

# Agent-state based policies in POMDPs: Beyond belief-state MDPs

Amit Sinha and Aditya Mahajan

**Abstract**—The traditional approach to POMDPs is to convert them into fully observed MDPs by considering a belief state as an information state. However, a belief-state based approach requires perfect knowledge of the system dynamics and is therefore not applicable in the learning setting where the system model is unknown. Various approaches to circumvent this limitation have been proposed in the literature. We present a unified treatment of some of these approaches by viewing them as models where the agent maintains a local recursively updateable “agent state” and chooses actions based on the agent state. We highlight the different classes of agent-state based policies and the various approaches that have been proposed in the literature to find good policies within each class. These include the designer’s approach to find optimal non-stationary agent-state based policies, policy search approaches to find a locally optimal stationary agent-state based policies, and the approximate information state to find approximately optimal stationary agent-state based policies. We then present how ideas from the approximate information state approach have been used to improve Q-learning and actor-critic algorithms for learning in POMDPs.

## I. INTRODUCTION

Partially observable Markov decision processes (POMDPs) are a widely used model for the optimal control of dynamical systems with partial state observation. They have been extensively studied across various research communities including systems and control, operations research, and artificial intelligence.

A key conceptual challenge for POMDPs is that the data available at the agent—the history of observations and actions—is increasing with time. The standard approach is to compress this increasing data into a finite dimensional statistic known as the belief state, which is the posterior density of the unobserved state conditioned on the history of observations and actions and, therefore, may be viewed as a generalization of non-linear filtering of controlled processes. The belief state is sufficient for evaluating the per-step reward, can be updated recursively, and is strategy independent. Therefore, we can write a dynamic programming decomposition using the belief state as an information state [3], [14], [49], [76]. Furthermore, the value function of the corresponding belief-state MDP has certain qualitative properties (it is piecewise linear and convex), which can be leveraged for efficient computational algorithms. The earliest such algorithm was the “one-pass” algorithm by Smallwood

The authors are with the Department of Electrical and Computer Engineering, McGill University, Montreal, QC, Canada. They are also affiliated with CIM, GERAD, Mila, and ILLS. Emails: amit.sinha@mail.mcgill.ca, aditya.mahajan@mcgill.ca

This research was supported in part by NSERC Alliance International Catalyst Grant ALLRP 580801-2022 and in part by a grant from Google’s Institutional Research Program in collaboration with Mila.

and Sondik [76]. Various efficient refinements of this algorithm have been presented in the literature, including linear-support algorithm [20], witness algorithm [14], incremental pruning [13], [93], point-based methods [66], [77], [79], and others. See [38], [46] for a unified overview of the numerical methods.

There are two limitations of the belief-state based approach. First, implementing a belief-state based policy is computationally challenging. Keeping track of the belief state requires non-linear filtering, which can be approximated via particle filtering, but is still computationally heavy and is therefore difficult to implement in embedded hardware in robotics and other applications. Second, the belief state is model dependent. Therefore, it cannot be used in model-free reinforcement learning (RL).

An alternative approach, which is more amenable to the learning setting, is to consider what we call the *agent state*. In particular, we relax the assumption that the agent can use the entire history of observations and actions (or a model-dependent compression of it like the belief state) to make a decision. Rather, we assume that the agent maintains a local state that does not depend on the model and makes a decision as a function of the agent state. Such an agent-state formulation has been proposed multiple times in the literature [29], [39], [54], [57], [78], [80]. Perhaps the simplest example of the agent state is an agent keeping track of a finite window of past observations and actions, which was first considered in [67], [88] and is commonly referred to as *frame stacking* in the RL literature [58]. It is argued in [29], [54] that in many of the popular implementations of RL algorithms, the agent state can be further decomposed into three parts: an algorithmic state which is data to be used by subsequent computations, a situational state, which is a summary of the agent’s current situation, and an epistemic state, which is the summary of the agent’s current knowledge of the environment. We do not pursue such a distinction here but, roughly speaking, our notion of agent state is similar to the situational state in [29], [54].

### Notation

We use uppercase letters to denote random variables (e.g.  $S, A$ , etc.), lowercase letters to denote their realizations (e.g.  $s, a$ , etc.) and calligraphic letters to denote sets (e.g.  $\mathcal{S}, \mathcal{A}; \text{ggh}$ , etc.). Subscripts (e.g.  $S_t, A_t$ , etc.) denote variables at time  $t$ .  $\Delta(\mathcal{S})$  denotes the space of probability measures on a set  $\mathcal{S}$ ;  $\mathbb{P}(\cdot)$  and  $\mathbb{E}[\cdot]$  denote the probability of an event and the expectation of a random variable, respectively; and  $\mathbb{1}$  denotes the indicator function.

## Organization

The rest of the paper is organized as follows. In [Sec. II](#) we present the mathematical model of POMDPs, introduce the agent-state framework, and present sufficient conditions for the agent state to be an information state. In [Sec. III](#) we discuss the different classes of agent-state policies and present the three approaches for finding optimal policies among different policy classes. In [Sec. IV](#) we discuss various RL approaches taken for learning agent-state based policies. Finally, we present concluding remarks and discussions in [Sec. V](#).

## II. THE POMDP MODEL AND AGENT-STATE BASED POLICIES

### A. System model

Consider a stochastic dynamical system with state  $S_t \in \mathcal{S}$ , input  $A_t \in \mathcal{A}$ , and output  $Y_t \in \mathcal{Y}$ . To simplify the discussion, we will assume that all sets are finite valued and ignore integrability and measurability issues, existence of suprema, etc. See the companion paper [26] for a nuanced discussion of these issues. The system operates in discrete time with the dynamics given as follows: for any time  $t \in \mathbb{N}$ , we have

$$\begin{aligned} \mathbb{P}(S_{t+1}, Y_{t+1} \mid S_{1:t}, Y_{1:t}, A_{1:t}) &= \mathbb{P}(S_{t+1}, Y_{t+1} \mid S_t, A_t) \\ &=: P(S_{t+1}, Y_{t+1} \mid S_t, A_t) \end{aligned}$$

where  $P$  is a probability transition matrix.<sup>1</sup> In addition, at each time the system yields a reward  $R_t = r(S_t, A_t)$ . We will assume that  $R_t \in [0, R_{\max}]$ .

There is an agent (also called a controller or a decision maker, depending on the research community) which observes the outputs of the system and chooses control actions as inputs to the system. In principle, this agent can be as sophisticated as we want and can, therefore, use the entire history of past observations and actions to choose its action, i.e.,<sup>2</sup>

$$A_t = \vec{\pi}_t(Y_{1:t}, A_{1:t-1}), \quad t \in \mathbb{N}$$

where  $\vec{\pi}_t$  is called the control law or decision rule at time  $t$ . We are using the vector accent to highlight the fact that the policy is history dependent.

We assume that the system runs for an infinite horizon and use  $\vec{\pi} = (\vec{\pi}_1, \vec{\pi}_2, \dots)$  to denote the control policy (or simply the policy).<sup>3</sup> Let  $\vec{\Pi}_{\text{ND}}$  denote the set of all history dependent

<sup>1</sup>For the simplicity of notation, we are using a slightly informal notation. Terms such as  $\mathbb{P}(S_{t+1}, Y_{t+1} \mid S_t, A_t)$  should either be viewed as the numerical value  $\mathbb{P}(S_{t+1} = s_{t+1}, Y_{t+1} = y_{t+1} \mid S_t = s_t, A_t = a_t)$  for specific realizations  $(s_{t+1}, y_{t+1}, s_t, a_t)$  of  $(S_{t+1}, Y_{t+1}, S_t, A_t)$  or as probability mass function  $\mathbb{P}(S_{t+1} = \cdot, Y_{t+1} = \cdot \mid S_t = s_t, A_t = a_t)$  or as a distribution-valued random variable  $\mathbb{P}(S_{t+1} = \cdot, Y_{t+1} = \cdot \mid S_t, A_t)$ . Typically, all of these interpretations are consistent. If a specific interpretation is needed, we will use a more elaborate notation as appropriate.

<sup>2</sup>Again, we are using a slightly informal notation. We can interpret the above either as  $a_t = \vec{\pi}_t(y_{1:t}, a_{1:t-1})$  for special realizations  $(y_{1:t}, a_{1:t})$  of  $(Y_{1:t}, A_{1:t})$  or as an equality between random variables  $A_t(\omega)$  and  $\vec{\pi}_t(Y_{1:t}(\omega), A_{1:t-1}(\omega))$ .

<sup>3</sup>We are making a deliberate choice of using bold  $\vec{\pi}$  to denote the policy to distinguish between control laws and control policies.

(indicated by the vector accent) non-stationary (i.e., time-varying) and deterministic policies.

We assume that the initial state  $S_1$  is distributed according to probability mass function  $\xi_1$ . Then, the performance of any policy  $\vec{\pi} \in \vec{\Pi}_{\text{ND}}$  is given by

$$J^{\vec{\pi}} := \mathbb{E}^{\vec{\pi}} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid S_1 \sim \xi_1 \right]$$

where  $\gamma \in (0, 1)$  is the discount factor. Let  $\vec{J}_{\text{ND}}^*$  (again, the vector accent highlights that we are optimizing over all history dependent policies) denote the optimal performance in  $\vec{\Pi}_{\text{ND}}$ , i.e.,

$$\vec{J}_{\text{ND}}^* := \sup_{\vec{\pi} \in \vec{\Pi}_{\text{ND}}} J^{\vec{\pi}}.$$

### B. Some remarks on the model

- 1) The system described above is referred to as a *partially observed* Markov decision process (POMDP) to highlight the fact that the agent sees *partial* observations of the state of the environment.
- 2) Since the per-step reward is uniformly bounded,  $J^{\vec{\pi}}$  is well defined. However, it is not immediately obvious that there exists an optimal policy  $\vec{\pi}^*$  such that  $\vec{J}_{\text{ND}}^* = J^{\vec{\pi}^*}$ .
- 3) There are several technical questions that need to be resolved carefully when the variables are continuous valued. We refer the reader to the companion paper [26] for a detailed discussion.
- 4) In the literature, it is often assumed that

$$\begin{aligned} \mathbb{P}(S_{t+1}, Y_{t+1} \mid S_t, A_t) \\ = \mathbb{P}(S_{t+1} \mid S_t, A_t) \mathbb{P}(Y_{t+1} \mid S_{t+1}, A_t) \end{aligned}$$

or sometimes even

$$\begin{aligned} \mathbb{P}(S_{t+1}, Y_{t+1} \mid S_t, A_t) \\ = \mathbb{P}(S_{t+1} \mid S_t, A_t) \mathbb{P}(Y_{t+1} \mid S_{t+1}). \end{aligned}$$

Such an assumption is not needed for the discussion presented in this paper.

