

# Structural Estimation of Partially Observable Markov Decision Processes

Yanling Chang, Alfredo Garcia, Zhide Wang, Lu Sun

**Abstract**—In many practical settings control decisions must be made under partial/imperfect information about the evolution of a relevant state variable. Partially Observable Markov Decision Processes (POMDPs) is a relatively well-developed framework for modeling and analyzing such problems. In this paper we consider the structural estimation of the primitives of a POMDP model based upon the observable history of the process. We analyze the structural properties of POMDP model with random rewards and specify conditions under which the model is identifiable without knowledge of the state dynamics. We consider a *soft* policy gradient algorithm to compute a maximum likelihood estimator and provide a finite-time characterization of convergence to a stationary point. We illustrate the estimation methodology with an application to optimal equipment replacement. In this context, replacement decisions must be made under partial/imperfect information on the true state (i.e. condition of the equipment). We use synthetic and real data to highlight the robustness of the proposed methodology and characterize the potential for misspecification when partial state observability is ignored.

## I. INTRODUCTION

POMDPs generalize Markov Decision Processes (MDPs) by taking into account partial state observability due to measurement noise and/or limited access to information. When the state is only partially observable (or “hidden”), optimal control policies may need to be based upon the complete history of implemented actions and recorded observations. For Bayes optimal control, the updated Bayes belief provides enough information to identify the optimal action. In this context, a POMDP can be transformed to a MDP in which the state is the updated Bayes belief distribution (see [1]–[4] and many others). POMDP-based models have been successfully applied in a variety of application domains (see e.g. [5]–[9]). In this paper we consider the task of estimating the primitives of a POMDP model (i.e. reward function, hidden transition and observation probabilities) based upon the observable histories of implemented actions and observables.

Under the assumption of complete state observability (for both the controller and the modeler), this problem has been widely studied in two strands of the literature where it is referred to as structural estimation of Markov decision processes (MDP) or alternatively as inverse reinforcement learning (IRL) (see literature review in Section II below). The partially observable case remains to be considered. The present paper addresses this important gap in the literature.

In the first part of the paper we introduce a POMDP model with random rewards. For every implemented action, the realization of rewards is observed by the controller but not by the modeler. We characterize optimal decisions in a POMDP model with privately observed random rewards by means of a *soft* Bellman equation. This result relies on showing that the updated Bayesian belief on the hidden state

is sufficient to define optimal dynamic choices. Here the term *soft* is used because optimal policies are characterized by the *softmax* function under a specific distributional assumption on random rewards. For a given choice of parameter estimates, the modeler can compute the evolution of latent Bayesian beliefs by the recursive application of Bayes rule for each sample path. Based upon this observation, we introduce a *soft* policy gradient algorithm in order to compute the maximum likelihood estimator. Since maximum likelihood is in general a non-convex function, we provide a finite-time characterization of convergence to a stationary point.

We show that the proposed model can be identified if the priori belief distribution, the cardinality of the system state space, the distribution of random i.i.d shocks, the discount factor, and the reward for a fixed reference action are given. We show the estimation is robust to the specification of the a priori belief distribution provided the data sequence is sufficiently long (to address the case where the priori belief distribution is unknown).

We test and validate the proposed methodology in an optimal engine replacement problem with both synthetic and real data. Using synthetic data, we show that our developed estimation procedure can recover the true model primitives. This experiment also numerically illustrates how model misspecification resulting from ignoring partial state observability can lead to models with poor fit. We further apply the model to a widely studied engine replacement real dataset. Compared to the results in [10], our new method can dramatically improve the data fit by 17.7% in terms of the log-likelihood. The model reveals a feature of route assignment behavior in the dataset which was hitherto ignored, i.e. buses with engines *believed* to be in worse condition exhibit less utilization (mileage) and higher maintenance costs.

The paper is organized as follows. Section II provides a literature review with emphasis on how this paper differs from the literature. Section III introduces a POMDP model with random reward perturbations. Section IV presents a methodology for structural estimation of POMDP and a policy gradient algorithm. The identification results are presented in Section V. Section VI provides an illustration to the application of the estimation method to optimal engine replacement problem using both synthetic dataset and the real dataset.

## II. LITERATURE REVIEW

Research on structural estimation of MDPs features efficient algorithms to address computational challenges in estimating the structural parameters such as determining the value function used in the likelihood function estimation [11]–[17].

In the computer science literature this problem has been studied under the label of inverse reinforcement learning

(IRL). A maximum entropy method proposed in [18] has been highly influential in this literature. Sample-based algorithms for implementing the maximum entropy method have scaled to scenarios with nonlinear reward functions (see e.g., [19], [20]). In [21], the authors extended the maximum entropy estimation method to a partially observable environment, assuming both the transition probabilities and observation probabilities are known with domain knowledge. These methods have also been used for apprenticeship learning where a robot learns from expert-based demonstrations [18]. To our best knowledge, none of these existing works have developed a general methodology to jointly estimate the reward structure and system dynamics based on observable trajectories of a POMDP process.

### III. A POMDP MODEL WITH RANDOM REWARDS

At each decision epoch  $t \geq 0$ , the value of state  $s_t \in S$  is not directly observable to the controller nor to the external modeler. However, both the controller and the external modeler are able to observe value of a random variable  $z_t \in Z$  correlated with the underlying state  $s_t$ . We assume the finite action, states and observations. If the hidden state is  $s_t$  and  $a_t \in A$  is implemented, the random reward accrued is  $r_{\theta_1}(z_t, s_t, a_t) + \epsilon_t(a_t)$  where  $\theta_1 \in \mathbb{R}^{p_1}$  for some  $p_1 \in \mathbb{N}_+$  and  $\epsilon_t(a_t)$  is a random variable. The realization of random reward is observed by the controller but *not* by the modeler. This asymmetry of information implies that from the point of view of the modeler, the controller's actions are not necessarily deterministic.

The system dynamics is described by probabilities  $P_{\theta_2}(z_{t+1}, \epsilon_{t+1}, s_{t+1}|z_t, \epsilon_t, s_t, a_t)$  where  $\theta_2 \in \mathbb{R}^{p_2}$  for some  $p_2 \in \mathbb{N}_+$ ; see Figure 1 for a schematic representation.

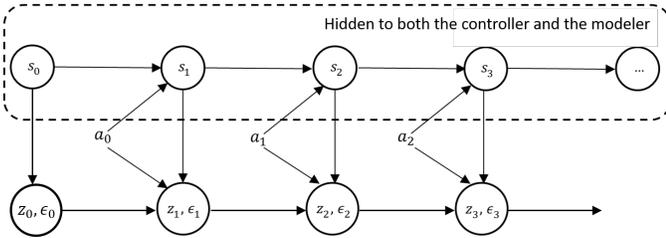


Fig. 1: Graphical illustration of the proposed POMDP model with random reward perturbations. At each stage,  $z_t$  is observed by both the controller and the modeler, and  $\epsilon_t$  is privately observed by the controller. The system state is  $s_t$  is hidden to both the controller and the modeler and  $s_t$  has its own hidden dynamics.

