

# Approximate information states for partially observable systems

Aditya Mahajan  
McGill University

Joint work with: Jayakumar Subramanian, Amit Sinha,  
Nima Akbarzadeh, Raihan Seraj, Erfan Seyedsalehi

DataAI Seminar  
11 Nov 2023

## Recent successes of RL

## Recent successes of RL



Alpha Go

## Recent successes of RL



Arcade games

## Recent successes of RL



Robotic grasping

## Recent successes of RL

- ▷ Algorithms based on comprehensive theory



Robotic grasping

## Recent successes of RL

- ▷ Algorithms based on comprehensive theory
- ▷ The theory is restricted almost exclusively to systems with **perfect state observations**.



Robotic grasping

## Recent successes of RL

- ▷ Algorithms based on comprehensive theory
- ▷ The theory is restricted almost exclusively to systems with **perfect state observations**.



Robotic grasping

## Many real-world applications are partially observed

- ▷ Healthcare
- ▷ Autonomous driving
- ▷ Finance (portfolio management)
- ▷ Retail and marketing

## Recent successes of RL

- ▷ Algorithms based on comprehensive theory
- ▷ The theory is restricted almost exclusively to systems with **perfect state observations**.



Robotic grasping

## Many real-world applications are partially observed

- ▷ Healthcare
- ▷ Autonomous driving
- ▷ Finance (portfolio management)
- ▷ Retail and marketing

How do we develop a theory for RL for partially observed systems?

# Outline



## Background

- ▷ Review of MDPs and RL
- ▷ Review of POMDPs
- ▷ Why is RL for POMDPs difficult?

# Outline



## Background

- ▷ Review of MDPs and RL
- ▷ Review of POMDPs
- ▷ Why is RL for POMDPs difficult?



## Approximate Planning for POMDPs

- ▷ Preliminaries on information state
- ▷ Approximate information state
- ▷ Approximation bounds

# Outline



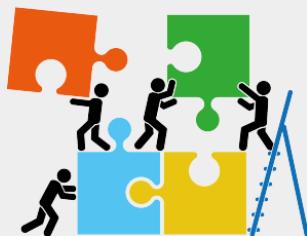
## Background

- ▷ Review of MDPs and RL
- ▷ Review of POMDPs
- ▷ Why is RL for POMDPs difficult?



## Approximate Planning for POMDPs

- ▷ Preliminaries on information state
- ▷ Approximate information state
- ▷ Approximation bounds



## RL for POMDPs

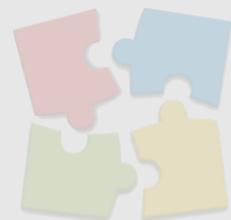
- ▷ From approximation bounds to RL
- ▷ Numerical experiments

# Outline



## Background

- ▷ Review of MDPs and RL
- ▷ Review of POMDPs
- ▷ Why is RL for POMDPs difficult?



## Approximate Planning for POMDPs

- ▷ Preliminaries on information state
- ▷ Approximate information state
- ▷ Approximation bounds

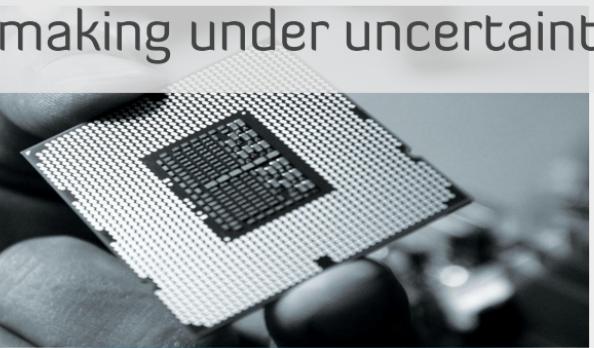


## RL for POMDPs

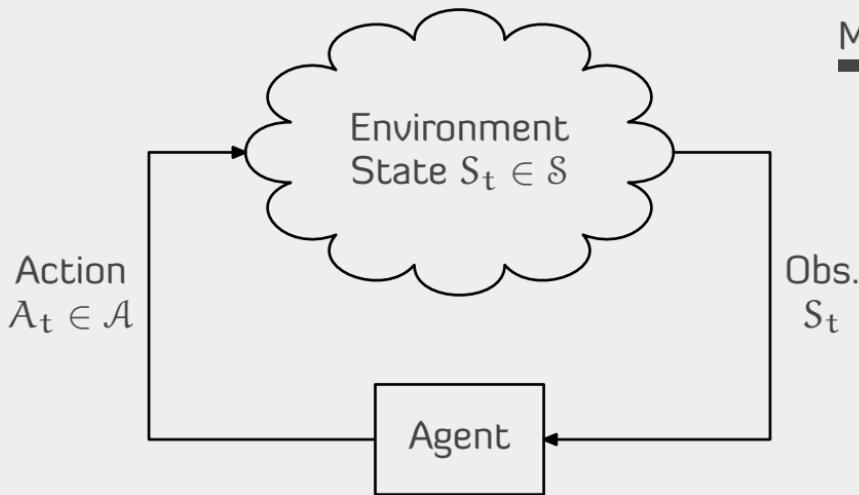
- ▷ From approximation bounds to RL
- ▷ Numerical experiments



**Common theme:** multi-stage  
decision making under uncertainty



# Review: Markov decision processes (MDPs)



## MDP: MARKOV DECISION PROCESS

Dynamics:  $\mathbb{P}(S_{t+1} | S_t, A_t)$

Observations:  $S_t$

Reward  $R_t = r(S_t, A_t)$ .

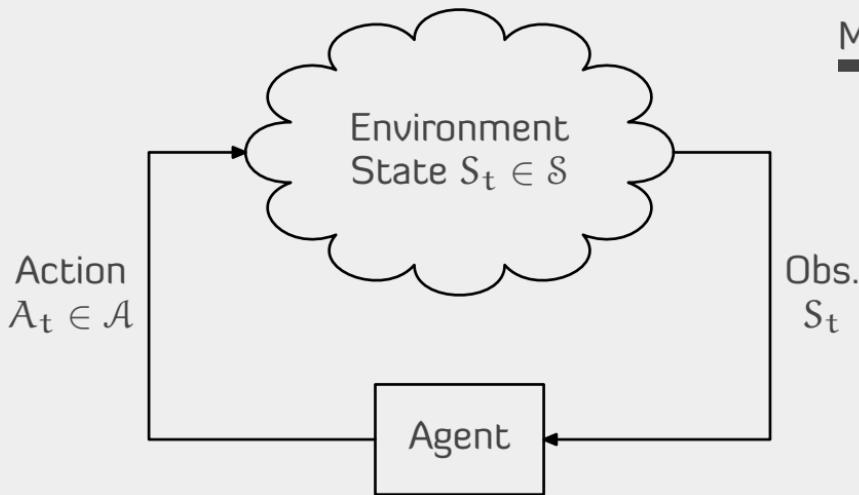
Action:  $A_t \sim \pi_t(S_{1:t}, A_{1:t-1})$ .

$\pi = (\pi_t)_{t \geq 1}$  is called a **policy**.

The objective is to choose a policy  $\pi$  to maximize:

$$J(\pi) := \mathbb{E}^\pi \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$$

# Review: Markov decision processes (MDPs)



## MDP: MARKOV DECISION PROCESS

Dynamics:  $\mathbb{P}(S_{t+1} | S_t, A_t)$

Observations:  $S_t$

Reward  $R_t = r(S_t, A_t)$ .

Action:  $A_t \sim \pi_t(S_{1:t}, A_{1:t-1})$ .

$\pi = (\pi_t)_{t \geq 1}$  is called a **policy**.

The objective is to choose a policy  $\pi$  to maximize:

$$\mathbb{E}_{\pi} \sum_{t=1}^{\infty} \gamma^t R_t$$

## Conceptual challenge

- ▷ Brute force search has an exponential complexity in time horizon.
- ▷ How to efficiently search an optimal policy?

# Review: Markov decision processes (MDPs)

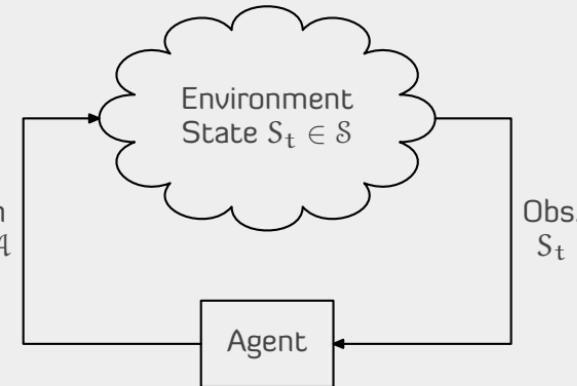
## Key simplifying ideas

### Principle of Irrelevant Information

#### Structure of optimal policy

There is no loss of optimality in choosing the action  $A_t$  as a function of the current state  $S_t$

█ Blackwell, "Memoryless strategies in finite-stage dynamic prog.," Annals Math. Stats, 1964.



# Review: Markov decision processes (MDPs)

## Key simplifying ideas

### Principle of Irrelevant Information

#### Structure of optimal policy

There is no loss of optimality in choosing the action  $A_t$  as a function of the current state  $S_t$

█ Blackwell, "Memoryless strategies in finite-stage dynamic prog.," Annals Math. Stats, 1964.

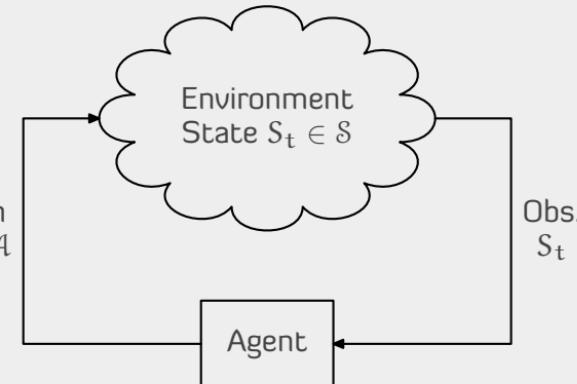
### Principle of Optimality

#### Dynamic Program

The optimal control policy is given a DP with state  $S_t$ :

$$V(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \int V(s') P(ds'|s, a) \right\}$$

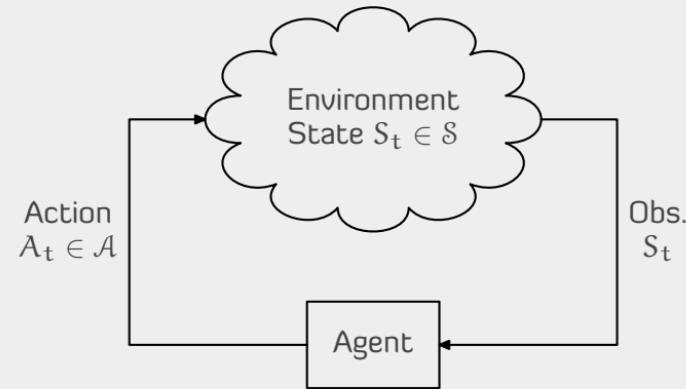
█ Bellman, "Dynamic Programming," 1957.



# Review: Reinforcement Learning (RL)

## The RL setting

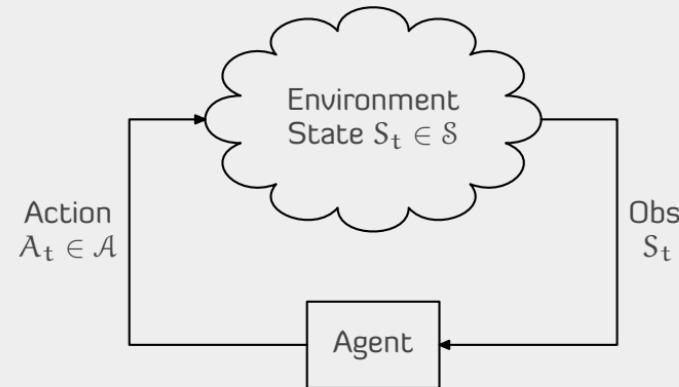
- ▷ Dynamics and reward functions are unknown.
- ▷ Agent can interact with the environment and observe states and rewards.
- ▷ Design an algorithm that asymptotically identifies an optimal policy.



