# Data-Driven
# Markov Decision Processes

**Wolfram Wiesemann**
Imperial College Business School

**Markov decision process**

Tuple $(\mathcal{S}, \mathcal{A}, q, p, r, \lambda)$ where

- $\mathcal{S} = \{1, \ldots, S\}$ is the (finite) state space;

- $\mathcal{A} = \{1, \ldots, A\}$ is the (finite) action space;

- $q = (q_1, \ldots, q_S) \in \Delta(\mathcal{S})$ is the initial state distribution;

- $p : \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$ is the transition kernel with elements $p(s' | s, a)$;

- $r : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ are the expected one-step rewards;

- $\lambda \in (0, 1)$ is the discount factor.
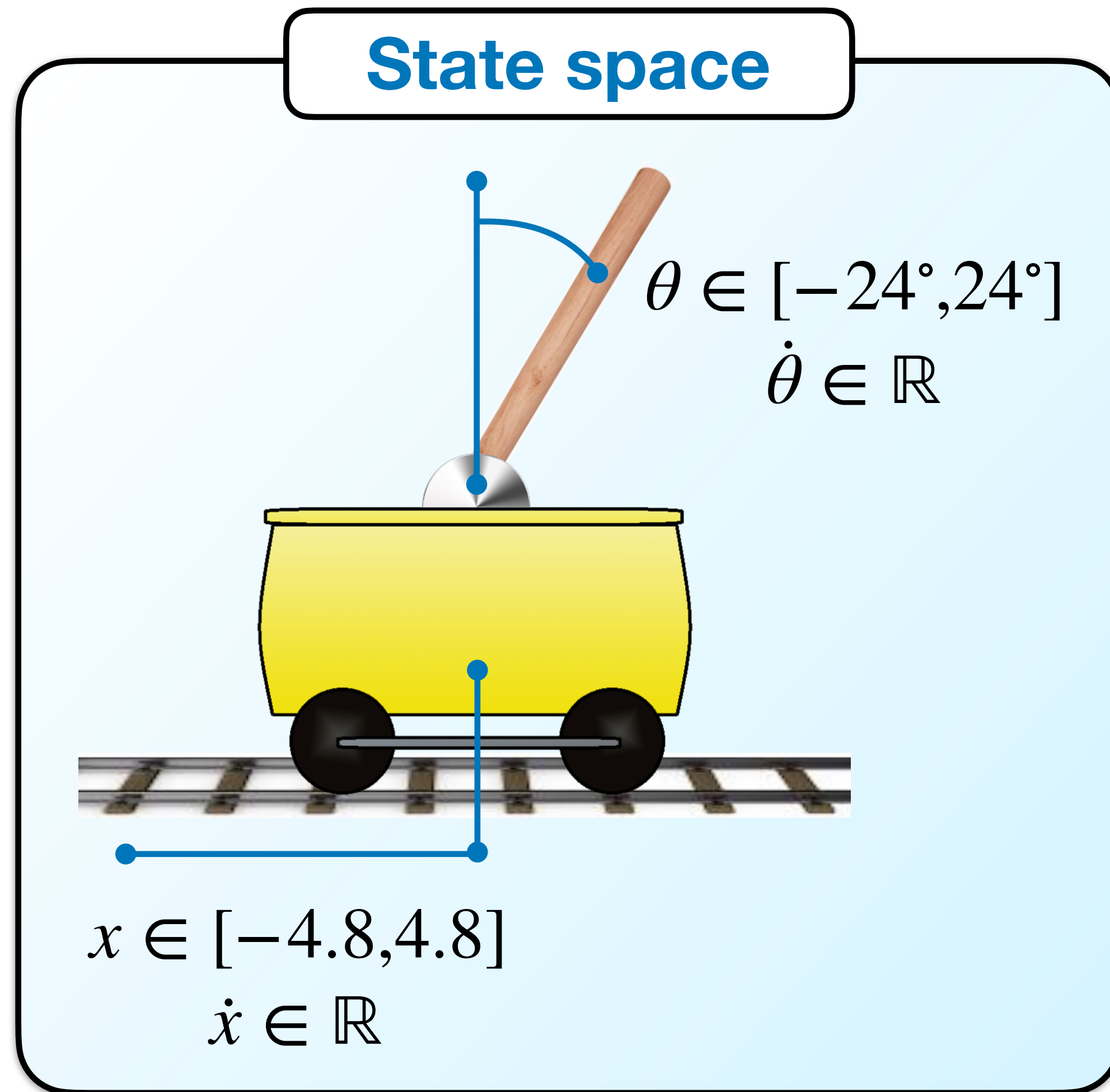
**Markov decision process**

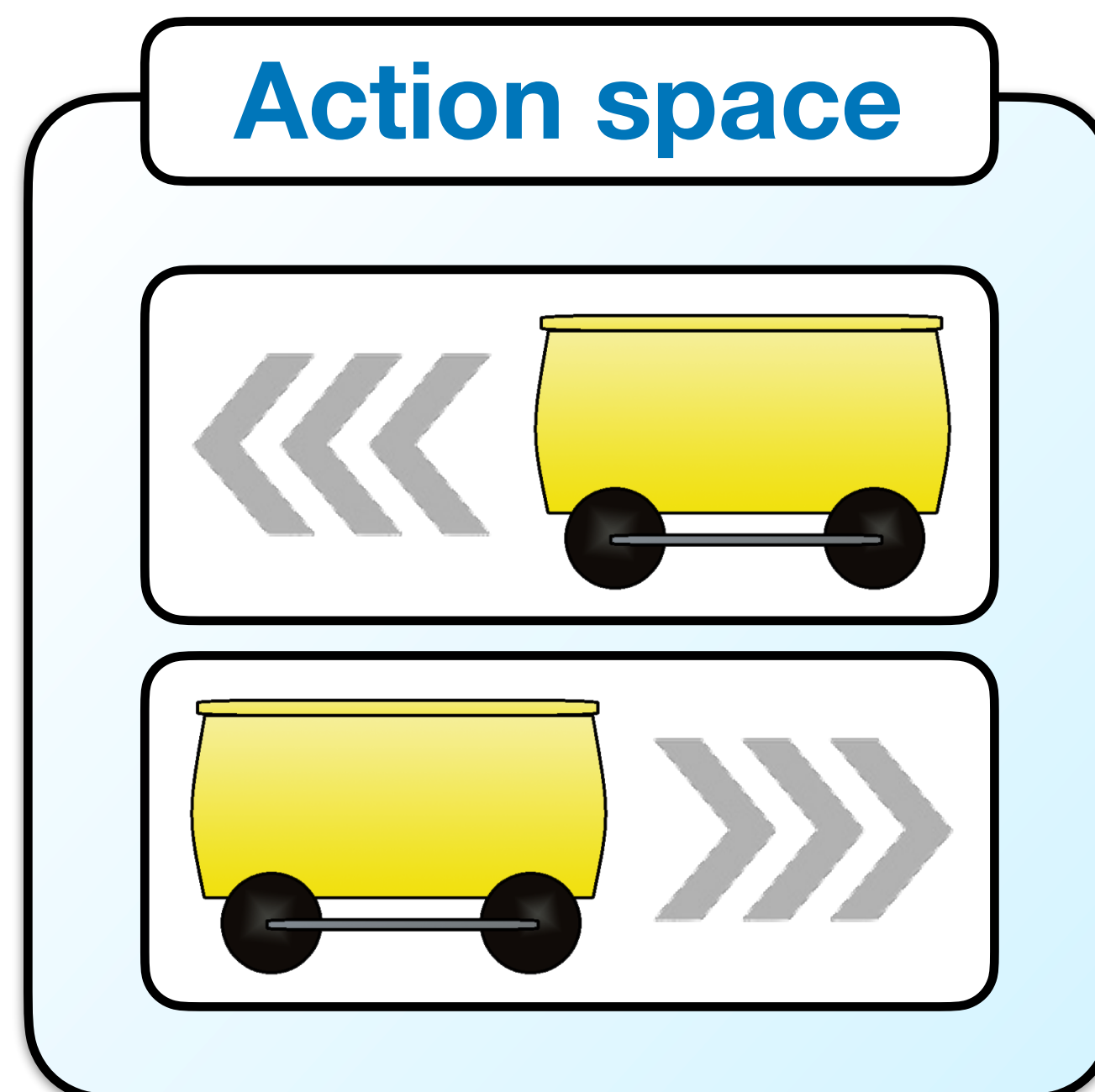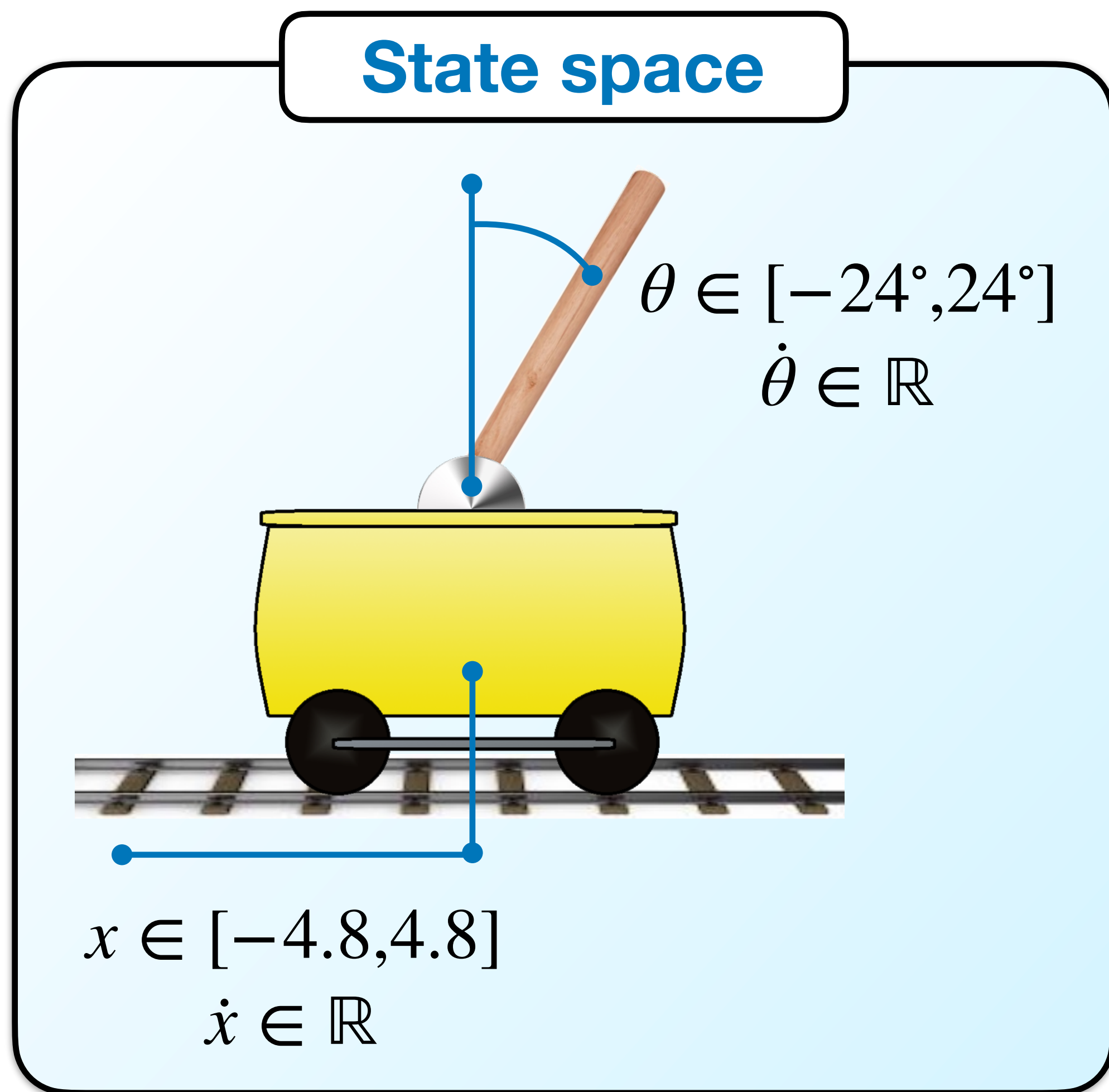Tuple $(\mathcal{S}, \mathcal{A}, q, p, r, \lambda)$ where

- $\mathcal{S} = \{1, \ldots, S\}$ is the (finite) state space;
- $\mathcal{A} = \{1, \ldots, A\}$ is the (finite) action space;
- $q = (q_1, \ldots, q_S) \in \Delta(\mathcal{S})$ is the initial state distribution;
- $p : \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$ is the transition kernel with elements $p(s' \,|\, s, a)$;
- $r : \mathcal{S} \times$
- $\lambda \in (0, 1)$