Let  $\zeta_t = \{z_t, \dots, z_0, a_{t-1}, \dots, a_0, x_0\}$  be the publicly received history of the dynamic decision process including all past and present revealed observations and all past actions at time  $t > 0$ , where  $x_0 = \{P(s_0), s_0 \in S\}$  is the prior belief distribution over  $S$ . The controller aims to maximize

$$E \left( \sum_{t=0}^{\infty} \beta^t [r(z_t, s_t, a_t) + \epsilon_t(a_t)] | x_0 \right)$$

We assume the following conditional independence (CI):

$$P_{\theta_2}(z_{t+1}, \epsilon_{t+1} | \zeta_t, \epsilon_t, a_t) = P(\epsilon_{t+1} | z_{t+1}) P_{\theta_2}(z_{t+1} | \zeta_t, a_t). \quad (\text{III.1})$$

Note that the process  $\{\zeta_t, \epsilon_t\}$  is not Markovian; however,  $z_{t+1}$  is a sufficient statistic for  $\epsilon_{t+1}$ , indicating  $\epsilon_t$  and  $\epsilon_{t+1}$  are independent given  $z_{t+1}$  and

$$P_{\theta_2}(z_{t+1}, \epsilon_{t+1}, s_{t+1} | z_t, \epsilon_t, s_t, a_t) = P(\epsilon_{t+1} | z_{t+1}) P_{\theta_2}(z_{t+1}, s_{t+1} | z_t, s_t, a_t). \quad (\text{III.2})$$

In addition, the conditional probability  $P_{\theta_2}(z_{t+1} | \zeta_t, a_t)$  does not depend on  $\epsilon_t$ . Intuitively, the CI assumption implies that the POMDP system dynamics  $P_{\theta_2}(z_{t+1}, s_{t+1} | z_t, s_t, a_t)$  is superimposed by a noise process  $\{\epsilon_t\}$ .

Let  $x_{t, \theta_2} = P_{\theta_2}(\cdot | \zeta_t) \in X \subset \mathbb{R}^{|S|}$  be the conditional probability distribution given history  $\zeta_t$  where  $X$  is the unit simplex. Given action  $a_t$  and observation  $z_t$ , the expected reward is:

$$r_{\theta_1}(z_t, x_{t, \theta_2}, a_t) = \sum_{s_t} x_{t, \theta_2}(s_t) r_{\theta_1}(z_t, s_t, a_t).$$

Under the CI assumption, the optimal decision process can be recursively formulated as:

$$U_{t, \theta}(\zeta_t, \epsilon_t) = \max_{a_t \in A} \left\{ r_{\theta_1}(z_t, x_{t, \theta_2}, a_t) + \epsilon_t(a_t) + \beta \sum_{z_{t+1}} \int P_{\theta_2}(z_{t+1} | \zeta_t, a_t) U_{t+1, \theta}(\zeta_{t+1}, \epsilon_{t+1}) d\mu(\epsilon_{t+1} | z_{t+1}) \right\}, \quad (\text{III.3})$$

where  $\mu(\epsilon_{t+1} | z_{t+1})$  is the cumulative probability distribution of the random perturbation vector  $\epsilon_{t+1}$  given the new observation  $z_{t+1}$ .

The objective for structural estimation is to identify  $\theta_1$  in reward  $r_{\theta_1}(z_t, s_t, a_t)$ , and  $\theta_2$  in dynamics  $P_{\theta_2}(z_{t+1}, s_{t+1} | z_t, s_t, a_t)$  from the publicly received histories  $\{\zeta_t^i\}_{i=1}^N$ .

Building upon the POMDP literature [1], [3], we show in Theorem 1 below that  $z_t$  and the updated Bayesian belief distribution  $x_{t, \theta_2}$  are sufficient to identify the optimal dynamic choices. To this end we mainly follow the notation of introduce the observation probabilities:

$$\sigma_{\theta_2}(z_{t+1}, z_t, x_{t, \theta_2}, a_t) \triangleq \sum_{s'} \sum_s x_{t, \theta_2}(s) P_{\theta_2}(z_{t+1}, s' | z_t, s, a_t),$$

and the belief update function:

$$\lambda_{\theta_2}(z_{t+1}, z_t, x_{t, \theta_2}, a_t) \triangleq \frac{x_{t, \theta_2} P_{\theta_2}(z_{t+1}, z_t, a_t)}{\sigma_{\theta_2}(z_{t+1}, z_t, x_{t, \theta_2}, a_t)}, \quad (\text{III.4})$$

assuming  $\sigma_{\theta_2}(z_{t+1}, z_t, x_{t, \theta_2}, a_t) \neq 0$ , where we denote the  $(s, s')$  element of the matrix  $[P_{\theta_2}(z_{t+1}, z_t, a_t)]_{s, s'} \triangleq P_{\theta_2}(z_{t+1}, s' | z_t, s, a_t)$ ,  $s, s' \in S$ , and

$$[x_{t, \theta_2} P_{\theta_2}(z_{t+1}, z_t, a_t)]_{s_{t+1}} \triangleq \sum_s x_{t, \theta_2}(s) P_{\theta_2}(z_{t+1}, s_{t+1} | z_t, s, a_t).$$

**Theorem 1.** Let  $\zeta_t$  denote a finite history with current observation  $z_t = z$  and updated belief  $x_{t, \theta_2} = x$ . Let  $V_{t, \theta}(z, x, \epsilon)$

be defined as follows:

$$V_{t,\theta}(z, x, \epsilon) = \max_{a \in A} \left\{ r_{\theta_1}(z, x, a) + \epsilon(a) + \beta \sum_{z'} \int \sigma_{\theta_2}(z', z, x, a) V_{t+1,\theta}(z', x', \epsilon') d\mu(\epsilon' | z') \right\}$$

where  $x' = \lambda_{\theta_2}(z', z, x, a)$ . It follows that  $V_{t,\theta}(z, x, \epsilon) = U_{t,\theta}(\xi_t, \epsilon)$ . Hence,  $(z, x)$  is a sufficient statistic for solving (III.3).

#### A. Soft Bellman Equation

We now state and prove the *soft* Bellman equation for the POMDP model with random reward. Let  $\mathcal{B}$  be the Banach space of bounded, Borel measurable functions  $Q : Z \times X \times A \rightarrow R$  under the supremum norm  $\|\cdot\|$ .

Define the *soft* Bellman operator  $H_\theta : \mathcal{B} \rightarrow \mathcal{B}$  by

$$[H_\theta Q](z, x, a) = r_\theta(z, x, a) + \beta \sum_{z'} \sigma_{\theta_2}(z', z, x, a) \int \max_{a \in A} \{Q(z', x', a) + \epsilon(a)\} d\mu(\epsilon | z'), \quad (\text{III.5})$$

where  $x' = \lambda_{\theta_2}(z', z, x, a)$ .