# Review: Reinforcement Learning (RL)

## The RL setting

- ▷ Dynamics and reward functions are unknown.
- ▷ Agent can interact with the environment and observe states and rewards.
- ▷ Design an algorithm that asymptotically identifies an optimal policy.



### Value based methods

Estimate the Q-function  $Q(s, a) = r(s, a) + \gamma \int V(s')P(ds'|s, a)$  using temporal difference learning (i.e., stochastic approximation).

[Watkins and Dayan, 1992; Tsitsiklis, 1994]

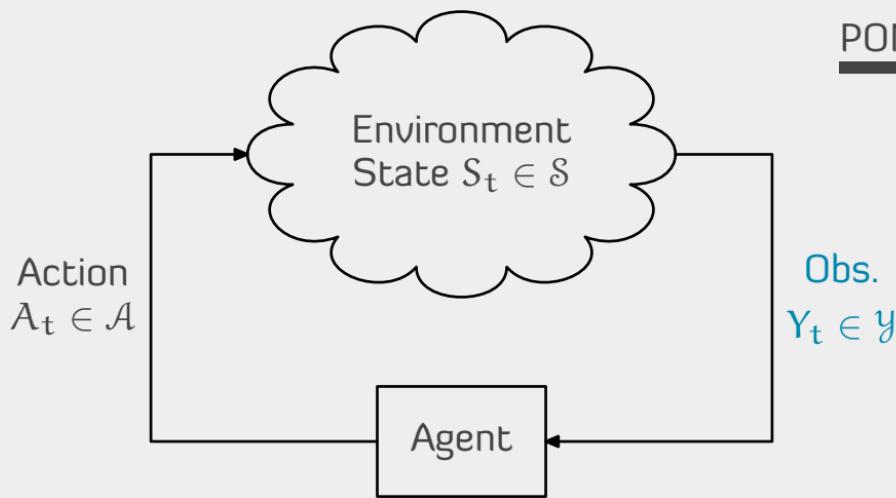
### Policy-based methods

Use parameterized policies  $\pi_\theta$ . Estimate  $\nabla_\theta V_\theta(s)$  using single trajectory gradient estimates (i.e., infinitesimal perturbation analysis).

[Sutton 2000, Marback and Tsitsiklis 2001], [Cao, 1985; Ho, 1987]

Why is learning difficult in partially  
observable environments?

# Review: Planning in partially observable environments



## POMDP: PARTIALLY OBSERVABLE MARKOV DECISION PROCESS

Dynamics:  $\mathbb{P}(S_{t+1} | S_t, A_t)$

Observations:  $\mathbb{P}(Y_t | S_t)$

Reward  $R_t = r(S_t, A_t)$ .

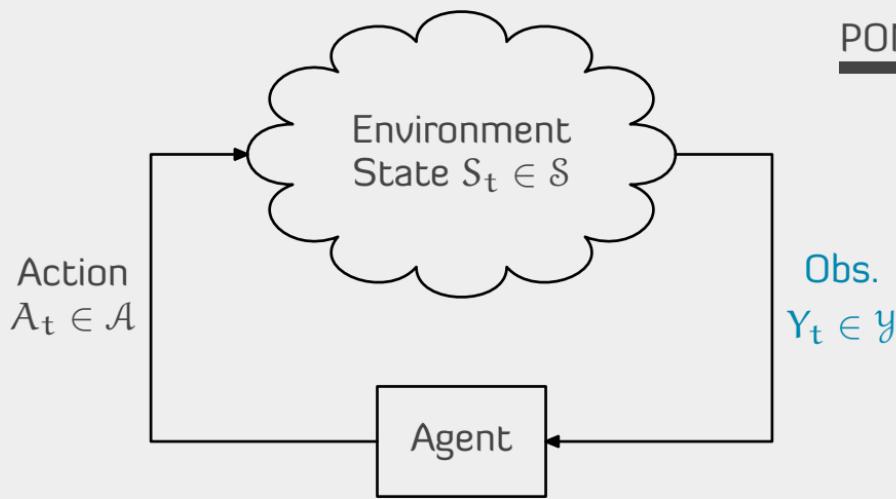
Action:  $A_t \sim \pi_t(Y_{1:t}, A_{1:t-1})$ .

$\pi = (\pi_t)_{t \geq 1}$  is called a **policy**.

The objective is to choose a policy  $\pi$  to maximize:

$$J(\pi) := \mathbb{E}^\pi \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$$

# Review: Planning in partially observable environments



## POMDP: PARTIALLY OBSERVABLE MARKOV DECISION PROCESS

Dynamics:  $\mathbb{P}(S_{t+1} | S_t, A_t)$

Observations:  $\mathbb{P}(Y_t | S_t)$

Reward  $R_t = r(S_t, A_t)$ .

Action:  $A_t \sim \pi_t(Y_{1:t}, A_{1:t-1})$ .

$\pi = (\pi_t)_{t \geq 1}$  is called a **policy**.

The objective is to choose a policy  $\pi$  to maximize:

### Conceptual challenge

- ▷ Action is a function of the history of observations and actions.
- ▷ The history is increasing in time. So, the search complexity increases exponentially in time.

# Review: Planning in partially observable environments

## Key simplifying idea

Define **belief state**  $B_t \in \Delta(\mathcal{S})$  as  $B_t(s) = \mathbb{P}(S_t = s | Y_{1:t}, A_{1:t-1})$ .

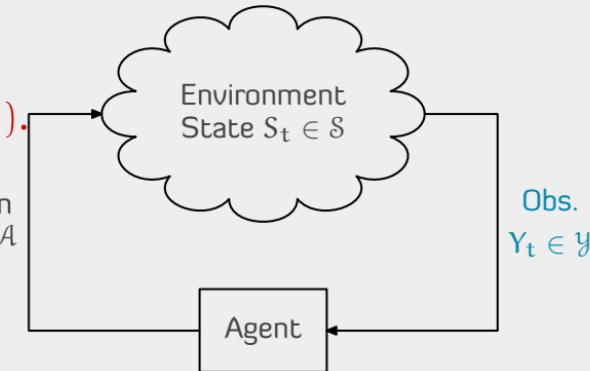
► Belief state updates in a state-like manner

$$B_{t+1} = \text{function}(B_t, Y_{t+1}, A_t).$$

► Belief state is sufficient to evaluate rewards

$$\mathbb{E}[R_t | Y_{1:t}, A_{1:t}] = \hat{r}(B_t, A_t).$$

Thus,  $\{B_t\}_{t \geq 1}$  is a **perfectly observed** controlled Markov process.



❑ Astrom, "Optimal control of Markov processes with incomplete information," JMAA 1965.

❑ Stratonovich, "Conditional Markov Processes," TVP 1960.

# Review: Planning in partially observable environments

## Key simplifying idea

Define **belief state**  $B_t \in \Delta(\mathcal{S})$  as  $B_t(s) = \mathbb{P}(S_t = s | Y_{1:t}, A_{1:t-1})$ .

► Belief state updates in a state-like manner

$$B_{t+1} = \text{function}(B_t, Y_{t+1}, A_t).$$

► Belief state is sufficient to evaluate rewards

$$\mathbb{E}[R_t | Y_{1:t}, A_{1:t}] = \hat{r}(B_t, A_t).$$

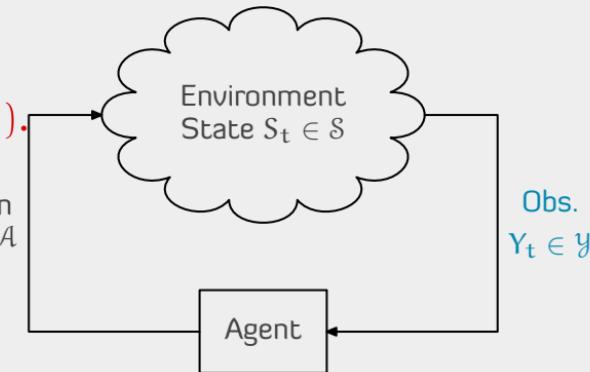
Thus,  $\{B_t\}_{t \geq 1}$  is a **perfectly observed** controlled Markov process. Therefore:

### Structure of optimal policy

There is no loss of optimality in choosing the action  $A_t$  as a function of the belief state  $B_t$

### Dynamic Program

The optimal control policy is given a DP with belief  $B_t$  as state.



# Implications of the POMDP modeling framework

## Implications for planning

- ▶ Allows the use of the MDP machinery for partially observed systems.
- ▶ Various exact and approximate algorithms to efficiently solve the DP.
  - Exact:** incremental pruning, witness algorithm, linear support algo
  - Approximate:** QMDP, point based methods, SARSOP, DESPOT, ...

# Implications of the POMDP modeling framework

## Implications for learning

- ▶ Allows the use of the MDP machinery for partially observed systems.
- ▶ The construction of the belief state depends on the system model.
- ▶ So, when the system model is unknown, we cannot construct the belief state and therefore cannot use standard RL algorithms.

# Implications of the POMDP modeling framework

## Implications for learning

- ▶ Allows the use of the MDP machinery for partially observed systems.
- ▶ The construction of the belief state depends on the system model.
- ▶ So, when the system model is unknown, we cannot construct the belief state and therefore cannot use standard RL algorithms.
- ▶ **On the theoretical side:**
  - ▶ Propose alternative methods: PSRs (predictive state representations), bisimulation metrics, . . .
  - ▶ Good theoretical guarantees, but difficult to scale.

# Implications of the POMDP modeling framework

## Implications for learning

- ▷ Allows the use of the MDP machinery for partially observed systems.
- ▷ The construction of the belief state depends on the system model.
- ▷ So, when the system model is unknown, we cannot construct the belief state and therefore cannot use standard RL algorithms.
- ▷ **On the theoretical side:**
  - ▷ Propose alternative methods: PSRs (predictive state representations), bisimulation metrics, . . .
  - ▷ Good theoretical guarantees, but difficult to scale.
- ▷ **On the practical side:**
  - ▷ Simply stack the previous  $k$  observations and treat it as a “state”.
  - ▷ Instead of a CNN, use an RNN to model policy and action-value fn.
  - ▷ Can be made to work but lose theoretical guarantees and insights.

**Our result: A theoretically grounded method  
for RL in partially observable models  
which has strong empirical performance  
for high-dimensional environments.**

- ▶ paper 1: JMLR, Feb 2022  
co-authors: J. Subramanian, A. Sinha, and R. Seraj.
- ▶ paper 2: arxiv:2306.05991, June 2024  
co-authors: E. SeyedSalehi, N. Akbarzadeh, A. Sinha

# Outline



## Background

- ▷ Review of MDPs and RL
- ▷ Review of POMDPs
- ▷ Why is RL for POMDPs difficult?