- 5) In the discussion above, we have restricted attention to deterministic policies. In principle, we could have also considered non-stationary history dependent stochastic policies  $\vec{\pi} = (\vec{\pi}_1, \vec{\pi}_2, \dots)$ , where  $\vec{\pi}_t: \mathcal{H}_t \rightarrow \Delta(\mathcal{A})$ . Let  $\vec{\Pi}_{\text{NS}}$  denote the set of all non-stationary history dependent stochastic policies. Define

$$\vec{J}_{\text{NS}}^* := \sup_{\vec{\pi} \in \vec{\Pi}_{\text{NS}}} J^{\vec{\pi}}.$$

By definition  $\vec{\Pi}_{\text{ND}} \subseteq \vec{\Pi}_{\text{NS}}$ ; hence,  $\vec{J}_{\text{ND}}^* \leq \vec{J}_{\text{NS}}^*$ . However, since the agent has perfect recall (i.e., remembers everything that it has seen and done in the past) it can be shown that **there is no loss of optimality in**

restricting attention to deterministic strategies,<sup>4</sup> i.e.,

$$\boxed{\vec{J}_{\text{ND}}^* = \vec{J}_{\text{NS}}^*} \quad (1)$$

### C. Agent-state based policies

We now describe the agent-state based approach, where it is assumed that instead of using the entire history of observations and actions to make a decision the agent maintains a local state (which we will refer to as the agent state and denote by  $Z_t \in \mathcal{Z}$ ) and makes a decision as a function of the agent state. The agent starts with an initial state  $Z_1 = \phi_0(Y_1)$  and recursively updates it as follows:<sup>5</sup>

$$Z_{t+1} = \phi(Z_t, Y_{t+1}, A_t), \quad \forall t \in \mathbb{N}. \quad (2)$$

We call  $\phi_0$  the *state-initialization function* and  $\phi$  the *state-update function*. The agent chooses an action either using a deterministic control law  $\pi_t: \mathcal{Z} \rightarrow \mathcal{A}$  as

$$A_t = \pi_t(Z_t)$$

or using a stochastic control law  $\pi_t: \mathcal{Z} \rightarrow \Delta(\mathcal{A})$  as

$$A_t \sim \pi_t(Z_t).$$

We call such a policy as an **agent-state based policy**.

A simple example of agent-state based policy is when the agent uses a finite window of past observations and actions, i.e.,  $Z_t = (Y_{t-n:t}, A_{t-n:t-1})$ . Such a model is sometimes called frame-stacking in the RL literature [58], [88]. Another example is when the controller is a finite state automaton [36], [67], [69]. It is not necessary for the agent state to be finite. In fact, the belief-state representation is a special case of agent-state model where the agent state belongs to a  $|\mathcal{S}|$ -dimensional simplex. However, for the convenience of notation, we would assume that  $\mathcal{Z}$  is finite. The discussion can be generalized to continuous valued  $\mathcal{Z}$  under appropriate technical assumptions.

There are three fundamental questions when working with an agent state:

- (Q1) When is there no loss of optimality in restricting attention to agent-state based policies?
- (Q2) For a given agent-state update function  $\phi$ , how do we find the optimal agent-state based policy?
- (Q3) For a given agent-state space  $\mathcal{Z}$ , what is the optimal agent-state update function and agent-state based policy?

<sup>4</sup>To explain the high-level idea of the result, we borrow the terminology of pure, mixed, and behavioral strategies used in game theory. A pure strategy is what we call deterministic policy; a behavioral strategy is what we call stochastic policy. A mixed strategy is a probability distribution over pure strategies where the agent picks a pure strategy at the beginning of the game according to specified probability distribution and then follows the pure strategy throughout the game. With this terminology, the result follows from two facts. First, since we are in a single agent unconstrained expectation maximization setting, mixed and pure strategies have the same performance. Second, since we are in a perfect recall setting, mixed and behavioral strategies have the same performance due to Kuhn's theorem (see e.g., [4]).

<sup>5</sup>In principle, the update can be stochastic and be given by  $Z_{t+1} = \phi(Z_t, Y_{t+1}, A_t, W_t)$  where  $\{W_t\}_{t \geq 1}$  is a sequence of i.i.d. (independent and identically distributed) randomizing variables that are independent of all other primitive random variables.

The short answer to (Q1) is simple: **if an agent state is an information state, then there is no loss of optimality in restricting attention to an agent-state policy**. Of course, this answer only makes sense if we define an information state, which we do in the next section. The answers to (Q2) and (Q3) are more difficult and depend on what we mean by "optimal". We will discuss the different conceptual approaches that have been used in the literature in Sec. III.

### D. Information state

We start with some notation. Let  $H_t = (Y_{1:t}, A_{1:t-1})$  denote the history of observations and actions of the agent up to time  $t$  and let  $\mathcal{H}_t$  denote the space of realization of all such histories. We can recursively unroll the agent-state update function (2) and define a sequence of functions  $\vec{\sigma} := (\vec{\sigma}_1, \vec{\sigma}_2, \dots)$ , where  $\vec{\sigma}_t: \mathcal{H}_t \rightarrow \mathcal{Z}$ , such that  $Z_t = \vec{\sigma}_t(H_t)$ . In particular,

$$\vec{\sigma}_1(H_1) = \phi_0(Y_1), \quad \vec{\sigma}_2(H_2) = \phi(\vec{\sigma}_1(H_1), Y_2, A_1), \quad \text{etc.} \quad (3)$$

We call  $\vec{\sigma}$  to be the *history compression function* corresponding to the agent-state initialization and update functions  $(\phi_0, \phi)$ .

The agent state is an information state if it satisfies the following two properties:

- (P1) Sufficient for performance evaluation.** There exists a function  $r_{\text{IS}}: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$  such that for any  $t$  and  $H_t$  and  $A_t$ , we have

$$\mathbb{E}[R_t | H_t, A_t] = r_{\text{IS}}(\vec{\sigma}_t(H_t), A_t).$$

- (P2) Sufficient for predicting itself.** There exists a controlled transition probability matrix  $P_{\text{IS}}: \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Z})$  such that for any  $t$  and  $H_t$  and  $A_t$ , we have

$$\mathbb{P}(Z_{t+1} | H_t, A_t) = P_{\text{IS}}(Z_{t+1} | \vec{\sigma}_t(H_t), A_t).$$

Recall that the agent state updates in a state-like manner (2). Thus,

$$\begin{aligned} &\mathbb{P}(Z_{t+1} | H_t, A_t) \\ &= \sum_{y_{t+1} \in \mathcal{Y}} \mathbb{P}(y_{t+1} | H_t, A_t) \mathbb{1}_{\{Z_{t+1} = \phi(\vec{\sigma}_t(H_t), y_{t+1}, A_t)\}} \end{aligned}$$

It was shown in [80] that the above relationship implies that a sufficient condition for (P2) is the following.

- (P2b) Sufficient for predicting output.** There exists a controlled transition probability matrix  $P_{\text{IS}}: \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Y})$  such that for any  $t$  and  $H_t$  and  $A_t$ , we have

$$\mathbb{P}(Y_{t+1} | H_t, A_t) = P_{\text{IS}}(Y_{t+1} | \vec{\sigma}_t(H_t), A_t).$$

The belief-state satisfies properties (P1) and (P2) and is, therefore, an information state. For specific models such as the linear quadratic Gaussian (LQG) control or certain classes of machine repair models, simpler and model-specific information states exist. We refer the reader to [80] for a detailed discussion of the history of information states and various other examples.

TABLE I: Abbreviations used for different dynamic programs

Abbreviation	Meaning
IS	Information-state based DP
DES	DP using designer's approach
PROD	DP based on product processes
AIS	Approximate information state based DP
ASQL	Agent-state based Q-learning

A key implication of an information state is that it always leads to a dynamic programming decomposition.<sup>6</sup> In particular, we have the following [80, Theorem 5 and 25].

**Theorem 1: Information-state based DP**

Suppose the agent state  $\{Z_t\}_{t \geq 1}$  is an information state, i.e., satisfies properties (P1) and (P2) with some  $(r_{\text{IS}}, P_{\text{IS}})$ . Define the following dynamic program<sup>a</sup>:

$$Q_{\text{IS}}^*(z, a) = r_{\text{IS}}(z, a) + \gamma \sum_{z' \in \mathcal{Z}} P_{\text{IS}}(z'|z, a) V_{\text{IS}}^*(z'), \quad (4a)$$

$$V_{\text{IS}}^*(z) = \max_{a \in \mathcal{A}} Q_{\text{IS}}^*(z, a). \quad (4b)$$

Let  $\pi_{\text{IS}}^*(z)$  denote any *arg max* of the right hand side of (4b). Then the policy  $\vec{\pi}_{\text{IS}}^* = (\vec{\pi}_{\text{IS},1}^*, \vec{\pi}_{\text{IS},2}^*, \dots)$  given by

$$\vec{\pi}_{\text{IS},t}^*(h_t) = \pi_{\text{IS}}^*(\vec{\sigma}_t(h_t))$$

is optimal, i.e.,  $J^{\vec{\pi}_{\text{IS}}^*} = \vec{J}_{\text{ND}}^*$ .

<sup>a</sup>For simplicity, we state the result for finite  $\mathcal{Z}$ . Similar decomposition holds for continuous  $\mathcal{Z}$ .

The classical belief-state based dynamic programming decomposition may be considered as an instance of [Theorem 1](#). For specific models which have a simpler information state such as LQG control, [Theorem 1](#) provides a simpler dynamic program than the standard belief-state based dynamic program.

In spite of its generality, [Theorem 1](#) is only applicable when the agent state satisfies properties (P1) and (P2), which is not always the case. A simple counterexample is when the agent uses a finite window of past observations with  $Z_t = (Y_{t-n:t}, A_{t-n:t-1})$ . In the companion paper [26], approximation results for this model are presented using ideas from filter stability. In the next section, we present other approaches that have been used in the literature for identifying optimal policies for models where properties (P1) and (P2) are not satisfied.

### III. OPTIMAL AGENT-STATE BASED POLICIES

In this section, we consider the different policy classes (deterministic/stochastic and stationary/non-stationary) of agent-state based policies. We present examples to show that non-stationary policies may outperform stationary ones and that

<sup>6</sup>We present multiple dynamic programming decompositions and differentiate between them via subscripts. A summary of these subscripts is shown in [Table I](#).

TABLE II: List of different policy classes

Policy class	Meaning
$\vec{\Pi}_{\text{ND}}$	history based non-stationary deterministic policies
$\vec{\Pi}_{\text{NS}}$	history based non-stationary stochastic policies
$\Pi_{\text{ND}}$	agent-state based non-stationary deterministic policies
$\Pi_{\text{NS}}$	agent-state based non-stationary stochastic policies
$\Pi_{\text{SD}}$	agent-state based stationary deterministic policies
$\Pi_{\text{SS}}$	agent-state based stationary stochastic policies

stochastic policies may outperform deterministic ones. We then present different schemes that have been used in the literature to identify optimal or sub-optimal policies within the different policy classes. In particular, the designer's approach for finding an optimal non-stationary agent-state based policy; policy evaluation and policy gradients for stationary agent-state based policies; and approximate information state approach for finding good stationary and deterministic agent-state based policies.

#### A. Policy classes

We start with a notation to denote (one-step) decision rules. Let

- $\mathcal{D}_{\text{D}} = [\mathcal{Z} \rightarrow \mathcal{A}]$  denote the family of all deterministic decision rules, and
- $\mathcal{D}_{\text{S}} = [\mathcal{Z} \rightarrow \Delta(\mathcal{A})]$  denote the family of all stochastic decision rules.

We define the following classes of agent-state based policies. Notation wise, we will use the absence of the vector accent to differentiate agent-state based policies from history-based policies. A summary of the different policy class used in this paper is shown in [Table II](#).

- $\Pi_{\text{ND}}$  denotes the class of non-stationary agent-state based deterministic policies  $\pi = (\pi_1, \pi_2, \dots)$  where  $\pi_t \in \mathcal{D}_{\text{D}}$ . Let

$$J_{\text{ND}}^* = \sup_{\pi \in \Pi_{\text{ND}}} J^\pi$$

denote the optimal performance within class  $\Pi_{\text{ND}}$ . A policy  $\pi_{\text{ND}}^* \in \Pi_{\text{ND}}$  is called optimal in policy class  $\Pi_{\text{ND}}$  if  $J^{\pi_{\text{ND}}^*} = J_{\text{ND}}^*$ .