**Theorem 2.** *Under CI and the mild regularity conditions listed in the Appendix,  $H_\theta : \mathcal{B} \rightarrow \mathcal{B}$  is a contraction mapping with modulus  $\beta$ . Hence,  $H_\theta$  has a unique fixed point  $Q_\theta$  (i.e.,  $Q_\theta = H_\theta Q_\theta$ ) and the optimal decision rule  $\delta_\theta$  is of the form:*

$$\delta_\theta(z, x, \epsilon) = \arg \max_{a \in A} \left\{ Q_\theta(z, x, a) + \epsilon(a) \right\}. \quad (\text{III.6})$$

Furthermore, with conditional choice probabilities

$$\pi_\theta(a | z, x) \triangleq P(a \in \delta_\theta(z, x, \epsilon))$$

it holds that  $\pi_\theta(a | z, x) = \frac{\partial \bar{V}_\theta}{\partial Q_\theta}$  where:

$$\begin{aligned} \bar{V}_\theta(z, x) &\triangleq \int \max_{a \in A} \{Q_\theta(z, x, a) + \epsilon(a)\} d\mu(\epsilon | z) \\ &= \sum_{a \in A} \pi_\theta(a | z, x) (Q_\theta(z, x, a) + E[\epsilon | a]) \\ &= E_{a \sim \pi_\theta(\cdot | z, x)} [Q_\theta(z, x, a) + E[\epsilon | a]]. \end{aligned} \quad (\text{III.7})$$

Finally when the distribution of  $\epsilon$  is standard Gumbel, the optimal policy takes a *softmax* form.

**Theorem 3.** *If the probability measure of  $\epsilon$  is multivariate extreme-value, i.e.,*

$$\mu(d\epsilon | z) = \prod_{a \in A} \exp\{-\epsilon(a) + \gamma\} \exp[-\exp\{-\epsilon(a) + \gamma\}], \quad (\text{III.8})$$

where  $\gamma > 0$  is the Euler constant. Then,

$$\pi_\theta(a | z, x) = \frac{\exp Q_\theta(z, x, a)}{\sum_{a' \in A} \exp Q_\theta(z, x, a')}, \quad (\text{III.9})$$

where

$$\begin{aligned} Q_\theta(z, x, a) &= r_{\theta_1}(z, x, a) + \beta \sum_{z'} \sigma_{\theta_2}(z', z, x, a) \bar{V}_\theta(z', x'), \\ x' &= \lambda_{\theta_2}(z', z, x, a), \end{aligned}$$

and

$$\begin{aligned} \bar{V}_\theta(z', x') &\triangleq \int V_\theta(z', x', \epsilon) d\mu(\epsilon | z') \\ &= \gamma + \log \left( \sum_{a'} \exp Q_\theta(z', x', a') \right). \end{aligned}$$

*Remark 1.* It can be easily verified that Theorems 1, 2 and 3 continue to hold for the case in which the controller is solving a finite horizon problem. Evidently, the results in this case require that the state-action function  $Q_{t,\theta}$  and the conditional choice probabilities  $\pi_{t,\theta}$  are time-dependent  $t$ .

#### IV. MAXIMUM LIKELIHOOD ESTIMATION

Given data corresponding to  $N \geq 1$  finite histories of pairs  $\{x_{0,i}, z_{t,i}, a_{t,i}, t = 1, \dots, T\}$  for  $i \in \{1, \dots, N\}$ , a sequence of trajectories for the belief  $\{x_{t,\theta_2,i} : t > 0\}$  can be recursively computed for a fixed value of  $\theta = (\theta_1, \theta_2)$  as follows:

$$\begin{aligned} x_{t+1,\theta_2,i} &= \lambda_{\theta_2}(z_{t+1,i}, z_{t,i}, x_{t,\theta_2,i}, a_{t,i}) \\ &= \frac{x_{t,\theta_2,i} P_{\theta_2}(z_{t+1,i}, z_{t,i}, a_{t,i})}{\sigma_{\theta_2}(z_{t+1,i}, z_{t,i}, x_{t,\theta_2,i}, a_{t,i})}. \end{aligned}$$

Thus, the log-likelihood can be written as:

$$\begin{aligned} \log \ell(\theta) &\triangleq \log \prod_{i=1}^N P(\zeta_{T,i} | x_{0,i}) x_{0,i} \\ &= \log \prod_{i=1}^N \prod_{t=0}^{T-1} P(z_{t+1,i} | \zeta_{t,i}, a_{t,i}) P(a_{t,i} | \zeta_{t,i}) x_{0,i} \\ &= \sum_{i=1}^N \sum_{t=0}^{T-1} [\log \sigma_{\theta_2}(z_{t+1,i}, z_{t,i}, x_{t,\theta_2,i}, a_{t,i}) \\ &\quad + \log \pi_\theta(a_{t,i} | z_{t,i}, x_{t,\theta_2,i})] + \sum_{i=1}^N \log x_{0,i}. \end{aligned} \quad (\text{IV.1})$$

That is, assuming the data is generated by a Bayesian agent controlling a partially observable Markov process, the external modeler can construct a POMDP model by finding the parameter values that maximize the log-likelihood in (IV.1).

#### A. A Soft Policy Gradient Algorithm

We now introduce a *Soft* policy gradient algorithm for approximately maximizing the log-likelihood expression in (IV.1). In what follows we shall assume the a priori distribution  $x_{0,i}$  is known (see Section V-A for uncertain/unknown  $x_{0,i}$ ). A two-stage estimator can be obtained by first solving for the value of  $\theta_2$  that maximizes the value of the first term on the right hand side in (IV.1) and then solve for the value of  $\hat{\theta}_1$  that maximizes the log of pseudo-likelihood  $\hat{\ell}(\theta)$  defined as:

$$\hat{\ell}(\theta_1) = \sum_{i=1}^N \sum_{t=0}^{T-1} \log \pi_{(\theta_1, \hat{\theta}_2)}(a_{t,i} | z_{t,i}, x_{t,\hat{\theta}_2,i}).$$

We simplify notation by using  $\pi_{\theta_1}$  and  $x_{t,i}$  to refer to  $\pi_{(\theta_1, \hat{\theta}_2)}$  and  $x_{t,\hat{\theta}_2,i}$  respectively. For a given value of  $\theta_1$ , consider the *soft* Bellman equation:

$$Q_{\theta_1}(z, x, a) = r_{\theta_1}(z, x, a) + \beta \sum_{z'} \sigma_{\hat{\theta}_2}(z', z, x, a) \bar{V}_{\theta_1}(z', x'),$$

where  $x' = \lambda_{\hat{\theta}_2}(z', z, x, a)$ .

After solving the *soft* Bellman equation for fixed  $\theta_1$  we can compute the gradient:

$$\begin{aligned} & \nabla_{\theta_1} \log \pi_{\theta_1}(a|z, x) \\ &= \nabla_{\theta_1} \log \left( \frac{\exp Q_{\theta_1}(z, x, a)}{\sum_{a' \in A} \exp Q_{\theta_1}(z, x, a')} \right) \\ &= \nabla_{\theta_1} Q_{\theta_1}(z, x, a) - \nabla_{\theta_1} \log \sum_{a'} \exp Q_{\theta_1}(z, x, a') \\ &= \nabla_{\theta_1} Q_{\theta_1}(z, x, a) - \nabla_{\theta_1} \bar{V}_{\theta_1}(z, x) \\ &= \nabla_{\theta_1} Q_{\theta_1}(z, x, a) - \sum_{a'} \pi_{\theta_1}(a'|z, x) \nabla_{\theta_1} Q_{\theta_1}(z, x, a'). \end{aligned}$$

The basic steps of a *soft* policy gradient algorithm are listed in Algorithm 1. Before analyzing the convergence of the *soft* policy gradient algorithm we state a preliminary result.