## Approximate Planning for POMDPs

- ▷ Preliminaries on information state
- ▷ Approximate information state
- ▷ Approximation bounds



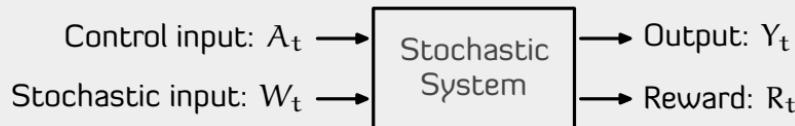
## RL for POMDPs

- ▷ From approximation bounds to RL
- ▷ Numerical experiments

# System model

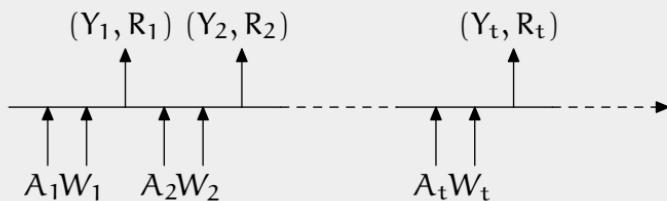
- ▶ In RL, unobserved state space may not be known
- ▶ So, we work directly with input-output model

# System model



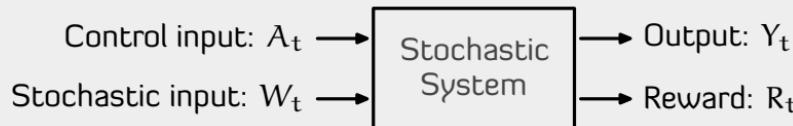
$$Y_t = f_t(A_{1:t}, W_{1:t}),$$

$$R_t = r_t(A_{1:t}, W_{1:t}).$$



- ▶ In RL, unobserved state space may not be known
- ▶ So, we work directly with input-output model

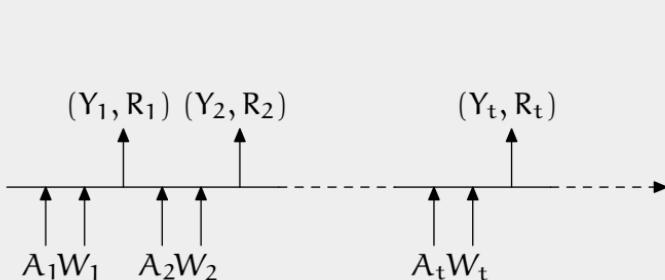
# System model



$$Y_t = f_t(A_{1:t}, W_{1:t}),$$

$$R_t = r_t(A_{1:t}, W_{1:t}).$$

- ▶  $H_t = (Y_{1:t-1}, A_{1:t-1})$  denotes the history of all data available to the agent at time  $t$ .
- ▶ Agent chooses an  $A_t \sim \pi_t(H_t)$ .
- ▶  $\pi = (\pi_1, \pi_2, \dots)$  denotes the control policy.



The objective is to choose a policy  $\pi$  to maximize:

$$J(\pi) := \mathbb{E}^\pi \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$$

- ▶ In RL, unobserved state space may not be known
- ▶ So, we work directly with input-output model

## Key solution concept: Information state

Informally, an information state is a compression of information which is sufficient for performance evaluation and predicting itself.

# Key solution concept: Information state

Informally, an information state is a compression of information which is sufficient for performance evaluation and predicting itself.

## Historical overview

- ▶ **Old concept.** May be viewed as a generalization of the notion of state (Nerode, 1958).
- ▶ Informal definitions given in Kwakernaak (1965), Bohlin (1970), Davis and Varaiya (1972), Kumar and Varaiya (1986) but no formal analysis.
- ▶ Related to but different from concepts such bisimulation, predictive state representations (PSR), and  $\varepsilon$ -machines.

# Information state: Intuition

## Sufficient Statistics

$$S \quad Y \quad A \quad R = r(S, A)$$

State    Obs.    Action    Reward

- ▷ Consider a **compression of data**  $Z = \sigma(Y)$
- ▷  $Z$  is a **sufficient statistic** for **evaluating reward** if

$$(P1) \quad \mathbb{E}[R | Y = y, A = a] = \mathbb{E}[R | Z = \sigma(y), A = a] =: \hat{r}(\sigma(y), a)$$

# Information state: Intuition

## Sufficient Statistics

$$S \quad Y \quad A \quad R = r(S, A)$$

State    Obs.    Action    Reward

- ▶ Consider a **compression of data**  $Z = \sigma(Y)$
  - ▶  $Z$  is a **sufficient statistic** for **evaluating reward** if
- (P1)  $\mathbb{E}[R | Y = y, A = a] = \mathbb{E}[R | Z = \sigma(y), A = a] =: \hat{r}(\sigma(y), a)$

## Self-predictive property

In a MDP/POMDP, a sequence of sufficient statistics  $\{Z_t = \sigma(Y_t)\}$  where  $Z_t$  is sufficient for  $R_t$  is not **sufficient for dynamic programming!** We also need

$$\begin{aligned} (P2) \quad & \mathbb{P}(Z_{t+1} = z_{t+1} | H_t = h_t, A_t = a_t) \\ & = \mathbb{P}(Z_{t+1} = z_{t+1} | Z_t = \sigma_t(H_t), A_t = a_t) \end{aligned}$$

## Information state: Definition

Given a state space  $\mathcal{Z}$ , an INFORMATION STATE GENERATOR is a tuple of

- ▶ history compression functions  $\{\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$
- ▶ reward function  $\hat{r}: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
- ▶ transition kernel  $\hat{P}: \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Z})$

which satisfies two properties:

# Information state: Definition

Given a state space  $\mathcal{Z}$ , an INFORMATION STATE GENERATOR is a tuple of

- ▶ history compression functions  $\{\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$
- ▶ reward function  $\hat{r}: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
- ▶ transition kernel  $\hat{P}: \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Z})$

which satisfies two properties:

**(P1) The reward function  $\hat{r}$  is sufficient for performance evaluation:**

$$\mathbb{E}[R_t | H_t = h_t, A_t = a_t] = \hat{r}(\sigma_t(h_t), a_t).$$

# Information state: Definition

Given a state space  $\mathcal{Z}$ , an INFORMATION STATE GENERATOR is a tuple of

- ▶ history compression functions  $\{\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$
- ▶ reward function  $\hat{r}: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$
- ▶ transition kernel  $\hat{P}: \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Z})$

which satisfies two properties:

**(P1) The reward function  $\hat{r}$  is sufficient for performance evaluation:**

$$\mathbb{E}[R_t | H_t = h_t, A_t = a_t] = \hat{r}(\sigma_t(h_t), a_t).$$

**(P2) The transition kernel  $\hat{P}$  is sufficient for predicting the info state:**

$$\mathbb{P}(Z_{t+1} \in B | H_t = h_t, A_t = a_t) = \hat{P}(B | \sigma_t(h_t), a_t).$$

## Information state: Key result

An information state **always** leads to a dynamic programming decomposition.

## Information state: Key result

An information state **always** leads to a dynamic programming decomposition.

Let  $\{Z_t\}_{t \geq 1}$  be **any** information state process. Let  $\hat{V}$  be the fixed point of:

$$\hat{V}(z) = \max_{a \in \mathcal{A}} \left\{ \hat{r}(z, a) + \gamma \int_Z \hat{V}(z_+) \hat{P}(dz_+ | z, a) \right\}$$

Let  $\pi^*(z)$  denote the arg max of the RHS.

**Then, the policy  $\pi = (\pi_t)_{t \geq 1}$  given by  $\pi_t(h_t) = \pi^*(\sigma_t(h_t))$  is optimal.**

# Examples of information state

Markov decision processes (MDP)

Current state  $S_t$  is an info state

POMDP

Belief state is an info state

# Examples of information state

## Markov decision processes (MDP)

Current state  $S_t$  is an info state

## MDP with delayed observations

$(S_{t-\delta+1}, A_{t-\delta+1:t-1})$  is an info state

## POMDP

Belief state is an info state

## POMDP with delayed observations

$(\mathbb{P}(S_{t-\delta}|Y_{1:t-\delta}, A_{1:t-\delta}), A_{t-\delta+1:t-1})$   
is info state

# Examples of information state

## Markov decision processes (MDP)

Current state  $S_t$  is an info state

## MDP with delayed observations

$(S_{t-\delta+1}, A_{t-\delta+1:t-1})$  is an info state

## POMDP

Belief state is an info state

## POMDP with delayed observations

$(\mathbb{P}(S_{t-\delta}|Y_{1:t-\delta}, A_{1:t-\delta}), A_{t-\delta+1:t-1})$   
is info state

## Linear Quadratic Gaussian (LQG)

The state estimate  $\mathbb{E}[S_t|H_t]$  is an info state

## Machine Maintenance

$(\tau, S_\tau^+)$  is info state,  
where  $\tau$  is the time of last maintenance

# And now to Approximate Information States ....

## Main idea

- ▶ Info state is defined in terms of two properties (P1) & (P2).
- ▶ An AIS is a process which satisfies these **approximately**

# And now to Approximate Information States ....

## Main idea

- ▶ Info state is defined in terms of two properties (P1) & (P2).
- ▶ An AIS is a process which satisfies these **approximately**
- ▶ Show that AIS always leads to approx. DP
- ▶ Recover (and improve upon) many existing results

## Approximate Information state: Definition

An  $(\varepsilon, \delta)$ -APPROXIMATE INFORMATION STATE (AIS) generator is a tuple  $(\sigma_t, \hat{r}, \hat{P})$  which approximately satisfies (P1) and (P2):

# Approximate Information state: Definition

An  $(\varepsilon, \delta)$ -APPROXIMATE INFORMATION STATE (AIS) generator is a tuple  $(\sigma_t, \hat{r}, \hat{P})$  which approximately satisfies (P1) and (P2):

**(AP1)  $\hat{r}$  is sufficient for approximate performance evaluation:**

$$|\mathbb{E}[R_t | H_t = h_t, A_t = a_t] - \hat{r}(\sigma_t(h_t), a_t)| \leq \varepsilon$$

# Approximate Information state: Definition

An  $(\varepsilon, \delta)$ -APPROXIMATE INFORMATION STATE (AIS) generator is a tuple  $(\sigma_t, \hat{r}, \hat{P})$  which approximately satisfies (P1) and (P2):

**(AP1)  $\hat{r}$  is sufficient for approximate performance evaluation:**

$$|\mathbb{E}[R_t | H_t = h_t, A_t = a_t] - \hat{r}(\sigma_t(h_t), a_t)| \leq \varepsilon$$

**(AP2)  $\hat{P}$  is sufficient for approximately predicting next AIS:**

$$d_{\mathfrak{F}}(\mathbb{P}(Z_{t+1} = \cdot | H_t = h_t, A_t = a_t), \hat{P}(\cdot | \sigma_t(h_t), a_t)) \leq \delta$$

# Approximate Information state: Definition

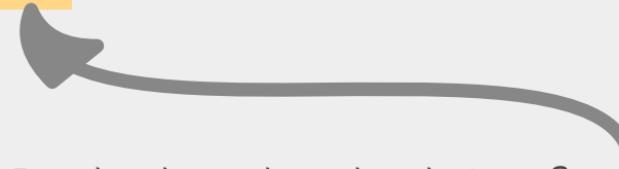
An  $(\varepsilon, \delta)$ -APPROXIMATE INFORMATION STATE (AIS) generator is a tuple  $(\sigma_t, \hat{r}, \hat{P})$  which approximately satisfies (P1) and (P2):

**(AP1)  $\hat{r}$  is sufficient for approximate performance evaluation:**

$$|\mathbb{E}[R_t | H_t = h_t, A_t = a_t] - \hat{r}(\sigma_t(h_t), a_t)| \leq \varepsilon$$

**(AP2)  $\hat{P}$  is sufficient for approximately predicting next AIS:**

$$d_{\mathfrak{F}}(\mathbb{P}(Z_{t+1} = \cdot | H_t = h_t, A_t = a_t), \hat{P}(\cdot | \sigma_t(h_t), a_t)) \leq \delta$$



Results depend on the choice of **metric on probability spaces**

## AIS based approximation bounds

Let  $V$  denote the optimal value and  $\hat{V}$  denote the fixed point of the following equations:

$$\hat{V}(z) = \max_{a \in \mathcal{A}} \left\{ \hat{r}(z, a) + \gamma \int_{\mathcal{Z}} \hat{V}(z_+) \hat{P}(dz_+ | z, a) \right\}$$

# AIS based approximation bounds

Let  $V$  denote the optimal value and  $\hat{V}$  denote the fixed point of the following equations:

$$\hat{V}(z) = \max_{a \in \mathcal{A}} \left\{ \hat{r}(z, a) + \gamma \int_{\mathcal{Z}} \hat{V}(z_+) \hat{P}(dz_+ | z, a) \right\}$$

