ECE276B: Planning & Learning in Robotics Lecture 3: Markov Decision Processes

Nikolay Atanasov

natanasov@ucsd.edu



JACOBS SCHOOL OF ENGINEERING Electrical and Computer Engineering

Outline

Markov Decision Processes

Open-Loop vs Closed-Loop Control

Partially Observable Models

Markov Chain

Markov Chain

Stochastic process defined by a tuple (\mathcal{X}, p_0, p_f) :

- X is a discrete or continuous space
- ▶ $p_0(\cdot)$ is a prior pdf defined on \mathcal{X}
- ▶ $p_f(\cdot | \mathbf{x})$ is a conditional pdf defined on \mathcal{X} for given $\mathbf{x} \in \mathcal{X}$ that specifies the stochastic process transitions
- When the state space is finite $\mathcal{X} := \{1, \dots, N\}$:
 - the prior pdf p_0 is represented by an $N \times 1$ vector with elements:

$$\mathbf{p}_{0,i} := \mathbb{P}(x_0 = i) = \mathbf{p}_0(i)$$

• the transition pdf p_f is represented by an $N \times N$ matrix with elements:

$$P_{ij} := \mathbb{P}(x_{t+1} = j \mid x_t = i) = p_f(j \mid x_t = i)$$

Example: Student Markov Chain



Markov Reward Process

Markov Reward Process

Markov chain with transition costs defined by a tuple $(\mathcal{X}, p_0, p_f, T, \ell, \mathfrak{q}, \gamma)$:

- X is a discrete or continuous space
- ▶ $p_0(\cdot)$ is a prior pdf defined on \mathcal{X}

▶ p_f(· | x) is a conditional pdf defined on X for given x ∈ X that specifies the stochastic process transitions

- T is a finite/infinite time horizon
- ▶ $\ell(\mathbf{x})$ is stage cost of state $\mathbf{x} \in \mathcal{X}$
- q(x) is terminal cost of being in state x at time T
- ▶ $\gamma \in [0, 1]$ is a discount factor

Example: Student Markov Reward Process



MRP Value Function

- ► Value function: the expected cumulative cost of an MRP starting from state x ∈ X at time t
- Finite-horizon MRP: trajectories terminate at fixed $T < \infty$

$$V_t(\mathbf{x}) := \mathbb{E}\left[\mathfrak{q}(\mathbf{x}_{\mathcal{T}}) + \sum_{ au=t}^{ au-1} \ell(\mathbf{x}_{ au}) \mid \mathbf{x}_t = \mathbf{x}
ight]$$

Infinite-horizon MRP:

- First-exit MRP: trajectories terminate at the first passage time $T = \min \{ t \in \mathbb{N} | \mathbf{x}_t \in \mathcal{T} \}$ to a terminal state $\mathbf{x}_t \in \mathcal{T} \subseteq \mathcal{X}$
- Discounted MRP: trajectories continue forever but stage costs are discounted by discount factor *γ* ∈ [0, 1):
 - γ close to 0 leads to myopic/greedy evaluation
 - γ close to 1 leads to nonmyopic/far-sighted evaluation
 - Mathematically convenient since discounting avoids infinite costs as $T
 ightarrow \infty$
- Average-cost MRP: trajectories continue forever and the value function is the expected average stage cost

Example: Student MRP Value Function



Example: Student MRP Value Function



Example: Student MRP Value Function



Markov Decision Process

Markov Decision Process

Markov Reward Process with controlled transitions defined by a tuple $(\mathcal{X}, \mathcal{U}, p_0, p_f, T, \ell, \mathfrak{q}, \gamma)$

- \mathcal{X} is a discrete or continuous state space
- $\blacktriangleright~\mathcal{U}$ is a discrete or continuous control space
- $p_0(\cdot)$ is a prior pdf defined on \mathcal{X}
- ▶ $p_f(\cdot | \mathbf{x}_t, \mathbf{u}_t)$ is a conditional pdf defined on \mathcal{X} for given $\mathbf{x}_t \in \mathcal{X}$ and $\mathbf{u}_t \in \mathcal{U}$ (matrices P^u with elements $P_{ij}^u := p_f(j | x_t = i, u_t = u)$ in finite-dim case)
- T is a finite or infinite time horizon
- $\blacktriangleright~\ell(x,u)$ is stage cost of applying control $u\in \mathcal{U}$ in state $x\in \mathcal{X}$
- $q(\mathbf{x})$ is terminal cost of being in state \mathbf{x} at time T
- ▶ $\gamma \in [0,1]$ is a discount factor

