

# Least Inferable Policies for Markov Decision Processes

Mustafa O. Karabag, Melkior Ornik, Ufuk Topcu

## Abstract

In a variety of applications, an agent's success depends on the knowledge that an adversarial observer has or can gather about the agent's decisions. It is therefore desirable for the agent to achieve a task while reducing the ability of an observer to infer the agent's policy. We consider the task of the agent as a reachability problem in a Markov decision process and study the synthesis of policies that minimize the observer's ability to infer the transition probabilities of the agent between the states of the Markov decision process. We introduce a metric that is based on the Fisher information as a proxy for the information leaked to the observer and using this metric formulate a problem that minimizes expected total information subject to the reachability constraint. We proceed to solve the problem using convex optimization methods. To verify the proposed method, we analyze the relationship between the expected total information and the estimation error of the observer, and show that, for a particular class of Markov decision processes, these two values are inversely proportional.

## I. INTRODUCTION

We consider a setting in which an agent is supposed to accomplish a task in a stochastic environment while an observer that is potentially adversarial tries to infer the characteristics of the agent's behavior. For a scenario where predictable behaviors may put the success of the task at risk, it is crucial for an agent to conceal its strategy. In this paper, we study the synthesis of policies that enable an agent to achieve its task while limiting the ability of the observer to infer.

We model the behavior of the agent by a Markov decision process (MDP). The agent follows a policy to achieve its objective, for example, reaching a set of target states with high probability. This policy determines the transition probabilities of the agent between the states of the MDP. The observer can observe the transitions of the agent at a subset of the states and, solely based on the observed transitions, infers the agent's transition probabilities at these states. As a counter objective, the agent aims to choose its policy such that it limits the ability of the observer to infer the transition probabilities in addition to achieving the agent's task with high probability.

A policy can limit the ability of the observer to infer by minimizing the amount of information on the transition probabilities that the observer can gather from each observed transition. We introduce a metric, *transition information*, to measure the amount of information that a single transition leaks to the observer. This metric is related to the Fisher information which measures the amount of information that a random variable has on a parameter [1]. An observer that is trying to estimate the parameter would have high expected estimation error if the random variable has low Fisher information on the parameter. The notion of transition information generalizes the Fisher information by providing a scalar value describing the information leaked for the agent's transition that is parametrized by the transition probabilities.

While the notion of transition information is appropriate for a single observed transition, we also need to consider the effect of the number of observed transitions on the ability of the observer to infer. A policy that solely minimizes the transition information for each observed state adjusts the transition probabilities to the successor states as close to each other as possible since the uniform distribution of the successor states minimizes the transition information for a state. However, this approach might increase the number of observed transitions, and the observer may be able to infer the transition probabilities due to high number of observed transitions. Hence a policy that minimizes the ability of the observer to infer the transition probabilities must also take into account the number of observed transitions and balance the number of observed transitions and the transition information of each observed transition.

We account for the two quantities of interest, the number of observed transitions and the transition information of each observed transition, through a unified notion of *expected total information* — the expected sum of transition informations over a path generated by the agent's policy. We propose to compute a policy that has the minimum expected total information subject to the constraint that the task of the agent is completed with high probability.

To the best of our knowledge, the proposed method is the first policy synthesis method that uses the Fisher information for planning in MDPs against an adversary. The method introduced in [2] uses the Fisher information for learning and control in unknown systems that are modeled by MDPs. However, in contrast to our approach, [2] aims to increase the information gathered from transitions. A-optimality criterion [3] for experiment design aims to minimize the total variance of estimators by minimizing the trace of the inverse Fisher information matrix. The transition information is the reciprocal of the trace of inverse Fisher information matrix. In contrast, by minimizing the transition information, we aim to maximize

M. O. Karabag is with the Department of Electrical and Computer Engineering, University of Texas at Austin. e-mail: karabag@utexas.edu

M. Ornik is with the Institute for Computational Engineering and Sciences, University of Texas at Austin. e-mail: mornik@ices.utexas.edu

U. Topcu is with the Department of Aerospace Engineering and Engineering Mechanics and the Institute for Computational Engineering and Sciences, University of Texas at Austin. e-mail: utopcu@utexas.edu

the total variance of estimators unlike A-optimality criterion. In terms of the use of Fisher information, the closest works to the method proposed in this paper are [4] and [5]. The methods introduced in [4] and [5] use the Fisher information to preserve privacy for database systems and smart meters, respectively, and they do not deal with MDPs. Planning in stochastic control settings in the presence of an adversary has been substantially explored previously; the works closest to our paper are [6]–[8]. The reference [6] provides a method for multi-agent perimeter patrolling scenarios and is not applicable to MDPs in general. Papers [7], [8] propose to randomize the policy of an agent by maximizing the entropy of an induced stochastic process. While, for an MDP, increasing the entropy of a process increases randomness of the paths, it does not necessarily limit the ability of an observer to infer the transition probabilities.

The rest of the paper is organized as follows. Section II provides necessary background on the proposed method. In Section III, the definition of information and the problem formulation are presented. Section IV includes the methodology to synthesize the policy that has minimum expected total information subject to a reachability constraint by convex optimization problems. In Section V, we show the relationship between considered problems and estimation errors of the observer. We present numerical examples in Section VI and conclude with suggestions for the future work in Section VII. We discuss some special cases of the proposed method in Appendix A and give the proofs for the technical results of this paper in Appendix B.

## II. PRELIMINARIES

In this section, we present some of the concepts and notation used in the rest of the paper.

We use  $[n]$  for the set  $\{1, \dots, n\}$ . For a finite set  $D$ , we denote the power set with  $2^D$  and cardinality with  $|D|$ .  $\mathbb{E}[\Theta]$  denotes the expectation of the random variable  $\Theta$  and  $\text{Var}(\Theta)$  denotes the variance of  $\Theta$  which is  $\mathbb{E}[(\Theta - \mathbb{E}[\Theta])(\Theta - \mathbb{E}[\Theta])^T]$ . We use  $\mathbb{1}_D$  for the indicator function of a set  $D$  where  $\mathbb{1}_D(x) = 1$  if  $x \in D$  and  $\mathbb{1}_D(x) = 0$  otherwise.

### A. Markov Decision Processes

A *Markov decision process* is a tuple  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$  where  $S$  is a finite set of states,  $\mathcal{A}$  is a finite set of actions,  $\mathcal{P} : S \times \mathcal{A} \times S \rightarrow [0, 1]$  is the transition probability function, and  $s_0$  is the initial state. We denote  $\mathcal{P}(s, a, q)$  by  $\mathcal{P}_{s,a,q}$ .  $\mathcal{A}(s)$  denotes the set of available actions at state  $s$  where  $\sum_{q \in S} \mathcal{P}_{s,a,q} = 1$  for all  $a \in \mathcal{A}(s)$ . We denote the successor states of state  $s$  by  $\text{Succ}(s)$  such that a state  $q \in S$  if and only if there exists an action  $a$  such that  $\mathcal{P}_{s,a,q} > 0$ . A state  $s$  is *absorbing* if it has only a single successor state that is itself, i.e.,  $\mathcal{P}_{s,a,s} = 1$  for all  $a \in \mathcal{A}(s)$ .

A *sub-MDP*  $(C, D)$  of  $\mathcal{M}$  is a pair where  $C \subseteq S$  is non-empty and  $D : C \rightarrow 2^{\mathcal{A}}$  is a function such that  $a \in D(s)$  only if  $\mathcal{P}_{s,a,q} = 0$  for all  $q \notin C$ . An *end component* is a sub-MDP  $(C, D)$  of  $\mathcal{M}$  such that the digraph induced by  $(C, D)$  is strongly connected. An end component  $(C, D)$  is *closed* if, for all  $s \in C$ ,  $\text{Succ}(s) \setminus C = \emptyset$ . A *maximal end component*  $(C, D)$  is an end component where there is no end component  $(C', D')$  such that  $(C, D) \neq (C', D')$ ,  $C \subseteq C'$ , and  $D \subseteq D'$ .

A *policy* is a sequence  $\pi = [\mu_0, \mu_1, \dots]$  where each  $\mu_t : S \times \mathcal{A} \rightarrow [0, 1]$  is a function such that  $\sum_{a \in \mathcal{A}(s)} \mu_t(s, a) = 1$  for every  $s \in S$ . A *stationary policy*  $\pi$  is a sequence  $\pi = [\mu, \mu, \dots]$  where  $\mu : S \times \mathcal{A} \rightarrow [0, 1]$  is a function such that  $\sum_{a \in \mathcal{A}(s)} \mu(s, a) = 1$  for every  $s \in S$ . We denote the set of all policies by  $\Pi(\mathcal{M})$  and the set of stationary policies by  $\Pi^{St}(\mathcal{M})$ . For a stationary policy  $\pi$ , we denote  $\mu(s, a)$  by  $\pi_{s,a}$ . A stationary policy  $\pi$  induces a *Markov chain*  $\mathcal{M}^\pi = (S, \mathcal{P}^\pi, s_0)$  from  $\mathcal{M}$  where  $\mathcal{P}^\pi : S \times S \rightarrow [0, 1]$  is the transition probability function such that

$$\mathcal{P}^\pi(s, q) := \sum_{a \in \mathcal{A}(s)} \pi_{s,a} \mathcal{P}_{s,a,q}$$

for all  $s, q \in S$ . We denote  $\mathcal{P}^\pi(s, q)$  by  $\mathcal{P}_{s,q}^\pi$ .

A *path*  $\xi = s_0 s_1 s_2 \dots$  is an infinite sequence of states under policy  $\pi = [\mu_0, \mu_1, \dots]$  such that  $\sum_{a \in \mathcal{A}(s_t)} \mathcal{P}_{s_t,a,s_{t+1}} \mu_t(s_t, a) > 0$  for all  $t \geq 0$ . The set of paths for  $\mathcal{M}$  under policy  $\pi$  is denoted by  $\text{Paths}(\mathcal{M}^\pi)$ .