## Value function approximation

The value function  $\hat{V}$  is approximately optimal, i.e.,

$$|V_t(h_t) - \hat{V}(\sigma_t(h_t))| \leq \alpha := \frac{\varepsilon + \gamma \rho_{\mathfrak{F}}(\hat{V}) \delta}{1 - \gamma}.$$

# AIS based approximation bounds

Let  $V$  denote the optimal value and  $\hat{V}$  denote the fixed point of the following equations:

$$\hat{V}(z) = \max_{a \in \mathcal{A}} \left\{ \hat{r}(z, a) + \gamma \int_{\mathcal{Z}} \hat{V}(z_+) \hat{P}(dz_+ | z, a) \right\}$$

Depends on metric

Value function  
approximation

The value function  $\hat{V}$  is approximately optimal, i.e.,

$$|V_t(h_t) - \hat{V}(\sigma_t(h_t))| \leq \alpha := \frac{\varepsilon + \gamma \rho_{\mathfrak{F}}(\hat{V}) \delta}{1 - \gamma}.$$

# AIS based approximation bounds

Let  $V$  denote the optimal value and  $\hat{V}$  denote the fixed point of the following equations:

$$\hat{V}(z) = \max_{a \in \mathcal{A}} \left\{ \hat{r}(z, a) + \gamma \int_{\mathcal{Z}} \hat{V}(z_+) \hat{P}(dz_+ | z, a) \right\}$$

Depends on metric

## Value function approximation

The value function  $\hat{V}$  is approximately optimal, i.e.,

$$|V_t(h_t) - \hat{V}(\sigma_t(h_t))| \leq \alpha := \frac{\varepsilon + \gamma \rho_{\mathfrak{F}}(\hat{V}) \delta}{1 - \gamma}.$$

## Policy approximation

Let  $\hat{\pi}^*: \mathcal{Z} \rightarrow \Delta(\mathcal{A})$  be an optimal policy for  $\hat{V}$ .

Then, the policy  $\pi = (\pi_1, \pi_2, \dots)$  where  $\pi_t = \hat{\pi}^* \circ \sigma_t$  is approx. optimal:

$$V_t(h_t) - V_t^\pi(h_t) \leq 2\alpha.$$

## Some remarks on AIS

- ▷ Two ways to interpret the results:
  - ▷ Given the information state space  $\mathcal{Z}$ , find the best compression  $\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}$
  - ▷ Given any compression function  $\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}$ , find the approximation error.

# Some remarks on AIS

- ▷ Two ways to interpret the results:
  - ▷ Given the information state space  $\mathcal{Z}$ , find the best compression  $\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}$
  - ▷ Given any compression function  $\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}$ , find the approximation error.
- ▷ Most of the existing literature on approximate DPs focuses on the first interpretation
- ▷ The second interpretation allows us to develop AIS-based RL algorithms

# Some remarks on AIS

- ▷ Two ways to interpret the results:
  - ▷ Given the information state space  $\mathcal{Z}$ , find the best compression  $\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}$
  - ▷ Given any compression function  $\sigma_t: \mathcal{H}_t \rightarrow \mathcal{Z}$ , find the approximation error.
- ▷ Most of the existing literature on approximate DPs focuses on the first interpretation
- ▷ The second interpretation allows us to develop AIS-based RL algorithms
- ▷ Results depend on the choice of metric on probability spaces.
- ▷ The bounds use what are known as **integral probability metrics (IPM)**, which include many commonly used metrics:
  - ▷ Total variation
  - ▷ Wasserstein distance
  - ▷ Maximum mean discrepancy (MMD)

## Examples of AIS

## Example 1: Robustness to model mismatch in MDPs

Real-world  
model  
 $(P, r)$

Simulation  
model  
 $(\hat{P}, \hat{r})$

What is the loss in performance if we choose a policy using the simulation model and use it in the real world?

# Example 1: Robustness to model mismatch in MDPs

Real-world  
model  
 $(P, r)$

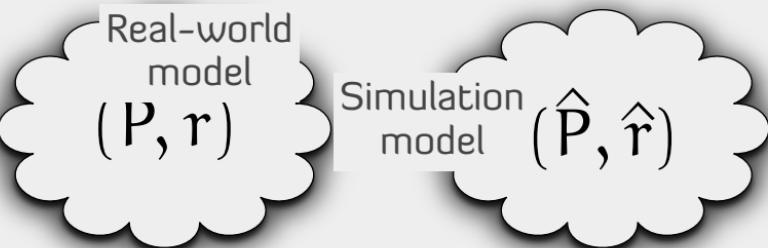
Simulation  
model  
 $(\hat{P}, \hat{r})$

What is the loss in performance if we choose a policy using the simulation model and use it in the real world?

## Model mismatch as an AIS

- (Identity,  $\hat{P}, \hat{r}$ ) is an  $(\varepsilon, \delta)$ -AIS with  $\varepsilon = \sup_{s, a} |r(s, a) - \hat{r}(s, a)|$  and  $\delta_{\mathcal{F}} = \sup_{s, a} d_{\mathcal{F}}(P(\cdot | s, a), \hat{P}(\cdot | s, a))$ .

# Example 1: Robustness to model mismatch in MDPs



■ Müller, "How does the value function of a Markov decision process depend on the transition probabilities?" MOR 1997.

## Model mismatch as an AIS

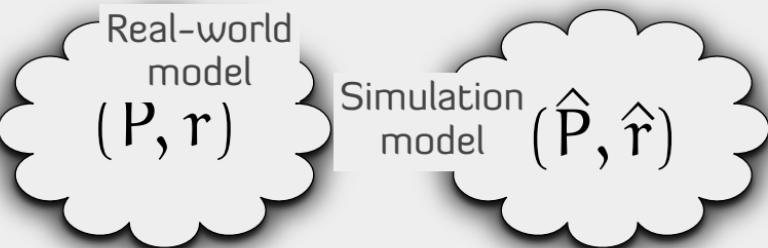
► (Identity,  $\hat{P}, \hat{r}$ ) is an  $(\varepsilon, \delta)$ -AIS with  $\varepsilon = \sup_{s, a} |r(s, a) - \hat{r}(s, a)|$  and  $\delta_{\mathfrak{F}} = \sup_{s, a} d_{\mathfrak{F}}(P(\cdot|s, a), \hat{P}(\cdot|s, a))$ .

$d_{\mathfrak{F}}$  is total variation

$$V(s) - V^\pi(s) \leq \frac{2\varepsilon}{1-\gamma} + \frac{\gamma\delta \text{span}(r)}{(1-\gamma)^2}$$

Recover bounds of Müller (1997).

# Example 1: Robustness to model mismatch in MDPs



- Müller, "How does the value function of a Markov decision process depend on the transition probabilities?" MOR 1997.
- Asadi, Misra, Littman, "Lipschitz continuity in model-based reinforcement learning," ICML 2018.

## Model mismatch as an AIS

► (Identity,  $\hat{P}, \hat{r}$ ) is an  $(\varepsilon, \delta)$ -AIS with  $\varepsilon = \sup_{s, a} |r(s, a) - \hat{r}(s, a)|$  and  $\delta_{\mathfrak{F}} = \sup_{s, a} d_{\mathfrak{F}}(P(\cdot|s, a), \hat{P}(\cdot|s, a))$ .

$d_{\mathfrak{F}}$  is total variation

$$V(s) - V^\pi(s) \leq \frac{2\varepsilon}{1-\gamma} + \frac{\gamma\delta \text{span}(r)}{(1-\gamma)^2}$$

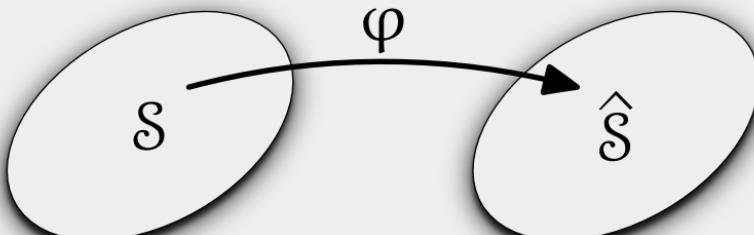
Recover bounds of Müller (1997).

$d_{\mathfrak{F}}$  is Wasserstein distance

$$V(s) - V^\pi(s) \leq \frac{2\varepsilon}{1-\gamma} + \frac{2\gamma\delta L_r}{(1-\gamma)(1-\gamma L_p)}$$

Recover bounds of Asadi, Misra, Littman (2018).

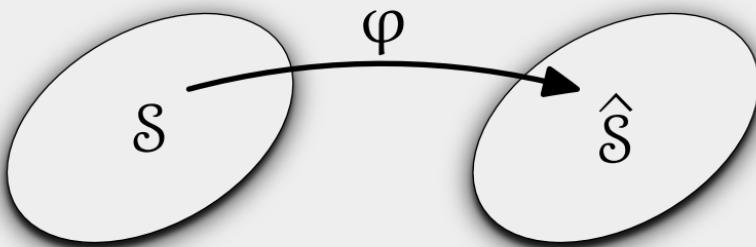
## Example 2: Feature abstraction in MDPs



$(\hat{P}, \hat{r})$  is determined from  $(P, r)$  using  $\varphi$

What is the loss in performance if we choose a policy using the abstract model and use it in the original model?

## Example 2: Feature abstraction in MDPs



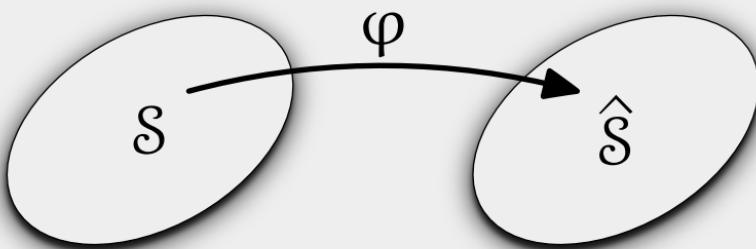
$(\hat{P}, \hat{r})$  is determined from  $(P, r)$  using  $\varphi$

What is the loss in performance if we choose a policy using the abstract model and use it in the original model?

### Feature abstraction as AIS

- $(\varphi, \hat{P}, \hat{r})$  is an  $(\varepsilon, \delta)$ -AIS with  $\varepsilon = \sup_{s, a} |r(s, a) - \hat{r}(\varphi(s), a)|$   
and  $\delta_{\mathfrak{F}} = \sup_{s, a} d_{\mathfrak{F}}(P(\varphi^{-1}(\cdot)|s, a), \hat{P}(\cdot|\varphi(s), a))$ .

## Example 2: Feature abstraction in MDPs



$(\hat{P}, \hat{r})$  is determined from  $(P, r)$  using  $\varphi$

### Feature abstraction as AIS

►  $(\varphi, \hat{P}, \hat{r})$  is an  $(\varepsilon, \delta)$ -AIS with  $\varepsilon = \sup_{s, a} |r(s, a) - \hat{r}(\varphi(s), a)|$

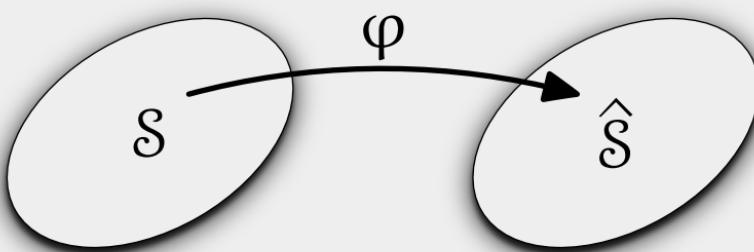
and  $\delta_{\mathfrak{F}} = \sup_{s, a} d_{\mathfrak{F}}(P(\varphi^{-1}(\cdot)|s, a), \hat{P}(\cdot|\varphi(s), a))$ .