Example: Markov Decision Process

A control u_t applied in state x_t determines the next state x_{t+1} and the stage cost l(x_t, u_t)



Example: Student Markov Decision Process



MDP Control Policy and Value Function

- Control policy: a function π that maps a time step t ∈ N and a state x ∈ X to a feasible control input u ∈ U
- Value function: expected cumulative cost of a policy π applied to an MDP with initial state x ∈ X at time t:
- Finite-horizon MDP: trajectories terminate at fixed $T < \infty$:

$$V_t^{\pi}(\mathbf{x}) := \mathbb{E}\left[\mathfrak{q}(\mathbf{x}_{\mathcal{T}}) + \sum_{\tau=t}^{T-1} \ell(\mathbf{x}_{\tau}, \pi_{\tau}(\mathbf{x}_{\tau})) \mid \mathbf{x}_t = \mathbf{x}\right]$$

▶ Infinite-horizon MDP: as $T \to \infty$, optimal policies become stationary, i.e., $\pi := \pi_0 \equiv \pi_1 \equiv \cdots$

- First-exit MDP: trajectories terminate at the first passage time $T = \min \{t \in \mathbb{N} | \mathbf{x}_t \in T\}$ to a terminal state $\mathbf{x}_t \in T \subseteq \mathcal{X}$
- ▶ Discounted MDP: trajectories continue forever but stage costs are discounted by a factor *γ* ∈ [0, 1)
- Average-cost MDP: trajectories continue forever and the value function is the expected average stage cost

Example: Value Function of Student MDP



Alternative Cost Formulations

Noise-dependent costs: stage costs ℓ' depend on motion noise \mathbf{w}_t :

$$V_0^{\pi}(\mathbf{x}) := \mathbb{E}_{\mathbf{w}_{0:T}, \mathbf{x}_{1:T}} \left[\mathfrak{q}(\mathbf{x}_T) + \sum_{t=0}^{T-1} \ell'(\mathbf{x}_t, \pi_t(\mathbf{x}_t), \mathbf{w}_t) \mid \mathbf{x}_0 = \mathbf{x} \right]$$

• Using the pdf $p_w(\cdot | \mathbf{x}_t, \mathbf{u}_t)$ of \mathbf{w}_t , this is equivalent to our formulation:

$$\ell(\mathbf{x}_t, \mathbf{u}_t) := \mathbb{E}_{\mathbf{w}_t | \mathbf{x}_t, \mathbf{u}_t} \left[\ell'(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t) \right] = \int \ell'(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t) \rho_w(\mathbf{w}_t \mid \mathbf{x}_t, \mathbf{u}_t) d\mathbf{w}_t$$

The expectation can be computed if p_w is known or approximated.

▶ Joint cost-state pdf: allow random costs ℓ' with joint pdf p(x', ℓ' | x, u). This is equivalent to our formulation as follows:

$$p_f(\mathbf{x}' \mid \mathbf{x}, \mathbf{u}) := \int p(\mathbf{x}', \ell' \mid \mathbf{x}, \mathbf{u}) d\ell'$$
$$\ell(\mathbf{x}, \mathbf{u}) := \mathbb{E}\left[\ell' \mid \mathbf{x}, \mathbf{u}\right] = \int \int \ell' p(\mathbf{x}', \ell' \mid, \mathbf{x}, \mathbf{u}) d\mathbf{x}' d\ell'$$

Alternative Motion-Model Formulations

- Fine-lag motion model: $\mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{u}_t, \mathbf{u}_{t-1}, \mathbf{w}_t)$
- Can be converted to the standard form via state augmentation

• Let $\mathbf{y}_t := \mathbf{x}_{t-1}$ and $\mathbf{s}_t := \mathbf{u}_{t-1}$ and define the augmented dynamics:

$$\tilde{\mathbf{x}}_{t+1} := \begin{bmatrix} \mathbf{x}_{t+1} \\ \mathbf{y}_{t+1} \\ \mathbf{s}_{t+1} \end{bmatrix} = \begin{bmatrix} f_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{u}_t, \mathbf{s}_t, \mathbf{w}_t) \\ \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} =: \tilde{f}_t(\tilde{\mathbf{x}}_t, \mathbf{u}_t, \mathbf{w}_t)$$

