# Single Agent Setup

Agent

Nature

# Single Agent Setup

$$x^+ \sim f(x, a, s)$$

Nature

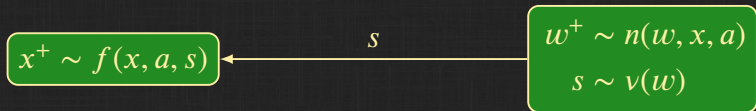
# Single Agent Setup



# Single Agent Setup



# Example: Asset Management



State holding  $x \in \{0..M\}$

Action sell one, hold, or buy one  $a \in \{-1, 0, 1\}$

Signal price  $p \in \{\text{Low, High}\}$

Stage cost  $p \cdot a$

Nature  $w$  represents market sentiment, political climate,  
other traders



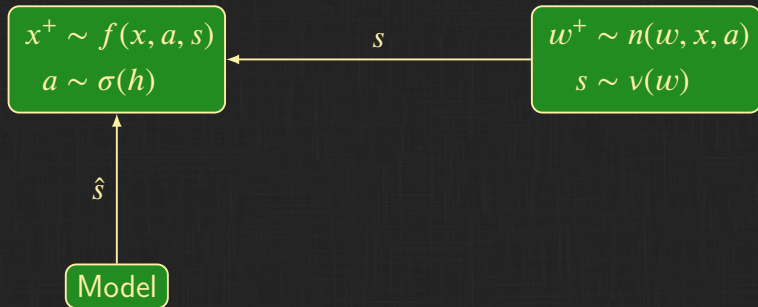
# Single Agent Setup



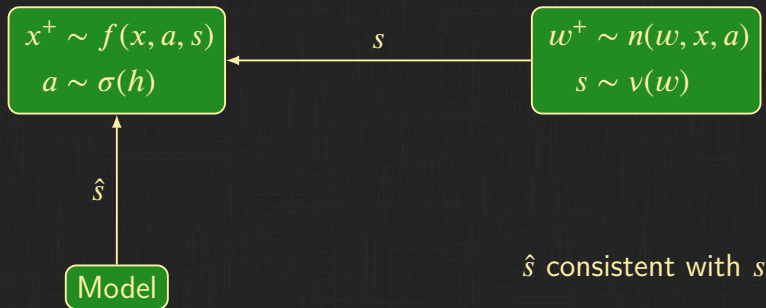
# Single Agent Setup



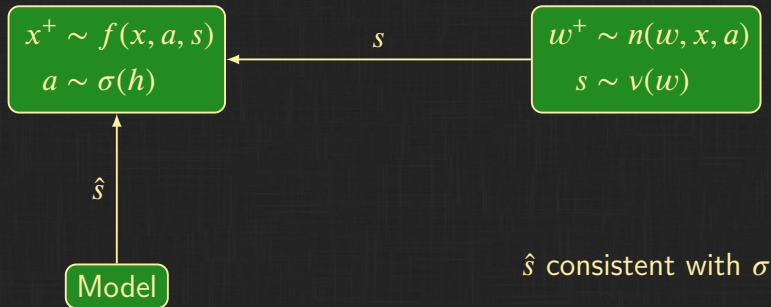
# Single Agent Setup



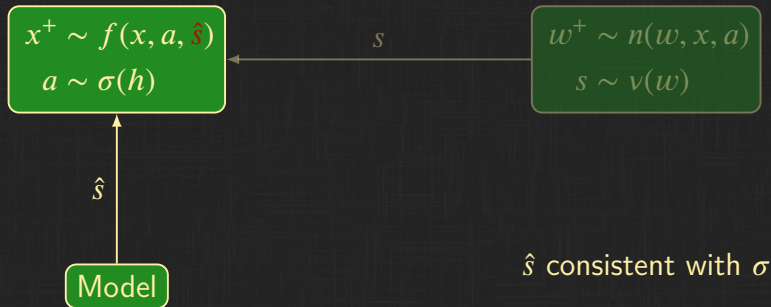
# Single Agent Setup



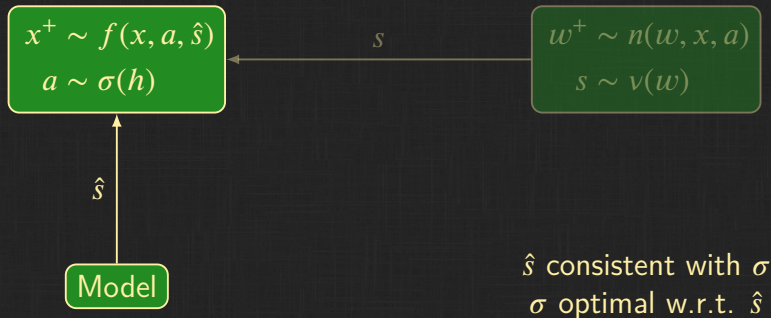
# Single Agent Setup



# Single Agent Setup



# Single Agent Setup



# Depth- $k$ Consistency

Consider a binary stochastic process  $s$

0100010001001010010110111010000111010101...