The reachability probability to the set  $B$  of states, i.e., the probability of reaching a state  $b \in B$  under policy  $\pi$ , is denoted by  $\Pr^\pi(\text{Reach}[B])$ .

The *expected state residence time* at state  $s$  is defined by

$$x_s^\pi := \sum_{t=0}^{\infty} \Pr(s_t = s | s_0),$$

where  $s_t$  is the state at time  $t$ . The expected state residence time  $x_s^\pi$  is also equal to  $\mathbb{E}[N_{s,\xi}]$  where  $N_{s,\xi}$  is the number of appearances of  $s$  in the random path  $\xi$  that is generated by the policy  $\pi$ . The *expected state-action residence time* at state  $s$  and action  $a$  is defined by

$$x_{s,a}^\pi := \sum_{t=0}^{\infty} \Pr(s_t = s | s_0) \mu_t(s_t, a).$$

The expected state-action residence time of a state and an action is the expected number of times that the action is taken at the state. For a stationary policy  $\pi \in \Pi^{St}(\mathcal{M})$ ,  $x_{s,a}^\pi = \pi_{s,a} x_s^\pi$ .

### B. The Fisher Information and the Cramér-Rao Bound

Let the random variable  $X$  represent the observed data from a random variable that is parametrized by  $\Theta \in \mathbb{R}^n$ . An *estimator* is a function  $\hat{\Theta} : X \rightarrow \mathbb{R}^n$  that estimates  $\Theta$  based on observed data. The estimator  $\hat{\Theta}$  is an *unbiased estimator* of  $\Theta$  if  $\mathbb{E}[\hat{\Theta}] - \Theta = 0$ .

The *precision* of a random variable is the reciprocal of the variance of the random variable. For an unbiased estimator, its precision is the reciprocal of the mean squared error (MSE) of the estimator.

The *Fisher information* [9]  $I_X(\theta)$  of a discrete random variable  $X$  parametrized by  $\theta \in \mathbb{R}$  is

$$I_X(\theta) := - \sum_{x \in \text{Supp}(X)} \frac{\partial^2 \log(\Pr(X = x|\theta))}{\partial \theta^2} \Pr(X = x|\theta).$$

An important property of the Fisher information is additivity, that is, when the samples are drawn from i.i.d. random variables, the Fisher information based on  $n$  samples  $I_{X^n}(\theta)$  satisfies  $I_{X^n}(\theta) = nI_X(\theta)$  where  $I_X(\theta)$  is the Fisher information of one sample.

The *Cramér-Rao inequality* [9] defines a relationship between the variance of an unbiased estimator of parameter  $\theta$  and the Fisher information on the parameter  $\theta$ . The inequality is stated as

$$\text{Var}(\hat{\theta}) \geq I_X(\theta)^{-1} \quad (1)$$

where  $\hat{\theta}$  is any unbiased estimator of  $\theta$ .

An unbiased estimator is *efficient* if it achieves the Cramér-Rao bound.

### III. PROBLEM STATEMENT

Consider an agent whose behavior is governed by a Markov decision process (MDP)  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$  where a stationary policy followed by the agent  $\pi$  implemented on this MDP induces a Markov chain. An adversary which we call *observer* observes the transitions and tries to infer the transition probabilities for a set  $W$  of states in the induced Markov chain. We assume that the observer can only observe the transitions at the states in  $W$  which we call *observed states*, and has no side information.

The problem we study is the synthesis of a policy for the agent with two objectives: (i) reach a set  $C_{reach}$  of states with probability higher than a given threshold  $0 \leq \nu_{reach} \leq 1$  and (ii) minimize the amount of information leaked to the observer.

For the first objective, we assume that the transitions of the agent after reaching  $C_{reach}$  are irrelevant, i.e., every state in  $C_{reach}$  is absorbing and is not observed.

For the second objective, we define the notion of *transition information* to measure the amount of information leaked to the observer due to a transition.

**Definition 1.** The *transition information* of a state  $s$  is defined by

$$\iota_s^\pi := \frac{1}{\sum_{q \in \text{Succ}(s)} I_Q(\mathcal{P}_{s,q}^\pi)^{-1}} \quad (2)$$

where  $Q$  is the random variable that is the successor state of state  $s$ .

We remark that the Fisher information and the transition information are analogous:

- The reciprocal of the Fisher information is a lower bound on the variance of an unbiased estimator for a single parameter.
- The reciprocal of the transition information is a lower bound on the variance of an unbiased estimator for many parameters.

For a state  $s$ , consider an unbiased estimator  $\hat{\mathcal{P}}_s$  of transition probabilities. The reciprocal of the transition information  $\iota_s^\pi$  is a lower bound on the variance of  $\hat{\mathcal{P}}_s$ :

$$\text{Var}(\hat{\mathcal{P}}_s) \geq \frac{1}{\iota_s^\pi}.$$

We use the transition information to define the *total information* of a path. The total information of a path  $\xi = s_0 s_1 s_2 \dots$  is defined as the sum of each observed transition's transition information such that

$$\iota_{W,\xi}^\pi := \sum_{t=0}^{\infty} \mathbb{1}_W(s_t) \iota_{s_t}^\pi.$$

We then state the synthesis problem formally as follows:

**Problem 1** (Synthesis of Minimum-Information Admissible Policies). *Given an MDP  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$ , a set  $C_{reach}$  of states, a probability threshold  $\nu_{reach}$ , and the set  $W$  of observed states, compute*

$$\inf_{\pi \in \Pi^{St}(\mathcal{M})} \mathbb{E}[\ell_{W, \xi}^{\pi}], \quad (3a)$$

$$s. t. \quad \Pr^{\pi}(\text{Reach}[C_{reach}]) \geq \nu_{reach} \quad (3b)$$

where  $\xi$  is a random path generated under policy  $\pi$ . If the optimal value is attainable, compute the optimal policy  $\pi^*$ .

Hereafter we call the policies that satisfy the reachability constraint *admissible policies* and an optimal policy for Problem 1 a *minimum-information admissible policy*.

**Example 1.** We explain the characteristics of a minimum-information admissible policy through the MDP given in Figure 1 with  $\nu_{reach} = 0$  for simplicity.

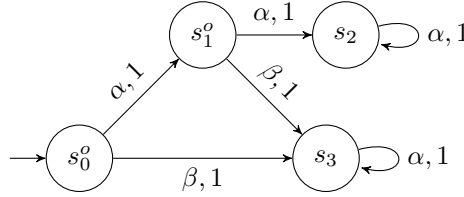


Figure 1: An MDP with 4 states. A label  $a, p$  of a transition refers to the transition that happens with probability  $p$  when action  $a$  is taken. The states marked with the superscript  $o$  are observed.

The goal of the agent is to find a policy that minimizes the expected total information. Consider the policy at state  $s_1$  and note that the policy decision at  $s_1$  does not affect information leaked from state  $s_0$  since it does not change the expected residence time at  $s_0$ . Hence we may only consider the information leaked from  $s_1$ . If the agent chooses a deterministic policy, the observer can estimate the transition probabilities with no error even after observing a single transition, which means infinite leaked information. Therefore, it is expected that the agent randomizes the transition probabilities. Formally, we explain the reasoning by the fact that the Fisher information is minimized for a  $Ber(p)$  random variable with  $p = 0.5$ . Similarly at state  $s_0$ , the agent randomizes the transition probabilities. However, unlike  $s_1$ , the policy at  $s_0$  affects the information leaked from  $s_1$ . As the agent decreases the probability of taking action  $\alpha$  at state  $s_0$ , the expected number of visits to state  $s_1$  decreases and consequently information leaked from  $s_1$  decreases. Hence, the agent must take the action  $\beta$  with a greater probability than the action  $\alpha$ . On the other hand, taking the action  $\beta$  with high probability increases the information leaked from  $s_0$ . We expect that, under this trade-off, the agent must choose a policy that takes both actions, but the action  $\beta$  more likely. Numerically, the optimal policy is  $\pi_{s_0, \alpha} = 0.38$ ,  $\pi_{s_0, \beta} = 0.62$ ,  $\pi_{s_1, \alpha} = 0.5$ , and  $\pi_{s_1, \beta} = 0.5$ .

**Remark 1.** Note that if the transition probabilities are not constant and change between observations, measurement of inference with a transition information is not meaningful since we assume underlying probability distribution is constant. To have a well-defined problem, we only focus on agents that have to follow stationary policies and we search the optimal policies only in the stationary policies.

#### IV. SYNTHESIS OF MINIMUM-INFORMATION ADMISSIBLE POLICIES

For an MDP  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$ , we aim to find a minimum-information admissible policy  $\pi$  that minimizes the expected total information of a path subject to the reachability constraint  $\Pr^{\pi}(\text{Reach}[C_{reach}]) \geq \nu_{reach}$  where the set of observed states is  $W$ . In this section, we represent the expected total transition information in terms of expected state-action residence times, show the existence of a minimum-information admissible policy, and give an optimization problem whose solution is a minimum-information admissible policy. We also show that the proposed optimization problem is convex in the expected state-action residence time parameters and hence can be solved using off-the-shelf convex optimization tools.

Note that the Fisher information for a parameter is well-defined if the regularity conditions are satisfied. These conditions require that the distributions depending on the parameter have a common support that is independent of the parameter [9]. For a random variable  $P \sim Ber(p)$ , the Fisher information is not defined when  $p = 0$  or  $p = 1$  since the probability distribution of  $P$  does not have a common support. However, such a case practically corresponds to infinite Fisher information, which means that the value of the parameter is estimated exactly even after a single observation. We assume that the Cramér-Rao lower bound is zero if the Fisher information is infinite.