---

**Algorithm 1** *Soft Policy Gradient Algorithm.*

---

Compute  $\hat{\theta}_2$  and  $x_{t,i} = x_{t,\hat{\theta}_2,i}$   $t = 1, \dots, T, i = 1, \dots, N$ ;

Initialize  $k = 0, \theta_1^0, \nabla_{\theta_1} \hat{\ell}(\theta_1^0), \epsilon$  and  $\rho$ ;

**while**  $\left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\| \geq \epsilon$  **do**

$k \leftarrow k + 1$ ;

Compute  $\nabla_{\theta_1} Q_{\theta_1^k}(a|z_{t,i}, x_{t,i}), a \in A$ ;

Compute  $\nabla_{\theta_1} \hat{\ell}(\theta_1^k) = \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta_1} \log \pi_{\theta_1^k}(a_{t,i}|z_{t,i}, x_{t,i})$ ;

Update parameters  $\theta_1^{k+1} = \theta_1^k + \rho \nabla_{\theta_1} \hat{\ell}(\theta_1^k)$ ;

**end**

---

**Lemma 1.** Assume  $r_{\theta_1}(z, x, a)$  is twice continuously differentiable in  $\theta_1 \in \mathbb{R}^{p_1}$  and

$$\begin{aligned} & \sup_{\theta_1} \|\nabla_{\theta_1} r_{\theta_1}(z, x, a)\| \leq L_{r,1} < \infty \\ & \sup_{\theta_1} \|\nabla_{\theta_1}^2 r_{\theta_1}(z, x, a)\| \leq L_{r,2} < \infty, \end{aligned}$$

$\forall(z, x, a) \in Z \times X \times A$ . Then,  $Q_{\theta_1}(z, x, a)$  and  $\bar{V}_{\theta_1}(z, x)$  are also twice continuously differentiable in  $\theta_1 \in \mathbb{R}^{p_1}$  and

$$\sup_{\theta_1} \|\nabla_{\theta_1}^2 Q_{\theta_1}(z, x, a)\| \leq L_Q, \quad \sup_{\theta_1} \|\nabla_{\theta_1}^2 \bar{V}_{\theta_1}(z, x)\| \leq L_{\bar{V}}$$

$\forall(z, x, a) \in Z \times X \times A$ , where

$$\begin{aligned} L_Q &:= \frac{1}{1-\beta} L_{r,2} + \frac{2\beta}{(1-\beta)^3} (L_{r,1})^2 \\ L_{\bar{V}} &:= \frac{1}{1-\beta} L_{r,2} + \frac{2}{(1-\beta)^3} (L_{r,1})^2. \end{aligned}$$

**Theorem 4.** Under the same assumptions of Lemma 1, the pseudo-log likelihood has Lipschitz continuous gradients with constant  $L := NT(L_Q + L_{\bar{V}})$ . With step size  $\rho < \frac{2}{L}$  it holds that

$$\min_{k \in \{1, \dots, K\}} \left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\|^2 \leq \frac{1}{\rho \left(1 - \frac{\rho L}{2}\right)} \frac{\hat{\ell}(\theta_1^*) - \hat{\ell}(\theta_1^0)}{K},$$

where  $\theta_1^*$  maximizes pseudo-log likelihood.

## V. MODEL IDENTIFICATION

The structure of the random reward POMDP model is defined by parameters :

$$\{r_{\theta_1}(Z, S, A), P_{\theta_2}(Z, S, A), \mu, \beta\},$$

where  $r_{\theta_1}(Z, S, A) \equiv \{r_{\theta_1}(z, s, a) : z \in Z, s \in S, a \in A\}$ ,  $P_{\theta_2}(Z, S, A) \equiv \{P_{\theta_2}(z', s'|z, s, a), z, z' \in Z, s, s' \in S, a \in A\}$ . For the rest of the section, we assume  $\mu$  and  $\beta$  are given and known. Then the conditional observation probabilities  $\{\sigma_{\theta_2}(z_{t+1}, z_t, x_t, \theta_2, a_t)\}$  and conditional choice probabilities  $\{\pi_{\theta}(a_t|z_t, x_t, \theta_2)\}$  are called the reduced form observation probabilities and choice probabilities under structure  $\theta = (\theta_1, \theta_2)$ . Under the true structure  $\theta^* = (\theta_1^*, \theta_2^*)$ , we must have

$$\forall(z_{t+1}, \zeta_t, a_t), \quad \underbrace{\hat{P}(z_{t+1}|\zeta_t, a_t)}_{\text{Data}} = \underbrace{\sigma_{\theta_2^*}(z_{t+1}, z_t, x_t, \theta_2^*, a_t)}_{\text{Model}}, \quad (\text{V.1})$$

$$\underbrace{\hat{P}(a_t|\zeta_t)}_{\text{Data}} = \underbrace{\pi_{\theta^*}(a_t|z_t, x_t, \theta_2^*)}_{\text{Model}}. \quad (\text{V.2})$$

where  $\hat{P}(z_{t+1}|\zeta_t, a_t)$  and  $\hat{P}(a_t|\zeta_t)$  are functions of the data.

[22] defines the observational equivalence and identification as follows.

**Definition 1** (Observational equivalence). Let  $\Theta$  be the set of structures  $\theta$ , and let  $\overset{\circ}{\iff}$  be observational equivalence.  $\forall \theta, \theta' \in \Theta$ ,  $\theta \overset{\circ}{\iff} \theta'$  if and only if

$$\sigma_{\theta_2}(z_{t+1}, z_t, x_t, \theta_2, a_t) = \sigma_{\theta_2'}(z_{t+1}, z_t, x_t, \theta_2', a_t)$$

and

$$\pi_{\theta}(a_t|z_t, x_t, \theta_2) = \pi_{\theta'}(a_t|z_t, x_t, \theta_2'), \forall z_{t+1}, z_t, a_t.$$

**Definition 2** (Identification). The model is identified if and only if  $\forall \theta, \theta' \in \Theta$ ,  $\theta \overset{\circ}{\iff} \theta'$  implies  $\theta = \theta'$ .

It is well known that the primitives of a MDP cannot be completely identified in general, and our POMDP model is not an exception. In addition, in the POMDP, the dynamics of the system under study cannot be directly observed as the system state is only partially observable. However, the next theorem shows that we could identify the hidden dynamics using two periods of data (including  $x_0$ ), assuming we know the cardinality of the state space.

**Theorem 5.** Assume  $|S|$  is known. The hidden dynamic  $P_{\theta_2}(Z, S, A)$  (not rank-1) can be uniquely identified from the first two periods of data (including  $x_0$ ).

Theorem 5 is crucial as it allows us to generalize the identification results in [12] and [22] for MDPs to POMDPs.

**Theorem 6.** The primitives of the POMDP model cannot be completely identified in general. However,  $\{r_{\theta_1}(Z, S, a), a \in A \setminus a^0\}$  and  $P_{\theta_2}(Z, S, A)$  can be uniquely identified from the data, given the initial belief  $x_0$ , the cardinality of the state space  $|S|$ , the discount factor  $\beta$ , the distribution of random shock  $\mu$ , and the rewards  $r_{\theta_1}(s, a^0), s \in S$  for a reference action  $a^0 \in A$  are all known.