- $\Pi_{\text{NS}}$  denotes the class of non-stationary agent-state based stochastic policies  $\pi = (\pi_1, \pi_2, \dots)$  where  $\pi_t \in \mathcal{D}_{\text{S}}$ . Let

$$J_{\text{NS}}^* = \sup_{\pi \in \Pi_{\text{NS}}} J^\pi$$

denote the optimal performance within class  $\Pi_{\text{NS}}$ . A policy  $\pi_{\text{NS}}^* \in \Pi_{\text{NS}}$  is called optimal in policy class  $\Pi_{\text{NS}}$  if  $J^{\pi_{\text{NS}}^*} = J_{\text{NS}}^*$ .

- $\Pi_{\text{SD}}$  denotes the class of stationary agent-state based deterministic policies  $\pi = (\pi, \pi, \dots)$  where  $\pi \in \mathcal{D}_{\text{D}}$ . Let

$$J_{\text{SD}}^* = \sup_{\pi \in \Pi_{\text{SD}}} J^\pi$$

denote the optimal performance within class  $\Pi_{\text{SD}}$ . A policy  $\pi_{\text{SD}}^* \in \Pi_{\text{SD}}$  is called optimal in policy class  $\Pi_{\text{SD}}$  if  $J^{\pi_{\text{SD}}^*} = J_{\text{SD}}^*$ .

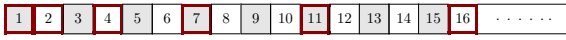


Fig. 1: The cells indicate the state of the environment. Cells with the same background color have the same observation. The cells with a thick red boundary correspond to elements of the set  $\mathcal{D}_0 := \{n(n+1)/2 + 1 : n \in \mathbb{N}\}$ , where the action 0 gives a reward of +1 and moves the state to the right, while the action 1 gives a reward of -1 and resets the state to 1. The cells with a thin black boundary correspond to elements of the set  $\mathcal{D}_1 = \mathbb{N} \setminus \mathcal{D}_0$ , where the action 1 gives the reward of +1 and moves the state to the right while the action 0 gives a reward of -1 and resets the state to 1.

- $\Pi_{SS}$  denotes the class of stationary agent-state based stochastic policies  $\pi = (\pi, \pi, \dots)$  where  $\pi \in \mathcal{D}_s$ . Let

$$J_{SS}^* = \sup_{\pi \in \Pi_{SS}} J^\pi$$

denote the optimal performance within class  $\Pi_{SS}$ . A policy  $\pi_{SS}^* \in \Pi_{SS}$  is called optimal in policy class  $\Pi_{SS}$  if  $J^{\pi_{SS}^*} = J_{SS}^*$ .

Note that  $\Pi_{ND}$  and  $\Pi_{NS}$  agent-state based policies are different from  $\bar{\Pi}_{ND}$  and  $\bar{\Pi}_{NS}$ , which are history based policies. Therefore,  $J_{ND}^*$  and  $J_{NS}^*$  are different from  $\bar{J}_{ND}^*$  and  $\bar{J}_{NS}^*$ .

When the agent state is an information state and satisfies (P1) and (P2), [Theorem 1](#) implies that all the performance functions  $J_{ND}^*$ ,  $J_{NS}^*$ ,  $J_{SD}^*$ ,  $J_{SS}^*$  are equal and are also equal to the performance of history-based policies  $\bar{J}_{ND}^*$  and  $\bar{J}_{NS}^*$ . However, when the agent state is not an information state, then the agent state does not satisfy the controlled Markov property, i.e.,

$$\mathbb{P}(Z_{t+1} | Z_{1:t}, A_{1:t}) \neq \mathbb{P}(Z_{t+1} | Z_t, A_t)$$

or the property of being sufficient for reward evaluation, i.e.,

$$\mathbb{E}[r(S_t, A_t) | Z_{1:t}, A_{1:t}] \neq \mathbb{E}[r(S_t, A_t) | Z_t, A_t].$$

Moreover, the agent does not have perfect recall (i.e., at time  $t$ , the agent does not have access to all the information that was available to the information that was available to it in the past).

The absence of the information state properties imply that (F1) non-stationary policies may outperform stationary policies and (F2) stationary stochastic policies may outperform stationary deterministic policies. These facts were first reported in [73]. We illustrate them via examples borrowed from [74].

(F1) *Non-stationary agent-state based policies may outperform stationary agent-state based policies:* Consider the POMDP described in [Fig. 1](#), where the system starts in state 1. Since the dynamics are deterministic, the agent can infer the current state from the history of past actions and can take the action to increment the current state and receive a per-step reward of +1. Thus,  $\bar{J}_{ND}^* = 1/(1-\gamma)$ . Furthermore, since the system is deterministic, the optimal policy can be implemented via an open-loop policy given by  $\{a_t^*\}_{t \geq 1}$ ,

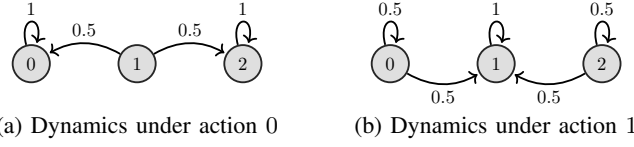


Fig. 2: A POMDP with  $\mathcal{S} = \{0, 1, 2\}$ ,  $\mathcal{A} = \{0, 1\}$  and  $\mathcal{Y} = \{0\}$ . The rewards functions are  $r(\cdot, 0) = [-1, 0, 2]$  and  $r(\cdot, 1) = -0.5$ .

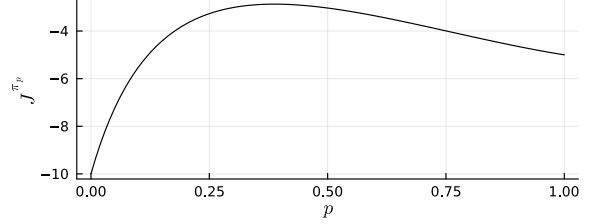


Fig. 3: Performance of stationary stochastic policies in the model of [Fig. 2](#).

where  $a_t^* = \mathbb{1}_{\{t \in \mathcal{D}_1\}}$ , which can be implemented irrespective of the information structure. Thus,  $J_{ND}^* = \bar{J}_{ND}^*$ .

When  $Z_t = Y_t$ ,  $\Pi_{SD}$  consists of four policies. A brute force evaluation shows that  $J_{SD}^* = (1 + \gamma - \gamma^2)/(1 - \gamma^3) < J_{ND}^*$ . See [75] for details. Thus, non-stationary agent-state policies may outperform stationary agent-state based policies.<sup>7</sup>

(F2) *Stationary stochastic agent-state policies may outperform stationary deterministic agent-state based policies:* Consider the POMDP shown in [Fig. 2](#), where the system starts in state 0. The agent gets no observations, i.e.,  $Y_t \equiv 0$ , and we consider agent-state policies with  $Z_t = Y_t \equiv 0$ . A policy  $\pi \in \Pi_{SS}$  is parameterized by a single parameter  $p \in [0, 1]$ , which indicates the probability of choosing action 1. We denote such a policy by  $\pi_p$ . Note that if  $p \in \{0, 1\}$ , then  $\pi_p \in \Pi_{SD}$ . Let  $(P_a, r_a)$  denote the probability transition matrix and reward function when  $a \in \mathcal{A}$  is chosen and let  $(P_p, r_p) = (1-p)(P_0, r_0) + p(P_1, r_1)$ . Then, the performance of policy  $\pi_p$  is given by  $J^{\pi_p} = [(1 - \gamma P_p)^{-1} r_p]_0$ . The performance for all  $p \in [0, 1]$  for  $\gamma = 0.9$  is shown in [Fig. 3](#), which shows that the best performance is achieved by the stochastic policy  $\pi_p$  with  $p \approx 0.39$ .

In summary, **the optimal agent-state based policy depends on the choice of the policy class**. The relationship between the different performance measures can be summarized as follows:

$$\begin{array}{ccccc}
 J_{SD}^* & \longrightarrow & J_{ND}^* & \longrightarrow & \bar{J}_{ND}^* \\
 \downarrow & & \updownarrow & & \updownarrow \\
 J_{SS}^* & \longrightarrow & J_{NS}^* & \longrightarrow & \bar{J}_{NS}^*
 \end{array} \tag{5}$$

where the arrows mean “less than or equal to”. The solid

<sup>7</sup>The conclusion continue to hold if we compare with the performance on stationary stochastic policies. Although it is not possible to evaluate  $J_{SS}^*$  in closed form, a brute force Monte Carlo evaluation shows that  $J_{SS}^* \approx J_{SD}^*$ . Therefore,  $J_{SS}^* < J_{ND}^*$ .

arrows follow from set inclusion relationships (that is, one policy class is a subset of the other); the dashed arrow relationships need to be established explicitly. The reason for the dashed arrow between  $\bar{J}_{\text{NS}}^*$  and  $\bar{J}_{\text{ND}}^*$  has already been presented in the discussion around (1). The reason for the dashed arrow between  $J_{\text{NS}}^*$  and  $J_{\text{ND}}^*$  will be presented later in (7).

Features (F1)–(F2) may appear to be surprising because they are not present in MDPs or POMDPs when using a belief state. In fact, they are absent when the agent state is an information state. However, when the agent state is not an information state, the system has what is called a *non-classical information structure*<sup>8</sup> and, therefore, the problem of finding the optimal agent-state based policy is a *decentralized* stochastic control problem (even though there is only one decision maker). This fact is well known in the decentralized control literature [55], [69] but perhaps not as well recognized in the POMDP literature.

The fact that the optimization problem at hand is a decentralized control problem means that we cannot directly use the standard results from Markov decision theory to find the optimal policy in a given policy class. In the remainder of this section, we summarize the approaches that have been taken in the literature to find the optimal or a good policy within a policy class.

### B. The designer's approach to find the optimal non-stationary policy

Since the problem of finding the optimal non-stationary agent-state policy is a decentralized stochastic control policy, it is possible to use solution techniques from decentralized stochastic control to find the optimal policy in class  $\Pi_{\text{ND}}$  and  $\Pi_{\text{NS}}$ . One such approach is the designer's approach, which was proposed in [55] and can be viewed as a refinement of an earlier approach known as the standard form proposed in [89]. The designer's approach was presented in [55] for a two-agent decentralized control system where each agent followed an agent-state based policy. We present a simplified version of this approach, restricted to the single agent POMDP setting.

For any policy  $\pi = (\pi_1, \pi_2, \dots) \in \Pi_{\text{NS}}$ , define  $\xi_t^\pi \in \Delta(\mathcal{S} \times \mathcal{Z})$  by

$$\xi_t^\pi(s, z) = \mathbb{P}^\pi(S_t = s, Z_t = z), \quad \forall s \in \mathcal{S}, z \in \mathcal{Z}.$$

Since  $\Pi_{\text{SD}}, \Pi_{\text{SS}}, \Pi_{\text{ND}} \subseteq \Pi_{\text{NS}}$ , the same definition holds for any policy in  $\Pi_{\text{SD}}, \Pi_{\text{SS}}, \Pi_{\text{ND}}$  as well. The key idea of the designer's approach is to show that the process  $\{\xi_t^\pi\}_{t \geq 1}$  is a controlled Markov process controlled by  $\{\pi_t\}_{t \geq 1}$ .

We first start with some definitions.