Consider a state  $w \in W$  whose successor state is denoted by the random variable  $Q$ . For each  $q \in \text{Succ}(w)$ , we have

$$I_Q(\mathcal{P}_{w,q}^{\pi}) = I_{\mathbb{1}_q(Q)}(\mathcal{P}_{w,q}^{\pi}) = \frac{1}{\mathcal{P}_{w,q}^{\pi}(1 - \mathcal{P}_{w,q}^{\pi})}$$

where  $\mathbb{1}_q(Q)$  is a  $Ber(\mathcal{P}_{w,q}^\pi)$  random variable. The transition information of a state  $w$  is a function

$$\iota_w : \{\mathcal{P}_w \in \mathbb{R}^{|Succ(w)|} : \sum_{q \in Succ(w)} \mathcal{P}_{w,q} = 1, \mathcal{P}_{w,q} \geq 0\} \rightarrow \mathbb{R} \cup \{\infty\}$$

and, under policy  $\pi$ , is equal to

$$\iota_w^\pi = \left( \sum_{q \in Succ(w)} \mathcal{P}_{w,q}^\pi (1 - \mathcal{P}_{w,q}^\pi) \right)^{-1}. \quad (4a)$$

**Remark 2.** The categorical random variable  $Q$  has the distribution  $\mathcal{P}_{w,q}^\pi$  where  $q \in Succ(w)$ . The covariance matrix  $\Sigma$  of  $Q$  has diagonal entries  $\mathcal{P}_{w,q}^\pi (1 - \mathcal{P}_{w,q}^\pi)$ . The transition information of state  $w$  given in (4a) is also equal to  $\text{tr}(\Sigma)^{-1}$ . Since  $Q$  is a categorical random variable, a sample mean estimator achieves the Cramér-Rao bound for a single transition. However, since the observed data consists of transitions from a path and the transitions are not independent in general, a sample mean estimator is not necessarily unbiased and efficient.

We now construct the optimization problem whose solution gives the expected state-action residence times for a minimum-information admissible policy. First, we rewrite (4a) as

$$\iota_w^\pi = \left( \sum_{q \in Succ(w)} \left( \sum_{a \in \mathcal{A}(w)} \frac{x_{w,a}^\pi}{\sum_{a' \in \mathcal{A}(w)} x_{w,a'}^\pi} \mathcal{P}_{w,a,q} \right) \left( 1 - \sum_{a' \in \mathcal{A}(w)} \frac{x_{w,a'}^\pi}{\sum_{a' \in \mathcal{A}(w)} x_{w,a'}^\pi} \mathcal{P}_{w,a,q} \right) \right)^{-1} \quad (5)$$

using the definitions of the induced Markov chain and expected state-action residence times.

We assume that the optimal value of Problem 1 is finite. If the optimal value is infinite any admissible policy is a minimum-information admissible policy.

**Proposition 1.** For an MDP  $\mathcal{M}$ , if  $\mathbb{E}[\iota_{W,\xi}^\pi]$  is finite where  $\xi$  is a path generated randomly under a policy  $\pi \in \Pi^{St}(\mathcal{M})$ , then

$$\mathbb{E}[\iota_{W,\xi}^\pi] = \sum_{w \in W} x_w^\pi \iota_w^\pi.$$

Note that the expected total information  $x_w^\pi \iota_w^\pi$  of a state  $w$  has some undefined points on the domain  $x_{w,a}^\pi \geq 0$  where  $a \in \mathcal{A}(w)$ . We define the function at such points as follows:

- If the expected state residence time is zero, i.e.,  $x_w^\pi = \sum_{a \in \mathcal{A}(w)} x_{w,a}^\pi = 0$ , then  $x_w^\pi \iota_w^\pi := 0$ . Since the state will never be visited, the observer cannot get information on the transition probabilities.
- If  $w$  deterministically transitions to one of the successor states and expected residence time is greater than zero, i.e., there exists a state  $q \in Succ(w)$  such that  $\sum_{a \in \mathcal{A}(w)} x_{w,a}^\pi \mathcal{P}_{w,a,q} > 0$  and  $\sum_{a \in \mathcal{A}(w)} x_{w,a}^\pi \mathcal{P}_{w,a,q'} = 0$  for all  $q' \in Succ(w) \setminus q$ , then  $x_w^\pi \iota_w^\pi := \infty$ . Since the observer can estimate the transition probabilities even after a single observation and there is a positive probability that the state will be visited, the expected total information is infinite.
- If the expected state residence time at  $w$  is infinite, i.e.,  $x_w^\pi = \sum_{a \in \mathcal{A}(w)} x_{w,a}^\pi = \infty$ , then  $x_w^\pi \iota_w^\pi := \infty$ . Since the observed distribution of transitions converges to the transition probabilities, the expected total information is infinite.

We represent the stationary policies of the agent with a set of constraints which use the expected state-action residence times. A stationary policy makes each state either recurrent or transient. We need to identify the states that can be reachable and recurrent. If a policy leaks finite information, a set of states can be reachable and recurrent if and only if they belong to an end component and are not observed since the recurrence of a reachable observed state results in infinite expected total information.

**Definition 2.** An *unobserved end component* (UEC) is a sub-MDP  $(C, D)$  such that the digraph induced by  $(C, D)$  is strongly connected and  $C \cap W = \emptyset$ . An *unobserved maximal end component* (UMEC)  $(C, D)$  is a UEC where  $C \subset S$  and there is no UEC  $(C', D')$  such that  $(C, D) \neq (C', D')$ ,  $C \subseteq C'$ , and  $D \subseteq D'$ .

We denote the set of states that belong to some UMEC by  $C_{end}$ . After reaching  $C_{end}$ , the agent can follow a stationary policy that always stays in the UMEC and leaks no more information. For example,  $s_2$  is a UMEC state in Figure 2. However, due to the reachability constraints the agent might need to follow a policy that leaves a UMEC. We disallow such cases and make the following assumption to ensure that the agent does not leave UMECs.

**Assumption 1.** All unobserved maximal end components are closed.

**Remark 3.** In the absence of Assumption 1, to find the optimal stationary policy, one needs to check every UEC to determine whether the agents needs to stay or leave the UEC. Such a check increases computational complexity of finding a minimum-information admissible policy. For clarity of presentation, we here adopt Assumption 1. In Appendix A, we investigate the more general problem without Assumption 1.

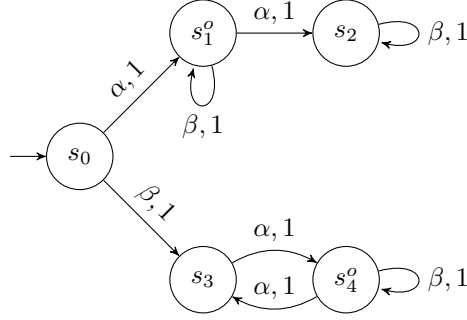


Figure 2: An MDP with 5 states. A label  $a, p$  of a transition refers to the transition that happens with probability  $p$  when action  $a$  is taken. The states marked with the superscript  $o$  are observed.

The optimal value of Problem 1 is

$$\inf \sum_{w \in W} x_w^\pi l_w^\pi \quad (6a)$$

$$\text{s. t. } x_s^\pi = \sum_{a \in \mathcal{A}(s)} x_{s,a}^\pi, \quad \forall s \in S \setminus C_{\text{end}} \quad (6b)$$

$$x_{s,a}^\pi \geq 0, \quad \forall s \in S \setminus C_{\text{end}}, \forall a \in \mathcal{A}(s) \quad (6c)$$

$$\sum_{a \in \mathcal{A}(s)} x_{s,a}^\pi - \sum_{q \in S} \sum_{a \in \mathcal{A}(q)} x_{q,a}^\pi \mathcal{P}_{q,a,s} = \mathbb{1}_{s_0}(s), \quad \forall s \in S \setminus C_{\text{end}} \quad (6d)$$

$$\sum_{q \in C_{\text{reach}}} \sum_{s \in S \setminus C_{\text{end}}} \sum_{a \in \mathcal{A}(s)} x_{s,a}^\pi \mathcal{P}_{s,a,q} + \mathbb{1}_{s_0}(q) \geq \nu_{\text{reach}}, \quad (6e)$$

where the decision variables are  $x_{s,a}^\pi$  for all  $s \in S \setminus C_{\text{end}}$  and  $a \in \mathcal{A}(s)$ . The objective function (6a) follows from Proposition 1 and the constraints (6b)-(6c) follow from definitions of expected residence times. The constraint (6d) is the flow equation indicating that the expected number of arrivals into a state, i.e., the inflow, is equal to the expected number of departures from the state, i.e., the outflow. These equations ensure that there exists a policy that gives the computed expected state-action residence times [10]. The reachability constraint in (3b) is equivalent to (6e).

Note that some stationary admissible policies are infeasible for the optimization problem given in (6). In detail, the stationary policies that eventually always stay in an end component and visit an observed state infinitely often are infeasible. For instance, consider a policy  $\pi$  such that  $\Pr^\pi(\text{Reach}[s_2]) = 0.5$  for the MDP given in Figure 2 with the reachability constraint  $\Pr(\text{Reach}[s_2] \geq 0.5)$ . While  $\pi$  leads to infinite expected total information and satisfies the reachability constraint, it is not feasible for the problem in (6). One can easily check the existence of a policy that satisfies the reachability constraint via model checking tools such as [11]. If there exists a policy that satisfies the task constraints, but the optimization problem given in (6) is infeasible, we can say that the minimum-information admissible policy leaks infinite information.

**Proposition 2.** *If there exists a policy  $\pi \in \Pi^{\text{St}}(\mathcal{M})$  that satisfies the reachability constraint given in (3b), then there exists a policy  $\pi^* \in \Pi^{\text{St}}(\mathcal{M})$  that attains the optimal value of the optimization problem given in (6).*

**Proposition 3.** *The optimization problem given in (6) is a convex optimization problem.*

**Remark 4.** *After a preprocessing step that has polynomial-time complexity in the size of  $\mathcal{M}$ , the optimization problem can be formulated as a conic optimization problem which can be solved using interior-point methods [12] in polynomial-time in the size of  $\mathcal{M}$ .*

After computing the optimal expected state-action residence times by the optimization problem in (6), a stationary, minimum-information admissible policy can be synthesized using the relationship  $x_{s,a}^\pi = \pi_{s,a} x_s^\pi$ .