# Depth- $k$ Consistency

Consider a binary stochastic process  $s$

0100010001001010010110111010000111010101...

- 0 characteristic:  $\mathbb{P}[s = 0], \mathbb{P}[s = 1]$
- 1 characteristic:  $\mathbb{P}[ss^+ = 00], \mathbb{P}[ss^+ = 10],$   
 $\mathbb{P}[ss^+ = 01], \mathbb{P}[ss^+ = 11]$
- ...
- $k$  characteristic: probability of strings of length  $k + 1$

# Depth- $k$ Consistency

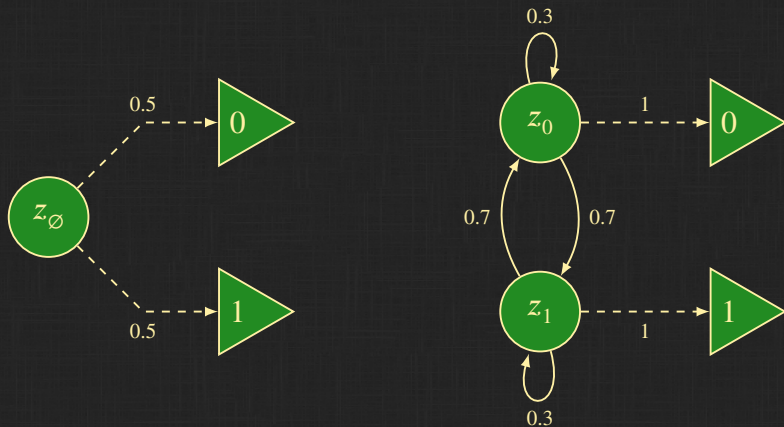
Consider a binary stochastic process  $s$

0100010001001010010110111010000111010101...

- 0 characteristic:  $\mathbb{P}[s = 0], \mathbb{P}[s = 1]$
- 1 characteristic:  $\mathbb{P}[ss^+ = 00], \mathbb{P}[ss^+ = 10],$   
 $\mathbb{P}[ss^+ = 01], \mathbb{P}[ss^+ = 11]$
- ...
- $k$  characteristic: probability of strings of length  $k + 1$

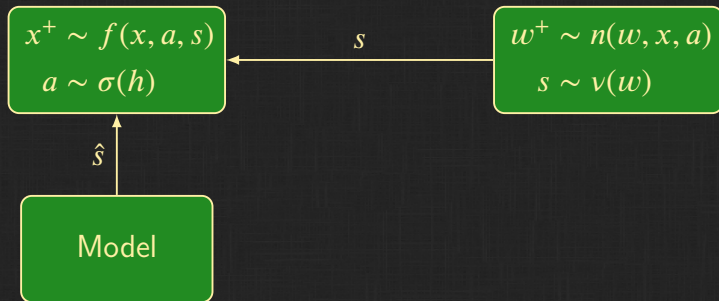
**Definition** Two processes  $s$  and  $s'$  are depth- $k$  consistent if they have the same  $k$  characteristic