- Define a function  $\phi_{\text{DES}}: \Delta(\mathcal{S} \times \mathcal{Z}) \times \mathcal{D}_s \rightarrow \Delta(\mathcal{S} \times \mathcal{Z})$  as follows: for any  $\xi \in \Delta(\mathcal{S} \times \mathcal{Z})$ ,  $\pi \in \mathcal{D}_s$  and for all

<sup>8</sup>See [56] for a general discussion of non-classical information structures.

$$s', z' \in \mathcal{S} \times \mathcal{Z},$$

$$\begin{aligned} & [\phi_{\text{DES}}(\xi, \pi)](s', z') \\ &= \sum_{\substack{(s, z, y', a) \in \\ \mathcal{S} \times \mathcal{Z} \times \mathcal{Y} \times \mathcal{A}}} \xi(s, z) \pi(a | z) P(y', s' | s, a) \mathbb{1}_{\{z' = \phi(z, y', a)\}} \end{aligned}$$

- Define the function  $r_{\text{DES}}: \Delta(\mathcal{S} \times \mathcal{Z}) \times \mathcal{D}_s \rightarrow \mathbb{R}$  as follows: for any  $\xi \in \Delta(\mathcal{S} \times \mathcal{Z})$  and  $\pi \in \mathcal{D}_s$ :

$$r_{\text{DES}}(\xi, \pi) := \sum_{\substack{(s, z, a) \in \\ \mathcal{S} \times \mathcal{Z} \times \mathcal{A}}} \xi(s, z) \pi(a | z) r(s, a).$$

Note that both  $\phi_{\text{DES}}$  and  $r_{\text{DES}}$  are bilinear functions.<sup>9</sup>

Then a simple application of Bayes rule and the definition of expectation imply the following [55].

**Proposition 1** *The process  $\{\xi_t^\pi\}_{t \geq 1}$  is a controlled Markov process controlled by  $\{\pi_t\}_{t \geq 1}$ , i.e.,*

- 1)  $\xi_{t+1}^\pi = \phi_{\text{DES}}(\pi_t, \xi_t^\pi)$ ;
- 2)  $\mathbb{E}^\pi[r(S, A)] = r_{\text{DES}}(\pi_t, \xi_t^\pi)$ .

The designer's approach is based on the following interpretation. Consider the system designer who wants to pick a policy  $\pi \in \Pi_{\text{NS}}$  (or  $\Pi_{\text{ND}}$ ). From the point of view of such a designer, the system is a completely unobserved input-output system, where the system designer chooses the control input  $\pi_t$ , receives a per-step reward  $r(S_t, A_t)$ , but does not observe anything. Proposition 1 implies that  $\xi_t^\pi$  is an information state for this system. Therefore, the system designer can use dynamic programming to identify its "control actions". Consequently, we have the following result [55].

### Theorem 2: DP using the designer's approach

Consider the following dynamic program: for all  $\xi \in \Delta(\mathcal{S} \times \mathcal{Z})$ ,

$$V_{\text{DES}}(\xi) = \max_{\pi \in \mathcal{D}_s} \{r_{\text{DES}}(\xi, \pi) + \gamma V_{\text{DES}}(\phi_{\text{DES}}(\xi, \pi))\}. \quad (6)$$

Let  $\psi_{\text{DES}}(\xi)$  denote any arg max of the right hand side of (6). Let  $\xi_1(s_1, a_1) := \mathbb{P}(S_1 = s_1, Z_1 = z_1)$  denote the initial distribution of the system and the agent state. Recursively define  $\{\xi_t^*\}_{t \geq 1}$  and  $\{\pi_t^*\}_{t \geq 1}$  as follows:  $\xi_1^* = \xi_1$  and for  $t \geq 1$ :

$$\pi_t^* = \psi_{\text{DES}}(\xi_t^*) \quad \text{and} \quad \xi_{t+1}^* = \phi_{\text{DES}}(\xi_t^*, \pi_t^*).$$

Then the policy  $\pi^* = (\pi_1^*, \pi_2^*, \dots) \in \Pi_{\text{NS}}$  is the optimal policy in  $\Pi_{\text{NS}}$ .

*Some remarks on the designer's approach*

- 1) The designer's approach presented here is adapted from the presentation in [55] for decentralized control

<sup>9</sup>Since  $\Delta(\mathcal{S} \times \mathcal{Z})$  and  $\mathcal{D}_s$  are bounded sets, linearity in the each component should be interpreted as follows: for any  $\xi^1, \xi^2 \in \Delta(\mathcal{S} \times \mathcal{Z})$ ,  $\pi \in \mathcal{D}_s$ , and  $\lambda \in [0, 1]$ , we have

$$\phi_{\text{DES}}(\lambda \xi^1 + (1 - \lambda) \xi^2, \pi) = \lambda \phi_{\text{DES}}(\xi^1, \pi) + (1 - \lambda) \phi_{\text{DES}}(\xi^2, \pi)$$

with a similar interpretation for  $r_{\text{DES}}$  and the  $\pi$  component of both functions.

problems with two agents, where the map  $\psi_{\text{DES}}: \Delta(\mathcal{S} \times \mathcal{Z}) \rightarrow \mathcal{D}_s$  was called a *meta-policy*. The general idea of using the probability distribution over all pertinent variables as an information state goes back to [89]. The method was re-discovered in [28], where it was called occupation MDP. A similar idea was presented in [5] for MDPs and POMDPs with memoryless policies.

- 2) In the argument above, we maximized over  $\pi \in \mathcal{D}_s$  in (6). If we instead maximize over  $\pi \in \mathcal{D}_D$ , we will obtain the optimal policy  $\pi^* \in \Pi_{\text{ND}}$ .
- 3) The probability distribution  $\xi_t^\pi$  may be viewed as the ‘‘belief’’ of the designer on the state sufficient for input-output mapping (i.e.,  $(S_t, Z_t)$ ) given the history of past observations and actions (i.e.,  $\pi_{1:t-1}$ , since the designer does not observe anything). Therefore, as stated in [55] and [28], we can follow the standard argument for POMDPs [76] to argue that the value function  $V_{\text{DES}}$  is convex (it is piecewise linear and convex for finite horizon models).
- 4) Since  $V_{\text{DES}}$  is convex and  $\phi_{\text{DES}}$  and  $r_{\text{DES}}$  are bilinear, it means that

$$Q_{\text{DES}}(\xi, \pi) := r_{\text{DES}}(\xi, \pi) + \gamma V_{\text{DES}}(\phi_{\text{DES}}(\xi, \pi))$$

is convex in  $\pi$ . Therefore, for any fixed  $\xi$ , the maximum value of the convex function  $Q_{\text{DES}}(\xi, \pi)$  over  $\pi \in \mathcal{D}_s$  (which is a convex polyhedron) is achieved at a vertex of  $\mathcal{D}_s$ . Thus, the arg max in (6) is a deterministic decision rule  $\pi \in \mathcal{D}_D$ . Therefore,

$$\boxed{J_{\text{ND}}^* = J_{\text{NS}}^*} \quad (7)$$

Consequently, **there is no loss of optimality in restricting attention to non-stationary deterministic (rather than stochastic) agent-state based policies.**

- 5) As far as we are aware, the result in (7) is new. A similar result is claimed in [5, Proposition 1] for memoryless policies in POMDPs (i.e.,  $Z_t = Y_t$ ), but for a different definition of optimality.
- 6) Note that the dynamics of  $\{\xi_t\}_{t \geq 1}$  in (6) are deterministic. Therefore, it is possible to use ideas from deterministic optimization to find the optimal *trajectory*  $\{\xi_t^*\}_{t \geq 1}$  and optimal *open-loop* ‘‘control actions’’  $\{\pi_t^*\}_{t \geq 1}$ .

### C. Policy search methods to find a locally optimal policy in $\Pi_{\text{SS}}$

We can use standard results from policy search methods for MDPs to identify a stationary policy that is locally optimal (within the class of stationary policies). The high-level idea is based on the fact that the *product process*  $\{(S_t, Z_t)\}_{t \geq 1}$  is a controlled Markov process with the controlled transition probability:

$$\begin{aligned} P_{\text{PROD}}(s', z' | s, z, a) \\ := \sum_{y' \in \mathcal{Y}} P(s', y' | s, a) \mathbb{1}_{\{z' = \phi(z, y', a)\}}. \end{aligned}$$

Now, for any  $\pi = (\pi, \pi, \dots) \in \Pi_{\text{SS}}$ , define:

$$P_{\text{PROD}}^\pi(s', z' | s, z) := \sum_{a \in \mathcal{A}} \pi(a|z) P_{\text{PROD}}(s', z' | s, z, a).$$

and

$$r_{\text{PROD}}^\pi(s, z) := \sum_{a \in \mathcal{A}} \pi(a|z) r(s, a).$$

Let  $V_{\text{PROD}}^\pi: \mathcal{S} \times \mathcal{Z} \rightarrow \mathbb{R}$  be the value function of the policy  $\pi$  in the  $(\mathcal{S}, \mathcal{Z})$  space, i.e., it is the solution of the following policy evaluation formula:

$$\begin{aligned} V_{\text{PROD}}^\pi(s, z) = r_{\text{PROD}}^\pi(s, z) \\ + \gamma \sum_{\substack{(s', z') \in \\ \mathcal{S} \times \mathcal{Z}}} P_{\text{PROD}}^\pi(s', z' | s, z) V_{\text{PROD}}^\pi(s', z'). \end{aligned} \quad (8)$$

Let  $Q_{\text{PROD}}^\pi: \mathcal{S} \times \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$  denote the action-value function of the policy  $\pi$ , i.e., for any  $(s, z, a) \in \mathcal{S} \times \mathcal{Z} \times \mathcal{A}$ ,

$$\begin{aligned} Q_{\text{PROD}}^\pi(s, z, a) = r(s, a) \\ + \gamma \sum_{\substack{(s', z') \in \\ \mathcal{S} \times \mathcal{Z}}} P_{\text{PROD}}(s', z' | s, z, a) V_{\text{PROD}}^\pi(s', z'). \end{aligned}$$

Finally, let  $d_{\text{PROD}}^\pi$  denote the *unnormalized* occupancy measure of the policy  $\pi$ , i.e.,

$$d_{\text{PROD}}^\pi(s, z, a) = \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{P}^\pi(S_t = s, Z_t = z, A_t = a).$$

Then, the performance of policy  $\pi \in \Pi_{\text{SS}}$  is given by

$$J^\pi = \sum_{(s, z) \in \mathcal{S} \times \mathcal{Z}} \xi_1(s, z) V_{\text{PROD}}^\pi(s, z),$$

where  $\xi_1(s, z) := \mathbb{P}(S_1 = s, Z_1 = z)$  denotes the initial distribution of the joint state of the environment and the agent. This shows that performance evaluation of any policy in  $\Pi_{\text{SS}}$  is straightforward.

Given the above performance evaluation formula, we can use policy-gradient based methods to find a good policy. In particular, suppose  $\pi = (\pi, \pi, \dots) \in \Pi_{\text{SS}}$  is a parameterized policy (e.g., softmax or mixture of Gaussians) with policy parameters  $\theta$ . We will use  $\pi_\theta$  to denote the policy with parameters  $\theta$  and  $\pi_\theta$  to denote the decision rule with parameters  $\theta$ . Then, by a straight forward modification of the policy gradient formula [12], [47], [48], [81], we get

$$\boxed{\nabla_\theta J^{\pi_\theta} = \sum_{\substack{(s, z, a) \in \\ \mathcal{S} \times \mathcal{Z} \times \mathcal{A}}} d_{\text{PROD}}^{\pi_\theta}(s, z, a) Q_{\text{PROD}}^{\pi_\theta}(s, z, a) \nabla_\theta \log \pi_\theta(a|z).} \quad (9)$$

Consequently, we can use any policy gradient based search algorithm (such as actor-critic and its variants; see [86]) to find locally optimal policy parameters, i.e., policy parameters  $\theta^*$  such that

$$\nabla_\theta J^{\pi_\theta} \Big|_{\theta=\theta^*} = 0.$$

Then the policy  $\pi_{\theta^*}$  is a locally optimal policy in class  $\Pi_{\text{SS}}$ .