## V. BOUNDS ON THE ESTIMATION ERROR

In this section, we consider estimators for the transition probabilities at the observed states and derive the bounds on the expected estimation error in terms of MSE. Define  $\sigma_w$  as the MSE of an unbiased estimator at a state  $w$ . We assume that, for the estimator at state  $w$ , the observed data are the whole path of the agent and the transition probabilities for the set  $S \setminus \{w\}$  of states are known.

**Proposition 4.** For an MDP  $\mathcal{M}$  and a policy  $\pi \in \Pi^{St}(\mathcal{M})$ ,

$$\sigma_w \geq \frac{\Pr^\pi(\text{Reach}[w])^2}{x_w^\pi \iota_w^\pi}$$

for every state  $w \in W$ .

**Corollary 5.** For an MDP  $\mathcal{M}$  and a policy  $\pi \in \Pi^{St}(\mathcal{M})$ , the total MSE  $\sum_{w \in W} \sigma_w$  satisfies

$$\sum_{w \in W} \sigma_w \geq \frac{\min_{w \in W} \Pr^\pi(\text{Reach}[w])^2 |W|^2}{\mathbb{E}[\iota_{W,\xi}^\pi]}.$$

Consequently, if  $\Pr^\pi(\text{Reach}[w]) = 1$  for every  $\pi \in \Pi^{St}(\mathcal{M})$  and for all  $w \in W$ , then  $\frac{|W|^2}{\mathbb{E}_\xi[\iota_{W,\xi}^\pi]}$  is a lower bound on the total MSE.

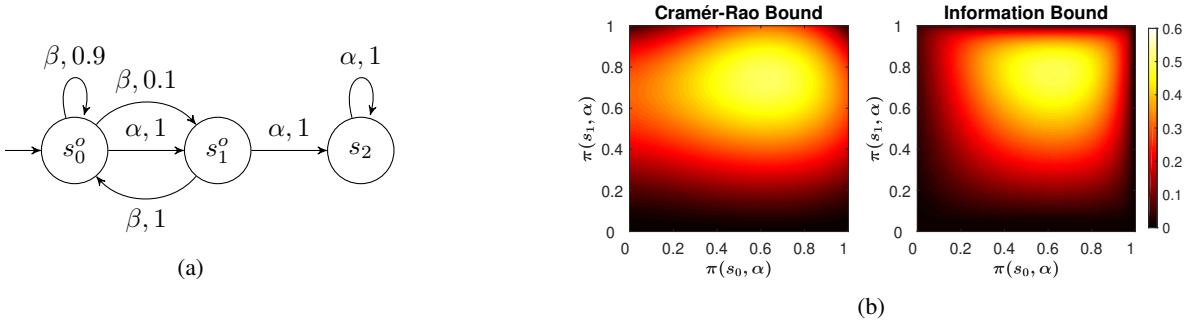


Figure 3: (a) An MDP with 3 states. A label  $a, p$  of a transition refers to the transition that happens with probability  $p$  when action  $a$  is taken. The states marked with the superscript  $o$  are observed. (b) The Cramér-Rao bound on the total MSE of the estimators and the error bound given in Corollary 5.

An example of the bound given in Corollary 5 is illustrated in Figure 3. Both of the observed states are visited under any stationary policy and the reciprocal of the expected total information is directly a lower bound on the total MSE of the estimators. One who wants to maximize the total MSE of the estimators may prefer to optimize over the expected total information instead of the Cramér-Rao bound since the Cramér-Rao is not a convex or concave function of the expected residence time parameters while the minimum-information admissible policy can be computed via a convex optimization problem.

## VI. NUMERICAL EXAMPLES

In this section, we illustrate the proposed method through two numerical examples. We solved the optimization problems using CVX toolbox [13] with MOSEK [14] on a computer with an Intel Core i7-8550u 1.8 GHz CPU and 8 GB of RAM.

### A. Partly Hidden Agent

In this example, we explain the characteristics of the minimum-information admissible policy through different scenarios.

The environment which is given in Figure 4 consists of 4 regions that are separated with walls and connected to each other with bridges. Each region is a  $20 \times 20$  grid world and each tile in these regions represents a state. Except for the reach state, the agent has 4 actions, namely, up, down, left, and right, at every state. When the agent takes an action the transition happens into the target direction with probability 0.8 and in the other directions uniformly randomly with probability 0.2. If a direction is out of the grid the transition probability to that direction is proportionally distributed to the other directions.

The initial state is the black top-left corner tile and the reach state is the green bottom-middle tile. The task of the agent is to reach the reach state with probability 1. While the agent is in the gray tiles, the observer cannot observe the transitions of the agent.

In the first scenario (see Figure 4a) all states are observed except the reach state and the bridge states. The agent completes the task with a low number of observed transitions (See Table I) with randomized transitions. Note that the randomization only happens between the states that are in the direction of the reach state since further randomization leads to more observations. When the unobserved regions are present in the environment (see Figure 4b), the policy generates paths that pass through the unobserved regions to reduce the number of observations. However, the unobserved regions are not always utilized. For example, in the top-right region if the agent is already away from the unobserved region, it directly goes to

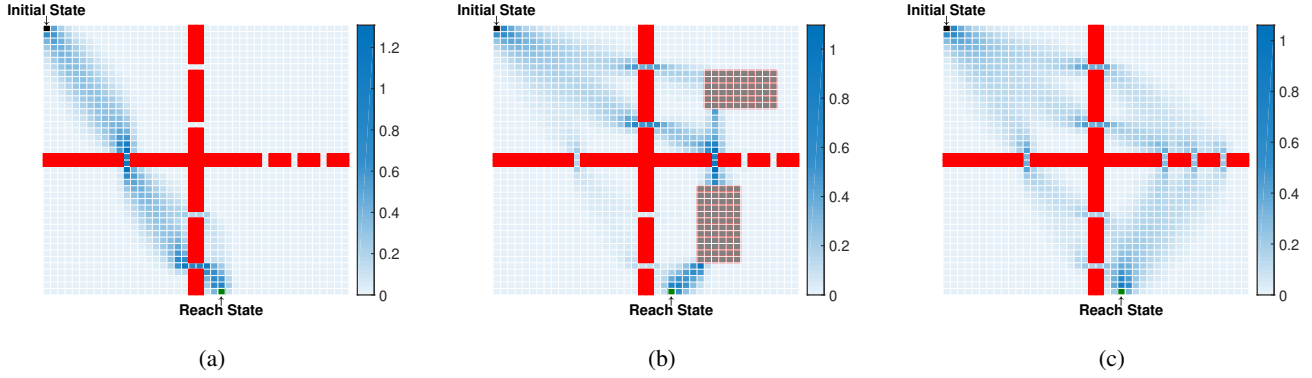


Figure 4: The heatmaps of expected state residence times for partly hidden agent example. For the scenario given in Figure 4b the environment has some unobserved regions while every state is observed for the scenario given in Figure 4b. The scenario given in Figure 4c considers exit information of in addition to the transition information.

the bottom-right region. Although, no information leaks in the unobserved regions, the agent leaks information during the process of reaching those states.

We remark that the minimum-information admissible policy minimizes only the information of transitions from the observed states. While this approach reduces the amount of leaked information in the local sense, i.e., the transitions between the states, the global behavior, i.e., the transitions between the regions, might be easily inferred. We observe such a phenomenon for the scenario given in Figure 4a; the agent leaves the regions using the same bridge. This behavior may be risky if there is an adversary that is interested in the information of which bridge is used. To avoid this behavior, we add a weighted penalty, *exit information*, for each region. The exit information of a region has the same form with the transition information and consists of the expected state residence times of the bridges. With the exit information (see Figure 4c) the agent randomizes its exit bridge from the regions compared to the initial case (see Figure 4a).

Table I: Numerical values for partly hidden agent example.

Scenario	Expected Total Information	Expected Number of Observations	Solving Time
Figure 4a	152.20	81.87	0.52
Figure 4b	146.00	78.89	0.38
Figure 4c	179.64	98.43	0.58

### B. Inference of Local Behavior

We explain the difference between the proposed method and the policy synthesis via entropy maximization through this example. The environment is a  $11 \times 11$  grid world given in Figure 5 where each tile represents a state. The black tile is the initial state, the green tile is the reach state, and the red tiles are the absorbing states. Except for the absorbing states and the reach state the agent can transition to 4 directions, namely, up, down, left, and right, at every state. When the agent takes an action, the transition happens in the target direction with probability 1. If a direction is out of the grid the action is not allowed. The task of the agent is to reach the reach state with probability 1.

We compare the policies in terms of their estimation error, which is calculated for different number of sample paths. The observer gets sample paths and estimates the transition probabilities at the observed states using a sample mean estimator. We measure the estimation error for a state by the mean squared error (MSE) between the observed and actual transition distributions at the observed states. The total error is the sum of MSE for each state. If there is no observation sample from a state, we set the MSE for that state. For the weighted MSE error, the weight of a state is ratio between the number of observations from the state and the total number of observations.

Maximizing the entropy of an MDP is equivalent to maximizing the entropy of the possible paths, and a high entropy value leads to unpredictable paths. Under the reachability constraint, the maximum entropy of the MDP given in Figure 5 is unbounded. For policy synthesis, we follow the procedure given in [8] and impose an upper bound on the expected total state residence time  $\Gamma$ . As the bound increases, the maximum entropy value of the MDP increases. We synthesize three policies that maximizes the entropy of MDP with different values for  $\Gamma = 15, 60$ , and  $120$ .