$d_{\mathfrak{F}}$  is total variation

$$V(s) - V^\pi(s) \leq \frac{2\varepsilon}{1-\gamma} + \frac{\gamma\delta_{\mathfrak{F}} \text{span}(r)}{(1-\gamma)^2}$$

Improve bounds of Abel et al. (2016)

## Example 2: Feature abstraction in MDPs



$(\hat{P}, \hat{r})$  is determined from  $(P, r)$  using  $\varphi$

### Feature abstraction as AIS

►  $(\varphi, \hat{P}, \hat{r})$  is an  $(\varepsilon, \delta)$ -AIS with  $\varepsilon = \sup_{s, a} |r(s, a) - \hat{r}(\varphi(s), a)|$

and  $\delta_{\mathfrak{F}} = \sup_{s, a} d_{\mathfrak{F}}(P(\varphi^{-1}(\cdot)|s, a), \hat{P}(\cdot|\varphi(s), a))$ .

$d_{\mathfrak{F}}$  is total variation

$$V(s) - V^\pi(s) \leq \frac{2\varepsilon}{1-\gamma} + \frac{\gamma\delta_{\mathfrak{F}} \text{span}(r)}{(1-\gamma)^2}$$

Improve bounds of Abel et al. (2016)

AIS for partially observed systems-(Mahajan)

■ Abel, Hershkowitz, Littman, "Near optimal behavior via approximate state abstraction," ICML 2016.

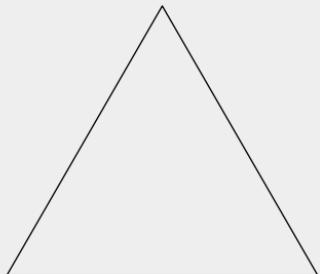
■ Gelada, Kumar, Buckman, Nachum, Bellemare, "DeepMDP: Learning continuous latent space models for representation learning," ICML 2019.

$d_{\mathfrak{F}}$  is Wasserstein distance

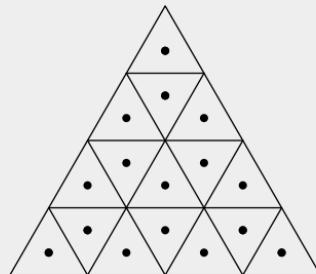
$$V(s) - V^\pi(s) \leq \frac{2\varepsilon}{1-\gamma} + \frac{2\gamma\delta_{\mathfrak{F}} \|\hat{V}\|_{\text{Lip}}}{(1-\gamma)^2}$$

Recover bounds of Gelada et al. (2019).

## Example 3: Belief approximation in POMDPs



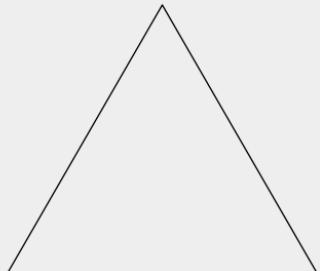
Belief space



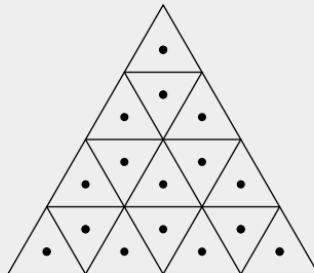
Quantized beliefs

What is the loss in performance if we choose a policy using the approximate beliefs and use it in the original model?

## Example 3: Belief approximation in POMDPs



Belief space



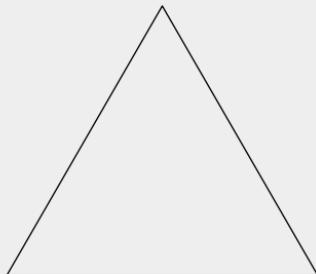
Quantized beliefs

What is the loss in performance if we choose a policy using the approximate beliefs and use it in the original model?

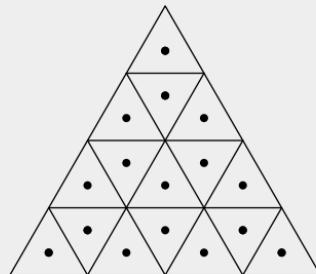
### Belief approximation in POMDPs

- Quantized cells of radius  $\varepsilon$  (in terms of total variation) are  $(\varepsilon\|r\|_\infty, 3\varepsilon)$ -AIS.

## Example 3: Belief approximation in POMDPs



Belief space



Quantized beliefs

■ Francois-Lavet, Rabusseau, Pineau, Ernst, Fonteneau, "On overfitting and asymptotic bias in batch reinforcement learning with partial observability," JAIR 2019.

### Belief approximation in POMDPs

- Quantized cells of radius  $\varepsilon$  (in terms of total variation) are  $(\varepsilon\|r\|_\infty, 3\varepsilon)$ -AIS.

$$V(s) - V^\pi(s) \leq \frac{2\varepsilon\|r\|_\infty}{1-\gamma} + \frac{6\gamma\varepsilon\|r\|_\infty}{(1-\gamma)^2}$$

**Improve** bounds of Francois Lavet et al. (2019) by a factor of  $1/(1-\gamma)$ .

Thus, the notion of AIS unifies many of the approximation results in the literature, both for MDPs and POMDPs.

# Outline



## Background

- ▷ Review of MDPs and RL
- ▷ Review of POMDPs
- ▷ Why is RL for POMDPs difficult?



## Approximate Planning for POMDPs

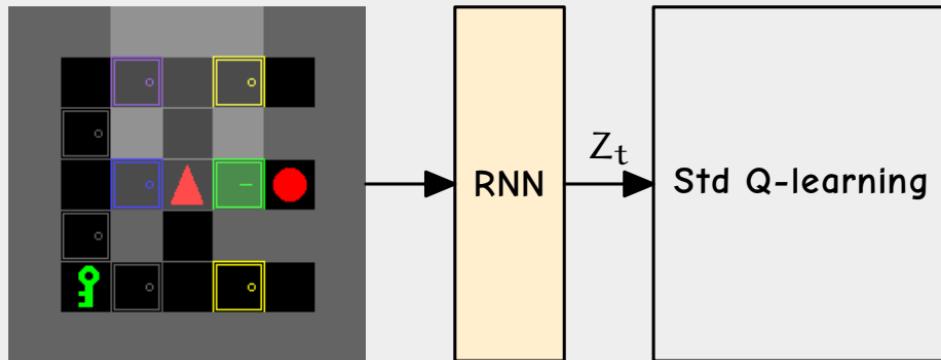
- ▷ Preliminaries on information state
- ▷ Approximate information state
- ▷ Approximation bounds



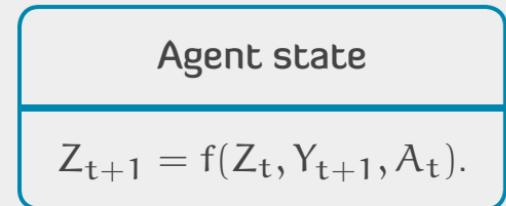
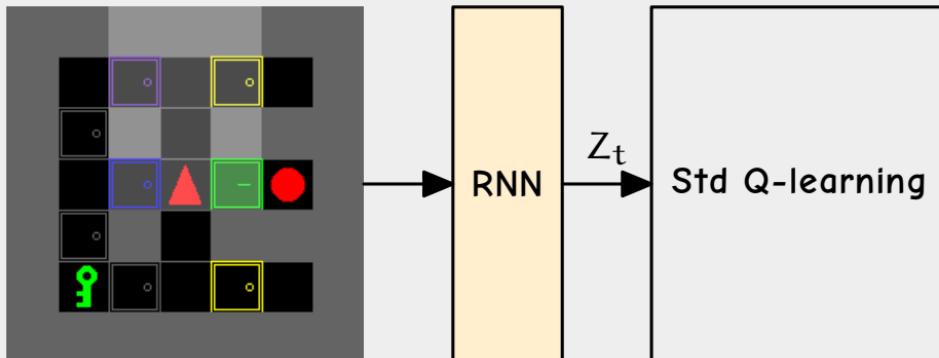
## RL for POMDPs

- ▷ From approximation bounds to RL
- ▷ Numerical experiments

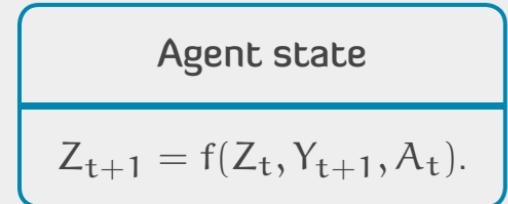
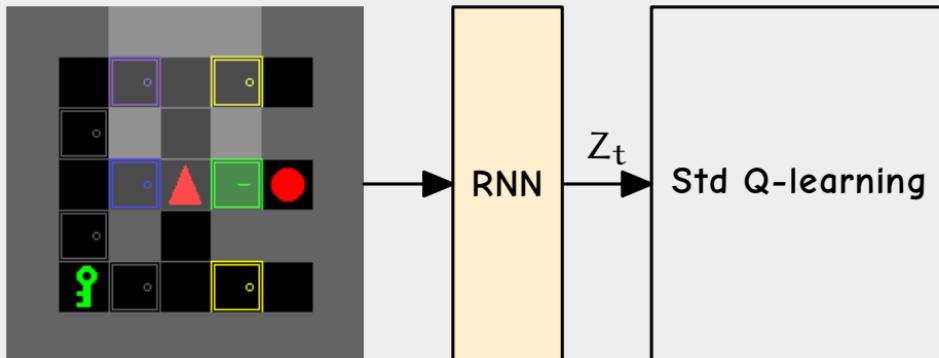
# Q-learning for POMDPs



# Q-learning for POMDPs



# Q-learning for POMDPs



Q-learning

$$\begin{aligned}\hat{Q}_{t+1}(z_t, a_t) &= \hat{Q}_t(z_t, a_t) \\ &+ \alpha_t(z_t, a_t) [R_t + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_t(z_{t+1}, a') - \hat{Q}_t(z_t, a_t)]\end{aligned}$$

What's the difficulty in analyzing  
Q-learning for POMDPs?

# Review: How does Q-learning work in MDPs?

## Dynamic programming decomposition

$$V(s) = \max_{a \in \mathcal{A}} Q(s, a)$$

$$Q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s')$$

# Review: How does Q-learning work in MDPs?

## Dynamic programming decomposition

$$V(s) = \max_{a \in \mathcal{A}} Q(s, a)$$

$$Q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s')$$

## Approximate via stochastic approximation

$$Q(s, a) \leftarrow Q(s, a)$$

$$+ \alpha [r(s, a) + \gamma \max_{a' \in \mathcal{A}} Q(s_+, a') - Q(s, a)]$$

# Review: How does Q-learning work in MDPs?

## Dynamic programming decomposition

$$V(s) = \max_{a \in \mathcal{A}} Q(s, a)$$

$$Q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s')$$

## Approximate via stochastic approximation

$$Q(s, a) \leftarrow Q(s, a)$$

$$+ \alpha [r(s, a) + \gamma \max_{a' \in \mathcal{A}} Q(s_+, a') - Q(s, a)]$$

unbiased sample

# Review: How does Q-learning work in MDPs?

## Dynamic programming decomposition

$$V(s) = \max_{a \in \mathcal{A}} Q(s, a)$$

$$Q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s')$$

## Approximate via stochastic approximation

$$Q(s, a) \leftarrow Q(s, a)$$

$$+ \alpha [r(s, a) + \gamma \max_{a' \in \mathcal{A}} Q(s', a') - Q(s, a)]$$

unbiased sample

## Why does Q-learning converge?

- ▶ Under appropriate technical conditions, SA tracks an ODE (Borkar 1997).
- ▶ **Since the Bellman operator is a contraction**, the ODE has a unique equilibrium point which is globally asymptotically stable (Borkar and Soumyanatha, 1997).