**Objective**

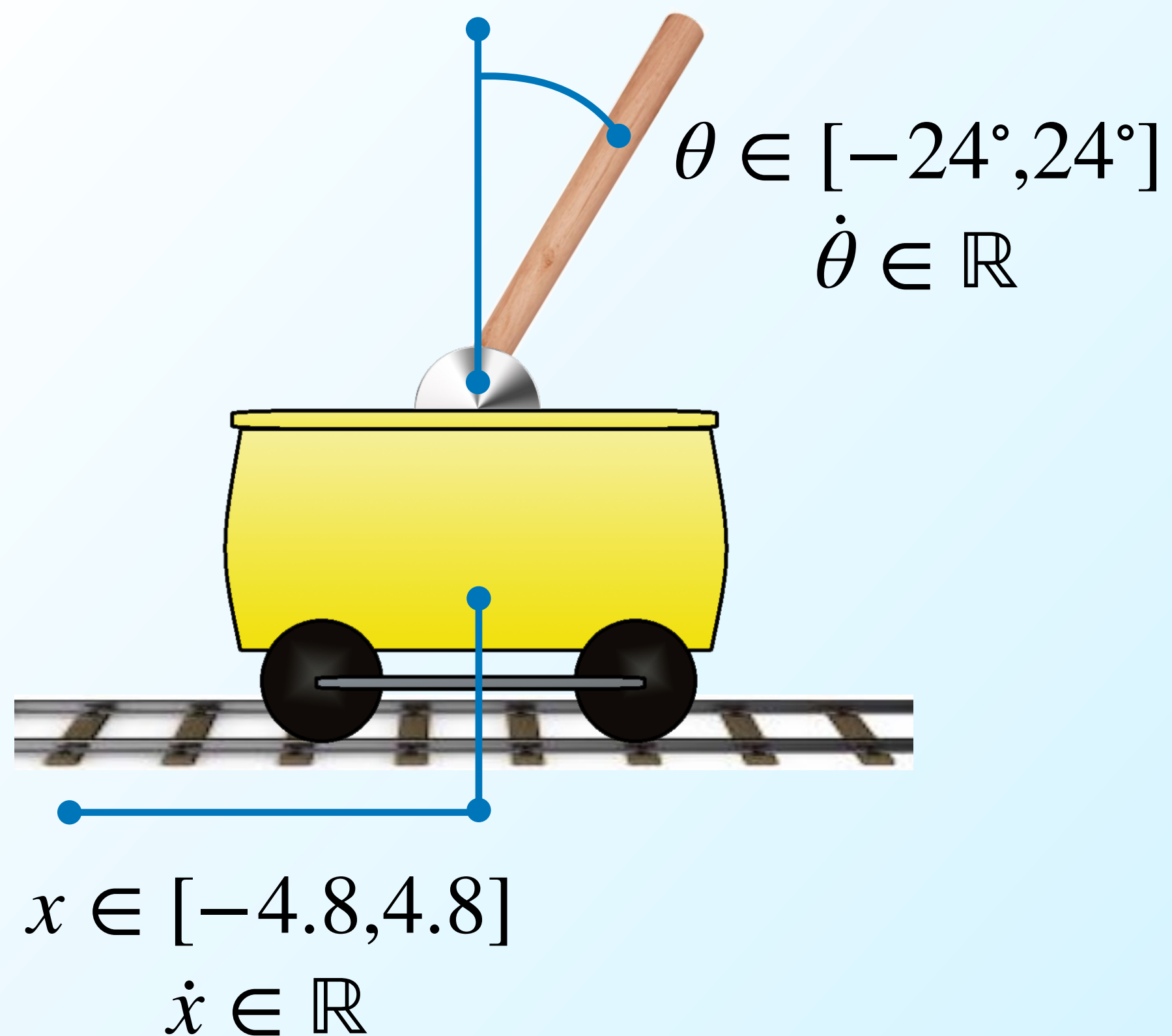find policy $\pi : \mathcal{S} \to \mathcal{A}$ that maximizes the expected total discounted rewards:

$$\underset{\pi \in \Pi}{\text{maximize}} \quad \mathbb{E}_p \left[ \sum_{t=1}^{\infty} \lambda^{t-1} \cdot r\big(s_t, \pi[s_t]\big) \right]$$

1

**State space**

$\theta \in [-24°, 24°]$

$\dot{\theta} \in \mathbb{R}$

$x \in [-4.8, 4.8]$

$\dot{x} \in \mathbb{R}$

Barto et al. (1983), *Neuronlike Adaptive Elements that can Solve Difficult Learning Control Problems*.

**State space**

$\theta \in [-24°, 24°]$
$\dot{\theta} \in \mathbb{R}$

$x \in [-4.8, 4.8]$
$\dot{x} \in \mathbb{R}$

**Action space**

Barto et al. (1983), *Neuronlike Adaptive Elements that can Solve Difficult Learning Control Problems*.

**State space**

$\theta \in [-24°, 24°]$
$\dot{\theta} \in \mathbb{R}$
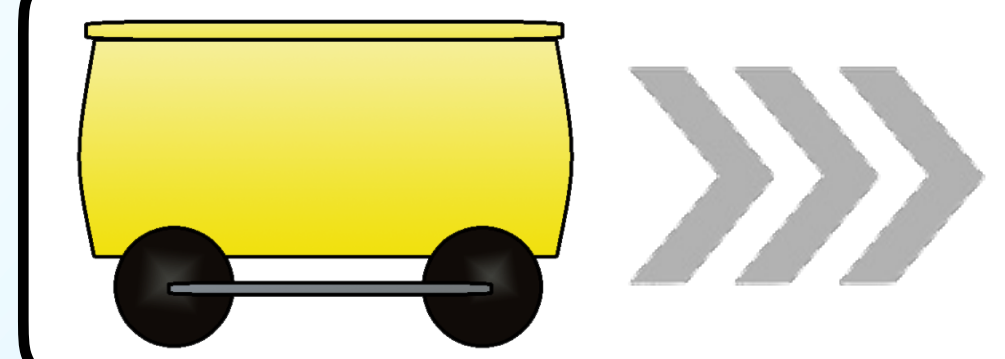
$x \in [-4.8, 4.8]$
$\dot{x} \in \mathbb{R}$

**Action space**

**Initial state**

$x, \dot{x}, \theta, \dot{\theta} \sim$
$\mathscr{U}[-0.05, 0.05]$

Barto et al. (1983), *Neuronlike Adaptive Elements that can Solve Difficult Learning Control Problems*.

**State space**

$\theta \in [-24°, 24°]$
$\dot{\theta} \in \mathbb{R}$

$x \in [-4.8, 4.8]$
$\dot{x} \in \mathbb{R}$

**Action space**

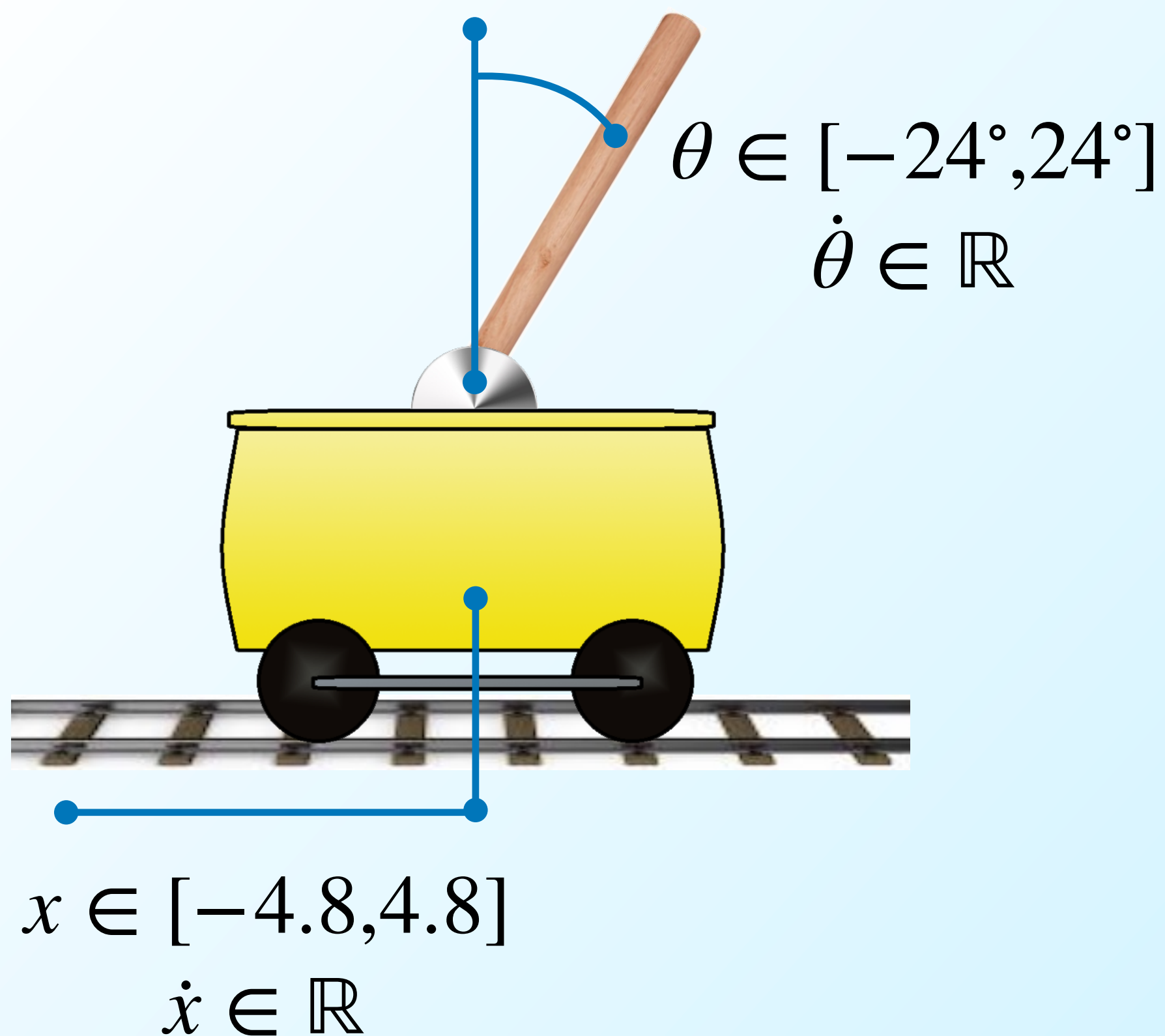**Initial state**

$x, \dot{x}, \theta, \dot{\theta} \sim$
$\mathcal{U}[-0.05, 0.05]$

**Transitions**

- deterministic via laws of mechanics
- terminate if
$x \notin [-2.4, 2.4]$
or $\theta \notin [-12°, 12°]$

Barto et al. (1983), *Neuronlike Adaptive Elements that can Solve Difficult Learning Control Problems*.

# Cart Pole Example

**State space**

$\theta \in [-24°, 24°]$
$\dot{\theta} \in \mathbb{R}$
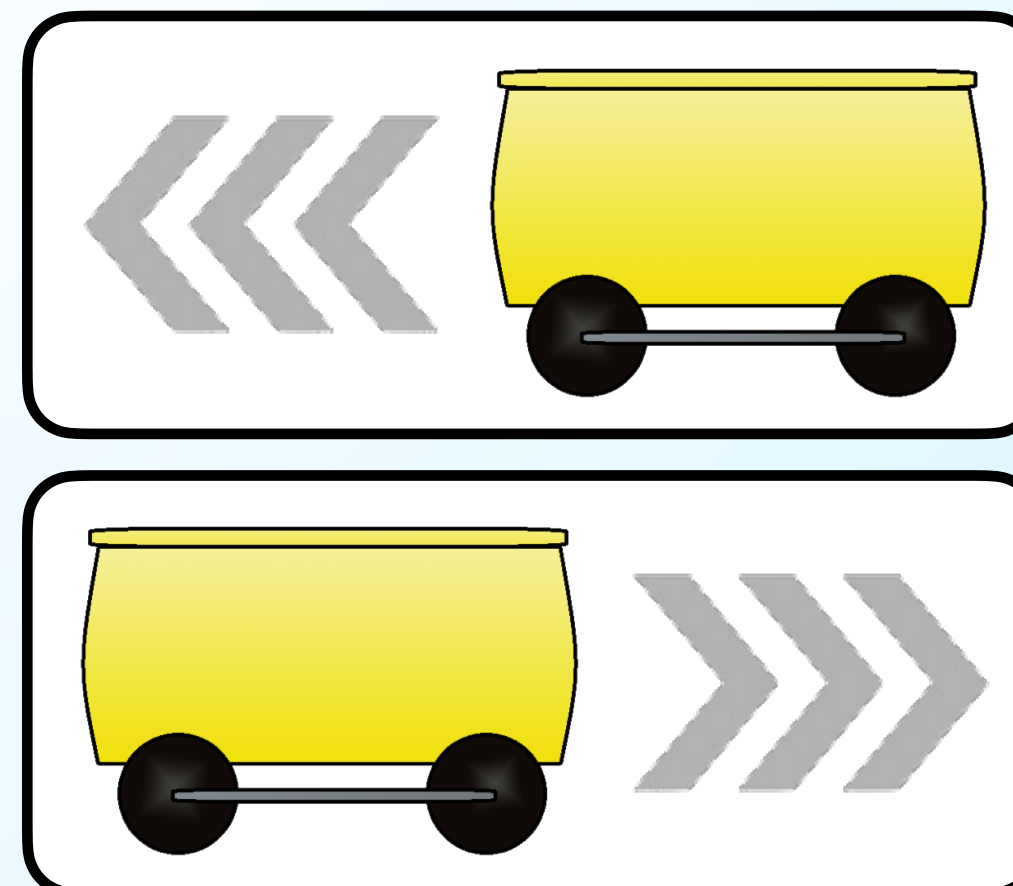
$x \in [-4.8, 4.8]$
$\dot{x} \in \mathbb{R}$

**Action space**

**Transitions**

- deterministic via laws of mechanics
- terminate if
  $x \notin [-2.4, 2.4]$
  or $\theta \notin [-12°, 12°]$

**Initial state**

$x, \dot{x}, \theta, \dot{\theta} \sim$
$\mathcal{U}[-0.05, 0.05]$

**Rewards**

+1/non-terminated time step

Barto et al. (1983), *Neuronlike Adaptive Elements that can Solve Difficult Learning Control Problems*.