For low values of  $\Gamma$  such as 15, the minimum-information admissible and the maximum-entropy policies show similar behavior. However, for the high values of  $\Gamma$ , the difference between the minimum-information admissible policy and the maximum-entropy policy becomes clear. The minimum-information admissible policy completes the task with a low number



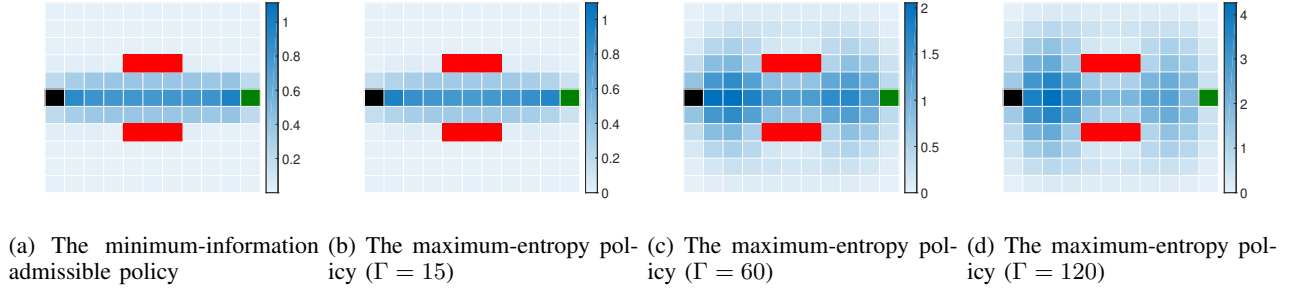


Figure 5: The expected state residence times for inference of local behavior example.

of non-informative observations. On the other hand, the maximum-entropy policy visits the observed states more to explore more paths and randomize the probabilities of paths. While the agent follows different paths, the expected residence times at the observed states increases and observer gets more samples. Although the policy is randomized and samples are less informative, transition probabilities are inferred due to the high number of observations. The result suggests that the unpredictability of the paths does not imply the limitation of inference for the transitions between states. Hence, the minimum-information admissible policy and the maximum-entropy policy serve different purposes.

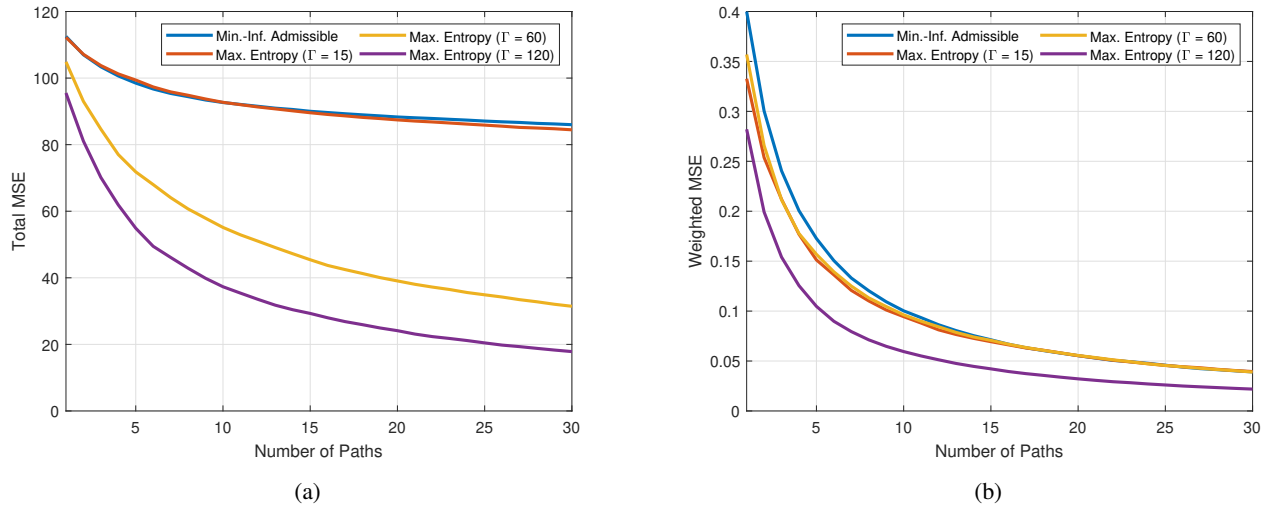


Figure 6: The expected estimation errors. The curves are averaged over 100 experiments.

## VII. CONCLUSION

We focus on policy synthesis for an agent whose behavior is inferred by an outside adversarial observer. Such an agent must as less informative observations as possible to the observer while completing its task. Based on this criterion, we introduced transition information which is based on the Fisher information and measures the amount of information leaked to the observer from a transition. Then, we formulated a problem that minimizes the expected total information leaked to the observer and showed the existence of such a policy. The significant feature of the proposed method is that it balances a possible trade-off between the number of observations and the informativeness of each observation.

The proposed method relies on the assumption that the agent follows a stationary policy on the observed states. A history dependent planning method may deceive the observer by actively changing the policy. We aim to remove this assumption and design an algorithm that takes the past transitions into account.

## ACKNOWLEDGEMENT

This work was supported in part by DARPA W911NF-16-1-0001.

## REFERENCES

- [1] B. R. Frieden. *Science from Fisher information: A Unification*. Cambridge University Press, 2004.
- [2] T. Alpcan and I. Shames. An information-based learning approach to dual control. *IEEE Transactions on Neural Networks and Learning Systems*, 26(11):2736–2748, 2015.

- [3] A. F. Emery and A. V. Nenarokomov. Optimal experiment design. *Measurement Science and Technology*, 9(6):864–876, 1998.
- [4] F. Farokhi and H. Sandberg. Optimal privacy-preserving policy using constrained additive noise to minimize the Fisher information. In *56th IEEE Conference on Decision and Control*, pages 2692–2697, 2017.
- [5] F. Farokhi and H. Sandberg. Fisher information as a measure of privacy: Preserving privacy of households with smart meters using batteries. *IEEE Transactions on Smart Grid*, 9(5):4726–4734, 2018.
- [6] N. Agmon, S. Kraus, and G. A. Kaminka. Multi-robot perimeter patrol in adversarial settings. In *IEEE International Conference on Robotics and Automation*, pages 2339–2345, 2008.
- [7] P. Paruchuri, M. Tambe, F. Ordóñez, and S. Kraus. Security in multiagent systems by policy randomization. In *Joint Conference on Autonomous Agents and Multiagent Systems*, pages 273–280, 2006.
- [8] Y. Savas, M. Ornik, M. Cubuktepe, and U. Topcu. Entropy maximization for Markov decision processes under temporal logic constraints. *arXiv preprint arXiv:1807.03223 [math.OC]*, 2018.
- [9] E. L. Lehmann and G. Casella. *Theory of Point Estimation*. Springer, 2nd edition, 1998.
- [10] K. Etessami, M. Kwiatkowska, M. Y. Vardi, and M. Yannakakis. Multi-objective model checking of Markov decision processes. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 50–65, 2007.
- [11] M. Kwiatkowska, G. Norman, and D. Parker. PRISM: Probabilistic symbolic model checker. In *International Conference on Modelling Techniques and Tools for Computer Performance Evaluation*, pages 200–204, 2002.
- [12] Yu. Nesterov and A. Nemirovskii. *Interior-Point Polynomial Algorithms in Convex Programming*. Society for Industrial and Applied Mathematics, 1994.
- [13] M. Grant and S. Boyd. CVX: Matlab software for disciplined convex programming, version 2.1. <http://cvxr.com/cvx>, 2014.
- [14] MOSEK ApS. The MOSEK optimization toolbox for matlab manual, version 8.1. <http://docs.mosek.com/8.1/toolbox/index.html>, 2017.
- [15] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [16] P. Zegers. Fisher information properties. *Entropy*, 17(7):4918–4939, 2015.

## APPENDIX A UNOBSERVED MAXIMAL END COMPONENTS

In Section IV, we said that after reaching an unobserved maximal end component (UMEC), the agent may leak no more information since there exists a stationary policy that always stays in the UMEC. However, such a policy may not be admissible due to the reachability constraint. In that case, the agent has to leave the UMEC.

Assumption 1 ensures that the agent cannot leave UMECs. Every policy stays in UMECs and hence the outflow from these states is zero. Thanks to this assumption, we only need to consider the policies where the agent stays in UMECs and synthesize the policy accordingly.

In the following subsections we investigate the cases where the assumption does not hold. Appendix A-A provides an exhaustive search algorithm to find the optimal stationary policy. Appendix A-B provides an algorithm that searches a different class of policies to find the optimal policy.

### A. Agents with Stationary Policies

Consider the MDPs given in Figure 7 where the reachability requirement is  $\Pr^\pi(\text{Reach}[s_4 \cup s_5]) \geq 0.5$ . For both MDPs information is leaked only at state  $s_3$  and it is proportional to the expected residence time at state  $s_3$ , i.e.,  $x_{s_3}^\pi$ . Note that also the reachability probability is equal to the expected residence time at state  $s_3$ , i.e.,  $x_{s_3}^\pi = \Pr^\pi(\text{Reach}[s_4 \cup s_5])$ .

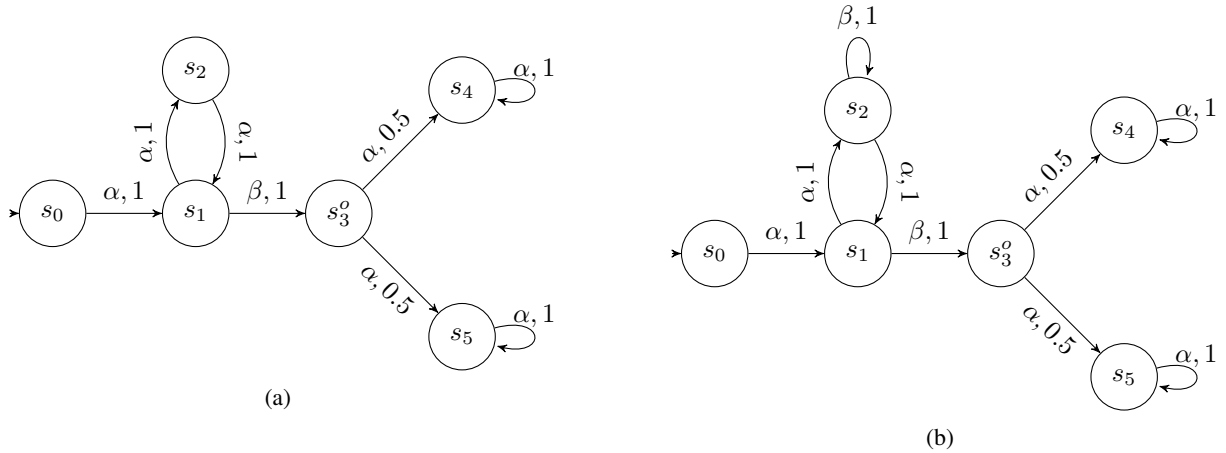


Figure 7: MDPs with 6 states. A label  $a, p$  of a transition refers to the transition that happens with probability  $p$  when action  $a$  is taken. The states marked with the superscript  $o$  are observed.