*Remark 2.* It is well known that the knowledge on the discount factor  $\beta$ , random shock  $\mu$ , and the reference reward function is necessary to uniquely identify the primitives of a MDP model [12]. The identification result of the POMDP in Theorem 6 only requires two additional mild conditions on the knowledge of the initial belief  $x_0$  and the cardinality of the state space, although the system dynamics is hidden and the state is partially observable. We examine the case where the initial belief  $x_0$  is unknown in Section V-A. It is an interesting future research question to examine whether or under what conditions the POMDP model is identifiable if  $|S|$  is unknown. In practice, the number of possible states can be obtained by domain knowledge for a particular application. For example, the possible stages of a cancer or system degradation are likely obtainable. A practitioner can also try possible values of  $|S|$  to examine which value can best explain the observed behaviors.

**Corollary 1.** *For  $T < \infty$ , the POMDP model can be uniquely identified from the data, if  $x_0, \mu, |S|$  and both the reward structure and the terminal value function  $Q_T$  in the reference action  $a^0 \in A$  are all known.*

#### A. Sensitivity to A Priori Distribution Specification

The requirement of knowing the agent's initial belief  $x_0$  seems to limit the potential use of the developed estimation approach. However, we now show that the effect of  $x_0$  on the estimation result is decreasing with increasing length of the history  $\zeta_T$ . Specifically, let  $\mathcal{D} : X \times X \mapsto \mathbb{R}^+$  be a metric on  $X$ , defined as

$$\mathcal{D}(x, x') \triangleq \max\{d(x, x'), d(x', x)\}, \quad (\text{V.3})$$

where

$$d(x, x') \triangleq 1 - \min \left\{ \frac{x(s)}{x'(s)} : s \in S, x'(s) > 0 \right\}, \forall x, x' \in X.$$

Define

$$\begin{aligned} \eta(P_{\theta_2}(z', z, a)) \\ = \max\{\mathcal{D}(\lambda_{\theta_2}(z', z, e_i, a), \lambda_{\theta_2}(z', z, e_j, a)) : i, j \in S\}, \end{aligned} \quad (\text{V.4})$$

where  $e_i \in X$  with 1 on its  $i$ th element. [23] and [24] showed that  $\forall x_1, x_2 \in X$ ,

$$\mathcal{D}(\lambda_{\theta_2}(z', z, x_1, a), \lambda_{\theta_2}(z', z, x_2, a)) \leq \eta(P_{\theta_2}(z', z, a)) < 1,$$

where  $\eta(P_{\theta_2}(z', z, a))$  is called a *contraction coefficient* (coefficient of ergodicity) for substochastic matrix  $P_{\theta_2}(z', z, a)$ . Thus, given a finite history,  $\{z_t, \dots, z_{t-M}, a_{t-1}, \dots, a_{t-M}\}$ , let  $\lambda^M$  be  $M$  applications of  $\lambda$  function for any  $x_M \in X$ , namely,

$$\begin{aligned} \lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}) &= \lambda_{\theta_2}(z_t, z_{t-1}, \lambda_{\theta_2}(\dots \\ &\quad \lambda_{\theta_2}(z_{t-M+1}, z_{t-M}, x_{t-M}, a_{t-M})), a_{t-1}), \end{aligned}$$

where

$$z_{t-M}^t = \{z_t, \dots, z_{t-M}\}, a_{t-M}^{t-1} = \{a_{t-1}, \dots, a_{t-M}\}$$

for short. Then

$$\begin{aligned} \mathcal{D}(\lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}), \lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x'_{t-M})) \\ \leq (\eta(P_{\theta_2}(z', z, a)))^M, \forall x_{t-M}, x'_{t-M} \in X \end{aligned}$$

(see Section 2.3 in [24]), showing that the effect of  $x_0$  decreases as  $M$  increases.

**Theorem 7.** *Assume  $|S|$  is known. A set of model primitives  $\theta$  can be obtained from the data (consistent with an unknown  $x_0 \in X$ ), given the discount factor  $\beta$ , the random shock distribution  $\mu$ , and the reward  $r_{\theta_1}(S, a^0)$ . The set of estimators will shrink to the singleton true value as  $M \rightarrow \infty$  (hence,  $T \rightarrow \infty$ ).*

Theorem 7 shows that if  $x_0$  is uncertain (to the modeler), we can still estimate a set of  $\theta_2$  (hence  $\theta_1$ ) by repeated applications of  $\lambda$  function and utilizing the entire data provided by the information sequence  $\zeta_T$  (not just the first two-period data), and by varying  $x_{t-M} \in X$ . The set of estimated  $\theta$ s will shrink to the singleton true value as  $M$  goes to infinity. In addition, in many applications, it is also possible to obtain some (or a small range of) belief points, i.e.,  $x_0 \in X' \subset X$ . For example, in the engine replacement example, it is acceptable to reason that the state of a newly replace engine is good. This information of  $X'$  can be very helpful in improving the accuracy of the estimates.

## VI. ILLUSTRATION: OPTIMAL EQUIPMENT REPLACEMENT

We illustrate the estimation methodology using both synthetic and real data for a bus engine replacement problem. POMDP approaches have been widely used in machine maintenance problems ([2], [25]–[27]), where maintenance and engine replacement decisions must be made based upon monthly inspection results. Cumulative mileage and other specialized tests only provide informative signals about the true underlying engine's state condition which is not readily observable.

The engine deterioration state  $s_t \in S = \{0, 1\}$ , where “0” is being the “good state” and “1” is being the “bad state”. The available actions are  $a_t = 1$  is for engine replacement and  $a_t = 0$  for regular maintenance. The model for hidden state dynamics is

$$P_{\theta_2}(s_{t+1}|s_t, a_t = 0) = \begin{pmatrix} \theta_{2,0} & 1 - \theta_{2,0} \\ 1 - \theta_{2,1} & \theta_{2,1} \end{pmatrix}, \quad (\text{VI.1})$$

and  $P_{\theta_2}(s_{t+1} = 0|s_t, a_t = 1) = 1$ . Per 2500-mile maintenance costs are parametrized by  $\theta_{1,0}$  (in good state) and  $\theta_{1,1}$  (in bad state). With a belief  $x_t \in (0, 1)$  of the engine being in good state and  $z_t$  cumulative mileage after  $t$  months, the expected (monthly) maintenance cost is of the form

$$r_{\theta_1}(z_t, x_t, a = 0) = -0.001[(\theta_{1,0}z_t)x_t + (\theta_{1,1}z_t)(1 - x_t)],$$

and replacement cost is  $r_{\theta_1}(z_t, x_t, a = 1) = -RC$ . The distribution of monthly mileage increments  $\Delta \in \{0, 1, 2, 3\}$  is parametrized as follows:

$$\begin{aligned} P_{\theta_3}(z_{t+1} = z_t + \Delta | z_t, s_t = 0, a_t = 0) \\ = \theta_{3,0,\Delta}, \quad \Delta \in \{0, 1, 2\}, \\ P_{\theta_3}(z_{t+1} = z_t + \Delta | z_t, s_t = 0, a_t = 0) \\ = 1 - \theta_{3,0,0} - \theta_{3,0,1} - \theta_{3,0,2}, \quad \Delta = 3. \end{aligned} \quad (\text{VI.2})$$

Similarly, we define  $P_{\theta_3}(z_{t+1} = z_t + \Delta | z_t, s_t = 1, a_t = 0) = \theta_{3,1,\Delta}$ ,  $\Delta \in \{0, 1, 2, 3\}$ . Furthermore, after a replacement, the mileage restarts from zero:  $P_{\theta_3}(z_{t+1} = 0 | z_t, s_t, a_t = 1) = 1$ .