Barto et al. (1983), *Neuronlike Adaptive Elements that can Solve Difficult Learning Control Problems*.

# STOCHASTICITY AND AMBIGUITY

**Two common sources of ambiguity:**

- **Modelling errors:** 32.67 secs/run



$15^4 = 50,625$ states

**Two common sources of ambiguity:**

- **Modelling errors:** 32.67 secs/run ➡ 2.45 secs/run



$$10^4 = 10,000 \text{ states}$$

**Two common sources of ambiguity:**

- **Modelling errors:** 32.67 secs/run ➡️ 2.45 secs/run
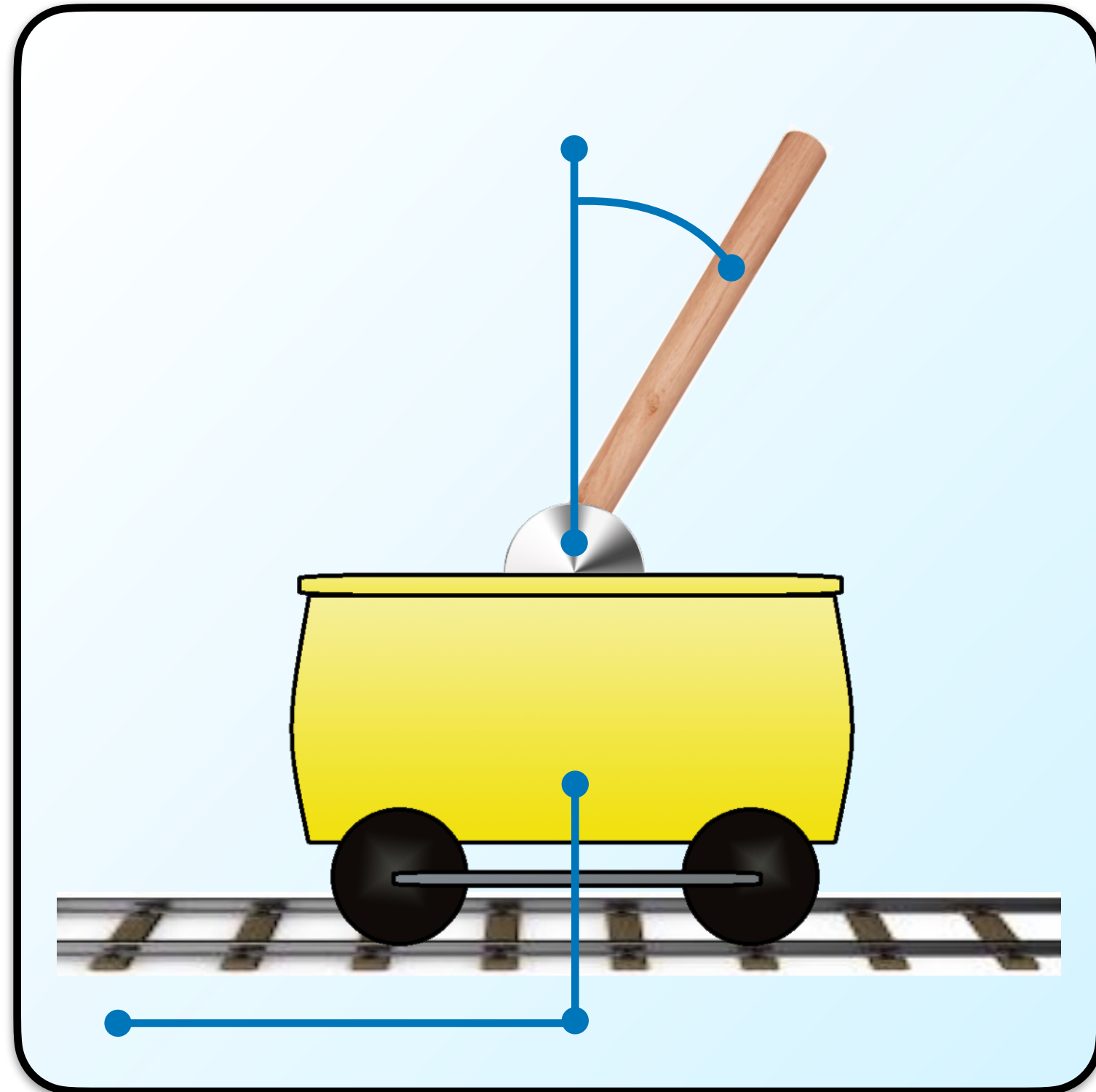
- **Estimation errors:** 32.67 secs/run

**Two common sources of ambiguity:**

- **Modelling errors**: 32.67 secs/run ➡ 2.45 secs/run
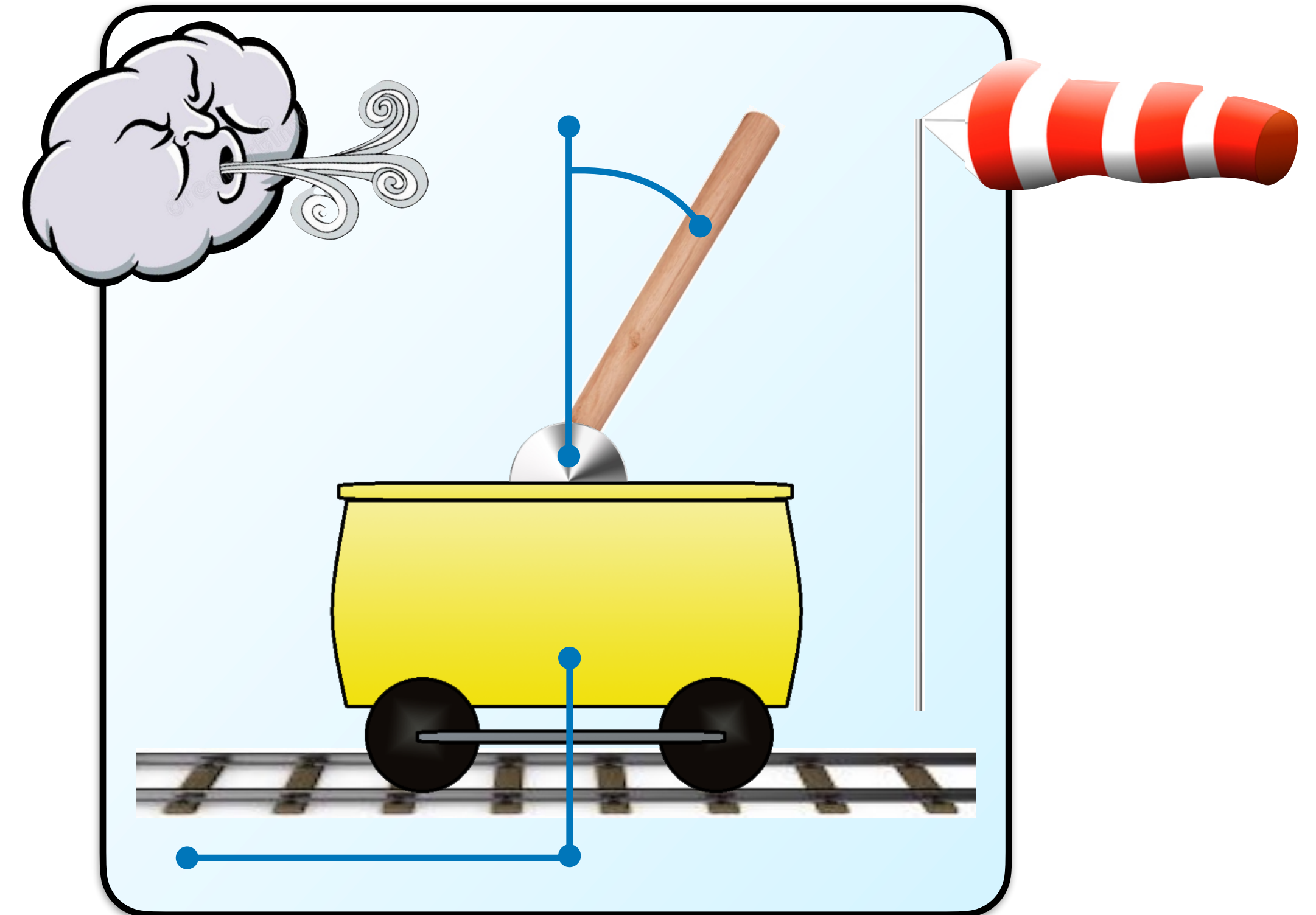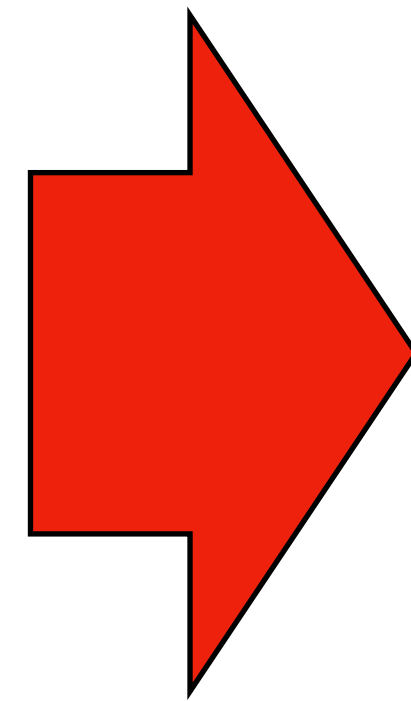
- **Estimation errors**: 32.67 secs/run ➡ 4.68 secs/run

**Two common sources of ambiguity:**

- **Modelling errors:** 32.67 secs/run ➡ 2.45 secs/run

- **Estimation errors:** 32.67 secs/run ➡ 4.68 secs/run

**Impact of ambiguity can be alleviated via robust optimization:**



**Robust MDP**

$$\underset{\pi \in \Pi}{\text{maximize}} \quad \underset{p \in \mathscr{P}}{\inf} \; \mathbb{E}_p \left[ \sum_{t=1}^{\infty} \lambda^{t-1} \cdot r\big(s_t, \pi[s_t]\big) \right]$$

*ambiguity set*

***Robust MDPs** admit interpretation as **regularized MDPs**!*

Derman et al. (2023), *Twice Regularized Markov Decision Processes: The Equivalence between Robustness and Regularization*.