One might naturally think that a policy that makes  $x_{s_3}^\pi = 0.5$  is a minimum-information admissible policy. However, we note that such a stationary policy might not exist since  $x_{s_3}^\pi$  and  $\Pr^\pi(\text{Reach}[s_4 \cup s_5])$  are not continuous functions of stationary policies. For the MDP given in Figure 7a, a stationary policy  $\pi$  has  $\Pr^\pi(\text{Reach}[s_4 \cup s_5]) = 1$  if  $\pi_{s_1, \beta} > 0$  and  $\Pr^\pi(\text{Reach}[s_4 \cup s_5]) = 0$  otherwise. Every policy  $\pi^*$  such that  $\pi^*(s_1, \beta) > 0$  is a stationary, minimum-information admissible policy. However, such a policy does not satisfy the reachability requirement with equality. For the MDP given

---

**Algorithm 1** Synthesis of a stationary, minimum-information admissible policy for MDPs with UMECs - Process 1

---

```
1: Input: An MDP  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$ , the set of observed states  $W$ , the set of states to be reached  $C_{reach}$ , and the reachability probability  $\nu_{reach}$ .
2: Output: A stationary, minimum-information admissible policy  $\pi^*$  for  $\mathcal{M}$ .
3:  $R := \emptyset$ .
4: Find every UMEC  $(C, D)$  and set  $R := R \cup C$ .
5:  $L := 2^R$ .
6:  $minval := \infty$ 
7: for all  $l \in L$  do
8:   if  $l$  is a union unobserved end component then
9:      $C_{end} := l$ 
10:    Solve (6) with  $C_{end}$ ,  $C_{reach}$ , and  $\nu_{reach}$ . Set the optimal value to  $val$  and set the solution to  $restimes$ .
11:    if  $val \leq minval$  then
12:       $minval := val$ ,  $minset := l$ ,  $minrestimes := restimes$ .
13:    end if
14:  end if
15: end for
16:  $C_{end} := minset$ .
17: Synthesize the minimum-information admissible  $\pi^*$  policy using Algorithm 2 with  $minrestimes$  and  $C_{end}$ .
18: return  $\pi^*$ 
```

---

---

**Algorithm 2** Synthesis of a stationary, minimum-information admissible policy for UMECs - Process 2

---

```
1: Input: An MDP  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$ , the expected residence times  $x_{s,a}^{\pi^*}$  for  $\mathcal{M}$ , and  $C_{end}$ .
2: Output: A stationary, minimum-information admissible policy  $\pi^*$  for  $\mathcal{M}$ .
3: Synthesize a policy  $\pi^{stay}$  such that for a state  $s \in C_{end}$ ,  $\sum_{q \in C_{end}} \sum_{a \in \mathcal{A}(s)} \pi_{s,a}^{stay} \mathcal{P}_{s,a,q} = 1$ .
4: for all  $s \in S$  do
5:   if  $\sum_{a \in \mathcal{A}(s)} x_{s,a}^{\pi^*} = 0$  then
6:     for all  $a \in \mathcal{A}(s)$  do
7:       Set  $\pi_{s,a}^*$  arbitrarily between 0 and 1 subject to  $\sum_{a' \in \mathcal{A}(s)} \pi_{s,a'}^* = 1$ .
8:     end for
9:   else if  $s \in C_{end}$  then
10:    for all  $a \in \mathcal{A}(s)$  do
11:       $\pi_{s,a}^* := \pi_{s,a}^{stay}$ .
12:    end for
13:   else
14:    for all  $a \in \mathcal{A}(s)$  do
15:       $\pi_{s,a}^* := \frac{x_{s,a}^{\pi^*}}{\sum_{a' \in \mathcal{A}(s)} x_{s,a'}^{\pi^*}}$ .
16:    end for
17:   end if
18: end for
19: return  $\pi^*$ 
```

---

in Figure 7b, it is possible to find a stationary policy that satisfies the reachability requirement with equality. The stationary policy  $\pi^*$  with  $\pi_{s_1,\alpha}^* = 0.5$ ,  $\pi_{s_1,\beta}^* = 0.5$ ,  $\pi_{s_2,\alpha}^* = 0$ , and  $\pi_{s_2,\beta}^* = 1$  is the minimum-information admissible policy.

By the MDPs given in Figure 7, we note that determining whether the optimal policy stays in a UEC is not trivial. To find the stationary, minimum-information admissible policy, we give an optimization algorithm that is based on exhaustive search of all unobserved end components.

**Definition 3.** A union unobserved end component is a sub-MDP  $(C, D)$  that is union of UECs  $(C_1, D_1), \dots, (C_N, D_N)$  such that  $C = C_1 \cup \dots \cup C_N$  and  $D(s) = D_1(s) \cup \dots \cup D_N(s)$  for every  $s$  in  $C$ .

Algorithm 1 takes a subset of UMEC states, checks whether this subset is a union unobserved end component (see Lines 7-8). If the subset is a union unobserved end component, it finds the optimal stationary policy that makes the agent stay in the union unobserved end component (see Lines 9 - 10). The algorithm outputs the minimum-information admissible policy after checking all subsets.

**Remark 5.** Note that the size of  $R$  is  $O(|S|)$  in Algorithm 1 and the size of  $L$  is  $O(2^{|S|})$ . Checking whether a set of states  $S$  is a union unobserved end component has  $O(|S|^3|A|)$  complexity. Hence, the exhaustive search given in Algorithm 1 increases the complexity by a factor of  $O(2^{|S|}|S|^3|A|)$ .

### B. Agents with Nonstationary Policies

In this section, we remove Assumption 1 and introduce an algorithm that avoids the exhaustive search given in Algorithm 1. The exhaustive search is required as a drawback of stationary policies. We extend the policy space of the agent to find the optimal policy with lower computational complexity by allowing the agent to pick a policy that might be nonstationary for the unobserved states. We call a policy  $\pi$  observation stationary if it is stationary at the observed states and define  $\Pi^{Obs\ St}(\mathcal{M})$  as set of the observation stationary policies of  $\mathcal{M}$ .

The new algorithm is based on the flow constraints that describe the policy space of the agent. Under Assumption 1, the flow constraints given in (6c) - (6e) disallow outflow from the observed maximal end components to the other states. We remove this assumption and allow outflow from UMECs.

To find the minimum-information admissible policy we first create a modified MDP. The modified MDP has two copies of UMECs that are connected to each other with an action called *switch*. For a UMEC, while the original copy is connected to the other states, the duplicate copy is closed. We use the duplicate copies to represent the cases where the agent decides to stay in the UMEC.

For MDP  $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, s_0)$ , we create a modified MDP  $\bar{\mathcal{M}} = (\bar{S}, \bar{\mathcal{A}}, \bar{\mathcal{P}}, s_0)$  as follows. Let  $C_{end}$  be the set of states that belong to some UMEC of  $\mathcal{M}$ . For each  $s \in C_{end}$ , we create a duplicate state  $\bar{s}$ . Let  $\bar{C}_{end}$  be the set of duplicate UMEC states. We define  $\bar{S} := S \cup \bar{C}_{end}$ . For all  $s \in S$ , we define  $\bar{\mathcal{A}}(s) := \mathcal{A}(s)$  and for all  $a \in \bar{\mathcal{A}}(s)$ ,  $q \in S$  we define  $\bar{\mathcal{P}}_{s,a,q} := \mathcal{P}_{s,a,q}$ . The duplicate state  $\bar{s}$  has the action  $a$  if and only if  $a \in \mathcal{A}(s)$  and  $\sum_{q \in C_{end}} \mathcal{P}_{s,a,q} = 1$ . For every  $\bar{s} \in \bar{C}_{end}$ ,  $\bar{q} \in \bar{C}_{end}$ , and  $a \in \bar{\mathcal{A}}(\bar{s})$ , we let  $\bar{\mathcal{P}}_{\bar{s},a,\bar{q}} = \mathcal{P}_{s,a,q}$ . For every state  $s \in C_{end}$ , we also add a new action *switch* to  $\bar{\mathcal{A}}(s)$  such that  $\bar{\mathcal{P}}_{s,switch,\bar{s}} = 1$ .

Note that by definition  $C_{reach}$  belongs to  $C_{end}$ . For the reachability constraint, we use the set of duplicate states  $\bar{C}_{reach}$ . For modified MDP  $\bar{\mathcal{M}}$  and  $\bar{C}_{end}$ , we find the expected residence times of a minimum-information admissible policy with the following optimization problem

$$\inf \sum_{w \in W} x_w^\pi l_w^\pi \quad (7a)$$

$$\text{s. t. } x_s^\pi = \sum_{a \in \mathcal{A}(s)} x_{s,a}^\pi, \quad \forall s \in \bar{S} \setminus \bar{C}_{end} \quad (7b)$$

$$x_{s,a}^\pi \geq 0, \quad \forall s \in \bar{S} \setminus \bar{C}_{end}, \forall a \in \bar{\mathcal{A}}(s) \quad (7c)$$

$$\sum_{a \in \bar{\mathcal{A}}(s)} x_{s,a}^\pi - \sum_{q \in S} \sum_{a \in \bar{\mathcal{A}}(q)} x_{q,a}^\pi \mathcal{P}_{q,a,s} = \mathbb{1}_{s_0}(s), \quad \forall s \in \bar{S} \setminus \bar{C}_{end} \quad (7d)$$

$$\sum_{q \in \bar{C}_{reach}} \sum_{s \in \bar{S} \setminus \bar{C}_{end}} \sum_{a \in \bar{\mathcal{A}}(s)} x_{s,a}^\pi \mathcal{P}_{s,a,q} + \mathbb{1}_{s_0}(q) \geq \nu_{reach}, \quad (7e)$$

and synthesize the optimal policy  $\bar{\pi}^*$ .