# Why is it difficult to analyze Q-learning for POMDPs?

## Q-learning

$$\begin{aligned}\hat{Q}_{t+1}(z_t, a_t) &= \hat{Q}_t(z_t, a_t) \\ &+ \alpha_t(z_t, a_t) [R_t + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_t(z_{t+1}, a') - \hat{Q}_t(z_t, a_t)]\end{aligned}$$

# Why is it difficult to analyze Q-learning for POMDPs?

## Q-learning

$$\begin{aligned}\hat{Q}_{t+1}(z_t, a_t) &= \hat{Q}_t(z_t, a_t) \\ &+ \alpha_t(z_t, a_t) [R_t + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_t(z_{t+1}, a') - \hat{Q}_t(z_t, a_t)]\end{aligned}$$

There is no “ubbiased sample from a DP”

We cannot just write a DP:

$$Q(z, a) = \mathbb{E}[R_t | Z_t = z, A_t = a] + \gamma \sum_{z' \in \mathcal{S}} P(z' | s, a) V(z')$$

# Why is it difficult to analyze Q-learning for POMDPs?

## Q-learning

$$\begin{aligned}\hat{Q}_{t+1}(z_t, a_t) &= \hat{Q}_t(z_t, a_t) \\ &+ \alpha_t(z_t, a_t) [R_t + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_t(z_{t+1}, a') - \hat{Q}_t(z_t, a_t)]\end{aligned}$$

There is no “ubbiased sample from a DP”

We cannot just write a DP:

$$Q(z, a) = \mathbb{E}[R_t | Z_t = z, A_t = a] + \gamma \sum_{z' \in \mathcal{S}} P(z' | s, a) V(z')$$

↑  
Not time-homogeneous

# Why is it difficult to analyze Q-learning for POMDPs?

## Q-learning

$$\begin{aligned}\hat{Q}_{t+1}(z_t, a_t) = & \hat{Q}_t(z_t, a_t) \\ & + \alpha_t(z_t, a_t) [R_t + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_t(z_{t+1}, a') - \hat{Q}_t(z_t, a_t)]\end{aligned}$$

There is no “ubbiased sample from a DP”

We cannot just write a DP:

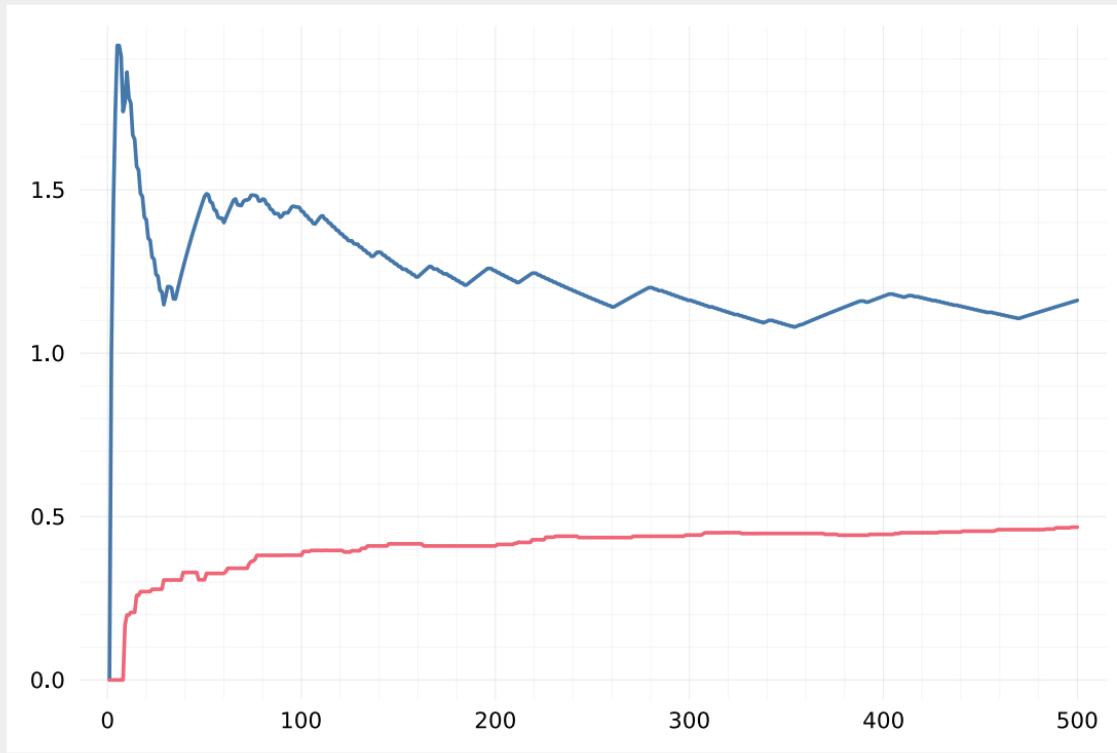
$$Q(z, a) = \mathbb{E}[R_t | Z_t = z, A_t = a] + \gamma \sum_{z' \in \mathcal{S}} P(z' | s, a) V(z')$$

Not time-homogeneous

Not-controlled Markov

# Implication: Exploration policy matters!

# Implication: Exploration policy matters!



# Characterizing convergence: Assumptions

## Assumptions

(A1) Restrict to tabular setting

(A2) The exploration policy  $\pi_{\text{exp}}: \mathcal{Z} \rightarrow \Delta(\mathcal{A})$  is such that the Markov chain  $\{(S_t, Y_t, Z_t, A_t)\}_{t \geq 1}$  has a unique stationary distribution  $\xi$ . Plus,  $\xi(s, y, z, a) > 0$ .

$$(A3) \quad \alpha_t(z, a) = \frac{\mathbb{1}\{Z_t = z, A_t = a\}}{\sum_{\tau=1}^t \mathbb{1}\{Z_\tau = z, A_\tau = a\}}$$

# Characterizing convergence: Guarantees

Define:

- $r_\xi(z, a) = \sum_{s \in \mathcal{S}} \xi(s | z, a) r(s, a).$
- $P_\xi(z' | z, a) = \sum_{s, s' \in \mathcal{S}} \sum_{y' \in \mathcal{Y}} \xi(s | z, a) P(s' | s, a) P(y' | s') \mathbb{1}\{z' = f(z, y', a)\}.$

# Characterizing convergence: Guarantees

Define:

- $r_\xi(z, a) = \sum_{s \in \mathcal{S}} \xi(s | z, a) r(s, a).$
- $P_\xi(z' | z, a) = \sum_{s, s' \in \mathcal{S}} \sum_{y' \in \mathcal{Y}} \xi(s | z, a) P(s' | s, a) P(y' | s') \mathbb{1}\{z' = f(z, y', a)\}.$

Let  $(V_\xi, Q_\xi)$  be the solution of the following DP:

$$V_\xi(z) = \max_{a \in \mathcal{A}} Q_\xi(z, a)$$

$$Q_\xi(z, a) = r_\xi(z, a) + \gamma \sum_{z' \in \mathcal{Z}} P_\xi(z' | z, a) V_\xi(z')$$

# Characterizing convergence: Guarantees

Define:

- $r_\xi(z, a) = \sum_{s \in \mathcal{S}} \xi(s | z, a) r(s, a).$
- $P_\xi(z' | z, a) = \sum_{s, s' \in \mathcal{S}} \sum_{y' \in \mathcal{Y}} \xi(s | z, a) P(s' | s, a) P(y' | s') \mathbb{1}\{z' = f(z, y', a)\}.$

Let  $(V_\xi, Q_\xi)$  be the solution of the following DP:

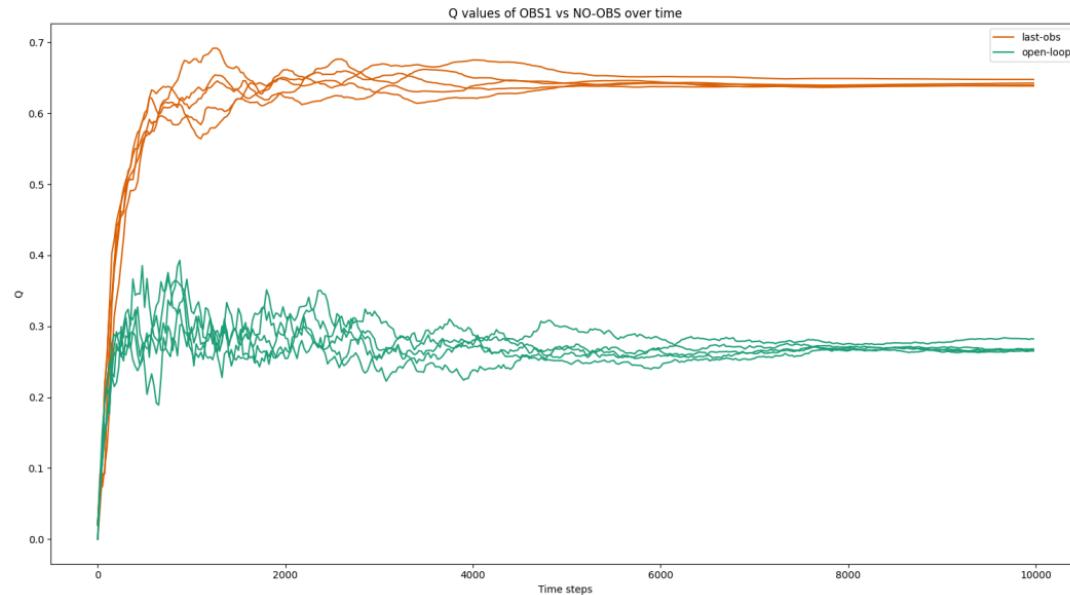
$$V_\xi(z) = \max_{a \in \mathcal{A}} Q_\xi(z, a)$$

$$Q_\xi(z, a) = r_\xi(z, a) + \gamma \sum_{z' \in \mathcal{Z}} P_\xi(z' | z, a) V_\xi(z')$$

Under (A1)–(A3),  $\hat{Q}_t \rightarrow Q_\xi$  a.s.

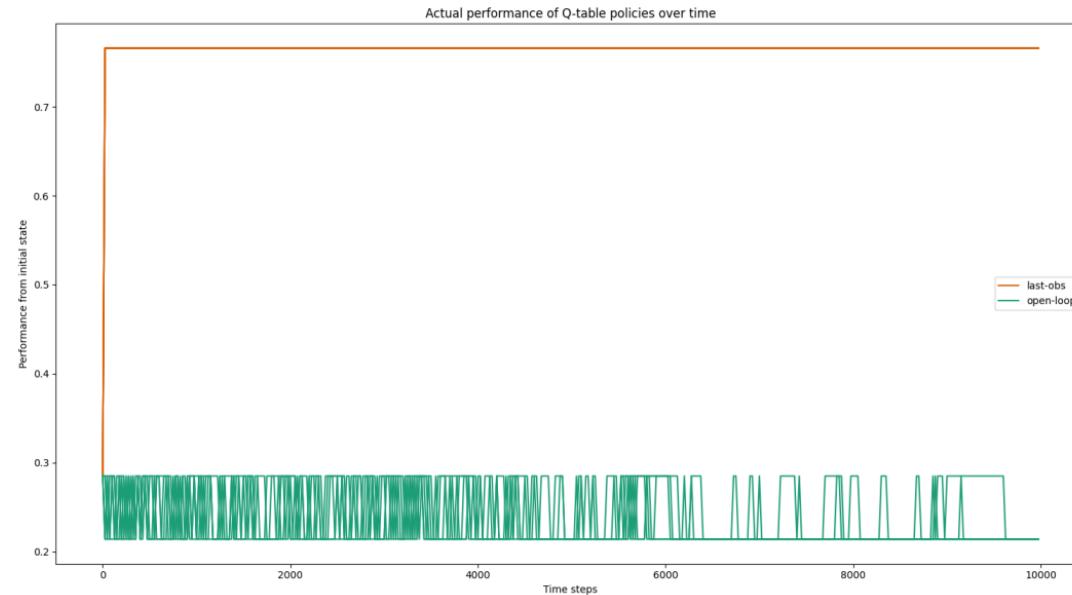
# Characterizing convergence is not enough!

## Convergence for two different update functions



# Characterizing convergence is not enough!

## Performance of greedy policy for two different update functions



# Quality of converged solution

# Quality of converged solution

$(r_\xi, P_\xi)$  is an AIS!