Some remarks on the policy gradient approach

- 1) The policy evaluation formula (8) using the product state  $(S_t, Z_t)$  was initially presented in [67] and has been rediscovered in slightly different forms multiple times [15], [36], [37], [53].
- 2) There is a rich history of policy based algorithms in Systems and Control, Operations Research, and Artificial Intelligence. We refer the reader to [8], [11] for detailed overviews.
- 3) Several other approaches for policy search in  $\Pi_{\text{NS}}$  have been proposed in the literature [36], [46], [49], [61].
- 4) Note that the formula in (9) is applicable when the system dynamics are known (and therefore the  $Q_{\text{PROD}}^\pi$  can be computed via policy evaluation formula described above). In the RL setting, the environment state  $S_t$  is not observed, so the above formula cannot be used directly.
- 5) For RL, the idea of using policy gradient methods for POMDPs was first presented in [45] and variations have been proposed in [1], [6], [40], [49], [57], [65], [92]. However, these policy gradient formula proposed in these papers is for the RL setting where the agent state  $S_t$  is not observed; therefore the exact expressions were different. We present an alternative approach in Sec. IV.
- 6) As pointed out in [45], the agent state in a POMDP can be viewed as a *feature* of the entire history of observations and actions. With this viewpoint, the policy gradient formula (9) may be viewed as a special case of actor-critic algorithms with features [47], [48].
- 7) Another idea that has been used in the RL setting is asymmetric actor critic [7], [75], where it is assumed that the critic has access to the environment state (which is the case in simulation environments). In such settings, it is possible to learn  $Q^\pi(s, z, a)$  via temporal difference learning. See [7], [75] for a discussion.

#### D. The approximate information state approach to find a good policy in $\Pi_{\text{SD}}$

In this section, we present an alternative approach to finding good policies in  $\Pi_{\text{SD}}$  called the *approximate information state* (AIS), which was originally developed in [80]. We first start with the intuition and then explain the technical results.

*Intuition:* One way to obtain a policy is to posit any dynamics and rewards  $(P_{\text{AIS}}, r_{\text{AIS}})$ , where  $P_{\text{AIS}}: \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  and  $r_{\text{AIS}}: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$ , and obtain a policy  $\bar{\pi}_{\text{AIS}}$  by solving the following dynamic program:

$$Q_{\text{AIS}}(z, a) = r_{\text{AIS}}(z, a) + \gamma \sum_{z' \in \mathcal{Z}} P_{\text{AIS}}(z'|z, a) V_{\text{AIS}}(z'), \quad (10a)$$

$$V_{\text{AIS}}(z) = \max_{a \in \mathcal{A}} Q_{\text{AIS}}(z, a). \quad (10b)$$

Let  $\pi_{\text{AIS}}(z)$  denote any arg max of the right hand side of (10b). Define the policy  $\bar{\pi}_{\text{AIS}} = (\bar{\pi}_{\text{AIS},1}, \bar{\pi}_{\text{AIS},2}, \dots)$  given

by<sup>10</sup>

$$\bar{\pi}_{\text{AIS},t}(h_t) = \pi_{\text{AIS}}(\bar{\sigma}_t(h_t)) \quad (11)$$

As shown in Theorem 1, if the posited dynamics and rewards  $(P_{\text{AIS}}, r_{\text{AIS}})$  satisfy properties (P1) and (P2), the policy  $\bar{\pi}_{\text{AIS}}$  is optimal. But what happens when properties (P1) and (P2) are not satisfied?

Let  $(\varepsilon, \delta)$  with  $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots)$  and  $\delta = (\delta_1, \delta_2, \dots)$  with  $\varepsilon_t, \delta_t \in \mathbb{R}_{\geq 0}$  be such that the following properties are satisfied:

**(AP1) Approximately sufficient for performance evaluation.** For any  $t$ , and  $H_t$  and  $A_t$  we have

$$|\mathbb{E}[R_t | H_t, A_t] - r_{\text{AIS}}(\bar{\sigma}_t(H_t), A_t)| \leq \varepsilon_t.$$

**(AP2) Sufficient for predicting itself.** For any  $t$ , and  $H_t$  and  $A_t$  we have

$$d_{\mathfrak{F}}(\mathbb{P}(Z_{t+1} | H_t, A_t), P_{\text{AIS}}(Z_{t+1} | \bar{\sigma}_t(H_t), A_t)) \leq \delta_t.$$

where  $d_{\mathfrak{F}}$  is a metric on  $\Delta(\mathcal{Z})$ , which we will make precise later.

Then the tuple  $(\bar{\sigma}, P_{\text{AIS}}, r_{\text{AIS}})$  is called an *approximate information state* (AIS).

Engineering intuition suggests that the sub-optimality of  $\bar{\pi}_{\text{AIS}}$ , i.e.,

$$\alpha := \bar{J}_{\text{ND}}^* - J^{\bar{\pi}_{\text{AIS}}}$$

should be a continuous function of  $(\varepsilon, \delta)$ . The results of [80] formalize this intuition.

*The choice of metric:* To formalize the above intuition, we need to choose a metric  $d_{\mathfrak{F}}$  on probability spaces. For our purposes, it is convenient to work with a class of metrics known as integral probability measures (IPMs) [60]. Let  $\mathcal{V}$  denote the family of all bounded (measurable) functions from  $\mathcal{Z}$  to  $\mathbb{R}$  and  $\mathcal{P}$  denote the set of all probability measures on  $\mathcal{Z}$  with finite mean.<sup>11</sup>

**Definition 1** Let  $\mathfrak{F}$  be a convex and balanced<sup>12</sup> subset of  $\mathcal{V}$ . Then, the IPM distance (w.r.t.  $\mathfrak{F}$ ) between two probability laws  $\nu_1, \nu_2 \in \mathcal{P}$  is given by

$$d_{\mathfrak{F}}(\nu_1, \nu_2) = \sup_{f \in \mathfrak{F}} \left| \int f d\nu_1 - \int f d\nu_2 \right|.$$

**Definition 2** In the setting of Definition 1, the Minkowski functional of any measurable function  $f \in \mathcal{V}$  is defined as

$$\rho_{\mathfrak{F}}(f) = \inf \left\{ \rho \in \mathbb{R}_{>0} : \frac{f}{\rho} \in \mathfrak{F} \right\}.$$

Note that if for every positive  $\rho$ ,  $f/\rho \notin \mathfrak{F}$ , then  $\rho_{\mathfrak{F}}(f) = \infty$ .

An immediate consequence of the above two definitions is that for any measurable function  $f \in \mathcal{V}$ ,

$$\left| \int f d\nu_1 - \int f d\nu_2 \right| \leq \rho_{\mathfrak{F}}(f) d_{\mathfrak{F}}(\nu_1, \nu_2). \quad (12)$$

<sup>10</sup>Recall that  $\bar{\sigma} = (\bar{\sigma}_1, \bar{\sigma}_2, \dots)$  is the history compression function corresponding to the agent-state update function  $\phi$  and is given by (3).

<sup>11</sup>When  $\mathcal{Z}$  is finite  $\mathcal{V} \equiv \mathbb{R}^{|\mathcal{Z}|}$  and  $\mathcal{P}$  is the set of all probability mass functions on  $\mathcal{Z}$ .

<sup>12</sup>A subset  $\mathfrak{F}$  of  $\mathcal{V}$  is balanced if for all  $f \in \mathfrak{F}$  and scalars  $a$  such that  $|a| \leq 1$ ,  $af \in \mathfrak{F}$ .



Many of the commonly used metrics on probability spaces are IPMs. For example

- **Total variation distance**, denoted by  $d_{TV}$ , corresponds to  $\mathfrak{F} = \mathfrak{F}_{TV} := \{f \in \mathcal{V} : \frac{1}{2} \text{span}(f) \leq 1\}$ , where  $\text{span}(f) = \sup(f) - \inf(f)$  [60], [84]. In this case  $\rho_{\mathfrak{F}}(f) = \frac{1}{2} \text{span}(f)$ .
- **Wasserstein distance**. Suppose  $(\mathcal{Z}, d_{\mathcal{Z}})$  is a metric space. Then the Wasserstein distance, denoted by  $d_{WAS}$ , corresponds to  $\mathfrak{F} = \mathfrak{F}_{WAS} := \{f \in \mathcal{V} : \text{Lip}(f) \leq 1\}$  where  $\text{Lip}(f)$  denotes the Lipschitz constant of a function  $f$  [83], [84]. In this case  $\rho_{\mathfrak{F}}(f) = \text{Lip}(f)$ .

Other distances such as Kantorovich distance, bounded Lipschitz, maximum mean discrepancy, are also instances of IPM. See [60], [80] for details.

We now state the main result of [80].

### Theorem 3: AIS based DP

Given a function class  $\mathfrak{F}$ , let  $(\vec{\sigma}, P_{AIS}, r_{AIS})$  be an  $(\varepsilon, \delta)$ -AIS with respect to  $d_{\mathfrak{F}}$ . Consider the DP (10). Then, the policy  $\vec{\pi}_{AIS}$  given by (11) satisfies

$$\vec{J}_{ND}^* - J^{\vec{\pi}_{AIS}} \leq \frac{2}{1-\gamma} [\varepsilon + \gamma \delta \rho_{\mathfrak{F}}(V_{AIS}^*)], \quad (13)$$

where

$$\varepsilon = (1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} \varepsilon_t, \quad \delta = (1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} \delta_t.$$

Some remarks on AIS

- 1) We have presented the simplest form of the AIS result. Several variations including finite horizon models, multi-agent systems, and others are presented in [80].
- 2) The results of Theorem 3 may be viewed as a generalization of approximation bounds for MDPs. The earliest such result is [59]. Several variations of such results for MDPs have been derived in the literature, including model approximation in MDPs [2], state abstraction in MDPs [2], [31], [52]. As shown in [80], many of these results are special instances of Theorem 3.
- 3) When defining information state, it was mentioned that (P2b) is a sufficient condition for (P2). Along similar lines, condition (AP2) can be relaxed to approximately predicting the next observation. See [80] for details.
- 4) One option for solving belief-state based POMDPs is by discretizing the belief space. Approximation bounds for belief discretization were derived in [30] and it was shown in [80] that such bounds may be viewed as a special instance of Theorem 3.
- 5) A heuristic algorithm to jointly learn a good AIS representation and a agent-state based policy is proposed in [91].
- 6) In addition to the AIS approach presented here, there are two popular approaches to determine good policies in  $\Pi_{SD}$ : bisimulation metrics [16], [17] and predictive state representations [10], [35], [41], [50], [51], [68].

These approaches are based on different approximation philosophies and are conceptually different.

- 7) The results of Theorem 3 are also applicable for worst case design. In particular, let the *worst-case distance*  $d_{\mathfrak{F}}(X_1, X_2)$  between two bounded random variables  $X_1$  and  $X_2$  be defined as the Hausdorff distance between their supports. It is then shown in [23]–[25] that under an additional smoothness assumption on the dynamics, one may define an AIS in terms of the worst-case distance and generalize the results of Theorem 3 to provide worst case guarantees for agent-state based policies.