# Depth- $k$ Consistency: Example



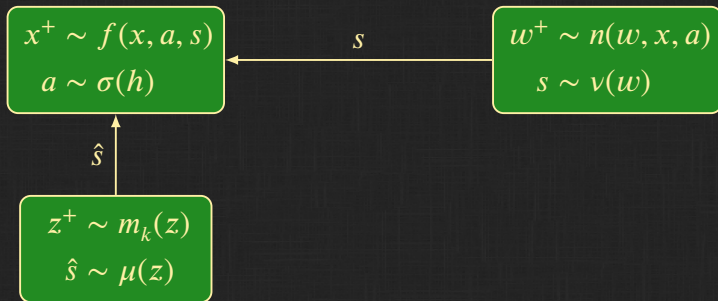
# Complete picture

Fix a depth  $k \in \mathbb{N}$



# Complete picture

Fix a depth  $k \in \mathbb{N}$

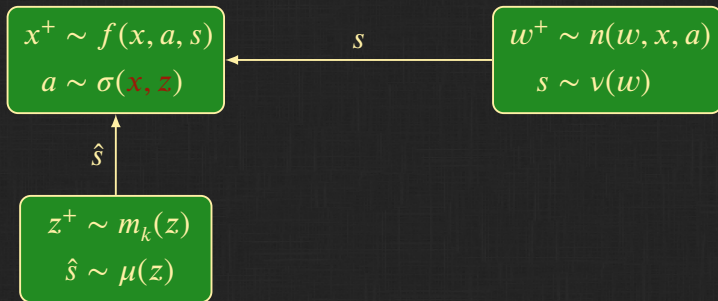


$z$  contains the last  $k$  observed signals

$$\mu(z = (s_1, s_2, \dots, s_k)) [s_{k+1}] = \mathbb{P}_\sigma [s^{t+1} = s_{k+1} \mid s^t = s_k, \dots, s^{t-k+1} = s_1]$$

# Complete picture

Fix a depth  $k \in \mathbb{N}$

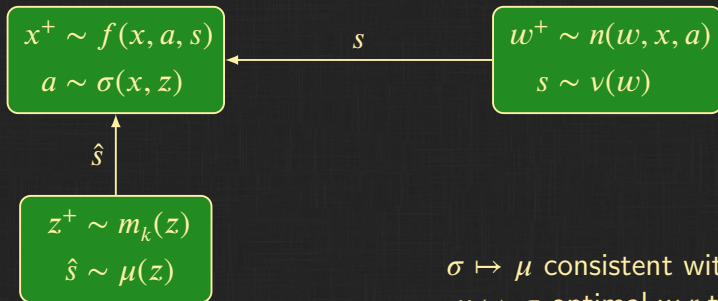


$z$  contains the last  $k$  observed signals

$$\mu(z = (s_1, s_2, \dots, s_k)) [s_{k+1}] = \mathbb{P}_\sigma [s^{t+1} = s_{k+1} \mid s^t = s_k, \dots, s^{t-k+1} = s_1]$$

# Complete picture

Fix a depth  $k \in \mathbb{N}$



$\sigma \mapsto \mu$  consistent with  $\sigma$   
 $\mu \mapsto \sigma$  optimal w.r.t.  $\mu$

$z$  contains the last  $k$  observed signals

$$\mu(z = (s_1, s_2, \dots, s_k)) [s_{k+1}] = \mathbb{P}_\sigma [s^{t+1} = s_{k+1} \mid s^t = s_k, \dots, s^{t-k+1} = s_1]$$

# Definition

$(\sigma, \mu)$  is an empirical-evidence optimum (EEO) for  $k$  iff

- $\sigma$  is optimal w.r.t.  $\mu$
- $\mu$  is depth- $k$  consistent with  $\sigma$



# Definition

$(\sigma, \mu)$  is an empirical-evidence optimum (EEO) for  $k$  iff

- $\sigma$  is optimal w.r.t.  $\mu$
- $\mu$  is depth- $k$  consistent with  $\sigma$

$(\sigma, \mu)$  is an  $\epsilon$  empirical-evidence optimum ( $\epsilon$  EEO) for  $k$  iff