# Quality of converged solution

$(r_\xi, P_\xi)$  is an AIS!

## Approximation guarantee

For any history  $h_t$

$$|V_t^*(h_t) - V_t^{\pi_\xi^* \circ \sigma_t}(h_t)| \leq (1 - \gamma)^{-1} [\varepsilon + \gamma \delta_{\mathfrak{F}} \rho_{\mathfrak{F}}(V_\xi)]$$

where

- $\varepsilon = \sup_{t \geq 1} \max_{h_t, a_t} |\mathbb{E}[r(S_t, a_t) | h_t, a_t] - r_\xi(\sigma_t(h_t), a_t)|$
- $\delta_{\mathfrak{F}} = \sup_{t \geq 1} \max_{h_t, a_t} d_{\mathfrak{F}} \left( \mathbb{P}(Z_{t+1} = \cdot | h_t, a_t), P_\xi(\cdot | \sigma_t(h_t), a_t) \right)$

Can we use this to improve Q-learning?

# Adding AIS-losses to Q-learning

## Main idea

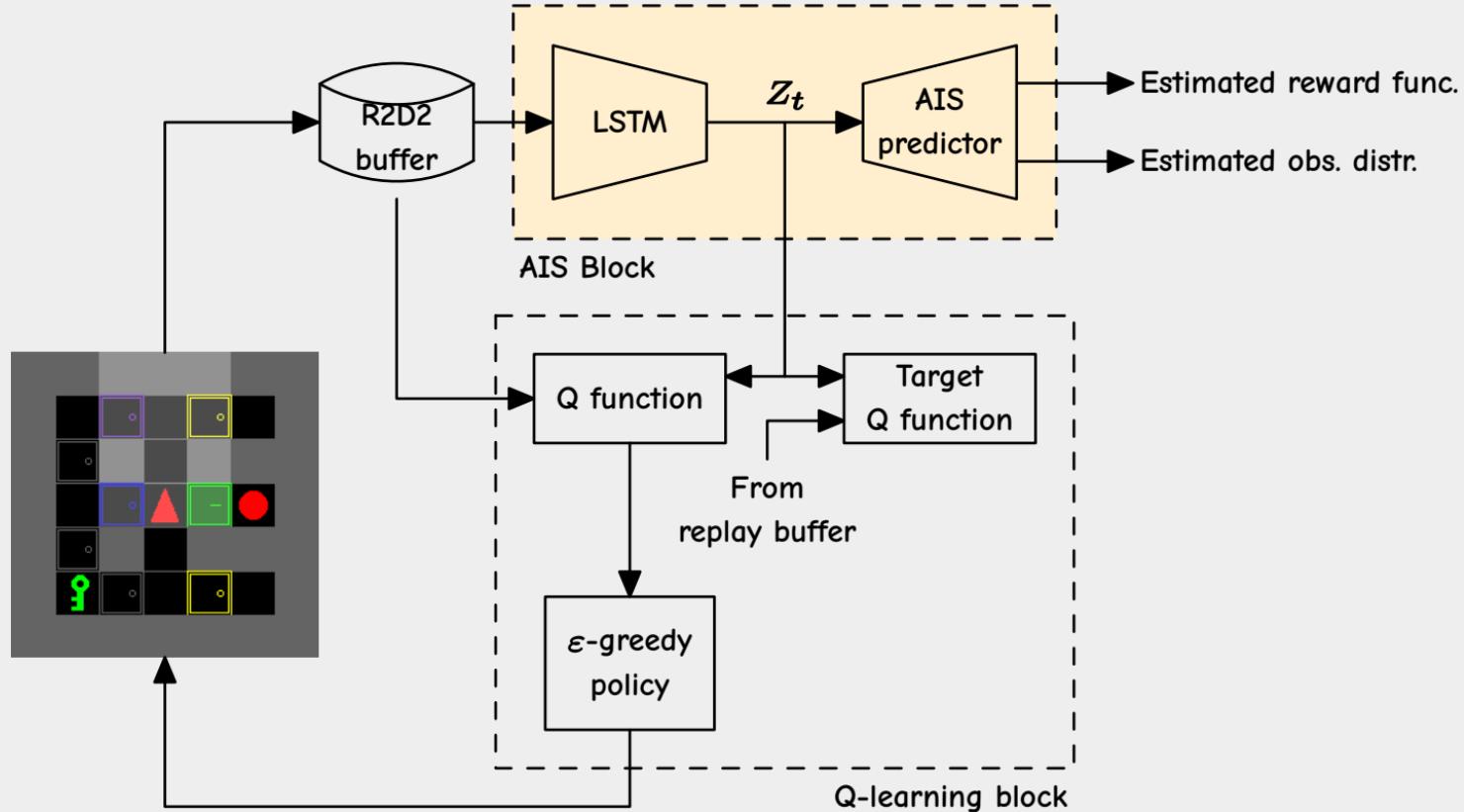
- ▶ AIS is defined in terms of two losses  $\varepsilon$  and  $\delta$ .
- ▶ Minimizing  $\varepsilon$  and  $\delta$  will minimize the AIS approximation loss.

# Adding AIS-losses to Q-learning

## Main idea

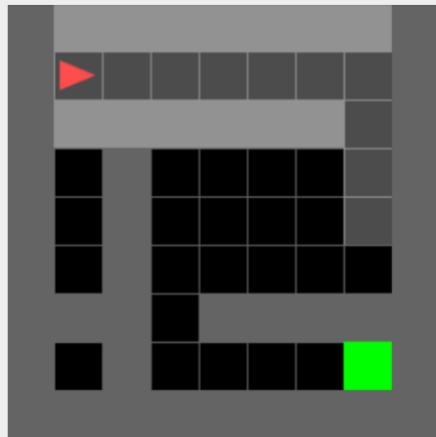
- ▶ AIS is defined in terms of two losses  $\varepsilon$  and  $\delta$ .
- ▶ Minimizing  $\varepsilon$  and  $\delta$  will minimize the AIS approximation loss.
- ▶ Use  $\lambda\varepsilon^2 + (1 - \lambda)\delta^2$  as surrogate loss for the AIS generator
- ▶ ... and combine it with standard Q-learning algorithm.

# Recurrent Q-learning with AIS losses (RQL-AIS)

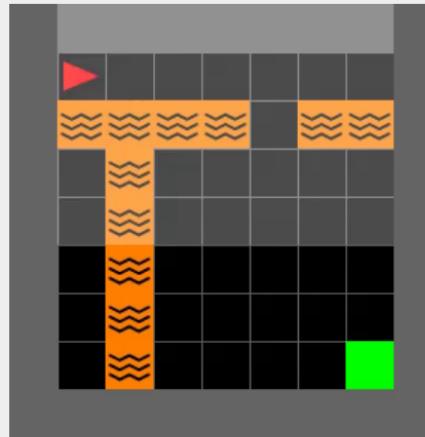


# Numerical Experiments

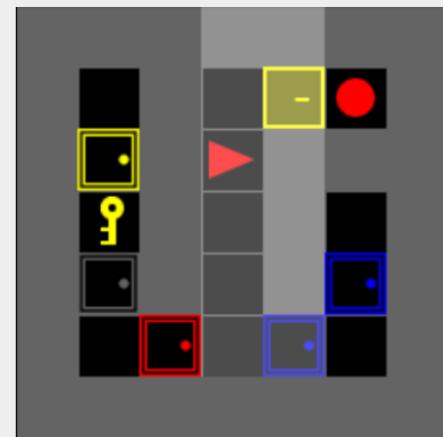
# MiniGrid Environments



Simple Crossing



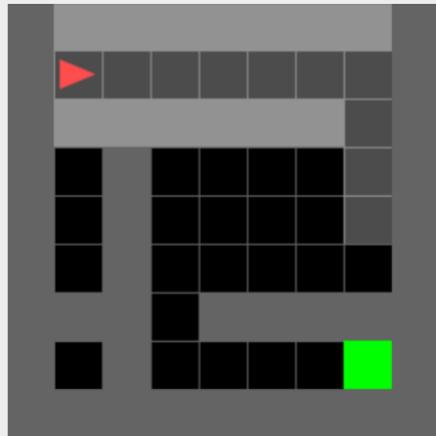
Lava Crossing



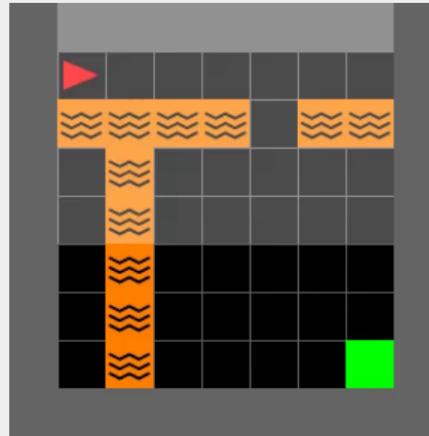
Key Corridor

- Features**
- ▶ Partially observable 2D grids. Agent has a view of a  $7 \times 7$  field in front of it. Observations are obstructed by walls.
  - ▶ Multiple entities (agents, walls, lava, boxes, doors, and keys)
  - ▶ Multiple actions (Move Forward, Turn Left, Turn Right, Open Door/Box, ...)

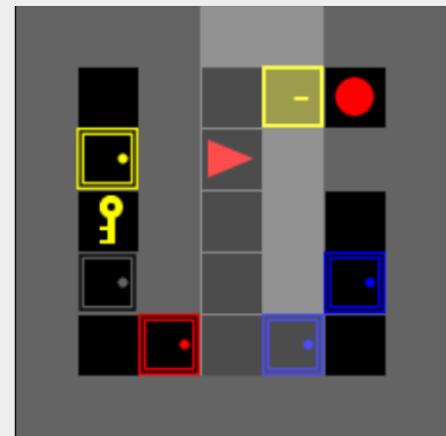
# MiniGrid Environments



Simple Crossing



Lava Crossing



Key Corridor

## Algorithms

R2D2

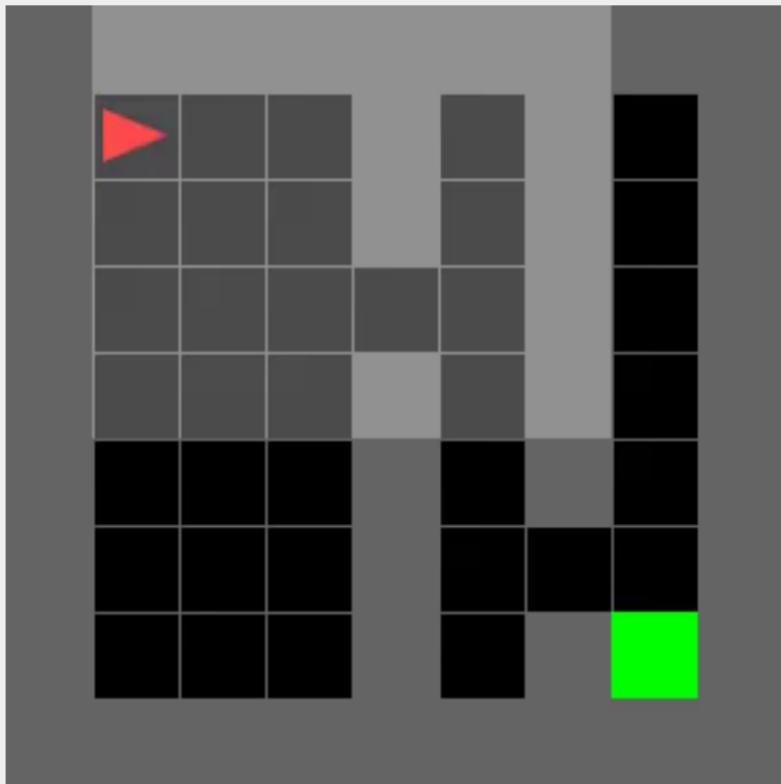
A competitive  
baseline for POMDPs

R2D2 + AIS

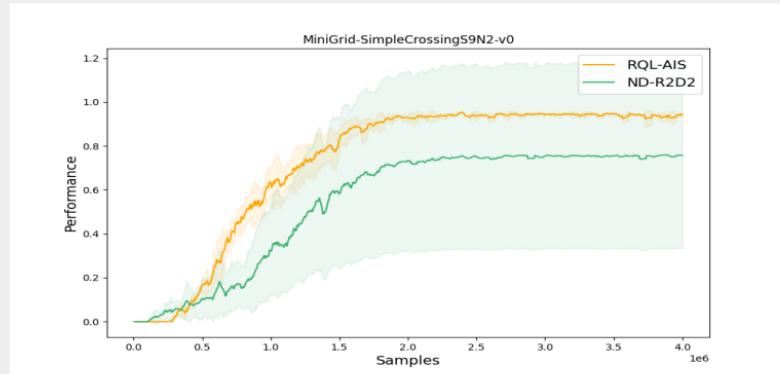
Add AIS losses to R2D2

AIS for partially observed systems-(Mahajan)