## IV. REINFORCEMENT LEARNING APPROACHES TO LEARN A POLICY IN $\Pi_{SD}$ OR $\Pi_{SS}$

As highlighted in the introduction, one of the motivations for considering agent-state based policies is that they are easier to use in the RL setting because the agent-state update does not depend on the system model. In this section, we review the different approaches for RL in POMDPs with agent state.

Before stating the results, we start with a remark that is not sufficiently highlighted in the literature. In the standard POMDP model, it is assumed that the actions are a function of the history of past observations and actions. In the learning setup, actions are also sometimes allowed to be a function of the history of rewards. Therefore, in the literature on RL for POMDPs, it is always implicitly assumed that the observations include the rewards as well. If that is not the case, we need to consider an agent-state update rule of the form  $Z_{t+1} = \phi(Z_t, Y_{t+1}, A_t, R_{t+1})$ , but that does not fundamentally change the nature of the results.

### A. Agent-state based Q-learning (ASQL)

One of the simplest RL algorithms is Q-learning [85], where the agent learns the Q-function by bootstrapping. For MDPs, it is well understood that the Q-learning iterates converge to the optimal Q-function, which is the solution of the dynamic programming equation.

Agent-state based Q-learning (ASQL) uses the same idea for POMDPs with agent state. In the simplest setting, the agent acts in the environment according to a behavior policy  $\mu: \mathcal{Z} \rightarrow \Delta(\mathcal{A})$  and generates the experience  $(Y_1, Z_1, A_1, R_1, Y_2, Z_2, A_2, R_2, \dots)$ . The agent uses the experience to recursively update Q-function  $\{Q_t\}_{t \geq 1}$ ,  $Q_t: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$  as follows:

$$Q_{t+1}(z, a) = Q_t(z, a) + \alpha_t(z, a) \left[ R_t + \gamma \max_{a' \in \mathcal{A}} Q_t(Z_{t+1}, a') - Q_t(z, a) \right] \quad (14)$$

where  $R_t = r(S_t, A_t)$  is the reward at time  $t$  and  $\{\alpha_t\}_{t \geq 1}$  are the learning rates which are chosen such that  $\alpha_t(z, a) = 0$  when  $(z, a) \neq (Z_t, A_t)$ .

As mentioned earlier, the agent state does not satisfy the Markov property. So the standard proof of Q-learning for MDPs [9], [40], [85] is not directly applicable. Nonetheless,

it is shown in [74] that ASQL converges under the following assumptions:

- (A1) The behavioral policy  $\mu$  is such that the Markov chain<sup>13</sup>  $\{(S_t, Y_t, Z_t, A_t)\}_{t \geq 1}$  is irreducible and aperiodic and, therefore, has a limiting distribution  $\zeta^\mu$ . Moreover, under the behavioral policy, all  $(z, a)$  are visited infinitely often, i.e.,  $\sum_{(s,y)} \zeta^\mu(s, y, z, a) > 0$ .
- (A2) The learning rates  $\{\alpha_t\}_{t \geq 1}$  satisfy the standard Robbins–Monro conditions, i.e., for all  $(z, a)$  we have that

$$\sum_{t \geq 1} \alpha_t(z, a) = \infty \quad \text{and} \quad \sum_{t \geq 1} \alpha_t(z, a)^2 < \infty, \text{ a.s.}$$

For the ease of notation, we will continue to use  $\zeta^\mu$  to denote the marginal and conditional distributions w.r.t.  $\zeta^\mu$ . In particular, for marginals we use  $\zeta^\mu(y, z, a)$  to denote  $\sum_{s \in \mathcal{S}} \zeta^\mu(s, y, z, a)$  and so on; for conditionals, we use  $\zeta^\mu(s|z, a)$  to denote  $\zeta^\mu(s, z, a)/\zeta^\mu(z, a)$  and so on.

We now state the main result of [74].<sup>14</sup>

#### Theorem 4: Convergence of ASQL

Under assumptions (A1) and (A2), the iterates  $\{Q_t\}_{t \geq 1}$  of ASQL (14) converge almost surely to a limit  $Q_{\text{ASQL}}^\mu$ , which is the fixed point of the following DP:

$$Q_{\text{ASQL}}^\mu(z, a) = r_{\text{ASQL}}^\mu(z, a) + \gamma \sum_{z' \in \mathcal{Z}} P_{\text{ASQL}}^\mu(z'|z, a) \max_{a' \in \mathcal{A}} Q_{\text{ASQL}}^\mu(z', a') \quad (15)$$

where

$$r_{\text{ASQL}}^\mu(z, a) = \sum_{s \in \mathcal{S}} \zeta^\mu(s|z, a) r(s, a),$$

$$P_{\text{ASQL}}^\mu(z'|z, a) = \sum_{(s, y') \in \mathcal{S} \times \mathcal{Y}} \mathbb{1}_{\{z' = \phi(z, y', a)\}} P(y'|s, a) \zeta^\mu(s|z, a).$$

Let  $\pi_{\text{ASQL}}^\mu \in \Pi_{\text{SD}}$  denote the policy

$$\pi_{\text{ASQL}}^\mu(z) := \arg \max_{a \in \mathcal{A}} Q_{\text{ASQL}}^\mu(z, a) \quad (16)$$

Theorem 4 implies that the greedy policy corresponding to  $\{Q_t\}_{t \geq 1}$  computed via (14) converges to  $\pi_{\text{ASQL}}^\mu$ . Define  $\bar{\pi}_{\text{ASQL}}^\mu = (\bar{\pi}_{\text{ASQL},1}^\mu, \bar{\pi}_{\text{ASQL},2}^\mu, \dots)$  to be the history dependent policy given by

$$\bar{\pi}_{\text{ASQL},t}^\mu(h_t) = \pi_{\text{ASQL}}^\mu(\sigma_t(h_t)). \quad (17)$$

A natural follow-up question is: How well does  $\bar{\pi}_{\text{ASQL}}^\mu$  perform? The answer to this question is surprisingly easy because we can view  $(\bar{\sigma}, P_{\text{ASQL}}^\mu, r_{\text{ASQL}}^\mu)$  as an AIS. Therefore, Theorem 3 provides an immediate bound on  $\bar{J}_{\text{ND}}^* - J^{\bar{\pi}_{\text{ASQL}}^\mu}$  as follows:

**Corollary 1** We have that

$$\bar{J}_{\text{ND}}^* - J^{\bar{\pi}_{\text{ASQL}}^\mu} \leq \frac{2}{1 - \gamma} [\varepsilon + \gamma \delta \rho_{\bar{\sigma}}(V_{\text{ASQL}}^\mu)]$$

<sup>13</sup>As argued in Sec. III-C, under any policy  $\mu \in \Pi_{\text{SS}}$  the product process  $\{(S_t, Z_t)\}_{t \geq 1}$  is a Markov chain. This directly implies that the process  $\{(S_t, Y_t, Z_t, A_t)\}_{t \geq 1}$  is a Markov chain as well.

<sup>14</sup>The result proved in [74] is a generalization of Theorem 4 to periodic agent-state based policies. Taking a period of 1 implies Theorem 4.

where  $V_{\text{ASQL}}^\mu(z) := \max_{a \in \mathcal{A}} Q_{\text{ASQL}}^\mu(z, a)$  and  $\varepsilon$  and  $\delta$  are defined in a manner similar to Sec. III-D

The results of Theorem 4 and Corollary 1 suggest that the performance of ASQL could be improved by choosing the state update function  $\phi$  so as to minimize the approximation losses  $\varepsilon$  and  $\delta$ . In particular, consider an implementation of Q-learning where a recurrent neural network (RNN) such as LSTM (long short term memory) or GRU (gated recurrent unit) is used as a history compression function  $\bar{\sigma}$ . Add an ‘‘AIS-block’’ which learns a generative model for  $(P_{\text{ASQL}}, r_{\text{ASQL}})$  by using the AIS loss  $\lambda \varepsilon^2 + (1 - \lambda) \delta^2$  as an auxiliary loss. This algorithm is called RQL-AIS in [71]. Note that we can generate RQL-AIS version of any implementation of recurrent Q-learning for POMDPs.

In [71], a version of RQL-AIS using R2D2 [42] as a base implementation is compared with R2D2 on the minigrid benchmark [21], [22]. The results of [71] show that adding an AIS block significantly improve the performance of Q-learning in the larger and more complicated minigrid environments.

#### Remarks on ASQL

- 1) A key feature of Theorem 4 is that the limit  $Q_{\text{ASQL}}^\mu$  depends on the behavioral policy  $\mu$ , unlike in MDPs where the limit is policy-independent (as long as all state-action pairs are visited infinitely often). This dependence arises because the agent state is not an information state and thus is not policy-independent. See [90] for a discussion on policy independence of information states.
- 2) The idea of adding AIS losses as an auxiliary loss for ASQL was first proposed in [80] for actor-critic algorithms. It is argued in [62] that in ASQL it is not useful to learn the reward function  $r_{\text{ASQL}}^\mu$  since the Q-function is already being learnt. Rather one should simply use  $\delta^2$  as an auxiliary loss.
- 3) Learning a dynamics model to minimize  $\delta^2$  is effectively the same as self-supervised representation learning, which has been used as a heuristic in different forms in MDPs and POMDPs in the literature. See [62] for a detailed discussion.
- 4) The convergence of the simplest version of ASQL was established in [73] for  $Z_t = Y_t$  under the assumption that the actions are chosen i.i.d. (and do not depend on  $Z_t$ ). In [64] it was established that  $Q_{\text{ASQL}}^\mu$  is the fixed point of (14), but convergence of  $\{Q_t\}_{t \geq 1}$  to  $Q_{\text{ASQL}}^\mu$  was not established. The convergence of ASQL when the agent state is a finite window of past observations and actions was established in [44]. The convergence result for a general agent state presented in Theorem 4 follows essentially the same proof argument.
- 5) ASQL is closely related to Q-learning for MDPs with state aggregation or quantization [43], [72], [82]. In particular, consider an MDP with large or continuous state space, where we partition the state space into bins, and treat each bin as an aggregated state.

This model is equivalent to a PODMP where the observation equals the index of the bin corresponding to the current state. Thus, Q-learning in MDPs with state aggregation is equivalent to ASQL with agent state equal to current observation. Therefore, there are considerable similarities between the convergence analysis and approximation bounds of these two setups.

- 6) ASQL may also be viewed as a MDP with non-Markovian state [18], [27]. As such there are considerable similarities between the convergence analysis and approximation bounds for these two setups.
- 7) The regret of an optimistic variant of ASQL is presented in [29], though the exact setting there is slightly different. See [29] for details.

### B. Agent-state based actor critic (ASAC)

ASQL always learns a deterministic policy (i.e., a policy belonging to  $\Pi_{\text{SD}}$ ). As highlighted by (5) and fact (F2), stochastic policies can perform better than deterministic policies. One approach to learn stochastic policies is to use the actor-critic class of algorithms.