**Remark 6.** The optimization problem given in (7) does not include the policies that always stay in  $C_{end}$ , in the feasible set. However, we remark that it does not effect the optimality of the solution since the value of such a policy can also be achieved by a policy that enters and always stays in  $\bar{C}_{end}$ .

We describe the policy  $\pi^*$  of the agent in the original MDP with Algorithm 3. We use a memory element *switched* that is *True* if and only if *switch* action is taken previously. We also synthesize a stationary policy  $\pi_{stay}$  for MDP  $\mathcal{M}$  that always stays in  $C_{end}$ .

Note that the resulting policy is not stationary for the original MDP  $\mathcal{M}$ . The agent remembers whether it switched to the stay mode in the past. However, it is stationary for all states in  $S \setminus C_{end}$ . The inference problem is still meaningful since the policy does not change over time for the observed states.

**Proposition 6.** For an MDP  $\mathcal{M}$ , the policy  $\pi^*$  that is synthesized via the optimization problem given in (7) and Algorithm 3, is a solution to the following problem

$$\begin{aligned} \min_{\pi \in \Pi^{Obs\ St}(\mathcal{M})} & \mathbb{E}[t_{W,\xi}^\pi] \\ \text{s. t.} & \Pr^\pi(\text{Reach}[C_{reach}]) \geq \nu_{reach} \end{aligned}$$

where  $\xi$  is a random path generated under policy  $\pi$ .

---

**Algorithm 3** Synthesis of a minimum-information admissible policy

---

```
1: Input: Current state  $s$  in  $\mathcal{M}$ ,  $switched$ ,  $\bar{\pi}^*$ ,  $\pi_{stay}$ , and  $C_{end}$ .
2: Output: The optimal policy  $\pi^*$  of  $\mathcal{M}$  and  $switched$ .
3: if  $switched$  then
4:    $\pi := \pi_{stay}$ .
5: else if  $s \notin C_{end}$  then
6:    $\pi^* := \bar{\pi}^*$ .
7: else
8:    $rnd := Unif[0, 1]$ .
9:   if  $rnd \leq \bar{\pi}_{s,switch}^*$  then
10:     $switched := True$ .
11:     $\pi^* = \pi_{stay}$ .
12:   else
13:    for all  $a \in \mathcal{A}(s)$  do
14:       $\pi_{s,a}^* = \frac{\bar{\pi}_{s,a}^*}{1 - \bar{\pi}_{s,switch}^*}$ .
15:    end for
16:   end if
17: end if
18: return  $\pi^*$  and  $switched$ 
```

---

## APPENDIX B

*Proof of Proposition 1.* We first consider two cases:

- $\iota_w^\pi = \infty$  for a reachable state  $w \in W$ .
- $\iota_w^\pi > 0$  for a state  $w \in W$  and  $w$  is recurrent under policy  $\pi$ .

Assume that the first case is possible. Since the path fragments of  $\mathcal{M}^\pi$  that end with  $w$  has a positive probability and  $\iota_w^\pi = \infty$ , the expected total information must be infinite. Thus, the first case is not possible. Assume that the second case is possible. Since the paths of  $\mathcal{M}^\pi$  that visit  $w$  infinitely often has a positive probability and  $\iota_w^\pi > 0$ , the expected total information must be infinite. The second case is also not possible. Hence, all observed states must be unreachable or must leak finite information and be transient.

The expected total information of a transient or unreachable state  $w \in W$  is

$$\mathbb{E}[\iota_{w,\xi}^\pi] = \sum_{n=0}^{\infty} \Pr(N_{w,\xi} = n) n \iota_w^\pi \quad (8a)$$

$$= \mathbb{E}[N_{w,\xi}] \iota_w^\pi \quad (8b)$$

$$= x_w^\pi \iota_w^\pi \quad (8c)$$

where  $N_{w,\xi}$  is the random variable that is the number of appearances of  $w$  in  $\xi$ .

The expected total information is

$$\mathbb{E}[\iota_{W,\xi}^\pi] = \sum_{w \in W} \mathbb{E}[\iota_{w,\xi}^\pi] \quad (9a)$$

$$= \sum_{w \in W} x_w^\pi \iota_w^\pi. \quad (9b)$$

■

*Sketch of Proof for Proposition 2.* If the optimal value of (6) is infinite then any policy that satisfies the reachability constraints is the optimal policy. Otherwise, let  $M$  be the optimal value of (6).  $\iota_s^\pi$  given in (5) is a lower semicontinuous function in the domain  $x_{s,a}^\pi \geq 0$  where  $a \in \mathcal{A}(s)$ . The objective function of (6) is a sum of lower semicontinuous functions and thus is a lower semicontinuous function in domain  $x_{s,a}^\pi \geq 0$  for all  $s \in S \setminus C_{end}$  and  $a \in \mathcal{A}(s)$ . For every  $x_{w,a}^\pi$ , that satisfies  $x_w^\pi \iota_w^\pi \leq M$ , is bounded where  $w \in W$ . Also every  $x_{s,a}^\pi$  is bounded since a state  $s \in S \setminus C_{end}$  must be transient. With the constraints (6c)-(6e) the feasible region is a compact set. Since a lower semicontinuous function attains its infimum on a compact set, we conclude that the proposition holds. ■

Before we proceed to the proof of Proposition 3, we give the following lemma that will be used in the proof.

**Lemma 7.** If  $f : V \rightarrow \mathbb{R}$  is a positive, concave function where  $V \subset \mathbb{R}^n$  is a convex set, then  $\frac{1}{f(x)}$  is a convex function on  $V$ .

*Proof of Lemma 7.* Since  $f$  and  $\log$  are concave functions,  $\log f$  is a concave function and consequently  $-\log f$  is a convex function on  $V$ . Finally,  $\exp(-\log f) = \frac{1}{f}$  is a convex function on  $V$ , due to convexity of  $\exp$  and  $-\log f$  on  $V$ . ■

*Proof of Proposition 3.* Let  $f_1 : Y_1 \rightarrow \mathbb{R}$  be a function such that  $f_1(p) = \sum_{i=1}^n p_i(1 - p_i)$ . Clearly  $f_1$  is a positive, concave function on the convex domain  $Y_1 = \{p \mid p_1, \dots, p_n \geq 0, \sum_{i=1}^n p_i = 1, \exists i, j \in [n], i \neq j, p_{i,j} > 0\}$ . Let  $f_2 : Y_1 \rightarrow \mathbb{R}$  be a function such that

$$f_2 := \frac{1}{f_1} = \frac{1}{\sum_{i=1}^n p_i(1 - p_i)}.$$

By Lemma 7,  $f_2$  is a convex function on the domain  $Y_1$ . A perspective function [15] of  $f_2$  is  $g_1 : Y_2 \rightarrow \mathbb{R}$  such that

$$\begin{aligned} g_1 \left( x, \sum_{i=1}^n x_i \right) &= \sum_{i=1}^n x_i f \left( \frac{x}{\sum_{i=1}^n x_i} \right) \\ &= \sum_{i=1}^n x_i f(p) \end{aligned}$$

where

$$p_i = \frac{x_i}{\sum_{i=1}^n x_i}.$$

Due to the convexity property of perspective functions [15],  $g_1$  is convex on

$$Y_2 = \left\{ \left( x, \sum_{i=1}^n x_i \right) \mid x_1, \dots, x_n \geq 0, \exists i, j \in [n], i \neq j, x_{i,j} > 0 \right\}$$

since  $f_2$  is convex on  $Y_1$ . We eliminate the redundant dimension  $\sum_{i=1}^n x_i$  and define  $g_2 : V_1 \rightarrow \mathbb{R}$  such that

$$g_2(x) = g_1 \left( x, \sum_{i=1}^n x_i \right).$$

$g_2$  is an affine transformation of  $g_1$  and convex on

$$V_1 = \{x \mid x_1, \dots, x_n \geq 0, \exists i, j \in [n], i \neq j, x_{i,j} > 0\}$$

We introduce  $V_0 = \{x \mid x_1 = \dots = x_n = 0\}$  and  $V_{det} = \{x \mid \exists i \in [n], \forall j \in [n], i \neq j, x_i > 0, x_j = 0\}$ . Note that  $V_0$ ,  $V_1$ , and  $V_{det}$  are disjoint sets.

Now we define  $g : V \rightarrow \mathbb{R} \cup \{\infty\}$  on  $V = V_0 \cup V_1 \cup V_{det}$  such that  $g(x) = 0$  if  $x \in V_0$ ,  $g(x) = g_1(x)$  if  $x \in V_1$ , and  $g(x) = \infty$  if  $x \in V_{det}$ .

Clearly  $g_3$  is convex on  $V_0$  and  $V_{det}$ . We check all possible combinations for convexity where  $\lambda \in [0, 1]$ :

- $\lambda g(v_1) + (1 - \lambda)g(v_2) \geq g(\lambda v_1 + (1 - \lambda)v_2)$  if  $v_1 \in V_{det}$  and  $v_2 \in V_0 \cup V_1$ ,
- $\lambda g(v_1) + (1 - \lambda)g(v_2) = g(\lambda Y_1 + (1 - \lambda)v_2)$  if  $v_1 \in V_0$  and  $v_2 \in V_1$ .

Hence  $g$  is convex on  $V$ .