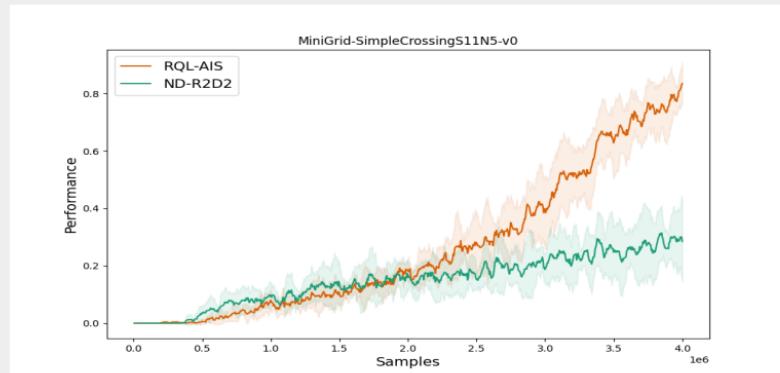
# Simple Crossing



AIS for partially observed systems-(Mahajan)

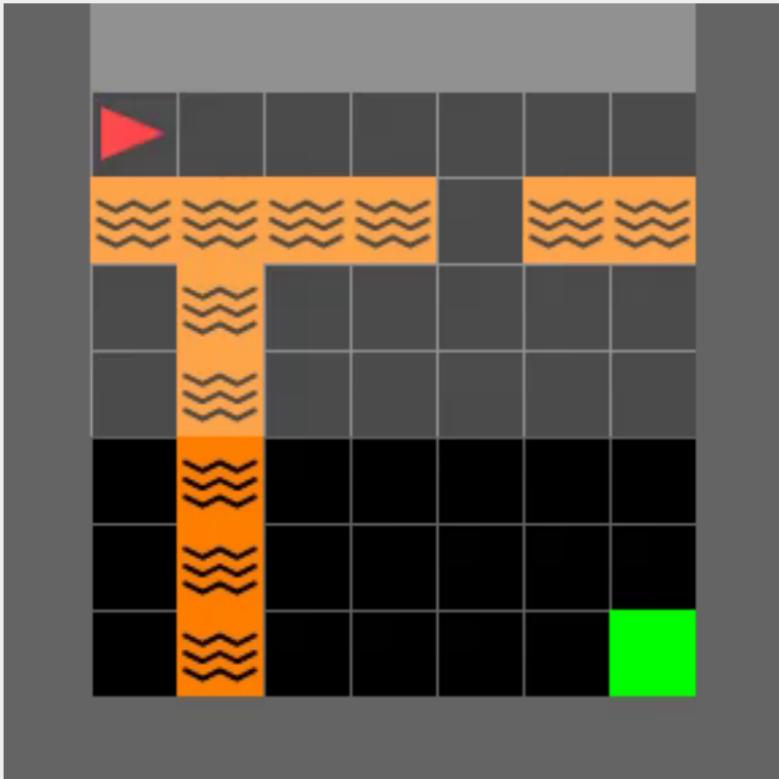


Simple Crossing S9N2

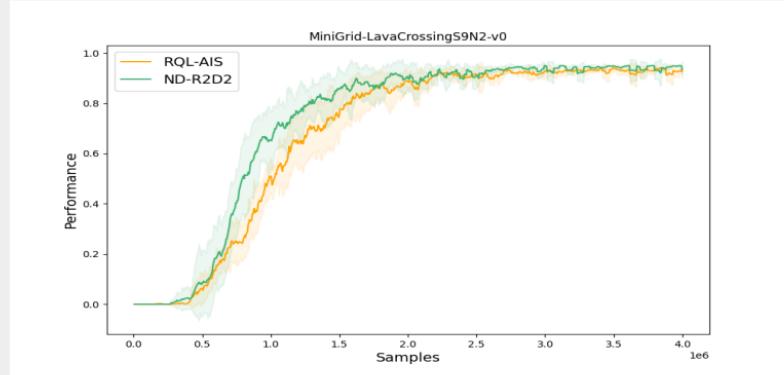


Simple Crossing S11N5

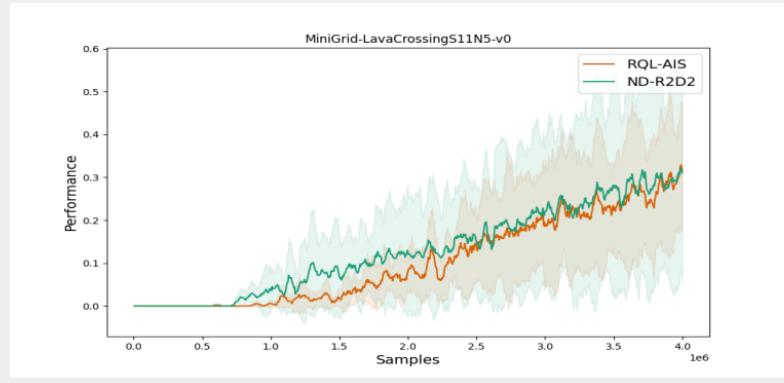
# Lava Crossing



AIS for partially observed systems-(Mahajan)

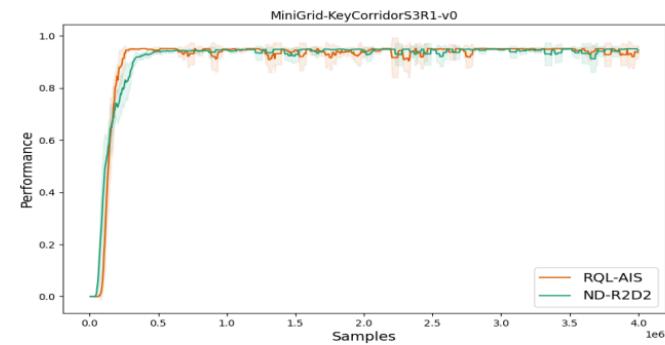
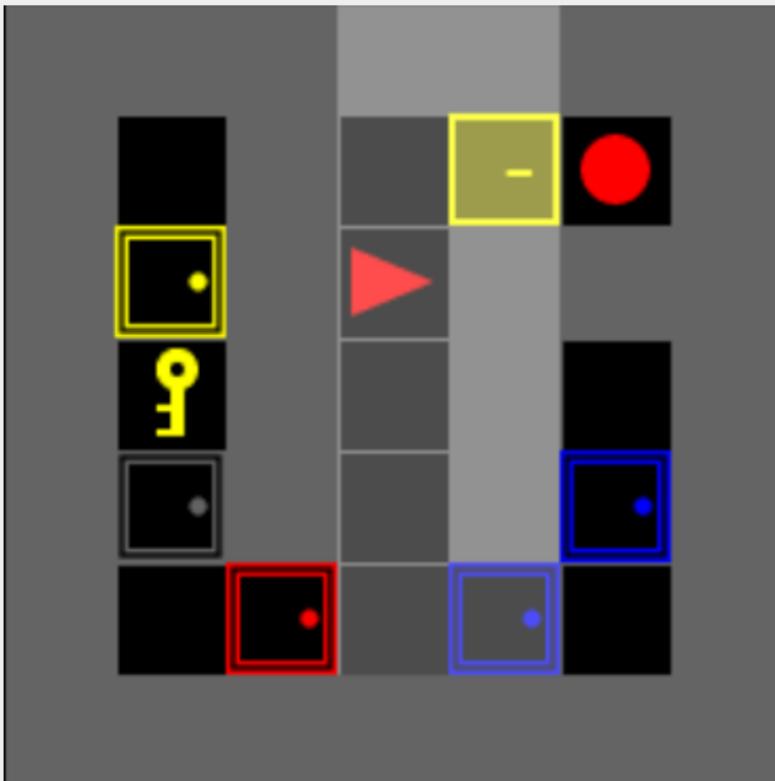


Lava Crossing S9N2

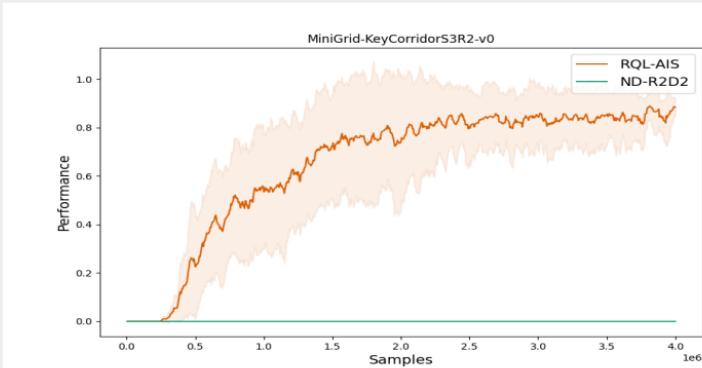


Lava Crossing S11N5

# Key Corridor



Key Corridor S3R1



Key Corridor S3R2

AIS for partially observed systems-(Mahajan)

## Similar improvements for AIS-enhanced Q-learning

Environment	RQL-AIS	ND-R2D2
SimpleCrossingS9N2	$0.944 \pm 0.007$	$0.757 \pm 0.423$
LavaCrossingS9N2	$0.926 \pm 0.014$	$0.934 \pm 0.034$
RedBlueDoors-8x8	$0.977 \pm 0.009$	$0.962 \pm 0.018$
MultiRoom-N2-S4	$0.790 \pm 0.049$	$0.839 \pm 0.010$
DoorKey-8x8	$0.942 \pm 0.038$	$0.371 \pm 0.508$
ObstructedMaze-1D	$0.916 \pm 0.020$	$0.000 \pm 0.000$
KeyCorridorS3R2	$0.885 \pm 0.038$	$0.000 \pm 0.000$
UnlockPickup	$0.517 \pm 0.474$	$0.000 \pm 0.000$

# Summary

A conceptually clean framework for approximate DP  
and online RL in partially observed systems

# Summary

A conceptually clean framework for approximate DP  
and online RL in partially observed systems

## Approximation results generalize to

- ▷ observation compression
- ▷ action quantization
- ▷ lifelong learning
- ▷ multi-agent teams
- ▷ Markov games

A conceptually clean framework for approximate DP and online RL in partially observed systems

## Approximation results generalize to

- ▷ observation compression
- ▷ action quantization
- ▷ lifelong learning
- ▷ multi-agent teams
- ▷ Markov games

## Ongoing work

- ▷ Other RL settings such as model based RL, offline RL, inverse RL.
- ▷ A building block for multi-agent RL.
- ▷ ...

- ▷ [email](mailto:aditya.mahajan@mcgill.ca): aditya.mahajan@mcgill.ca
- ▷ [web](http://cim.mcgill.ca/~adityam): http://cim.mcgill.ca/~adityam

# Thank you

- ▷ [paper 1](#): JMLR, Feb 2022  
[code](https://github.com/info-structures/ais): <https://github.com/info-structures/ais>
- ▷ [paper 2](#): JMLR, Feb 2022  
[code](https://github.com/info-structures/RQL-ais): <https://github.com/info-structures/RQL-ais>