An agent-state based actor critic (ASAC) algorithm was proposed in [80]. We present it with a slightly different perspective here. Consider any policy  $\pi \in \Pi_{\text{SS}}$ . Motivated by the result of Theorem 4, we can run a temporal difference learning algorithm to learn a Q-function corresponding to  $\pi$  as follows:

$$Q_{t+1}^\pi(z, a) = Q_t^\pi(z, a) + \alpha_t(z, a) \left[ R_t + \gamma Q_t^\pi(Z_{t+1}, A_{t+1}) - Q_t^\pi(z, a) \right] \quad (18)$$

where  $R_t = r(S_t, A_t)$  as before. Theorem 4 implies that under conditions similar to (A1) and (A2), the iteration (18) converges almost surely to a limit  $Q_{\text{ASAC}}^\pi$ , which is the fixed point of the following DP:

$$Q_{\text{ASAC}}^\pi(z, a) = r_{\text{ASAC}}^\pi(z, a) + \gamma \sum_{(z', a') \in \mathcal{Z} \times \mathcal{A}} \pi(a'|z') P_{\text{ASAC}}^\pi(z'|z, a) Q_{\text{ASAC}}^\pi(z', a') \quad (19)$$

where

$$r_{\text{ASAC}}^\pi(z, a) = \sum_{s \in \mathcal{S}} \zeta^\pi(s, a) r(s, a),$$

$$P_{\text{ASAC}}^\pi(z'|z, a) = \sum_{(s, y') \in \mathcal{S} \times \mathcal{Y}} \mathbb{1}_{\{z' = \phi(z, y', a)\}} P(y'|s, a) \zeta^\pi(s|z, a).$$

Define  $J_{\text{ASAC}}^\pi = \sum_{(z, a) \in \mathcal{Z} \times \mathcal{A}} \xi_1(z, a) Q_{\text{ASAC}}^\pi(z, a)$ . Moreover, define  $\vec{\pi} = (\vec{\pi}_1, \vec{\pi}_2, \dots)$  to be the history dependent policy given by

$$\vec{\pi}_t(h_t) = \pi(\vec{\sigma}_t(h_t)).$$

Then, by an argument similar to Theorem 3 we have that (see [80])

$$|J^{\vec{\pi}} - J_{\text{ASAC}}^\pi| \leq \frac{1}{1 - \gamma} [\varepsilon + \gamma \delta \rho_{\vec{\sigma}}(V_{\text{ASAC}}^\pi)]$$

where  $V_{\text{ASAC}}^\pi(z) = \max_{a \in \mathcal{A}} Q_{\text{ASAC}}^\pi(z, a)$  and  $\varepsilon$  and  $\delta$  are defined in a manner similar to Sec. III-D.

The above discussion suggests the following algorithm. Consider an implementation of an actor-critic algorithm where an RNN such as LSTM or GRU is used as a history compression function  $\vec{\sigma}$ . Similar to RQL-AIS, add an ‘‘AIS block’’ which learns a generative model for  $(P_{\text{ASAC}}, r_{\text{ASAC}})$  by using  $\lambda \varepsilon^2 + (1 - \lambda) \delta^2$  as an auxiliary loss. This algorithm is called PORL-AIS in [80]. Note that we can generate a PORL-AIS version of any implementation of recurrent actor-critic for POMDPs.

In [80], a version of PORL-AIS is compared with LSTM-based actor critic on the minigrid benchmark [21], [22]. The results of [80] show that adding an AIS block significantly improves the performance of actor-critic algorithms in the larger and more complicated minigrid environments.

### Remarks of ASAC

- 1) It is also possible to add AIS losses to actor only algorithms such as G(PO)MDP [8]. See [80] for details.
- 2) We have argued that the TD iterations for  $\{Q_t\}_{t \geq 1}$  converge to a limit. Therefore, for a parameterized policy  $\pi_\theta$ , the gradient computed via the policy gradient formula

$$\frac{1}{T} \sum_{t=1}^T \nabla_\theta \log \pi_\theta(A_t | Z_t) Q_t^{\pi_\theta}(Z_t, A_t)$$

asymptotically converges to

$$\sum_{(z, a) \in \mathcal{Z} \times \mathcal{A}} \zeta^\pi(z, a) \log \pi_\theta(a|z) Q_{\text{ASQL}}^\pi(z, a)$$

which is an unbiased estimate of  $\nabla_\theta J_{\text{ASAC}}^{\pi_\theta}$  rather than  $\nabla_\theta J^{\pi_\theta}$ . Therefore, PORL-AIS does not have the local convergence guarantees of actor critic for MDPs (but neither do other actor-critic algorithms for POMDPs!).

- 3) AIS error bounds for asymmetric actor critic [7] are developed in [75] but the numerical experiments presented there suggest that adding an AIS block to asymmetric actor critic algorithms does not provide any measurable performance improvement over the vanilla asymmetric actor critic algorithm.

## V. CONCLUSION

In this tutorial, we have summarized some of the approaches used for planning and learning agent-state based policies in POMDPs. Our survey is not exhaustive. There are other alternative approaches which we have not discussed (e.g., predictive state representations [10], [35], [41], [51], bisimulation [16], [17], world models [32]–[34], [70], [87], decision transformers [19], [63]). Nonetheless, it is our hope that this paper will present the reader with a unified perspective on a subset of the literature. We conclude by emphasizing that this is a rapidly developing research area and many of the fundamental questions are still open.

### ACKNOWLEDGMENT

The authors are grateful to Demosthenis Teneketzis (University of Michigan), Jayakumar Subramanian (Adobe Inc.),

and Matthieu Geist (Coher) for discussions and collaborations which helped the authors develop the viewpoint presented in this tutorial. They also thank Charalambos Charalambous (University of Cyprus), Vijay Subramanian (University of Michigan), and Serdar Yüksel (Queen’s University), and Tianwei Ni (University of Montreal) for helpful feedback..

## REFERENCES

- [1] D. Aberdeen and J. Baxter, “Scaling internal-state policy-gradient methods for POMDPs,” in *International Conference on Machine Learning*. San Francisco, CA, USA: Morgan Kaufman, 2002, pp. 3–10.
- [2] K. Asadi, D. Misra, and M. Littman, “Lipschitz continuity in model-based reinforcement learning,” in *International Conference on Machine Learning*, vol. 80. PMLR, Jul. 2018, pp. 264–273.
- [3] K. Åström, “Optimal control of Markov processes with incomplete state information,” *Journal of Mathematical Analysis and Applications*, vol. 10, no. 1, pp. 174–205, Feb. 1965.
- [4] R. J. Aumann, *Mixed and behavior strategies in infinite extensive games*. Princeton University Princeton, 1961.
- [5] J. A. Bagnell, S. Kakade, A. Y. Ng, and J. Schneider, “Policy search by dynamic programming,” in *Conference on Neural Information Processing Systems*. Cambridge, MA, USA: MIT Press, 2003, p. 831–838.
- [6] L. Baird and A. Moore, “Gradient descent for general reinforcement learning,” in *Conference on Neural Information Processing Systems*. Cambridge, MA, USA: MIT Press, 1998, p. 968–974.
- [7] A. Baisero and C. Amato, “Unbiased asymmetric reinforcement learning under partial observability,” in *International Conference on Autonomous Agents and Multiagent Systems*. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2022, p. 44–52.
- [8] J. Baxter and P. L. Bartlett, “Infinite-horizon policy-gradient estimation,” *Journal of Artificial Intelligence Research*, vol. 15, pp. 319–350, 2001.
- [9] D. Bertsekas and J. N. Tsitsiklis, *Neuro-dynamic programming*. Athena Scientific, 1996.
- [10] B. Boots, S. M. Siddiqi, and G. J. Gordon, “Closing the learning-planning loop with predictive state representations,” *The International Journal of Robotics Research*, vol. 30, no. 7, pp. 954–966, 2011.
- [11] X.-R. Cao, “From perturbation analysis to Markov decision processes and reinforcement learning,” *Discrete Event Dynamic Systems*, vol. 13, no. 1/2, pp. 9–39, 2003.
- [12] —, *Stochastic Learning and Optimization*. Springer, 2007.
- [13] A. Cassandra, M. L. Littman, and N. L. Zhang, “Incremental pruning: A simple, fast, exact method for partially observable Markov decision processes,” in *Uncertainty in Artificial Intelligence*, 1997.
- [14] A. R. Cassandra, L. P. Kaelbling, and M. L. Littman, “Acting optimally in partially observable stochastic domains,” in *AAAI Conference on Artificial Intelligence*, vol. 94, 1994, pp. 1023–1028.
- [15] A. R. Cassandra, “Exact and approximate algorithms for partially observable markov decision processes,” Ph.D. dissertation, Brown University, 1998.
- [16] P. S. Castro, T. Kastner, P. Panangaden, and M. Rowland, “Mico: Improved representations via sampling-based state similarity for Markov decision processes,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 30 113–30 126, 2021.
- [17] P. S. Castro, P. Panangaden, and D. Precup, “Equivalence relations in fully and partially observable Markov decision processes,” in *International Joint Conference on Artificial Intelligence*, 2009, pp. 1653–1658.
- [18] S. Chandak, P. Shah, V. S. Borkar, and P. Dodhia, “Reinforcement learning in non-Markovian environments,” *Systems & Control Letters*, vol. 185, p. 105751, 2024.
- [19] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, “Decision transformer: Reinforcement learning via sequence modeling,” in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 15 084–15 097.
- [20] H.-T. Cheng, “Algorithms for partially observable Markov decision processes,” Ph.D. dissertation, University of British Columbia, Vancouver, BC, 1988.
- [21] M. Chevalier-Boisvert, B. Dai, M. Towers, R. Perez-Vicente, L. Willems, S. Lahlou, S. Pal, P. S. Castro, and J. Terry, “Minigrid & Miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks,” in *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp. 73 383–73 394.
- [22] M. Chevalier-Boisvert, L. Willems, and S. Pal, “Minimalistic gridworld environment for gymnasium,” <https://github.com/Farama-Foundation/MiniGrid>, 2018.
- [23] A. Dave, I. Faros, N. Venkatesh, and A. A. Malikopoulos, “Worst-case control and learning using partial observations over an infinite time horizon,” in *IEEE Conference on Decision and Control*. IEEE, 2023, pp. 6020–6025.
- [24] A. Dave, N. Venkatesh, and A. A. Malikopoulos, “Approximate information states for worst-case control of uncertain systems,” in *Conference on Decision and Control*. IEEE, 2022, pp. 4945–4950.
- [25] —, “Approximate information states for worst-case control and learning in uncertain systems,” 2024, arXiv preprint, arXiv:2301.05089.
- [26] A. Devran Kera and S. Yüksel, “Partially observed optimal stochastic control: Regularity, existence, approximations, and learning,” in *submitted to IEEE Conference on Decision and Control*. IEEE, Dec. 2024.
- [27] —, “Q-learning for stochastic control under general information structures and non-markovian environments,” *Transactions on Machine Learning Research*, 2024.
- [28] J. S. Dibangoye, C. Amato, O. Buffet, and F. Charpillet, “Optimally solving Dec-POMDPs as continuous-state MDPs,” *Journal of Artificial Intelligence Research*, vol. 55, pp. 443–497, Feb. 2016.
- [29] S. Dong, B. V. Roy, and Z. Zhou, “Simple agent, complex environment: Efficient reinforcement learning with agent states,” *Journal of Machine Learning Research*, vol. 23, no. 255, pp. 1–54, 2022.
- [30] V. Francois-Lavet, G. Rabusseau, J. Pineau, D. Ernst, and R. Fonteneau, “On overfitting and asymptotic bias in batch reinforcement learning with partial observability,” *Journal of Artificial Intelligence Research*, vol. 65, pp. 1–30, 2019.
- [31] C. Gelada, S. Kumar, J. Buckman, O. Nachum, and M. G. Bellemare, “DeepMDP: Learning continuous latent space models for representation learning,” in *International Conference on Machine Learning*, 2019, pp. 2170–2179.
- [32] D. Ha and J. Schmidhuber, “Recurrent world models facilitate policy evolution,” in *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc., 2018.
- [33] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi, “Dream to control: Learning behaviors by latent imagination,” in *International Conference on Learning Representations*, 2020.
- [34] D. Hafner, T. P. Lillicrap, M. Norouzi, and J. Ba, “Mastering atari with discrete world models,” in *International Conference on Learning Representations*, 2021.
- [35] W. Hamilton, M. M. Fard, and J. Pineau, “Efficient learning and planning with compressed predictive states,” *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3395–3439, 2014.
- [36] E. A. Hansen, “Solving POMDPs by searching in policy space,” in *Uncertainty in Artificial Intelligence*, San Francisco, CA, USA, 1998, p. 211–219.
- [37] M. Hauskrecht, “Planning and control in stochastic domains with imperfect information,” Ph.D. dissertation, Massachusetts Institute of Technology, 1997.
- [38] —, “Value-function approximations for partially observable markov decision processes,” *Journal of artificial intelligence research*, vol. 13, pp. 33–94, 2000.
- [39] M. Hutter, *Universal artificial intelligence: Sequential decisions based on algorithmic probability*. Springer Science & Business Media, 2005.
- [40] T. Jaakkola, S. Singh, and M. Jordan, “Reinforcement learning algorithm for partially observable Markov decision problems,” in *Advances in Neural Information Processing Systems*, G. Tesauro, D. Touretzky, and T. Leen, Eds., vol. 7. MIT Press, 1994, pp. 345–352.
- [41] N. Jiang, A. Kulesza, and S. P. Singh, “Improving predictive state representations via gradient descent,” in *AAAI Conference on Artificial Intelligence*, 2016, pp. 1709–1715.
- [42] S. Kapturovski, G. Ostrovski, W. Dabney, J. Quan, and R. Munos, “Recurrent experience replay in distributed reinforcement learning,” in *International Conference on Learning Representations*, 2019.