This procedure works for an arbitrary number of time lags but the dimension of the state space grows and increases the computational burden exponentially ("curse of dimensionality")

Alternative Motion-Model Formulations

System dynamics: $\mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t)$

Correlated Disturbance: w_t correlated across time (colored noise):

$$\mathbf{y}_{t+1} = A_t \mathbf{y}_t + \mathbf{\xi}_t$$

 $\mathbf{w}_t = C_t \mathbf{y}_{t+1}$

where A_t , C_t are known and ξ_t are independent random variables

• Augmented state: $\tilde{\mathbf{x}}_t := (\mathbf{x}_t, \mathbf{y}_t)$ with dynamics:

$$\tilde{\mathbf{x}}_{t+1} = \begin{bmatrix} \mathbf{x}_{t+1} \\ \mathbf{y}_{t+1} \end{bmatrix} = \begin{bmatrix} f_t(\mathbf{x}_t, \mathbf{u}_t, C_t(A_t\mathbf{y}_t + \boldsymbol{\xi}_t)) \\ A_t\mathbf{y}_t + \boldsymbol{\xi}_t \end{bmatrix} =: \tilde{f}_t(\tilde{\mathbf{x}}_t, \mathbf{u}_t, \boldsymbol{\xi}_t)$$

State estimator: y_t must be observed at time t, which can be done using a state estimator

MDP Notation and Terminology (Summary)

$t \in \{0, \dots, T\}$ $\mathbf{x} \in \mathcal{X}$ $\mathbf{u} \in \mathcal{U}$	discrete time discrete/continuous state discrete/continuous control
$p_0(\mathbf{x}) \ p_f(\mathbf{x}' \mid \mathbf{x}, \mathbf{u})$	prior probability density function defined on ${\cal X}$ transition/motion model
$\ell(\mathbf{x},\mathbf{u})$ $\mathfrak{q}(\mathbf{x})$	stage cost of choosing control \boldsymbol{u} in state \boldsymbol{x} terminal cost at state \boldsymbol{x}
$\pi_t(\mathbf{x}) \ V_t^{\pi}(\mathbf{x})$	control policy: function from state x at time <i>t</i> to control u value function: expected cumulative cost of starting at state x at time <i>t</i> and acting according to π
$\pi_t^*(\mathbf{x}) \ V_t^*(\mathbf{x})$	optimal control policy optimal value function

MDP Finite-horizon Optimal Control (Summary)

Finite-horizon Optimal Control

The finite-horizon optimal control problem in an MDP $(\mathcal{X}, \mathcal{U}, p_0, p_f, T, \ell, \mathfrak{q}, \gamma)$ with initial state **x** at time *t* is:

$$\min_{\pi_{t:\tau-1}} V_t^{\pi}(\mathbf{x}) := \mathbb{E}_{\mathbf{x}_{t+1:\tau}} \left[\gamma^{T-t} \mathfrak{q}(\mathbf{x}_T) + \sum_{\tau=t}^{T-1} \gamma^{\tau-t} \ell(\mathbf{x}_{\tau}, \pi_{\tau}(\mathbf{x}_{\tau})) \ \middle| \ \mathbf{x}_t = \mathbf{x} \right]$$
s.t. $\mathbf{x}_{\tau+1} \sim p_f(\cdot \mid \mathbf{x}_{\tau}, \pi_{\tau}(\mathbf{x}_{\tau})), \qquad \tau = t, \dots, T-1$
 $\mathbf{x}_{\tau} \in \mathcal{X}, \ \pi_{\tau}(\mathbf{x}_{\tau}) \in \mathcal{U}$

Outline

Markov Decision Processes

Open-Loop vs Closed-Loop Control

Partially Observable Models

Open-Loop vs Closed-Loop Control

- Open-loop policy: control inputs u_{0:T-1} are determined at once at time 0 as a function π_t(x₀) of x₀ and do not change online depending on x_t
- Closed-loop policy: control inputs are determined "just-in-time" as a function π_t(x_t) of the current state x_t
- Open-loop control is a special case of closed-loop control that disregards the state x_t and, hence, never gives better performance
- In the absence of noise (deterministic problem) and in special linear quadratic Gaussian (LQG) cases, open-loop and closed-loop control have the same performance
- Open-loop control is computationally much cheaper than closed-loop control. Consider a discrete-space example with |X| = 10 states, |U| = 10 control inputs, planning horizon T = 4, and given x₀:
 - There are $|\mathcal{U}|^{T} = 10^{4}$ open-loop strategies

▶ There are $|\mathcal{U}|(|\mathcal{U}|^{|\mathcal{X}|})^{T-1} = |\mathcal{U}|^{|\mathcal{X}|(T-1)+1} = 10^{31}$ closed-loop strategies

Open-loop feedback control (OLFC) recomputes a new open-loop sequence u_{t:T-1} online, whenever a new state x_t is available. OLFC is guaranteed to perform better than open-loop control and is computationally more efficient than closed-loop control.