**Modelling errors:** 32.67 secs/run ➡ 2.45 secs/run ➡ 15.77 secs/run

5

**Estimation errors:** 32.67 secs/run ➡ 4.68 secs/run ➡ 15.76 secs/run

**Structural ambiguity set**



$\mathscr{P}^0$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Structural ambiguity set**

**Historical sample**

$\mathscr{P}(\mathscr{H}_n)$

$\mathscr{P}^0$

$\cap$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Structural ambiguity set** $\cap$ **Historical sample** $=$ **Out-of-sample guarantee**

$\mathscr{P}(\mathscr{H}_n)$

$\mathscr{P}^0$

$\mathscr{P}(\mathscr{H}_n)$

$\mathscr{P}(\mathscr{H}_n)$

$\mathscr{P}_n$

$\mathscr{P}^0$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Structural ambiguity set**

$$\mathscr{P}^0 \subseteq \{p : \mathcal{S} \times \mathscr{A} \to \Delta(\mathcal{S})\}$$

$p^0 \in \text{rel int } \mathscr{P}^0$

$\mathscr{P}^0$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Structural ambiguity set**

$$\mathscr{P}^0 \subseteq \{p : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})\}$$

$p^0 \in \text{rel int } \mathscr{P}^0$

**Possible transitions**

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

Structural ambiguity set

$$\mathscr{P}^0 \subseteq \{p : \mathscr{S} \times \mathscr{A} \to \Delta(\mathscr{S})\}$$

$$p^0 \in \text{rel int } \mathscr{P}^0$$

**Possible transitions**

**Equal probabilities**

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Historical sample**

historical policy $\pi^0$
(stationary, randomized)

state-action history
$\mathscr{H}_n = (s_1, a_1, \ldots, s_n, a_n) \in (\mathscr{S} \times \mathscr{A})^n$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Historical sample**

historical policy $\pi^0$
(stationary, randomized)

state-action history
$$\mathscr{H}_n = (s_1, a_1, \ldots, s_n, a_n) \in (\mathscr{S} \times \mathscr{A})^n$$

**Likelihood, given history**

$$\mathscr{L}_n(p) \;=\; q(s_1) \cdot \pi^0(a_n \,|\, s_n) \cdot \prod_{t=1}^{n-1} \left[ \pi^0(a_t \,|\, s_t) \cdot p(s_{t+1} \,|\, s_t, a_t) \right]$$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Historical sample**

$$\mathscr{P}(\mathscr{H}_n) = \left\{ p : \log \mathscr{L}_n(p) \geq \log \mathscr{L}_n(p^\star) - \delta \right\}$$

historical policy $\pi^0$
(stationary, randomized)

state-action history
$$\mathscr{H}_n = (s_1, a_1, \ldots, s_n, a_n) \in (\mathscr{S} \times \mathscr{A})^n$$

**Likelihood, given history**

$$\mathscr{L}_n(p) \;=\; q(s_1) \cdot \pi^0(a_n \,|\, s_n) \cdot \prod_{t=1}^{n-1} \left[ \pi^0(a_t \,|\, s_t) \cdot p(s_{t+1} \,|\, s_t, a_t) \right]$$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Theorem**

**Assumption:** Historical policy $\pi^0$ visits every $s \in \mathcal{S}$ infinite often as $n \longrightarrow \infty$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

9

**Theorem**

$$\mathscr{P}_n = \mathscr{P}^0 \cap \mathscr{P}(\mathscr{H}_n) \text{ with } \delta = (1 - \beta)\text{-quantile}$$
$$\text{of } \chi^2\text{-distribution with } \kappa \text{ degrees of freedom}$$

**Assumption:** Historical policy $\pi^0$ visits every $s \in \mathcal{S}$ infinite often as $n \longrightarrow \infty$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Theorem**

$$\mathscr{P}_n = \mathscr{P}^0 \cap \mathscr{P}(\mathscr{H}_n) \text{ with } \delta = (1 - \beta)\text{-quantile}$$
$$\text{of } \chi^2\text{-distribution with } \kappa \text{ degrees of freedom}$$

**Assumption:** Historical policy $\pi^0$ visits every $s \in \mathscr{S}$ infinite often as $n \longrightarrow \infty$

**1**

$$\lim_{n \longrightarrow \infty} \mathbb{P}\left[p^0 \in \mathscr{P}_n\right] = 1 - \beta$$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Theorem**

$$\mathscr{P}_n = \mathscr{P}^0 \cap \mathscr{P}(\mathscr{H}_n) \text{ with } \delta = (1 - \beta)\text{-quantile}$$
$$\text{of } \chi^2\text{-distribution with } \kappa \text{ degrees of freedom}$$

**Assumption:** Historical policy $\pi^0$ visits every $(s, a)$ infinite often as $n \longrightarrow \infty$

**1**
$$\lim_{n \longrightarrow \infty} \mathbb{P}\left[p^0 \in \mathscr{P}_n\right] = 1 - \beta$$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

## Theorem

$$\mathscr{P}_n = \mathscr{P}^0 \cap \mathscr{P}(\mathscr{H}_n) \text{ with } \delta = (1 - \beta)\text{-quantile}$$
$$\text{of } \chi^2\text{-distribution with } \kappa \text{ degrees of freedom}$$

**Assumption:** Historical policy $\pi^0$ visits every $(s, a)$ infinite often as $n \longrightarrow \infty$

**1**
$$\lim_{n \longrightarrow \infty} \mathbb{P}\left[p^0 \in \mathscr{P}_n\right] = 1 - \beta$$

**2**
$$\operatorname*{plim}_{n \longrightarrow \infty}\left[\sqrt{n} \cdot d^{\mathsf{H}}(\mathscr{P}_n, \{p^0\})\right] = 0$$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

**Theorem**

$$\mathscr{P}_n = \mathscr{P}^0 \cap \mathscr{P}(\mathscr{H}_n) \text{ with } \delta = (1-\beta)\text{-quantile}$$
$$\text{of } \chi^2\text{-distribution with } \kappa \text{ degrees of freedom}$$

**Assumption:** Historical policy $\pi^0$ visits every $(s, a)$ infinite often as $n \longrightarrow \infty$

**1** $\quad \lim_{n \longrightarrow \infty} \mathbb{P}\left[p^0 \in \mathscr{P}_n\right] = 1 - \beta$

**2** $\quad \mathrm{plim}_{n \longrightarrow \infty} \left[\sqrt{n} \cdot d^{\mathsf{H}}(\mathscr{P}_n, \{p^0\})\right] = 0$



$\mathscr{P}_{100}$ $\quad$ $\mathscr{P}_{300}$ $\quad$ $\mathscr{P}_{500}$ $\quad$ $p^0$

$\mathscr{P}_{400}$

$\mathscr{P}_{200}$

Wiesemann et al. (2013), *Robust Markov Decision Processes*.

$$\mathcal{P} \subseteq \Delta^{S \times A}$$

**General (non-rectangular) ambiguity sets**

👎 Optimal policy can be randomized & history-dependent

👎 Bellman optimality principle violated; NP-hard

$$\mathcal{P} \subseteq \Delta^{S \times A}$$

**General (non-rectangular) ambiguity sets**

👎 Optimal policy can be randomized & history-dependent

👎 Bellman optimality principle violated; NP-hard

**(s,a)-rectangular ambiguity sets**

$$\mathscr{P} = \prod_{(s,a) \in \mathscr{S} \times \mathscr{A}} \mathscr{P}_{s,a} \quad \text{with} \quad \mathscr{P}_{s,a} \subseteq \Delta(\mathscr{S})$$

$$\mathcal{P} = \times \, \mathcal{P}_{sa} \qquad\qquad \mathcal{P} = \times \, \mathcal{P}_{s} \qquad\qquad \mathcal{P} \subseteq \Delta^{S \times A}$$

$$\mathcal{P} \subseteq \Delta^{S \times A}$$

**General (non-rectangular) ambiguity sets**

👎 Optimal policy can be randomized & history-dependent

👎 Bellman optimality principle violated; NP-hard



**(s,a)-rectangular ambiguity sets**

👍 Optimal policy stationary and deterministic

👍 Bellman optimality principle holds

$$\mathcal{P} = \times \mathcal{P}_{sa} \qquad \mathcal{P} = \times \mathcal{P}_s \qquad \mathcal{P} \subseteq \Delta^{S \times A}$$

General (non-rectangular) ambiguity sets

**Example**



Action 1

Action 2

for some unknown $\xi \in [0,1]$

Bellman optimality principle holds

$\mathcal{P} \subseteq \Delta^{S}$

endent

10

$\mathcal{P} = \times \mathcal{P}_{sa}$        $\mathcal{P} = \times \mathcal{P}_{s}$        $\mathcal{P} \subseteq \Delta^{S \times A}$

**General (non-rectangular) ambiguity sets**

**Example**



**Action 1**

**Action 2**

for some unknown $\xi_1, \xi_2 \in [0,1]$

Bellman optimality principle holds

$\mathcal{P} = \times \mathcal{P}_{sa}$          $\mathcal{P} = \times \mathcal{P}_s$          $\mathcal{P} \subseteq \Delta^{S \times A}$

$\mathcal{P} \subseteq \Delta^{S}$

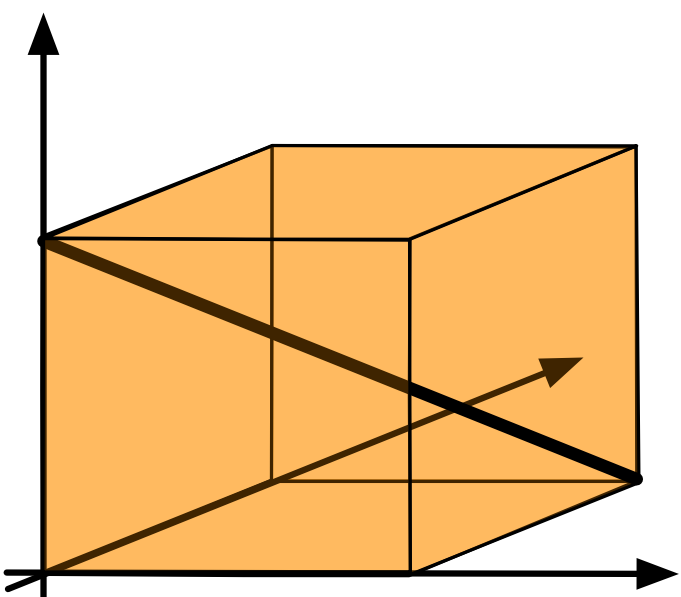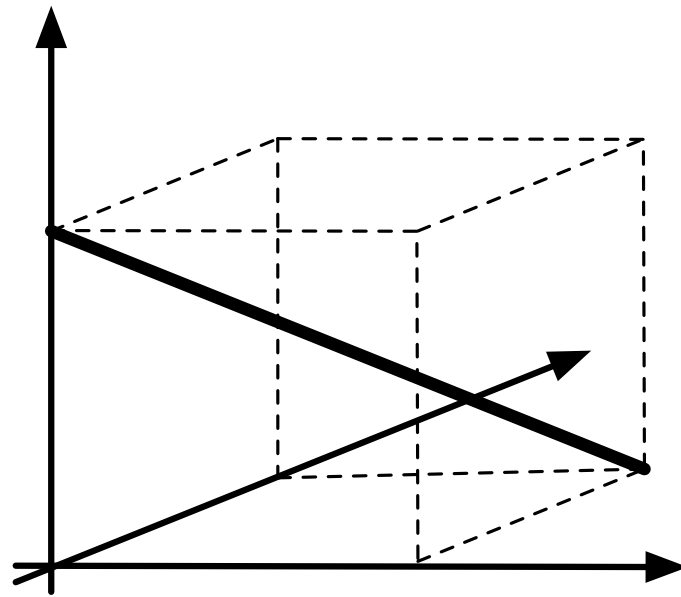**General (non-rectangular) ambiguity sets**

👎 Optimal policy can be randomized & history-dependent

👎 Bellman optimality principle violated; NP-hard

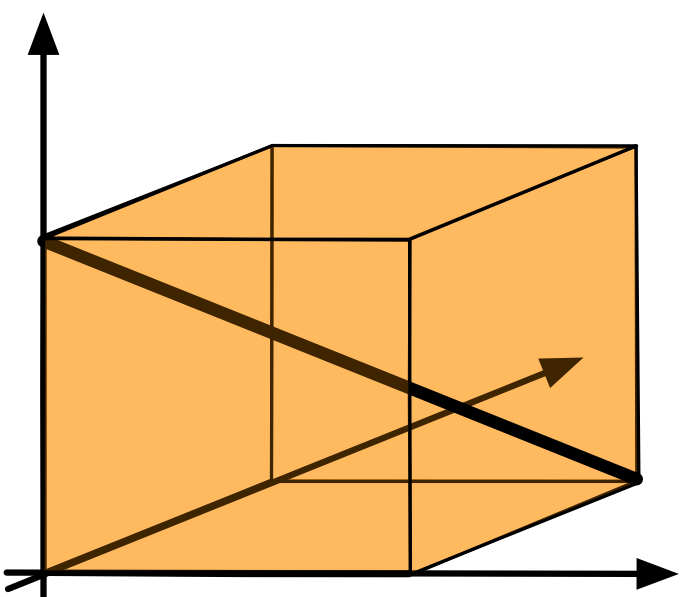$$\mathcal{P} \subseteq \Delta^{S \times A}$$

**s-rectangular ambiguity sets**

$$\mathscr{P} = \prod_{s \in \mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s \subseteq [\Delta(\mathscr{S})]^A$$

**(s,a)-rectangular ambiguity sets**

$$\mathcal{P} \subseteq \Delta^{S \times}$$

👍 Optimal policy stationary and deterministic

👍 Bellman optimality principle holds

$$\mathcal{P} = \times \, \mathcal{P}_{sa} \qquad \mathcal{P} = \times \, \mathcal{P}_s \qquad \mathcal{P} \subseteq \Delta^{S \times A}$$