- $\sigma$  is  $\epsilon$  optimal w.r.t.  $\mu$
- $\mu$  is depth- $k$  consistent with  $\sigma$

# Existence Result

## Theorem

For all  $k$  and  $\epsilon$ , there exists an  $\epsilon$  EEO for  $k$

# Existence Result

## Theorem

For all  $k$  and  $\epsilon$ , there exists an  $\epsilon$  EEO for  $k$

## Proof sketch

Prove continuity of  $\sigma \mapsto \mu \mapsto \sigma$

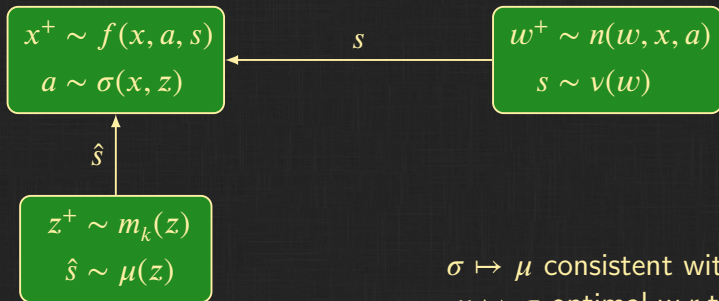
$$\sigma : \mathcal{X} \times \mathcal{L} \rightarrow \Delta(\mathcal{A})$$

$\sigma$  parametrized over a simplex (convex and compact)

Apply Brouwer's fixed point theorem

# Complete picture

Fix a depth  $k \in \mathbb{N}$

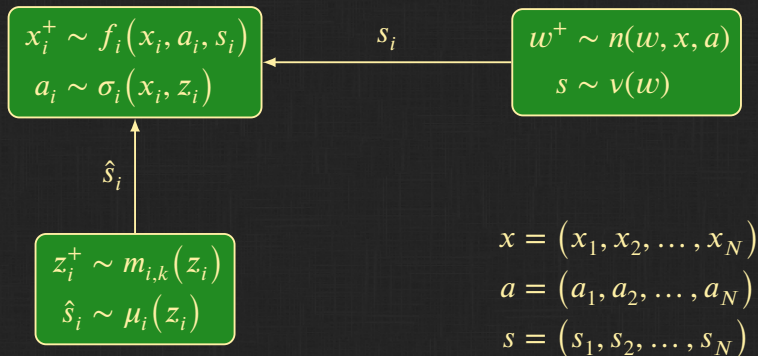


$\sigma \mapsto \mu$  consistent with  $\sigma$   
 $\mu \mapsto \sigma$  optimal w.r.t.  $\mu$

$z$  contains the last  $k$  observed signals

$$\mu(z = (s_1, s_2, \dots, s_k)) [s_{k+1}] = \mathbb{P}_\sigma [s^{t+1} = s_{k+1} \mid s^t = s_k, \dots, s^{t-k+1} = s_1]$$

# Multiagent Setting



# Empirical-evidence Equilibrium

$(\sigma, \mu)$  is an empirical-evidence equilibrium (EEE) for  $K = (k_1, k_2, \dots, k_N)$  iff

- for all  $i$ ,  $\sigma_i$  is optimal w.r.t.  $\mu_i$
- for all  $i$ ,  $\mu_i$  is depth- $k_i$  consistent with  $\sigma$

# Empirical-evidence Equilibrium

$(\sigma, \mu)$  is an empirical-evidence equilibrium (EEE) for  $K = (k_1, k_2, \dots, k_N)$  iff

- for all  $i$ ,  $\sigma_i$  is optimal w.r.t.  $\mu_i$
- for all  $i$ ,  $\mu_i$  is depth- $k_i$  consistent with  $\sigma$

## Theorem

For all  $K$  and  $\epsilon$ , there exists an  $\epsilon$  EEE for  $K$

# Open Questions

- endogenous model depending on action
- large number of agents
- large  $k$
- relating EEE to other concepts (MFE, optimum)
- offline computation
- online learning using empirical evidence



# Open Questions

- endogenous model depending on action
- large number of agents
- large  $k$
- relating EEE to other concepts (MFE, optimum)
- offline computation
- online learning using empirical evidence