Example: Chess Strategy Optimization

- Objective: come up with a strategy that maximizes the chances of winning a 2 game chess match
- Possible outcomes:
 - Win/Lose: 1 point for the winner, 0 for the loser
 - Draw: 0.5 points for each player
 - If the score is equal after 2 games, the players continue playing until one wins (sudden death)
- Playing styles:
 - **Timid**: draw with probability p_d and lose with probability $(1 p_d)$
 - **Bold**: win with probability p_w and lose with probability $(1 p_w)$
 - Assumption: $p_d > p_w$

Chess Match Model

- **State** x_t: 2-D vector with our and the opponent's score after the *t*-th game
- ▶ **Control** $u_t \in U = {\text{timid, bold}}$
- **Noise** *w_t*: score of the next game
- Since timid play does not make sense during the sudden death stage, the planning horizon is T = 2
- We can construct a time-dependent motion model P^u_{ijt} for t ∈ {0,1} (shown on the next slide)

• **Cost**: minimize loss probability:
$$-P_{win} = \mathbb{E}_{\mathbf{x}_{1:2}}\left[q(\mathbf{x}_2) + \sum_{t=0}^{1} \ell(\mathbf{x}_t, u_t)\right]$$
, where

$$\ell(\mathbf{x}, u) = 0 \text{ and } q(\mathbf{x}) = \begin{cases} -1 & \text{if } \mathbf{x} = \left(\frac{3}{2}, \frac{1}{2}\right) \text{ or } (2, 0) \\ -p_w & \text{if } \mathbf{x} = (1, 1) \\ 0 & \text{if } \mathbf{x} = \left(\frac{1}{2}, \frac{3}{2}\right) \text{ or } (0, 2) \end{cases}$$

Chess Transition Probabilities







Bold Play





Open-Loop Chess Strategy

► There are 4 possible open-loop policies:

1. timid-timid: $P_{win} = p_d^2 p_w$ 2. bold-bold: $P_{win} = p_w^2 + p_w(1 - p_w)p_w + (1 - p_w)p_w p_w = p_w^2(3 - 2p_w)$ 3. bold-timid: $P_{win} = p_w p_d + p_w(1 - p_d)p_w$

4. timid-bold:
$$P_{win} = p_d p_w + (1 - p_d) p_w^2$$

▶ Since $p_d^2 p_w \le p_d p_w \le p_d p_w + (1 - p_d) p_w^2$, timid-timid is not optimal

The best achievable winning probability is:

$$P_{win}^{*} = \max\{\overline{p_{w}^{2}(3-2p_{w})}, \overline{p_{d}p_{w}+(1-p_{d})p_{w}^{2}}\}$$
$$= p_{w}^{2} + p_{w}(1-p_{w})\max\{2p_{w}, p_{d}\}$$

▶ If
$$p_w \le 0.5$$
, then $P_{win}^* \le 0.5$
▶ For $p_w = 0.45$ and $p_d = 0.9$, $P_{win}^* = 0.43$
▶ For $p_w = 0.5$ and $p_d = 1.0$, $P_{win}^* = 0.5$

If p_d > 2p_w, bold-timid and timid-bold are optimal open-loop policies; otherwise bold-bold is optimal

Closed-Loop Chess Strategy

- There are 16 closed-loop policies
- Consider one option: play timid if and only if ahead (it will turn out that this is optimal)



The probability of winning is:

$$P_{win} = p_d p_w + p_w ((1 - p_d) p_w + p_w (1 - p_w)) = p_w^2 (2 - p_w) + p_w (1 - p_w) p_d$$

▶ In the closed-loop case, we can achieve P_{win} larger than 0.5 even when p_w is less than 0.5:

For
$$p_w = 0.45$$
 and $p_d = 0.9$, $P_{win} = 0.5$

For $p_w = 0.5$ and $p_d = 1.0$, $P_{win} = 0.625$

Outline

Markov Decision Processes

Open-Loop vs Closed-Loop Control

Partially Observable Models

Hidden Markov Model

Hidden Markov Model

Markov Chain with partially observable states defined by tuple $(\mathcal{X}, \mathcal{Z}, p_0, p_f, p_h)$

- \mathcal{X} is a discrete or continuous state space
- $\blacktriangleright \ \mathcal{Z}$ is a discrete or continuous observation space
- $p_0(\cdot)$ is a prior pdf defined on \mathcal{X}
- ▶ $p_f(\cdot | \mathbf{x}_t)$ is a conditional pdf defined on \mathcal{X} for given $\mathbf{x}_t \in \mathcal{X}$ (matrix P with $P_{ij} = p_f(j | x_t = i)$ in finite-dim case)
- *p_h*(· | **x**_t) is a conditional pdf defined on Z for given **x**_t ∈ X (matrix O with O_{ij} := *p_h*(*j* | *x_t* = *i*) in finite-dim case)

Partially Observable Markov Decision Process

Partially Observable Markov Decision Process

Markov Decision Process with partially observable states defined by tuple $(\mathcal{X}, \mathcal{U}, \mathcal{Z}, p_0, p_f, p_h, T, \ell, \mathfrak{q}, \gamma)$

- \mathcal{X} is a discrete or continuous state space
- \blacktriangleright \mathcal{U} is a discrete or continuous control space
- $\blacktriangleright \ \mathcal{Z}$ is a discrete or continuous observation space
- $p_0(\cdot)$ is a prior pdf defined on \mathcal{X}
- ▶ $p_f(\cdot | \mathbf{x}_t, \mathbf{u}_t)$ is a conditional pdf defined on \mathcal{X} for given $\mathbf{x}_t \in \mathcal{X}$ and $\mathbf{u}_t \in \mathcal{U}$ (matrices P^u with elements $P_{ij}^u = p_f(j | x_t = i, u_t = u)$ in finite-dim case)
- *p_h*(· | **x**_t) is a conditional pdf defined on Z for given **x**_t ∈ X (matrix O with O_{ij} := p_h(j | x_t = i) in finite-dim case)
- T is a finite/infinite time horizon
- ▶ $\ell(\mathbf{x}, \mathbf{u})$ is stage cost of applying control $\mathbf{u} \in \mathcal{U}$ in state $\mathbf{x} \in \mathcal{X}$
- q(x) is terminal cost of being in state x at time T
- ▶ $\gamma \in [0,1]$ is a discount factor

Comparison of Markov Models

	observed	partially observed
uncontrolled	Markov Chain/MRP	НММ
controlled	MDP	POMDP

- Markov Chain + Partial Observability = HMM
- Markov Chain + Control = MDP
- Markov Chain + Partial Observability + Control = HMM + Control = MDP + Partial Observability = POMDP

Bayes Filter

- A probabilistic inference technique for summarizing information $\mathbf{i}_t := (\mathbf{z}_{0:t}, \mathbf{u}_{0:t-1})$ about a partially observable state \mathbf{x}_t
- ► The Bayes filter keeps track of: $\frac{p_{t|t}(\mathbf{x}_t) := p(\mathbf{x}_t \mid \mathbf{z}_{0:t}, \mathbf{u}_{0:t-1})}{p_{t+1|t}(\mathbf{x}_{t+1}) := p(\mathbf{x}_{t+1} \mid \mathbf{z}_{0:t}, \mathbf{u}_{0:t})}$
- Derived using total probability, conditional probability, and Bayes rule based on the motion and observation models of the system
- Motion model: $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t) \sim p_f(\cdot \mid \mathbf{x}_t, \mathbf{u}_t)$
- Observation model: $\mathbf{z}_t = h(\mathbf{x}_t, \mathbf{v}_t) \sim p_h(\cdot \mid \mathbf{x}_t)$
- Bayes filter: consists of predict and update steps:

$$p_{t+1|t+1}(\mathbf{x}_{t+1}) = \underbrace{\frac{1}{p(\mathbf{z}_{t+1}|\mathbf{z}_{0:t}, \mathbf{u}_{0:t})} p_{h}(\mathbf{z}_{t+1} \mid \mathbf{x}_{t+1})}_{\mathbf{Update}} \underbrace{\int p_{f}(\mathbf{x}_{t+1} \mid \mathbf{x}_{t}, \mathbf{u}_{t}) p_{t|t}(\mathbf{x}_{t}) d\mathbf{x}_{t}}_{\mathbf{Update}}$$