### A. Synthetic Dataset

We first show that the developed method can recover the model parameters using synthetic data. In addition, we show how misspecification errors can arise if an MDP-model (where the state is cumulative mileage) is used to fit data generated by a Bayesian agent in a partially observable environment.

**Setup.** To generate synthetic data, we simulate data with ground truth parameters in Table I. We simulate 3000 buses for 100 decision epochs. Specifically, we randomly generate initial belief  $x_0, s_0, z_0$ . The selected action is sampled from  $\pi_\theta(\cdot|x_t, z_t)$  on the basis of current belief  $x_t$  and current mileage  $z_t$ . Once the action is selected, the system state evolves and generates new mileage  $z_{t+1}$  according to system dynamics (VI.1)-(VI.2). The belief is updated by  $\lambda$  function based upon  $z_{t+1}, z_t, x_t, a_t$ . In this process, only  $z_t$  and  $a_t$  are recorded and both  $x_t$  (except  $x_0$ ) and  $s_t$  are discarded (see Algorithm 2).

---

#### Algorithm 2 Synthetic Data Generation from the POMDP Model.

---

```

record = empty holder
while i < 3000 do
  randomly generate  $x_0, z_0, s_0$ ;
  sample  $a_0$  from  $\pi_\theta(\cdot|x_0, z_0)$ ;
  record(i) = [ $x_0, z_0, a_0$ ];
  while t < 100 do
     $z_{t+1}, s_{t+1}$  sampled from  $P_{\theta_2, \theta_3}(s, z|s_t, z_t, a_t)$ ;
     $x_{t+1} = \lambda(z_{t+1}, z_t, x_t, a_t)$ ;
    sample  $a_{t+1}$  from  $\pi_\theta(\cdot|x_{t+1}, z_{t+1})$ ;
    record=record+[ $z_{t+1}, a_{t+1}$ ]; // record only  $z, a$ 
    t ← t + 1;
  end
  i ← i + 1;
end

```

---

**Estimation Results.** The estimation results for the POMDP model are presented in Table I. It shows that our algorithm can identify the model parameters accurately with maximal element-wise deviation of 0.006 in dynamics and 0.012 in reward. In addition, the prior knowledge on initial belief  $x_0$  does not influence the estimation result in a significant way (and only improve log-likelihood by 0.06%). Theorem 7 shows that if  $x_0$  is unknown, the resulting estimates may deviate from their true values. However, the difference between the estimated results and the true values quickly diminishes as the number of decision epochs increases. Fig. 2 clearly illustrates this fact (where the estimation deviation is caused by varying  $x_0 \in X$ ). For example, with only 8-period of data, the estimated dynamics are within 0.1 range of the its true value in 2-norm.

**Model Misspecification.** When the data is generated by a POMDP, using existing MDP-based models will lead to misspecification errors. To see this, we apply the model in [10] to the same synthetic dataset and present the estimation results in Table II. Unsurprisingly, the mis-specification error manifests

itself by a significant drop in log-likelihood ( $\frac{301750-262973}{301750} = 12.9\%$ ). In general, the modeling options for a given dataset include MDPs (possibly high-order MDPs), POMDP, or other non-Markovian processes. A central question for the modeler is to select an appropriate model which in the end may not necessarily be Markovian. In this regard, we note that [28] has recently developed a model selection procedure for testing the Markov assumption. The developed estimation approach for POMDPs can be an appealing alternative when the Markov assumption in sufficiently high order models is still rejected.

### B. Real Dataset

To illustrate the application of the developed methodology, we now revisit a subset of dataset reported in [10]. Specifically, Group 4 consisting of buses with 1975 GMC engines. Evidence of positive serial correlation in mileage increments is quite strong as the Durbin-Watson statistic is less than 1.13 for all buses except one with a value of 1.32. Thus, we fit the data by the POMDP-based model as in Section VI-A; see Table III.

**Discussion.** We compared our estimation results (from the POMDP model) with the results found in [10] (from a MDP-based model). In terms of log-likelihood, our POMDP-based model outperformed the MDP-based model by  $\frac{4495-3819}{3819} = 17.7\%$  (see Table III and Table IV).

Compared to the MDP model displayed in Table IV, the POMDP model also captures a feature of engine utilization: the distribution of mileage increments for engines considered in bad state is dominated (in the first-order stochastic sense) by the distribution of mileage increments of engines in good state. Furthermore, we find that the marginal operation costs ( $\theta_{1,1}$ ) for buses in *bad* state is significantly higher than those ( $\theta_{1,0}$ ) in *good* state (at least about two times, taking the standard deviation into consideration).

To gauge the economic interpretation of this result, we follow the same scaling procedure as in [10] to get the dollar estimates for  $\theta_{1,0}$  and  $\theta_{1,1}$ , respectively. All estimates are scaled with respect to reported (average) replacement cost (in 1985 US dollars) which for group 4 is \$7513 (see Table III in [10], p. 1005). The *perceived* average monthly maintenance costs increases  $\frac{7513}{RC}\theta_{1,0} = \$0.231$  per 2500 miles in good state and  $\frac{7513}{RC}\theta_{1,1} = \$1$  in bad state. That is, an engine with 300K miles in *good* condition has (according to the model) a monthly maintenance cost of  $(300/2.5) \times 0.231 = \$27.6$  whereas an engine with 300K miles in *bad* condition has monthly maintenance cost of  $(300/2.5) \times 1.00 = \$120$ .

## VII. CONCLUSIONS

In this paper, we developed a novel estimation method to recover the primitives of a POMDP model based on observable trajectories of the process. First, we provide a characterization of optimal decisions in a POMDP model with random rewards by mean of *soft* Bellman equation. We then developed a soft policy gradient algorithm to obtain the maximum likelihood estimator. We also show that the proposed model can be identified if the a priori belief distribution, the cardinality of the system state space, the distribution of random i.i.d shocks,

TABLE I. PARAMETER ESTIMATES AND LOG-LIKELIHOOD OF POMDP-BASED MODEL.