- [43] A. D. Kara, N. Saldi, and S. Yüksel, “Q-learning for MDPs with general spaces: Convergence and near optimality via quantization under weak continuity,” *Journal of Machine Learning Research*, 2023.
- [44] A. D. Kara and S. Yüksel, “Convergence of finite memory Q learning for POMDPs and near optimality of learned policies under filter stability,” *Mathematics of Operations Research*, Nov. 2022.
- [45] H. Kimura, K. Miyazaki, and S. Kobayashi, “Reinforcement learning in POMDPs with function approximation,” in *International Conference on Machine Learning*, 1997, pp. 152–160.
- [46] M. J. Kochenderfer, T. A. Wheeler, and K. H. Wray, *Algorithms for decision making*. MIT press, 2022.
- [47] V. Konda and J. Tsitsiklis, “Actor-critic algorithms,” in *Advances in neural information processing systems*, vol. 12, 1999.
- [48] V. R. Konda and J. N. Tsitsiklis, “On actor-critic algorithms,” *SIAM Journal on Control and Optimization*, vol. 42, no. 4, pp. 1143–1166, 2003.
- [49] V. Krishnamurthy, *Partially observed Markov decision processes*. Cambridge university press, 2016.
- [50] A. Kulesza, N. Jiang, and S. Singh, “Low-rank spectral learning with weighted loss functions,” in *Artificial Intelligence and Statistics*, 2015, pp. 517–525.
- [51] A. Kulesza, N. Jiang, and S. P. Singh, “Spectral learning of predictive state representations with insufficient statistics,” in *AAAI Conference on Artificial Intelligence*, 2015, pp. 2715–2721.
- [52] L. Li, T. J. Walsh, and M. L. Littman, “Towards a unified theory of state abstraction for MDPs,” in *ISAAC*, 2006.
- [53] M. L. Littman, “Algorithms for sequential decision-making,” Ph.D. dissertation, Brown University, 1996.
- [54] X. Lu, B. Van Roy, V. Dwaracherla, M. Ibrahim, I. Osband, Z. Wen *et al.*, “Reinforcement learning, bit by bit,” *Foundations and Trends in Machine Learning*, vol. 16, no. 6, pp. 733–865, 2023.
- [55] A. Mahajan, “Sequential decomposition of sequential dynamic teams: applications to real-time communication and networked control systems,” Ph.D. dissertation, University of Michigan, Ann Arbor, MI, 2008.
- [56] A. Mahajan, N. C. Martins, M. C. Rotkowitz, and S. Yüksel, “Information structures in optimal decentralized control,” in *IEEE Conference on Decision and Control*. IEEE, Dec. 2012.
- [57] N. Meuleau, L. Peshkin, K.-E. Kim, and L. P. Kaelbling, “Learning finite-state controllers for partially observable environments,” in *Conference on Uncertainty in Artificial Intelligence*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1999, p. 427–436.
- [58] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [59] A. Müller, “How does the value function of a Markov decision process depend on the transition probabilities?” *Math. Oper. Res.*, vol. 22, no. 4, pp. 872–885, 1997.
- [60] —, “Integral probability metrics and their generating classes of functions,” *Advances in Applied Probability*, vol. 29, no. 2, pp. 429–443, 1997.
- [61] A. Y. Ng and M. I. Jordan, “PEGASUS: A policy search method for large mdps and pomdps,” in *Conference on Uncertainty in Artificial Intelligence*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2000, p. 406–415.
- [62] T. Ni, B. Eysenbach, E. SeyedSalehi, M. Ma, C. Gehring, A. Mahajan, and P.-L. Bacon, “Bridging state and history representations: Understanding self-predictive RL,” in *International Conference on Learning Representations*, 2024.
- [63] T. Ni, M. Ma, B. Eysenbach, and P.-L. Bacon, “When do transformers shine in RL? decoupling memory from credit assignment,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: <https://openreview.net/forum?id=APGXBNkt6h>
- [64] T. J. Perkins and M. D. Pendrith, “On the existence of fixed points for q-learning and sarsa in partially observable domains,” in *International Conference on Machine Learning*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2002, p. 490–497.
- [65] L. Peshkin, N. Meuleau, and L. P. Kaelbling, “Learning policies with external memory,” in *International Conference on Machine Learning*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1999, p. 307–314.
- [66] J. Pineau, G. Gordon, S. Thrun *et al.*, “Point-based value iteration: An anytime algorithm for POMDPs,” in *International Joint Conference on Artificial Intelligence*, vol. 3, 2003, pp. 1025–1032.
- [67] L. K. Platzman, “Finite memory estimation and control of finite probabilistic systems.” Ph.D. dissertation, Massachusetts Institute of Technology, 1977.
- [68] M. Rosencrantz, G. Gordon, and S. Thrun, “Learning low dimensional predictive representations,” in *International Conference on Machine Learning*, 2004.
- [69] J. Sandell, Nils R., “Control of finite-state, finite-memory stochastic systems,” Ph.D. dissertation, Massachusetts Institute of Technology, Cambridge, MA, 1974.
- [70] J. Schmidhuber, *Making the world differentiable: on using self supervised fully recurrent neural networks for dynamic reinforcement learning and planning in non-stationary environments*. Inst. für Informatik, 1990, vol. 126.
- [71] E. SeyedSalehi, N. Akbarzadeh, A. Sinha, and A. Mahajan, “pproximate information state based convergence analysis of recurrent Q-learning,” in *European conference on reinforcement learning*, 2023. [Online]. Available: <https://arxiv.org/abs/2306.05991>
- [72] S. Singh, T. Jaakkola, and M. Jordan, “Reinforcement learning with soft state aggregation,” in *Advances in Neural Information Processing Systems*, G. Tesauro, D. Touretzky, and T. Leen, Eds., vol. 7. MIT Press, 1994, pp. 361–368.
- [73] S. P. Singh, T. Jaakkola, and M. I. Jordan, “Learning without state-estimation in partially observable markovian decision processes,” in *Machine Learning*. Elsevier, 1994, pp. 284–292.
- [74] A. Sinha, M. Geist, and A. Mahajan, “Title withheld due to double blind review restrictions,” in *under review*, 2024.
- [75] A. Sinha and A. Mahajan, “Asymmetric actor-critic with approximate information state,” in *IEEE Conference on Decision and Control*, 2023, pp. 7810–7816.
- [76] R. D. Smallwood and E. J. Sondik, “The optimal control of partially observable Markov processes over a finite horizon,” *Operations Research*, vol. 21, no. 5, pp. 1071–1088, Oct. 1973.
- [77] T. Smith and R. Simmons, “Heuristic search value iteration for POMDPs,” in *Conference on Uncertainty in Artificial Intelligence*, Arlington, Virginia, USA, 2004, p. 520–527.
- [78] E. J. Sondik, “The optimal control of partially observable markov processes over the infinite horizon: Discounted costs,” *Operations research*, vol. 26, no. 2, pp. 282–304, 1978.
- [79] M. T. Spaan and N. Vlassis, “Perseus: Randomized point-based value iteration for POMDPs,” *Journal of Artificial Intelligence Research*, vol. 24, pp. 195–220, 2005.
- [80] J. Subramanian, A. Sinha, R. Seraj, and A. Mahajan, “Approximate information state for approximate planning and reinforcement learning in partially observed systems,” *Journal of Machine Learning Research*, vol. 23, no. 12, pp. 1–83, 2022.
- [81] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, “Policy gradient methods for reinforcement learning with function approximation,” in *Advances in Neural Information Processing Systems*, vol. 12, 1999.
- [82] J. N. Tsitsiklis and B. Van Roy, “Feature-based methods for large scale dynamic programming,” *Machine Learning*, vol. 22, no. 1, pp. 59–94, 1996.
- [83] L. N. Vaserstein, “Markov processes over denumerable products of spaces, describing large systems of automata,” *Problemy Peredachi Informatsii*, vol. 5, no. 3, pp. 64–72, 1969.
- [84] C. Villani *et al.*, *Optimal transport: old and new*. Springer, 2008, vol. 338.
- [85] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, pp. 279–292, 1992.
- [86] L. Weng, “Policy gradient algorithms,” 2018. [Online]. Available: <https://lilianweng.github.io/posts/2018-04-08-policy-gradient/>
- [87] P. Werbos, “Beyond regression: New tools for prediction and analysis in the behavioral sciences,” *PhD thesis, Committee on Applied Mathematics, Harvard University, Cambridge, MA*, 1974.
- [88] C. C. White III and W. T. Scherer, “Finite-memory suboptimal design for partially observed markov decision processes,” *Operations Research*, vol. 42, no. 3, pp. 439–455, 1994.
- [89] H. S. Witsenhausen, “A standard form for sequential stochastic control,” *Mathematical Systems Theory*, vol. 7, no. 1, pp. 5–11, Mar. 1973.
- [90] —, “On policy independence of conditional expectations,” *Information and Control*, vol. 28, no. 1, pp. 65–75, 1975.
- [91] L. Yang, K. Zhang, A. Amice, Y. Li, and R. Tedrake, “Discrete approximate information states in partially observable environments,” in *American Control Conference*. IEEE, 2022, pp. 1406–1413.

- [92] H. Yu, "Approximate solution methods for partially observable Markov and semi-Markov decision processes," Ph.D. dissertation, Massachusetts Institute of Technology, 2006.
- [93] N. Zhang and W. Liu, "Planning in stochastic domains: Problem characteristics and approximation," Hong Kong Univeristy of Science and Technology, Tech. Rep. HKUST-CS96-31, 1996.