Bayes Filter Example



Equivalence of POMDPs and MDPs

- A POMDP (X, U, Z, p₀, p_f, p_h, T, ℓ, q, γ) is equivalent to an MDP (P(X), U, p₀, p_ψ, T, ℓ, q̄, γ) such that:
 - **State space**: $\mathcal{P}(\mathcal{X})$ is the **continuous** space of pdfs over \mathcal{X}
 - If \mathcal{X} is continuous, then $\mathcal{P}(\mathcal{X}) := \{ p : \mathcal{X} \to \mathbb{R}_{\geq 0} \mid \int p(\mathbf{x}) d\mathbf{x} = 1 \}$

• If
$$|\mathcal{X}| = N$$
, then $\mathcal{P}(\mathcal{X}) := \{\mathbf{p} \in [0, 1]^N \mid \mathbf{1}^\top \mathbf{p} = 1\}$

- lnitial state: $p_0 \in \mathcal{P}(\mathcal{X})$
- Motion model: the Bayes filter p_{t+1|t+1} = ψ(p_{t|t}, u_t, z_{t+1}) acts as a motion model for p_{t|t} with motion noise given by the observations z_{t+1} with density:

$$\eta(\mathbf{z} \mid \boldsymbol{p}_{t|t}, \mathbf{u}_t) := \int \int p_h(\mathbf{z} \mid \mathbf{x}_{t+1}) p_f(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{u}_t) p_{t|t}(\mathbf{x}_t) d\mathbf{x}_t d\mathbf{x}_{t+1}$$

Cost: the equivalent MDP stage and terminal cost functions are the expected values of the POMDP stage and terminal costs:

$$\bar{\ell}(p,\mathbf{u}) := \int \ell(\mathbf{x},\mathbf{u})p(\mathbf{x})d\mathbf{x} \qquad \bar{\mathfrak{q}}(p) := \int \mathfrak{q}(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

POMDP Finite-horizon Optimal Control

▶ POMDP $(\mathcal{X}, \mathcal{U}, \mathcal{Z}, p_0, p_f, p_h, T, \ell, \mathfrak{q}, \gamma)$:

$$\min_{\boldsymbol{\pi}_{0:T-1}} \mathbb{E} \left[\gamma^{T} \boldsymbol{\mathfrak{q}}(\mathbf{x}_{T}) + \sum_{t=0}^{T-1} \gamma^{t} \ell(\mathbf{x}_{t}, \mathbf{u}_{t}) \right]$$
s.t. $\mathbf{x}_{t+1} \sim p_{f}(\cdot \mid \mathbf{x}_{t}, \mathbf{u}_{t}), \quad t = 0, \dots, T-1$
 $\mathbf{z}_{t+1} \sim p_{h}(\cdot \mid \mathbf{x}_{t}), \quad t = 0, \dots, T-1$
 $\mathbf{u}_{t} \sim \pi_{t}(\cdot \mid \mathbf{i}_{t}), \quad t = 0, \dots, T-1$
 $\mathbf{x}_{0} \sim p_{0}(\cdot)$

• Equivalent MDP $(\mathcal{P}(\mathcal{X}), \mathcal{U}, p_0, p_{\psi}, T, \overline{\ell}, \overline{\mathfrak{q}}, \gamma)$ with state $p_{t|t}$:

$$\min_{\pi_{0:T-1}} V_0^{\pi}(p_0) = \mathbb{E} \left[\gamma^T \bar{\mathfrak{q}}(p_{T|T}) + \sum_{t=0}^{T-1} \gamma^t \bar{\ell}(p_{t|t}, \mathbf{u}_t) \right]$$

s.t. $p_{t+1|t+1} = \psi(p_{t|t}, \mathbf{u}_t, \mathbf{z}_{t+1}), \ t = 0, \dots, T-1$
 $\mathbf{z}_{t+1} \sim \eta(\cdot \mid p_{t|t}, \mathbf{u}_t), \qquad t = 0, \dots, T-1$
 $u_t \sim \pi_t(\cdot \mid p_{t|t}), \qquad t = 0, \dots, T-1$

Due to the equivalence between POMDPs and MDPs, we will focus exclusively on MDPs