Parameter	$\theta_{3,0,0}$	$\theta_{3,0,1}$	$\theta_{3,0,2}$	$\theta_{3,1,0}$	$\theta_{3,1,1}$	$\theta_{3,1,2}$	$\theta_{2,0}$	$\theta_{2,1}$	$\theta_{1,0}$	$\theta_{1,1}$	$RC$
True value	0.039	0.333	0.590	0.181	0.757	0.061	0.949	0.988	0.2	1.2	9.243
$x_0$ known	0.038	0.327	0.596	0.181	0.754	0.064	0.950	0.987	0.2	1.2	9.231
	log-Likelihood: -262,973										
$x_0$ unknown	0.038	0.327	0.596	0.181	0.754	0.064	0.950	0.987	0.2	1.2	9.230
	log-Likelihood: -262,814										

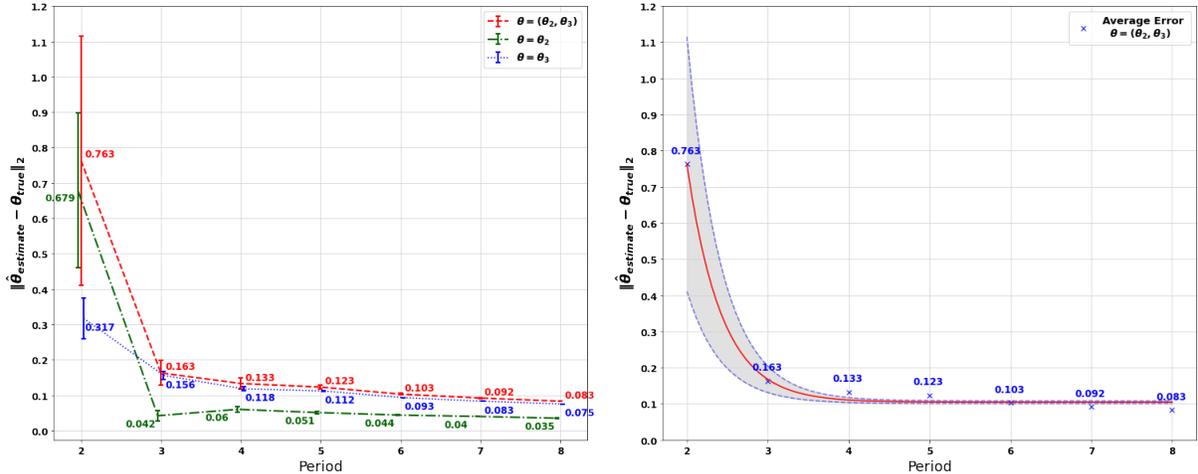


Fig. 2: Estimation results are affected by the prior knowledge of  $x_0$ . Without knowing  $x_0$ , more periods of data are needed for more accurate estimation. The unknown  $x_0$  induced deviation vanishes quickly as the number of data periods grows.

TABLE II. PARAMETER ESTIMATES AND LOG-LIKELIHOOD PROVIDED BY THE MDP MODEL FOR THE SYNTHETIC DATA.

Parameter	$\theta_{3,0}$	$\theta_{3,1}$	$\theta_{3,2}$	$\theta_{3,3}$	$\theta_1$	$RC$
MDP Model	0.128	0.601	0.257	0.015	1.1	9.811
	log-Likelihood: -301,750					

the discount factor, and the reward for a fixed reference action are given. Moreover, we show the estimation is robust to the specification of the a priori belief distribution provided the data sequence is sufficiently long (to address the case where the priori belief distribution is unknown). Finally we provide a numerical illustration with an application to optimal equipment replacement. With synthetic data, we show that highly accurate estimation of the true model primitives can be obtained despite having no prior knowledge of the underlying system dynamics. We also compared our POMDP approach to an MDP approach using a real data on the bus engine replacement problem. Our POMDP approach significantly improved the log-likelihood function and also revealed economically meaningful features that are conflated in the MDP model.

As this research represents a first effort on developing estimation methods for partially observable systems, future research directions are numerous. For example, computational

challenges of our model are obviously not trivial. POMDPs suffer from the well-known curse of dimensionality, and observations in many real applications can be high dimensional. Thus, a research direction is to address computational challenges of high dimensional hidden state models. This could be done for example via projection methods [29] or variational inference [30].

## VIII. APPENDIX

### A. Soft Bellman Equation for POMDPs

*Proof.* Proof of Theorem 1. The proof is by induction.

Assume  $U_{t+1,\theta}(\zeta_{t+1}, \epsilon_{t+1}) = V_{t+1,\theta}(z_{t+1}, x_{t+1}, \epsilon_{t+1})$ , then

$$\begin{aligned}
 & U_{t,\theta}(\zeta_t, \epsilon_t) \\
 &= \max_{a_t \in A} \left\{ \sum_{s_t} P_{\theta_2}(s_t | \zeta_t) r_{\theta_1}(z_t, s_t, a_t) + \epsilon_t(a_t) + \beta \sum_{z_{t+1}} \int \dots \right.
 \end{aligned}$$

TABLE III. PARAMETER ESTIMATES AND LOG-LIKELIHOOD PROVIDED BY THE POMDP-MODEL (STANDARD ERRORS OBTAINED BY BOOTSTRAPPING METHOD ARE IN PARENTHESES)

Parameter	$\theta_{3,0,0}$	$\theta_{3,0,1}$	$\theta_{3,0,2}$	$\theta_{3,1,0}$	$\theta_{3,1,1}$	$\theta_{3,1,2}$	$\theta_{2,0}$	$\theta_{2,1}$	$\theta_{1,0}$	$\theta_{1,1}$	$RC$
Good State	0.039 (.005)	0.335 (.018)	0.588 (.018)	*	*	*	0.949 (.004)	*	0.3 (.3)	*	9.738 (1.052)
Bad State	*	*	*	0.182 (.008)	0.757 (.008)	0.061 (.006)	*	0.988 (.002)	*	1.3 (.2)	
log-Likelihood	-3819										

TABLE IV. PARAMETER ESTIMATES AND LOG-LIKELIHOOD WITH MDP MODEL

Parameter	$\theta_{3,0}$	$\theta_{3,1}$	$\theta_{3,2}$	$\theta_{3,3}$	$\theta_1$	$RC$
MDP Model [10] p. 1022 (Standard errors in parentheses)	0.119 (0.005)	0.576 (0.008)	0.287 (0.007)	0.016 (0.002)	1.2 (0.3)	10.90 (1.581)
log-Likelihood	-4495					

$$\begin{aligned}
& \left. P_{\theta_2}(z_{t+1}|\zeta_t, a_t) V_{t+1, \theta}(z_{t+1}, x_{t+1}, \epsilon_{t+1}) d\mu(\epsilon_{t+1}|z_{t+1}) \right\} \\
= & \max_{a_t \in A} \left\{ r_{\theta_1}(z_t, x_t, a_t) + \epsilon_t(a_t) + \beta \sum_{z_{t+1}} \int \dots \right. \\
& \left. \sigma_{\theta_2}(z_{t+1}, z_t, x_t, a_t) V_{t+1, \theta}(z_{t+1}, \lambda(z_{t+1}, z_t, x_t, a_t), \epsilon_{t+1}) \dots \right. \\
& \left. d\mu(\epsilon_{t+1}|z_{t+1}) \right\} \\
= & V_{t, \theta}(z_t, x_t, \epsilon_t)
\end{aligned}$$

where the second equality follows from  $\sigma_{\theta_2}(z_{t+1}, z_t, x_t, a_t) = P_{\theta_2}(z_{t+1}|\zeta_t, a_t)$ .  $\square$

*Proof.* Proof of Theorem 2. Assume the following regularity conditions:

**R.1** (Bounded Upper Semicontinuous) For each  $a \in A$ ,  $r_{\theta}(z, x, a)$  is upper semicontinuous in  $z$  and  $x$  with bounded expectation and

$$h_{\theta}(z, x) \triangleq \sum_{t=1}^{\infty} \beta^t h_{t, \theta}(z, x) < \infty,$$

$$h_{1, \theta}(z, x) = \max_{a \in A} \sum_{z' \in Z} \sigma_{\theta_2}(z', z, x, a) \int \max_{a' \in A} \{r_{\theta}(z', x', a') + \epsilon'(a')\} d\mu(\epsilon'|z'),$$

$$h_{t, \theta}(z, x) = \max_{a \in A} \sum_{z' \in Z} \sigma_{\theta_2}(z', z, x, a) h_{t-1}(z', x');$$

where  $x' = \lambda_{\theta_2}(z', z, x, a)$ .