# Open Questions

- endogenous model depending on action
- large number of agents
- large  $k$
- relating EEE to other concepts (MFE, optimum)
- offline computation
- online learning using empirical evidence

# Open Questions

- endogenous model depending on action
- large number of agents
- large  $k$
- relating EEE to other concepts (MFE, optimum)
- offline computation
- online learning using empirical evidence

# Example: Asset Management

State holdings  $x_i \in \{0..M\}$

Action sell one, hold, or buy one  $a_i \in \{-1, 0, 1\}$

Signal price  $p \in \{\text{Low}, \text{High}\}$

Dynamic  $x_i^+ = x_i + a_i$

Stage cost  $p \cdot a_i$

Nature market trend  $b \in \{\text{Bull}, \text{Bear}\}$

$w = (b, p)$

Nature is a *sticky bear*

# Example: Asset Management

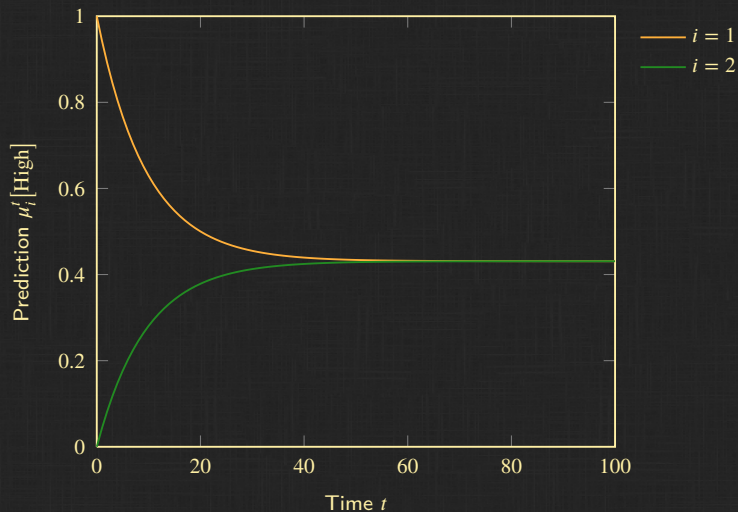
0. Pick arbitrary models  $\mu$
1. Design strategies  $\sigma$  optimal w.r.t. models  $\mu$
2. Formulate consistent models  $\mu_{\text{upd}}$ , then, back to 1.

Depth-0 consistency:

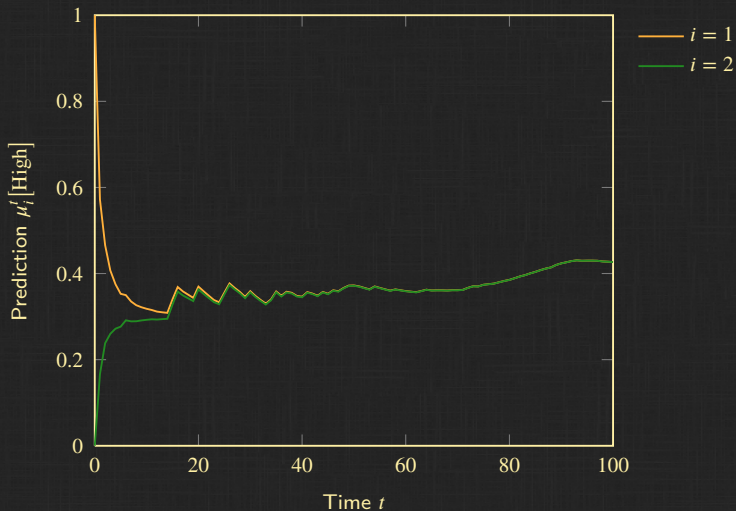
- $\mu_1 = 1$
- $\mu_2 = 0$

$$\mu_i^{t+1} = (1 - \alpha)\mu_i^t + \alpha\left(\mu_{i,\text{upd}}^t - \mu_i^t\right)$$

# Learning Results: Offline



# Learning Results: Online



# Empirical-evidence Equilibria

- Introduce
- Contrast
- Compute