Now we represent the objective function of (6) using  $g$ . Without loss of generality assume that the successor states of state  $s$  are  $q_1, \dots, q_{|Succ(s)|}$  and the actions at state  $s$  are  $a_1, \dots, a_{|A(s)|}$ . Note that

$$x_s^\pi \iota_s^\pi = g(Px)$$

where  $x = [x(s, a_1), \dots, x(s, a_{|A(s)|})]^T$  and  $P$  is a  $|Succ(s)| \times |A(s)|$  matrix with  $(i, j)$ -th entry  $\mathcal{P}_{s, a_j, q_i}$ .

Since  $x_s^\pi \iota_s^\pi$  is an affine mapping of  $g$ ,  $x_s^\pi \iota_s^\pi$  is convex on  $T' = \{x \in \mathbb{R}^{|A(s)|} \mid \sum_{a \in A(s)} x_{s,a}^\pi \mathcal{P}_{s,a,q} \geq 0\}$  and consequently on  $T = \{x \in \mathbb{R}^{|A(s)|} \mid x_{s,a}^\pi \geq 0\} \subseteq T'$ .

The objective function (6a) is a sum of convex functions and the constraints in (6) are linear. Therefore, we conclude that the optimization problem is convex. ■

*Proof of Proposition 4.* Due to the stochasticity of MDP, we might encounter the cases where the observer has no observation from a state and hence no sample for estimation. For such cases, denote  $\sigma_{w,0}$  for the MSE when there is no sample for estimation. Denote  $\sigma_{w,+}$  for the MSE when there is at least one sample for estimation.

The MSE of the  $q$ -th element  $(\sigma_w)_q$  is the estimation error for transition probability to the successor state  $q$  such that  $\sum_{q \in Succ(s)} (\sigma_w)_q = \sigma_w$ .

Denote the result of the successor state at time  $t$  for the random path  $\xi$  by  $R_{t,\xi}$  where by definition  $R_{-1,\xi} = s_0$  and  $N_{w,\xi}$  for the number of times that state  $w$  appears in  $\xi$ . We have

$$(\sigma_w)_q = \Pr(N_{w,\xi} = 0|\pi)(\sigma_{w,0})_q + \Pr(N_{w,\xi} > 0|\pi)(\sigma_{w,+})_q \quad (10a)$$

$$\geq \Pr^\pi(\text{Reach}[w])(\sigma_{w,+})_q \quad (10b)$$

$$\geq \frac{\Pr^\pi(\text{Reach}[w])}{I_{\xi|N_{w,\xi}>0}(\pi_{w,q})} \quad (10c)$$

$$= \frac{\Pr^\pi(\text{Reach}[w])}{\sum_{t=0}^{\infty} I_{R_{t,\xi}|R_{t-1,\xi},\dots,R_{0,\xi},N_{w,\xi}>0}(\mathcal{P}_{w,q}^\pi)} \quad (10d)$$

$$= \frac{\Pr^\pi(\text{Reach}[w])}{\sum_{t=0}^{\infty} I_{R_{t,\xi}|R_{t-1,\xi},N_{w,\xi}>0}(\mathcal{P}_{w,q}^\pi)} \quad (10e)$$

$$= \frac{\Pr^\pi(\text{Reach}[w])\mathcal{P}_{w,q}^\pi(1 - \mathcal{P}_{w,q}^\pi)}{\sum_{t=0}^{\infty} \Pr(R_{t-1} = w|N_{w,\xi} > 0)} \quad (10f)$$

$$= \frac{\Pr^\pi(\text{Reach}[w])^2\mathcal{P}_{w,q}^\pi(1 - \mathcal{P}_{w,q}^\pi)}{x_w^\pi} \quad (10g)$$

where (10c) is due to Cramér-Rao bound, (10d) is due to chain rule of the Fisher information [16], and (10e) is due to Markovian property of paths.

The MSE at state  $w$  is bounded such that

$$\sigma_w \geq \sum_{q \in \text{Succ}(s)} \frac{\Pr^\pi(\text{Reach}[w])^2\mathcal{P}_{w,q}^\pi(1 - \mathcal{P}_{w,q}^\pi)}{x_w^\pi} \quad (11a)$$

$$= \frac{\Pr^\pi(\text{Reach}[w])^2}{x_w^\pi \iota_w^\pi}. \quad (11b)$$

■

*Proof of Corollary 5.* Total MSE at state  $w$  is

$$\sigma_w \geq \frac{\Pr^\pi(\text{Reach}[w])^2}{x_w^\pi \iota_w^\pi}. \quad (12a)$$

The total MSE for the set of states  $W$  is

$$\sum_{w \in W} \sigma_w \geq \sum_{w \in W} \frac{\Pr^\pi(\text{Reach}[w])^2}{x_w^\pi \iota_w^\pi} \quad (13a)$$

$$\geq \sum_{w \in W} \frac{\min_{w' \in W} \Pr^\pi(\text{Reach}[w'])^2}{x_w^\pi \iota_w^\pi} \quad (13b)$$

$$\geq \frac{\min_{w' \in W} \Pr^\pi(\text{Reach}[w'])^2 |W|^2}{\sum_{w \in W} x_w^\pi \iota_w^\pi} \quad (13c)$$

$$= \frac{\min_{w' \in W} \Pr^\pi(\text{Reach}[w'])^2 |W|^2}{\mathbb{E}_\xi[\iota_{W,\xi}^\pi]}. \quad (13d)$$

■

*Sketch of Proof for Proposition 6.* The proof steps are as follows,

- show that a stationary policy  $\bar{\pi}^*$  is optimal for modified MDP  $\bar{\mathcal{M}}$  among all policies in  $\Pi^{\text{Obs.St}}(\bar{\mathcal{M}})$ ,
- show that the minimum-information admissible of  $\mathcal{M}$  is not lower than  $\bar{\mathcal{M}}$ ,
- show that the expected total informations are equal for  $\pi^*$  of  $\mathcal{M}$  and  $\bar{\pi}^*$  of  $\bar{\mathcal{M}}$ .

Consider the minimum-information admissible policy  $\bar{\pi}^*$  for  $\bar{\mathcal{M}}$ . For every state  $s \in \bar{S}$ , we identify whether  $\bar{\pi}^*$  makes  $s$  recurrent or transient.

Let  $(C, D)$  be an original UMEC of  $\bar{\mathcal{M}}$  and  $\pi^1$  be a policy such that  $\Pr^{\pi^1}(s_t \in C \text{ eventually always}) > 0$ . We claim that the expected total information under policy  $\pi^1$  can also be achieved by a policy  $\pi^2$  such that  $\Pr^{\pi^2}(s_t \in C \text{ infinitely often}) = 0$  since upon deciding to stay in  $C$  the agent can first take *switch* action and then take the same actions in the duplicate UMEC  $\bar{C}$ . Note that staying in  $C$  or  $\bar{C}$  does not affect the total information since both UMECs leak no information. Hence, we only look for policies that makes  $C_{\text{end}}$  transient and  $\bar{C}_{\text{end}}$  recurrent.

Let  $(C, D)$  be an end component of  $\bar{\mathcal{M}}$  such that  $C \cap C_{end} = \emptyset$  and there exists  $w \in W$  and  $w \in C$ . We claim that a policy that always stays in  $C$  visits an observed state infinitely often and leaks infinite information. If it does not, then there must exist a state  $s \in C$  such that  $s$  is recurrent and  $s \notin W$ . Such a state  $s$  must belong to a UMEC, but by construction it is not possible. Hence if there exists a policy that leaks finite information every state  $s \in C$  must be transient. Also note that a state that does not belong to an end component must be transient by definition.

We partition  $\bar{S}$  into two sets: transient states  $\bar{S} \setminus \bar{C}_{end}$  and recurrent states  $\bar{C}_{end}$ . Under any policy  $\pi \in \Pi^{Obs.St}(\bar{\mathcal{M}})$  that makes the  $\bar{C}_{end}$  recurrent and  $\bar{S} \setminus \bar{C}_{end}$  transient, we have the flow equation

$$\sum_{a \in \mathcal{A}(s)} x_{s,a}^\pi - \sum_{q \in \bar{S} \setminus \bar{C}_{end}} \sum_{a \in \mathcal{A}(q)} x_{q,a}^\pi \mathcal{P}_{q,a,s} = \mathbf{1}_{s_0}(s), \quad \forall s \in \bar{S} \setminus \bar{C}_{end}.$$

Since we optimize over the observation stationary policies, the Proposition 1 still holds. The optimization problem given in (7) finds the state-action residence times of the optimal policy subject to the flow equation and the reachability constraint. The stationary policy synthesized via (3) yields to the optimal expected residence times and hence is optimal.

Let

$$v^* = \inf_{\pi \in \Pi^{Obs.St}(\mathcal{M})} \mathbb{E}[\iota_{W,\xi}^\pi] \quad (15a)$$

$$\text{s. t. } \Pr^\pi(\text{Reach}[C_{reach}]) \geq \nu_{reach} \quad (15b)$$

and

$$\bar{v}^* = \inf_{\bar{\pi} \in \Pi^{Obs.St}(\bar{\mathcal{M}})} \mathbb{E}[\iota_{W,\xi}^{\bar{\pi}}] \quad (16a)$$

$$\text{s. t. } \Pr^{\bar{\pi}}(\text{Reach}[C_{reach}]) \geq \nu_{reach}. \quad (16b)$$

Since every  $\pi \in \Pi(\mathcal{M})$  is also realizable for  $\bar{\mathcal{M}}$  with the same expected total information and reachability probabilities, we have  $v^* \geq \bar{v}^*$ .

Finally, we note that  $\pi^*$  of  $\mathcal{M}$  and  $\bar{\pi}$  of  $\bar{\mathcal{M}}$  yield to the same expected total information  $\bar{v}^*$  since the expected residence times and the policies are the same at the observed states for both policies. Consequently,  $\pi^*$  is a minimum-information admissible policy of  $\mathcal{M}$ . ■