**General (non-rectangular) ambiguity sets**

👎 Optimal policy can be randomized & history-dependent

👎 Bellman optimality principle violated; NP-hard

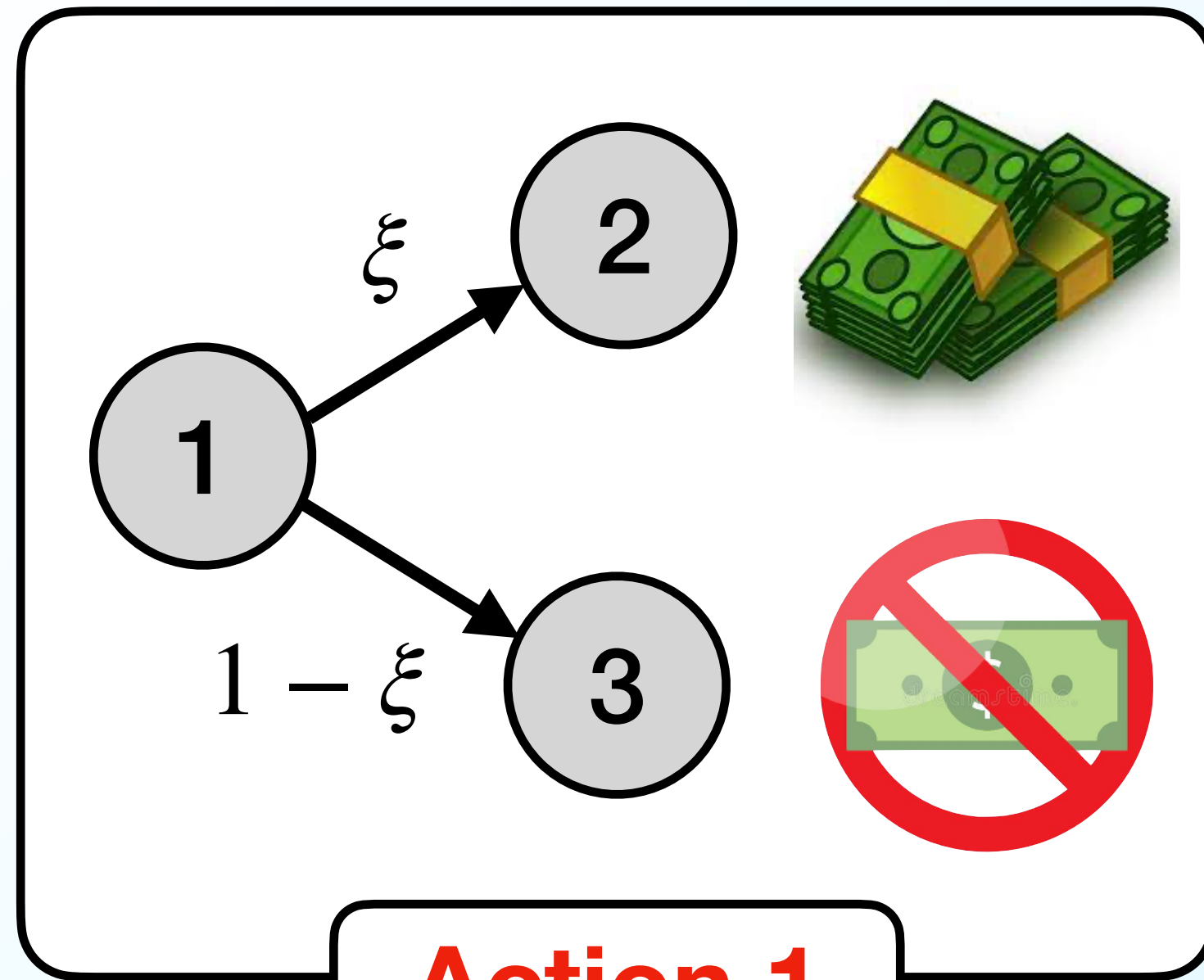$\mathcal{P} \subseteq \Delta^{S \times A}$

**s-rectangular ambiguity sets**

👍 Optimal policy stationary but can be *randomized*

👍 Bellman optimality principle holds

$\mathcal{P} = \times \mathcal{P}$

$\mathcal{P} \subseteq \Delta^{S \times}$

**(s,a)-rectangular ambiguity sets**

👍 Optimal policy stationary and deterministic

👍 Bellman optimality principle holds

$$\mathcal{P} = \times \mathcal{P}_{sa} \qquad \mathcal{P} = \times \mathcal{P}_s \qquad \mathcal{P} \subseteq \Delta^{S \times A}$$

**General (non-rectangular) ambiguity sets**

**Example**

Action 1

Action 2

for some unknown $\xi \in [0,1]$

Bellman optimality principle holds

$\mathcal{P} = \bigtimes \mathcal{P}_{sa}$          $\mathcal{P} = \bigtimes \mathcal{P}_s$          $\mathcal{P} \subseteq \Delta^{S \times A}$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Classical (non-robust) Bellman equations**

$$v^\star(s) = \max_{a \in \mathscr{A}} \left\{ r(s, a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \,|\, s, a) \cdot v^\star(s') \right\}$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Robust Bellman equations**

$$v^\star(s) = \max_{\pi \in \Delta(\mathcal{A})} \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathcal{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathcal{S}} p(s'|s,a) \cdot v^\star(s') \right] \right\}$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Robust Bellman operator**

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathscr{A})} \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \mid s, a) \cdot v(s') \right] \right\}$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Robust Bellman operator**

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathscr{A})} \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \,|\, s, a) \cdot v(s') \right] \right\}$$

**Distance-constrained *s*-rectangular ambiguity set**

$$\mathscr{P} = \prod_{s \in \mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s = \left\{ p(\,\cdot\,|\,s,\,\cdot\,) \,:\, \sum_{a \in \mathscr{A}} d\left[ p(\,\cdot\,|\,s,a), p^0(\,\cdot\,|\,s,a) \right] \leq \kappa \right\}$$



$$p^0(\,\cdot\,|\,s,a)$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathcal{A})} \min_{p \in \mathcal{P}_s} \left\{ \sum_{a \in \mathcal{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathcal{S}} p(s'|s,a) \cdot v(s') \right] \right\}$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathscr{A})} \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right] \right\}$$

$$= \min_{p \in \mathscr{P}_s} \max_{\pi \in \Delta(\mathscr{A})} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right] \right\}$$

**Minimax theorem:** exchange order of min and max

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$[\mathfrak{B}v](s) = \max_{\pi\in\Delta(\mathscr{A})}\min_{p\in\mathscr{P}_s}\left\{\sum_{a\in\mathscr{A}}\pi(a)\cdot\left[r(s,a)+\lambda\sum_{s'\in\mathscr{S}}p(s'\mid s,a)\cdot v(s')\right]\right\}$$

$$= \min_{p\in\mathscr{P}_s}\max_{\pi\in\Delta(\mathscr{A})}\left\{\sum_{a\in\mathscr{A}}\pi(a)\cdot\left[r(s,a)+\lambda\sum_{s'\in\mathscr{S}}p(s'\mid s,a)\cdot v(s')\right]\right\}$$

$$= \min_{p\in\mathscr{P}_s}\max_{a\in\mathscr{A}}\left\{r(s,a)+\lambda\sum_{s'\in\mathscr{S}}p(s'\mid s,a)\cdot v(s')\right\}$$

**Linearity:** we only need to consider ext $\Delta(\mathscr{A})=\mathscr{A}$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathscr{A})} \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right] \right\}$$

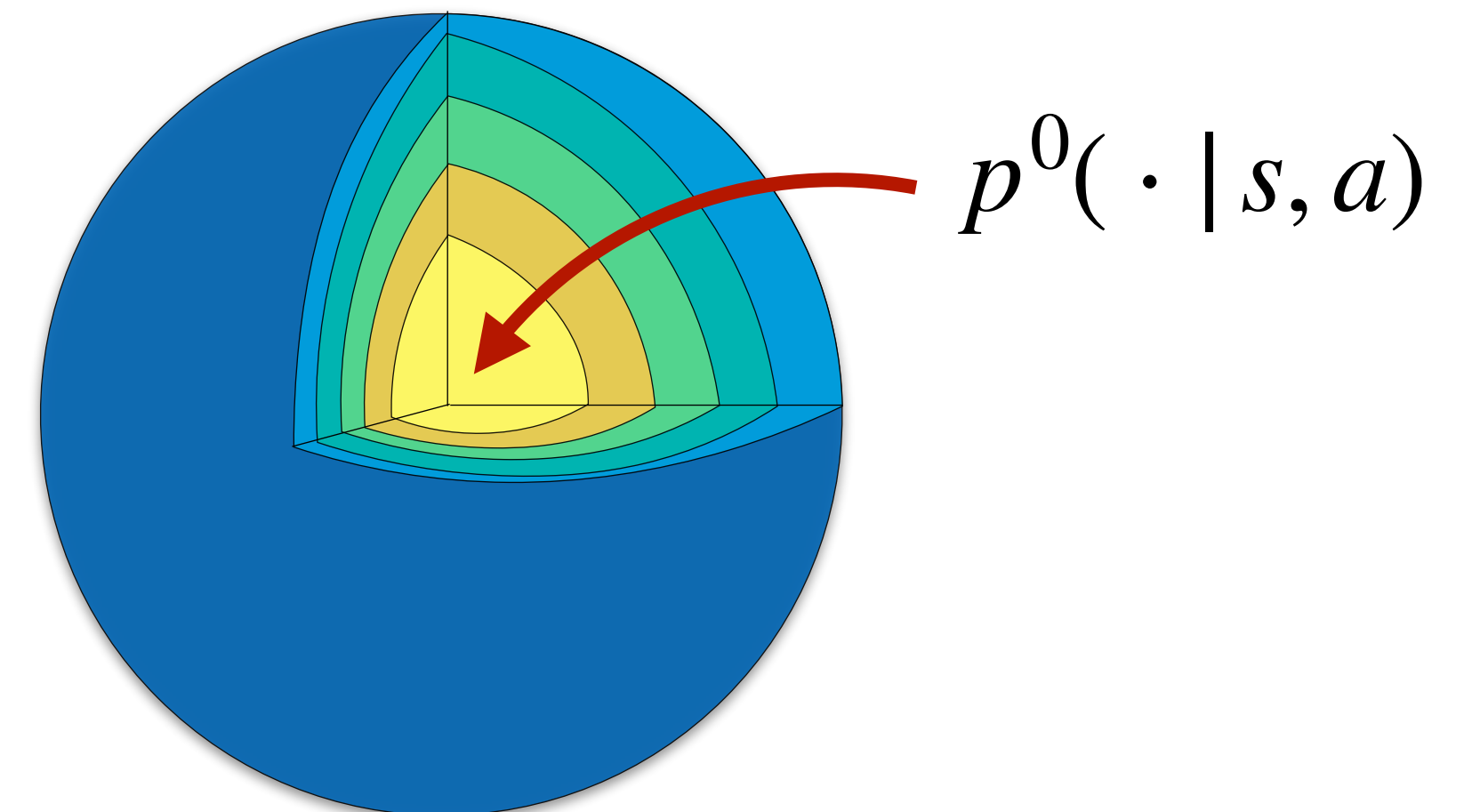$$= \min_{p \in \mathscr{P}_s} \max_{\pi \in \Delta(\mathscr{A})} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right] \right\}$$

$$\min_{p \in \mathscr{P}_s} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta \; ?$$

**Bisection search:**



Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$\min_{p \in \mathscr{P}_s} \quad \max_{a \in \mathscr{A}} \left\{ r(s, a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \,|\, s, a) \cdot v(s') \right\} \quad \leq \quad \theta \; ?$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$\min_{p \in \boxed{\mathscr{P}_s}} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta \, ?$$

$$\min_{p \in [\Delta(\mathscr{S})]^A} \left\{ \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} : \boxed{\sum_{a \in \mathscr{A}} d\left[ p(\cdot|s,a), p^0(\cdot|s,a) \right] \leq \kappa} \right\} \leq \theta$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$\min_{p \in \mathscr{P}_s} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta \; ?$$

$$\min_{p \in [\Delta(\mathscr{S})]^A} \left\{ \underbrace{\max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\}}_{f(p)} : \underbrace{\sum_{a \in \mathscr{A}} d\left[ p(\cdot|s,a), p^0(\cdot|s,a) \right] \leq \kappa}_{g(p)} \right\} \leq \theta$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$\min_{p \in \mathscr{P}_s} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta \ ?$$

$$\min_{p \in [\Delta(\mathscr{S})]^A} \left\{ \underbrace{\max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\}}_{f(p)} : \underbrace{\sum_{a \in \mathscr{A}} d\left[ p(\cdot|s,a), p^0(\cdot|s,a) \right] \leq \kappa}_{g(p)} \right\} \leq \theta$$

$$\min_{p \in [\Delta(\mathscr{S})]^A} \left\{ \underbrace{\sum_{a \in \mathscr{A}} d\left[ p(\cdot|s,a), p^0(\cdot|s,a) \right]}_{g(p)} : \underbrace{\max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta}_{f(p)} \right\} \leq \kappa$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$\min_{p \in \mathscr{P}_s} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta \, ?$$

$$\min_{p \in [\Delta(\mathscr{S})]^A} \left\{ \underbrace{\sum_{a \in \mathscr{A}} d\left[ p(\cdot|s,a), p^0(\cdot|s,a) \right]}_{g\,(p)} : \underbrace{\max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'|s,a) \cdot v(s') \right\} \leq \theta}_{f\,(p)} \right\} \leq \kappa$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

$$\min_{p \in \mathscr{P}_s} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \mid s,a) \cdot v(s') \right\} \leq \theta \ ?$$

$$\min_{p \in [\Delta(\mathscr{S})]^A} \left\{ \underbrace{\sum_{a \in \mathscr{A}} d\left[p(\cdot \mid s,a), p^0(\cdot \mid s,a)\right]}_{g\,(p)} : \underbrace{\max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \mid s,a) \cdot v(s') \right\} \leq \theta}_{f\,(p)} \right\} \leq \kappa$$

$$\iff \sum_{a \in \mathscr{A}} \min_{p_a \in \Delta(\mathscr{S})} \left\{ d\left[p(\cdot \mid s,a), p^0(\cdot \mid s,a)\right] \ : \ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s' \mid s,a) \cdot v(s') \leq \theta \right\} \leq \kappa$$

**Separability:** of both objective and constraints in $a \in \mathscr{A}$

Ho et al. (2023), *Robust Phi-Divergence MDPs.*

$$\min_{p \in \mathscr{P}_s} \max_{a \in \mathscr{A}} \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'\,|\,s,a) \cdot v(s') \right\} \;\leq\; \theta \,?$$

$$\sum_{a \in \mathscr{A}} \min_{p_a \in \Delta(\mathscr{S})} \left\{ d\left[p(\,\cdot\,|\,s,a), p^0(\,\cdot\,|\,s,a)\right] \;:\; r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'\,|\,s,a) \cdot v(s') \leq \theta \right\} \;\leq\; \kappa$$

Ho et al. (2023), *Robust Phi-Divergence MDPs.*

$$\min_{p \in \mathscr{P}_s} \; \max_{a \in \mathscr{A}} \; \left\{ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'\,|\,s,a) \cdot v(s') \right\} \;\; \leq \;\; \theta \;?$$

$$\sum_{a \in \mathscr{A}} \min_{p_a \in \Delta(\mathscr{S})} \left\{ d\left[p(\,\cdot\,|\,s,a), p^0(\,\cdot\,|\,s,a)\right] \;:\; r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'\,|\,s,a) \cdot v(s') \leq \theta \right\} \;\; \leq \;\; \kappa$$

$$\Longleftrightarrow \sum_{a \in \mathscr{A}} \mathfrak{P}(p^0; \lambda v, \theta - r(s\,|\,a)) \;\; \leq \;\; \kappa$$

with $\mathfrak{P}(p^0; b, \beta) = \begin{bmatrix} \underset{p}{\text{minimize}} & d\left[p, p^0\right] \\ \text{subject to} & \sum_{s' \in \mathscr{S}} b_{s'} \cdot p_{s'} \leq \beta \\ & p \in \Delta(\mathscr{S}) \end{bmatrix}$



$p^0$

$p^\star$

$$\sum_{s' \in \mathscr{S}} b_{s'} \cdot p_{s'} \leq \beta$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Distance-constrained *s*-rectangular ambiguity set**

$$\mathscr{P} = \prod_{s\in\mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s = \left\{ p(\,\cdot\,|s,\,\cdot\,) : \sum_{a\in\mathscr{A}} d\left[p(\,\cdot\,|s,a), p^0(\,\cdot\,|s,a)\right] \leq \kappa \right\}$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Distance-constrained *s*-rectangular ambiguity set**

$$\mathscr{P} = \prod_{s \in \mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s = \left\{ p(\,\cdot\,|\,s,\,\cdot\,) \; : \; \sum_{a \in \mathscr{A}} d\left[p(\,\cdot\,|\,s,a), p^0(\,\cdot\,|\,s,a)\right] \leq \kappa \right\}$$

**Robust Bellman operator**

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathscr{A})} \; \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'\,|\,s,a) \cdot v(s') \right] \right\}$$

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Distance-constrained *s*-rectangular ambiguity set**

$$\mathscr{P} = \prod_{s \in \mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s = \left\{ p(\,\cdot\,|\,s,\,\cdot\,) \; : \; \sum_{a \in \mathscr{A}} d\left[ p(\,\cdot\,|\,s,a), \, p^0(\,\cdot\,|\,s,a) \right] \leq \kappa \right\}$$

**Robust Bellman operator**

$$[\mathfrak{B}v](s) = \max_{\pi \in \Delta(\mathscr{A})} \; \min_{p \in \mathscr{P}_s} \left\{ \sum_{a \in \mathscr{A}} \pi(a) \cdot \left[ r(s,a) + \lambda \sum_{s' \in \mathscr{S}} p(s'\,|\,s,a) \cdot v(s') \right] \right\}$$

**Projection problem**

$$\mathfrak{P}(p^0; b, \beta) = \left[ \begin{array}{ll} \underset{p}{\text{minimize}} & d\left[p, p^0\right] \\ \text{subject to} & \sum_{s' \in \mathscr{S}} b_{s'} \cdot p_{s'} \leq \beta \\ & p \in \Delta(\mathscr{S}) \end{array} \right]$$

15

**Distance-constrained *s*-rectangular ambiguity set**

$$\mathscr{P} = \prod_{s \in \mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s = \left\{ p(\,\cdot\,|s,\,\cdot\,) \; : \; \sum_{a \in \mathscr{A}} d \left[ p(\,\cdot\,|s,a), p^0(\,\cdot\,|s,a) \right] \leq \kappa \right\}$$

**Theorem**

Assume $\mathfrak{P}$ can be computed exactly in time $\mathscr{O}(h(S))$.
Then $\mathfrak{B}$ can be computed to accuracy $\epsilon > 0$ in time
$\mathscr{O}(AS \cdot h(S) \cdot \log[\overline{R}/\epsilon])$.

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Distance-constrained *s*-rectangular ambiguity set**

$$\mathscr{P} = \prod_{s \in \mathscr{S}} \mathscr{P}_s \quad \text{with} \quad \mathscr{P}_s = \left\{ p(\cdot \,|\, s, \cdot) \,:\, \sum_{a \in \mathscr{A}} d\left[ p(\cdot \,|\, s, a), p^0(\cdot \,|\, s, a) \right] \leq \kappa \right\}$$

**Theorem**

Assume $\mathfrak{P}$ can be computed exactly in time $\mathcal{O}(h(S))$.
Then $\mathfrak{B}$ can be computed to accuracy $\epsilon > 0$ in time
$\mathcal{O}(AS \cdot h(S) \cdot \log[\overline{R}/\epsilon])$.

Assume $\mathfrak{P}$ can be computed to any accuracy $\delta > 0$
in time $\mathcal{O}(h(\delta))$. Then $\mathfrak{B}$ can be computed to accuracy
$\epsilon > 0$ in time $\mathcal{O}(AS \cdot h(\epsilon\kappa/[2A\overline{R} + A\epsilon]) \cdot \log[\overline{R}/\epsilon])$.

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

| Divergence | $d_a(\,\cdot\,,p^0)$ | Ours | Previous |
|---|---|---|---|
| KL-Divergence | $\sum_{s'\in\mathcal{S}} p(s'\,\|\,s,a)\cdot\log\left(\dfrac{p(s'\,\|\,s,a)}{p^0(s'\,\|\,s,a)}\right)$ | $\mathcal{O}(S^2\cdot A\log A)$ | $\mathcal{O}(\ell^2\cdot S^2\cdot A)$ |
| Burg Entropy | $\sum_{s'\in\mathcal{S}} p^0(s'\,\|\,s,a)\cdot\log\left(\dfrac{p^0(s'\,\|\,s,a)}{p(s'\,\|\,s,a)}\right)$ | $\mathcal{O}(S^2\cdot A\log A)$ | (none) |
| Variation Distance | $\sum_{s'\in\mathcal{S}} |p(s'\,\|\,s,a)-p^0(s'\,\|\,s,a)|$ | $\mathcal{O}(S^2\log S\cdot A)$ | $\mathcal{O}(S^2\log S\cdot A)$ |
| $\chi^2$-Distance | $\sum_{s'\in\mathcal{S}}\dfrac{\left[p(s'\,\|\,s,a)-p^0(s'\,\|\,s,a)\right]^2}{p^0(s'\,\|\,s,a)}$ | $\mathcal{O}(S^2\log S\cdot A)$ | $\mathcal{O}(S^{4.5}\cdot A)$ |

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

| Divergence | $d_a(\,\cdot\,,p^0)$ | Ours | Previous |
|---|---|---|---|
| KL-Divergence | $\displaystyle\sum_{s'\in\mathcal{S}} p(s'\,\vert\,s,a)\cdot\log\left(\frac{p(s'\,\vert\,s,a)}{p^0(s'\,\vert\,s,a)}\right)$ | $\mathcal{O}(S^2\cdot A\log A)$ | $\mathcal{O}(\ell^2\cdot S^2\cdot A)$ |
| Burg Entropy | $\displaystyle\sum_{s'\in\mathcal{S}} p^0(s'\,\vert\,s,a)\cdot\log\left(\frac{p^0(s'\,\vert\,s,a)}{p(s'\,\vert\,s,a)}\right)$ | $\mathcal{O}(S^2\cdot A\log A)$ | (none) |
| Variation Distance | $\displaystyle\sum_{s'\in\mathcal{S}} \vert p(s'\,\vert\,s,a) - p^0(s'\,\vert\,s,a)\vert$ | $\mathcal{O}(S^2\log S\cdot A)$ | $\mathcal{O}(S^2\log S\cdot A)$ |
| $\chi^2$-Distance | $\displaystyle\sum_{s'\in\mathcal{S}} \frac{\left[p(s'\,\vert\,s,a) - p^0(s'\,\vert\,s,a)\right]^2}{p^0(s'\,\vert\,s,a)}$ | $\mathcal{O}(S^2\log S\cdot A)$ | $\mathcal{O}(S^{4.5}\cdot A)$ |

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

**Projection problem**

**Bellman operator**

Ho et al. (2023), *Robust Phi-Divergence MDPs*.

[1] WW, D. Kuhn, B. Rustem, Robust Markov Decision Processes, *Mathematics of Operations Research* 38(1):153-183, 2013.

[2] C. Ho, M. Petrik, WW, Fast Bellman Updates for Robust MDPs, *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018.

[3] J. C. D'Aeth, WW et al. Optimal National Prioritization Policies for Hospital Care During the SARS-CoV-2 Pandemic, *Nature Computational Science* 1(8):521-531, 2021.

[4] J. C. D'Aeth, WW et al. Optimal Hospital Care Scheduling During the SARS-CoV-2 Pandemic, *Management Science* (Online First), 2023.

[5] C. Ho, M. Petrik, WW, Partial Policy Iteration for L1-Robust Markov Decision Processes, *The Journal of Machine Learning Research* 22(1):12612-12657, 2021.

[6] C. Ho, M. Petrik, WW, Robust Phi-Divergence MDPs, *Advances in Neural Information Processing Systems 35 (NeurIPS Proceedings)*, 2022.

**ww@imperial.ac.uk**

# — BACKUP —

# LARGE-SCALE PROBLEMS

House of Commons Research Briefing: *NHS Key Statistics, England, March 2023*.

**MDP model** of an **individual patient**:

D'Aeth et al. (2021), *Optimal National Prioritization Policies for Hospital Care During the SARS-CoV-2 Pandemic*.

**MDP model** of an **individual patient**:



**Wait**

**MDP model** of an **individual patient**:

**MDP model** of an **individual patient**:



**Wait**

**Admit**

**Treat**

19

**MDP model** of an **individual patient**:



Wait   Admit   Treat   Release

19

**Markov decision process**

Tuple $(\mathcal{S}, \mathcal{A}, q, p, r, T)$ where

- $\mathcal{S} = \prod_{i=1}^{n} \mathcal{S}_i$ with $\mathcal{S}_i = \{1, \ldots, S_i\}$ is the (finite) state space;

- $\mathcal{A} = \prod_{i=1}^{n} \mathcal{A}_i$ with $\mathcal{A}_i = \{1, \ldots, A_i\}$ is the (finite) action space;



$(\mathcal{S}_1, \mathcal{A}_1)$

X

$\vdots$

X

$(\mathcal{S}_n, \mathcal{A}_n)$

**Markov decision process**

Tuple $(\mathcal{S}, \mathcal{A}, q, p, r, T)$ where

- $\mathcal{S} = \prod_{i=1}^{n} \mathcal{S}_i$ with $\mathcal{S}_i = \{1, \ldots, S_i\}$ is the (finite) state space;

- $\mathcal{A} = \prod_{i=1}^{n} \mathcal{A}_i$ with $\mathcal{A}_i = \{1, \ldots, A_i\}$ is the (finite) action space;

- $q(s) = \prod_{i=1}^{n} q_i(s_i)$ is the initial state distribution;

- $p_t(s' \,|\, s, a) = \prod_{i=1}^{n} p_{ti}(s'_i \,|\, s_i, a_i)$ is the transition kernel;

- $r_t(s, a) = \sum_{i=1}^{n} r_{ti}(s_i, a_i)$ are the expected one-step rewards;

- $T \in \mathbb{N}$ is the (finite) time horizon

**Weakly coupled** **Markov decision process**

Tuple $(\mathcal{S}, \mathcal{A}, q, p, r, T)$ where

- $\mathcal{S} = \prod_{i=1}^{n} \mathcal{S}_i$ with $\mathcal{S}_i = \{1, \ldots, S_i\}$ is the (finite) state space;

- $\mathcal{A} = \prod_{i=1}^{n} \mathcal{A}_i$ with $\mathcal{A}_i = \{1, \ldots, A_i\}$ is the (finite) action space;

- $q(s) = \prod_{i=1}^{n} q_i(s_i)$ is the initial state distribution;

- $p_t(s' | s, a) = \prod_{i=1}^{n} p_{ti}(s_i' | s_i, a_i)$ is the transition kernel;

- $r_t(s, a) = \sum_{i=1}^{n} r_{ti}(s_i, a_i)$ are the expected one-step rewards;

- $T \in \mathbb{N}$ is the (finite) time horizon

**and**

$\quad a \in \mathcal{A}$ admissible only if $\sum_{i=1}^{n} c_{tli}(s_i, a_i) \leq b_{tl}$ for all $l \in \mathcal{L}$

20

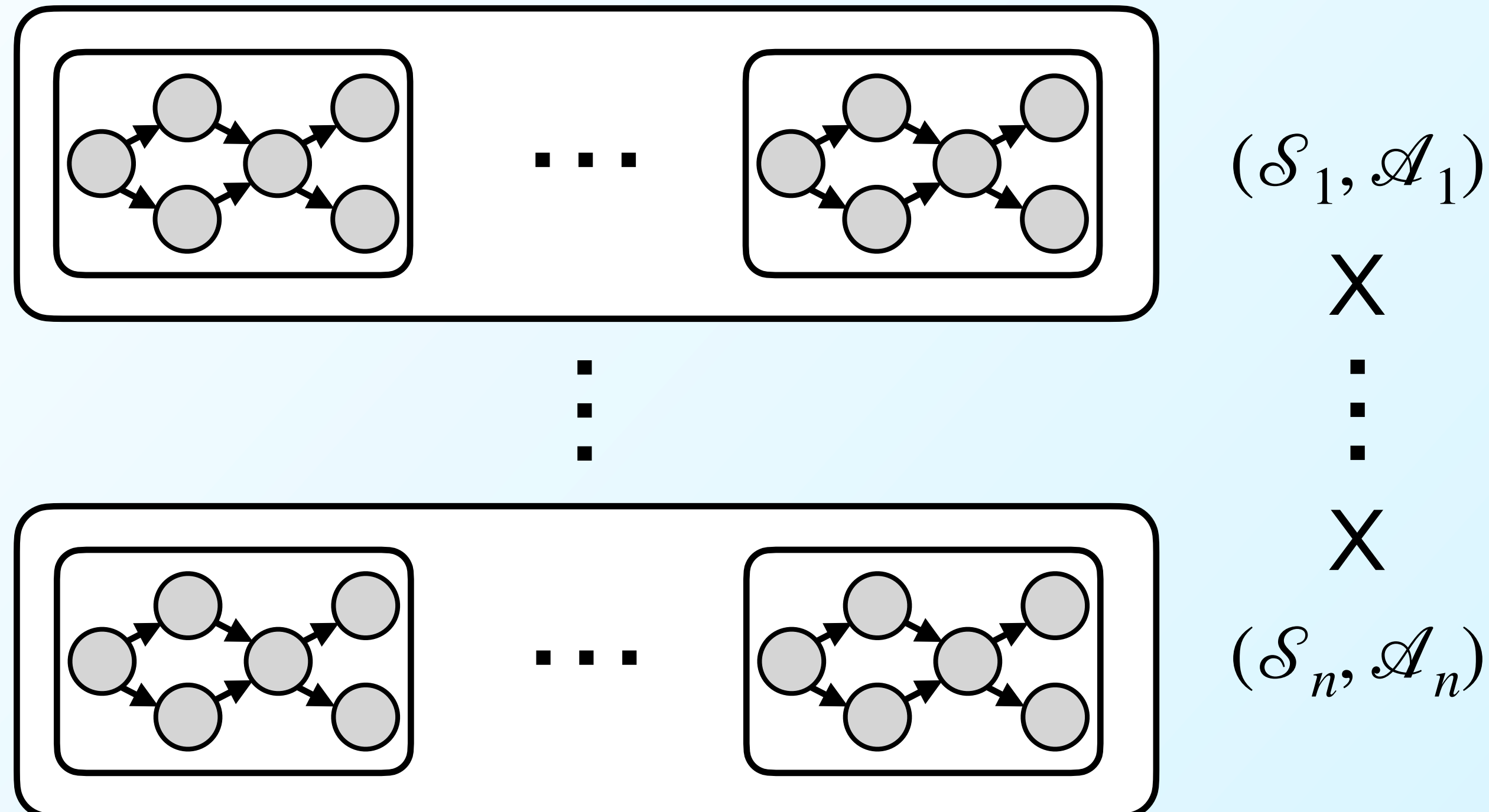**Weakly coupled** **Markov decision process**

Tuple $(\mathcal{S}, \mathcal{A}, q, p, r, T)$ where

- $\mathcal{S} = \prod_{i=1}^{n} \mathcal{S}_i$ with $\mathcal{S}_i = \{1, \ldots, S_i\}$ is the (finite) state space;

- $\mathcal{A} = \prod_{i=1}^{n} \mathcal{A}_i$ with $\mathcal{A}_i = \{1, \ldots, A_i\}$ is the (finite) action space;

- $q(s) = \prod_{i=1}^{n} q_i(s_i)$ is the initial state distribution;

- $p_t(s' | s, a) = \prod_{i=1}^{n} p_{ti}(\ldots)$ is the transition kernel;

- $r_t(s, \ldots)$ is the rewards;

- $T \in \ldots$

**and**

$a \in \mathcal{S}$ $\ldots$ $\in \mathcal{L}$

**Objective**

find policy $\pi = \mathcal{S} \rightarrow \mathcal{A}$ that

maximizes the expected total rewards:

$$\underset{\pi \in \Pi}{\text{maximize}} \quad \mathbb{E}_p \left[ \sum_{t=1}^{T} r\big(s_t, \pi_t[s_t]\big) \right]$$

20

**Fluid Linear Program**

$$\begin{array}{ll}
\underset{\sigma,\, \pi \geq 0}{\text{maximize}} & \displaystyle\sum_{t=1}^{T}\sum_{i=1}^{n}\sum_{s_i\in\mathcal{S}_i}\sum_{a_i\in\mathcal{A}_i} r_{ti}(s_i, a_i)\cdot \pi_{ti}(s_i, a_i)
\end{array}$$

$$\text{subject to} \quad \sigma_{1i}(s_i) = q_i(s_i) \qquad\qquad\qquad\qquad\qquad\quad \forall i,\, \forall s_i \in \mathcal{S}_i$$

$$\sigma_{t+1,i}(s_i') = \sum_{s_i\in\mathcal{S}_i}\sum_{a_i\in\mathcal{A}_i} p_{ti}(s_i' \mid s_i, a_i)\cdot \pi_{ti}(s_i, a_i) \quad \forall i,\, \forall s_i' \in \mathcal{S}_i,\, \forall t$$

$$\sum_{i=1}^{n}\sum_{s_i\in\mathcal{S}_i}\sum_{a_i\in\mathcal{A}_i} c_{tli}(s_i, a_i)\cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l,\, \forall t$$

$$\sum_{a_i\in\mathcal{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad\qquad\qquad\quad \forall i,\, \forall s \in \mathcal{S}_i,\, \forall t$$

**Fluid Linear Program**

maximize
$\sigma, \pi \geq 0$

$$\sum^{T} \sum^{n} \sum \sum r_{ti}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i)$$

$\sigma_{ti}(s_i)$: % of MDP *i* that is in
state $s_i$ in stage *t*

subject to $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \forall i, \forall s_i \in \mathcal{S}_i$

$$\sigma_{t+1,i}(s_i') = \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} p_{ti}(s_i' \,|\, s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \qquad \forall i, \forall s_i' \in \mathcal{S}_i, \forall t$$

$$\sum_{i=1}^{n} \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} c_{tli}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l, \forall t$$

$$\sum_{a_i \in \mathcal{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad \forall i, \forall s \in \mathcal{S}_i, \forall t$$

21

**Fluid Linear Program**

$$\text{maximize} \quad \sigma, \pi \geq 0 \quad \sum_{t}^{T} \sum_{i}^{n} \sum \sum r_{ti}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i)$$

$\sigma_{ti}(s_i)$: % of MDP $i$ that is in state $s_i$ in stage $t$

subject to

$\pi_{ti}(s_i, a_i)$: % of MDP from $\sigma_{ti}(s_i)$ that we apply action $a_i$ to

$$\sum_{i=1} \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} c_{tli}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l, \forall t$$

$$\sum_{a_i \in \mathcal{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad \forall i, \forall s \in \mathcal{S}_i, \forall t$$

$$\forall i, \forall s_i \in \mathcal{S}_i$$

$$\forall i, \forall s_i' \in \mathcal{S}_i, \forall t$$

21

**Fluid Linear Program**

$$\begin{array}{ll} \underset{\sigma, \pi \geq 0}{\text{maximize}} & \displaystyle\sum_{t=1}^{T} \sum_{i=1}^{n} \sum_{s_i \in \mathscr{S}_i} \sum_{a_i \in \mathscr{A}_i} r_{ti}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \end{array}$$

**objective function:**
maximize rewards

subject to

$$\sigma_{1i}(s_i) = q_i(s_i) \qquad \forall i, \forall s_i \in \mathscr{S}_i$$

$$\sigma_{t+1,i}(s_i') = \sum_{s_i \in \mathscr{S}_i} \sum_{a_i \in \mathscr{A}_i} p_{ti}(s_i' \mid s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \qquad \forall i, \forall s_i' \in \mathscr{S}_i, \forall t$$

$$\sum_{i=1}^{n} \sum_{s_i \in \mathscr{S}_i} \sum_{a_i \in \mathscr{A}_i} c_{tli}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l, \forall t$$

$$\sum_{a_i \in \mathscr{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad \forall i, \forall s \in \mathscr{S}_i, \forall t$$

**Fluid Linear Program**

$$
\underset{\sigma, \pi \geq 0}{\text{maximize}} \quad \sum_{t=1}^{T} \sum_{i=1}^{n} \sum_{s_i \in \mathscr{S}_i} \sum_{a_i \in \mathscr{A}_i} r_{ti}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i)
$$

**initial states:** must follow $q$

subject to

$$
\sigma_{1i}(s_i) = q_i(s_i) \qquad\qquad \forall i, \forall s_i \in \mathscr{S}_i
$$

$$
\sigma_{t+1,i}(s_i') = \sum_{s_i \in \mathscr{S}_i} \sum_{a_i \in \mathscr{A}_i} p_{ti}(s_i' \mid s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \qquad \forall i, \forall s_i' \in \mathscr{S}_i, \forall t
$$

$$
\sum_{i=1}^{n} \sum_{s_i \in \mathscr{S}_i} \sum_{a_i \in \mathscr{A}_i} c_{tli}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l, \forall t
$$

$$
\sum_{a_i \in \mathscr{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad\qquad \forall i, \forall s \in \mathscr{S}_i, \forall t
$$

21

**Fluid Linear Program**

maximize
$\sigma, \pi \geq 0$

$$\sum_{t=1}^{T}\sum_{i=1}^{n}\sum_{s_i\in\mathcal{S}_i}\sum_{a_i\in\mathcal{A}_i} r_{ti}(s_i,a_i)\cdot\pi_{ti}(s_i,a_i)$$

**transitions:**

must follow $p$

subject to

$$\sigma_{1i}(s_i) = q_i(s_i) \qquad \forall i, \forall s_i \in \mathcal{S}_i$$

$$\sigma_{t+1,i}(s_i') = \sum_{s_i\in\mathcal{S}_i}\sum_{a_i\in\mathcal{A}_i} p_{ti}(s_i'\,|\,s_i,a_i)\cdot\pi_{ti}(s_i,a_i) \qquad \forall i, \forall s_i' \in \mathcal{S}_i, \forall t$$

$$\sum_{i=1}^{n}\sum_{s_i\in\mathcal{S}_i}\sum_{a_i\in\mathcal{A}_i} c_{tli}(s_i,a_i)\cdot\pi_{ti}(s_i,a_i) \leq b_{tl} \qquad \forall l, \forall t$$

$$\sum_{a_i\in\mathcal{A}_i} \pi_{ti}(s_i,a_i) = \sigma_{ti}(s_i) \qquad \forall i, \forall s \in \mathcal{S}_i, \forall t$$

**Fluid Linear Program**

$$\underset{\sigma, \pi \geq 0}{\text{maximize}} \quad \sum_{t=1}^{T} \sum_{i=1}^{n} \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} r_{ti}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i)$$

**resources:**
budgets must be kept

subject to

$$\sigma_{1i}(s_i) = q_i(s_i) \qquad\qquad \forall i, \forall s_i \in \mathcal{S}_i$$

$$\sigma_{t+1,i}(s_i') = \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} p_{ti}(s_i' \,|\, s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \quad \forall i, \forall s_i' \in \mathcal{S}_i, \forall t$$

$$\sum_{i=1}^{n} \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} c_{tli}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l, \forall t$$

$$\sum_{a_i \in \mathcal{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad\qquad \forall i, \forall s \in \mathcal{S}_i, \forall t$$

21

**Fluid Linear Program**

$$\text{maximize}_{\sigma, \pi \geq 0} \quad \sum_{t=1}^{T} \sum_{i=1}^{n} \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} r_{ti}(s_i, a_i) \cdot \pi_{ti}(s_i,$$

**"flow preservation":**
we cannot "drop" MDPs

subject to

$$\sigma_{1i}(s_i) = q_i(s_i) \qquad \forall i, \forall s_i \in \mathcal{S}_i$$

$$\sigma_{t+1,i}(s_i') = \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} p_{ti}(s_i' \mid s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \qquad \forall i, \forall s_i' \in \mathcal{S}_i, \forall t$$

$$\sum_{i=1}^{n} \sum_{s_i \in \mathcal{S}_i} \sum_{a_i \in \mathcal{A}_i} c_{tli}(s_i, a_i) \cdot \pi_{ti}(s_i, a_i) \leq b_{tl} \qquad \forall l, \forall t$$

$$\sum_{a_i \in \mathcal{A}_i} \pi_{ti}(s_i, a_i) = \sigma_{ti}(s_i) \qquad \forall i, \forall s \in \mathcal{S}_i, \forall t$$

21

**Observation**

The fluid LP constitutes a relaxation of the weakly coupled MDP.

**Observation**

The fluid LP constitutes a relaxation of the weakly coupled MDP.

**Randomized policy**

For each MDP $i$, take action $a_i$ in state $s_i$ with probability $\dfrac{\pi_{ti}(s_i, a_i)}{\sigma_{ti}(s_i)}$ at time $t$.

**Observation**

The fluid LP constitutes a relaxation of the weakly coupled MDP.

**Randomized policy**

For each MDP $i$, take action $a_i$ in state $s_i$ with probability $\dfrac{\pi_{ti}(s_i, a_i)}{\sigma_{ti}(s_i)}$ at time $t$.

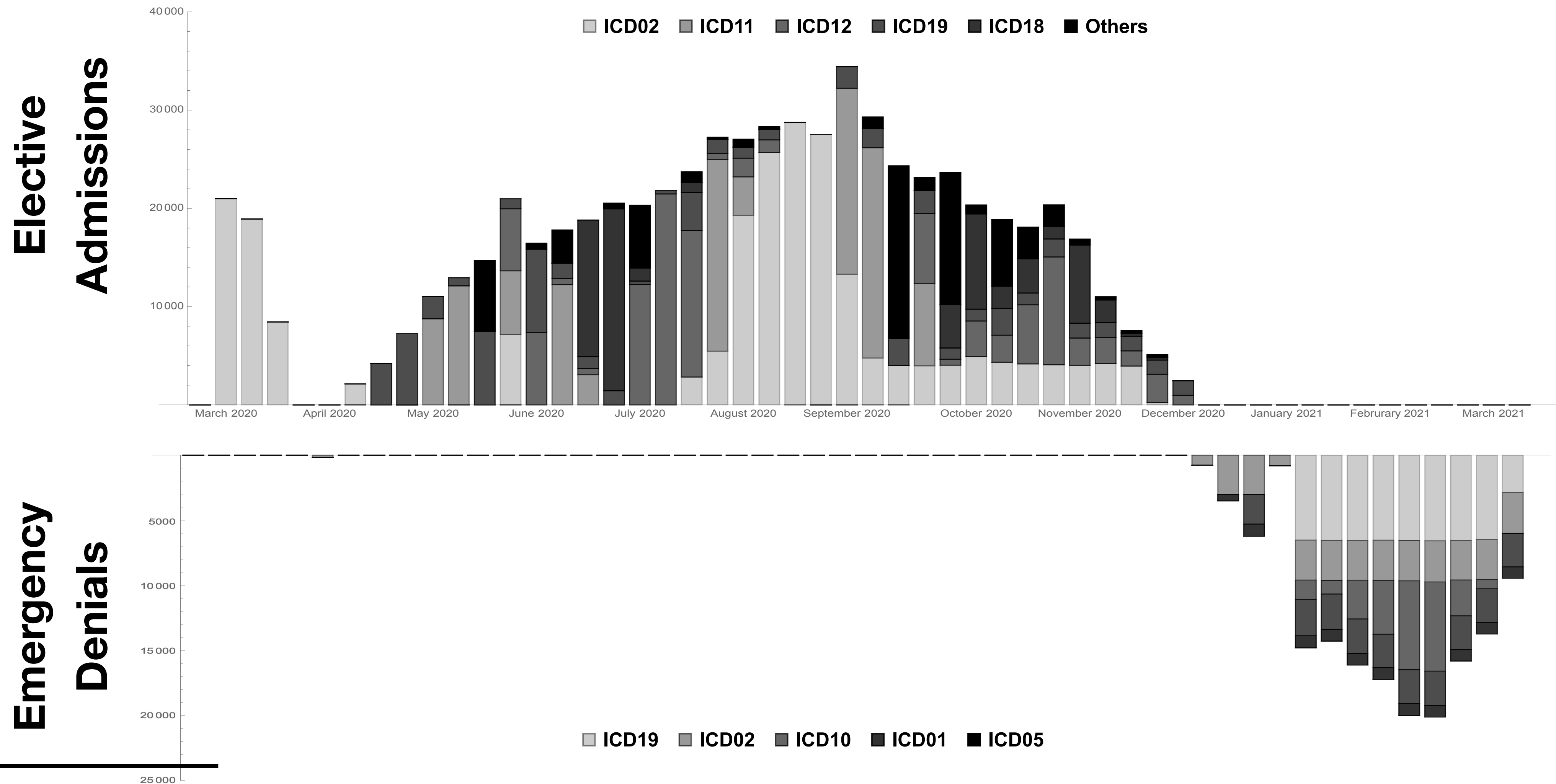**Performance guarantee**

For suitably adapted $b_{tl}$, the randomized policy is guaranteed to be feasible in the weakly coupled MDP. Moreover, the relative optimality gap for large MDPs is:

$$T \cdot \sqrt{\frac{\log n}{n}} + \frac{T^2 L}{n^2} \xrightarrow[n \to \infty]{} 0$$

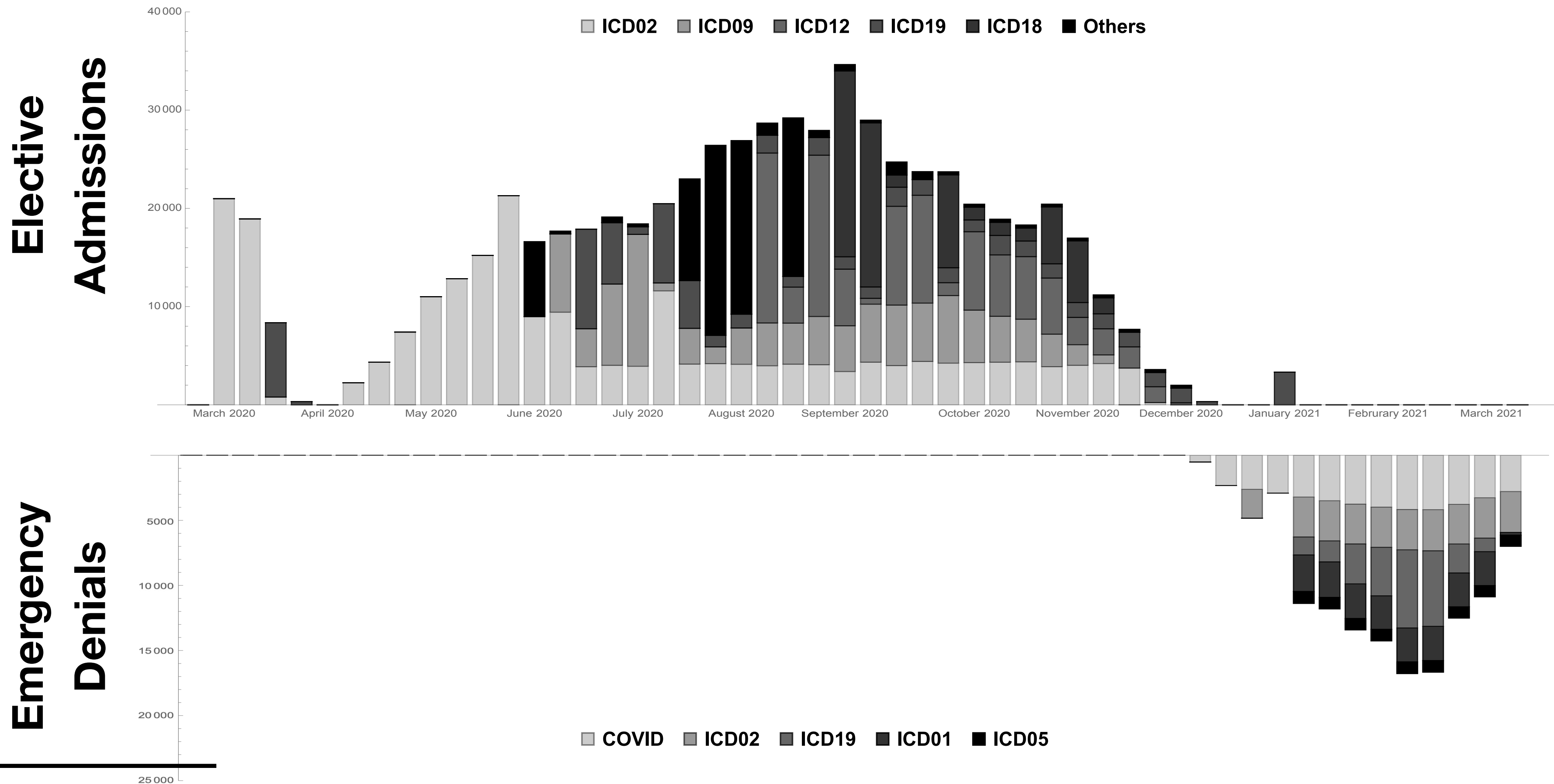D'Aeth et al. (2021), *Optimal National Prioritization Policies for Hospital Care During the SARS-CoV-2 Pandemic*.

D'Aeth et al. (2021), *Optimal National Prioritization Policies for Hospital Care During the SARS-CoV-2 Pandemic*.

23

**Years of Life Gained by Optimized Schedule**



* 720k YLG  (+8.7%)
* 22.1% less emergencies
* up to 53.5% more electives

injury, poisoning, external causes

digestive diseases

respiratory diseases

circulatory diseases

neoplasms (cancer)

infectious and parasitic diseases

Years of Life Gained

D'Aeth et al. (2021), *Optimal National Prioritization Policies for Hospital Care During the SARS-CoV-2 Pandemic*.

24

Years of Life Gained by Optimized Schedule

**Randomized Policy:**
- G&A  +0.05%
- CC  +1.56%

Others
ICD19 — *injury, poisoning, external causes*
ICD18
ICD14
ICD13
ICD11 — *digestive diseases*
ICD10 — *respiratory diseases*
ICD09 — *circulatory diseases*
ICD06
ICD05
ICD04
ICD02 — *neoplasms (cancer)*
ICD01 — *infectious and parasitic diseases*
COVID

−200 000    −100 000    100 000    200 000    300 000

Years of Life Gained

D'Aeth et al. (2021), *Optimal National Prioritization Policies for Hospital Care During the SARS-CoV-2 Pandemic*.