**R.2** (Weakly Continuous) The stochastic kernel

$$\sigma_{\theta_2}(\cdot, z, x, a) = \{\sigma_{\theta_2}(z', z, x, a)\}_{z' \in |Z|}$$

is weakly continuous in  $Z \times X \times A$ ;

**R.3** (Bounded Expectation) The reward function  $r_{\theta} \in \mathcal{B}$  and

for each  $Q \in \mathcal{B}$ ,  $E_{\theta}Q \in \mathcal{B}$ , where

$$\begin{aligned}
[E_{\theta}Q](z, x, a) &= \sum_{z' \in Z} \sigma_{\theta_2}(z', z, x, a) \\
&\quad \times \int \max_{a \in A} \{Q(z', x', a) + \epsilon(a)\} d\mu(\epsilon|z')
\end{aligned}$$

where  $x' = \lambda_{\theta_2}(z', z, x, a)$ .

Under these regularity conditions,  $H_{\theta} : \mathcal{B} \rightarrow \mathcal{B}$  is well defined (see related discussion in [11]).  $\forall Q, Q' \in \mathcal{B}$ ,  $\forall a$ , we have

$$\begin{aligned}
& \|H_{\theta}Q - H_{\theta}Q'\| \\
& \leq \beta \sum_{z'} \sigma_{\theta_2}(z', z, x, a) \int |\max_{a \in A} \{Q(z', x', a) + \epsilon(a)\} \\
& \quad - \max_{a \in A} \{Q'(z', x', a) + \epsilon(a)\}| d\mu(\epsilon|z') \\
& \leq \beta \sum_{z'} \sigma_{\theta_2}(z', z, x, a) \int \max_{a \in A} \{|Q(z', x', a) \\
& \quad - Q'(z', x', a)\}| d\mu(\epsilon|z') \\
& \leq \beta \|Q - Q'\|.
\end{aligned}$$

Hence,  $H_{\theta}$  is a contraction mapping and  $Q_{\theta}$  is the fixed point of  $H_{\theta}$ . Note that the controlled process  $\{z_{t+1}, z_t, x_t, a_t\}$  is Markovian because the conditional probability of  $a_t$  is  $P(a_t|z_t, x_t)$ , the conditional probability of  $x_{t+1}$  is provided by  $\lambda(z_{t+1}, z_t, x_t, a_t)$ , and the conditional probability of  $z_{t+1}$  is given by  $\sigma(z_{t+1}, z_t, x_t, a_t)$ . In addition,

$$\left\| \frac{\partial(\max_{a \in A} [Q_{\theta}(z, x, a) + \epsilon(a)])}{\partial Q_{\theta}(z, x, a)} \right\| \leq 1$$

for almost all  $\epsilon$ , by the Lebesgue dominated convergence

theorem,

$$\begin{aligned}\frac{\partial \bar{V}_\theta}{\partial Q_\theta} &= \int \left( \frac{\partial (\max_{a \in A} [Q_\theta(z, x, a) + \epsilon(a)])}{\partial Q_\theta(z, x, a)} \right) d\mu(\epsilon|z) \\ &= \int I\{a = \arg \max_{a \in A} [Q_\theta(z, x, a) + \epsilon(a)]\} d\mu(\epsilon|z) \\ &= \pi_\theta(a|z, x).\end{aligned}$$

□

*Proof.* Proof of Theorem 3. The result follows by Theorem 2 and [31] that

$$E_\epsilon[\max_a Q_\theta(z, x, a) + \epsilon_a] = \gamma + \ln \left( \sum_{a \in A} \exp(Q_\theta(z, x, a)) \right).$$

□

### B. Convergence of Soft Policy Gradient Algorithm

*Proof.* Proof of Lemma 1.

For fixed  $\theta_1 \in \mathbb{R}^{p_1}$ , consider the mapping  $G_{\theta_1}^1 : \mathcal{B} \mapsto \mathcal{B}$  defined as:

$$\begin{aligned}[G_{\theta_1}^1 g](z, x, a) &= \nabla_{\theta_1} r_{\theta_1}(z, x, a) \\ &\quad + \beta \sum_{z'} \sigma_{\hat{\theta}_2}(z', z, x, a) \sum_{a'} \pi_{\theta_1}(a'|z', x') g(z', x', a'),\end{aligned}$$

where  $x' = \lambda_{\hat{\theta}_2}(z', z, x, a)$ . It follows that  $G_{\theta_1}^1$  is a contraction map with unique fixed point  $\nabla_{\theta_1} Q_{\theta_1}$  and  $\|\nabla_{\theta_1} Q_{\theta_1}(z, x, a)\| \leq \frac{1}{1-\beta} L_{r,1}$ . Since

$$\nabla_{\theta_1} \bar{V}_{\theta_1}(z', x') = \sum_{a'} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1} Q_{\theta_1}(z', x', a')$$

it follows that  $\nabla_{\theta_1} \bar{V}_{\theta_1}$  also exists and  $\|\nabla_{\theta_1} \bar{V}_{\theta_1}\| \leq \frac{1}{1-\beta} L_{r,1}$ .

Consider the mapping  $G_{\theta_1}^2 : \mathcal{B} \mapsto \mathcal{B}$  defined as follows:

$$\begin{aligned}[G_{\theta_1}^2 g](z, x, a) &= \nabla_{\theta_1}^2 r_{\theta_1}(z, x, a) + \beta \sum_{z'} \sigma_{\hat{\theta}_2}(z', z, x, a) \\ &\quad \times \sum_{a'} \nabla_{\theta_1} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1} Q_{\theta_1}(z', x', a') \\ &\quad + \beta \sum_{z'} \sigma_{\hat{\theta}_2}(z', z, x, a) \sum_{a'} \pi_{\theta_1}(a'|z', x') g(z', x', a'),\end{aligned}$$

where  $x' = \lambda_{\hat{\theta}_2}(z', z, x, a)$ . Thus,  $G_{\theta_1}^2$  is a contraction map with unique fixed point  $\nabla_{\theta_1}^2 Q_{\theta_1}$ . Since

$$\begin{aligned}\nabla_{\theta_1}^2 \bar{V}_{\theta_1}(z', x') &= \sum_{a'} \nabla_{\theta_1} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1} Q_{\theta_1}(z', x', a') \\ &\quad + \sum_{a'} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1}^2 Q_{\theta_1}(z', x', a'),\end{aligned}\quad (\text{VIII.1})$$

it follows that  $\nabla_{\theta_1}^2 \bar{V}_{\theta_1}$  exists.

From the fixed point characterization of  $\nabla_{\theta_1}^2 Q_{\theta_1}$  we obtain:

$$\begin{aligned}\|\nabla_{\theta_1}^2 Q_{\theta_1}(z, x, a)\| &\leq \frac{1}{1-\beta} \left( \|\nabla_{\theta_1}^2 r_{\theta_1}(z, x, a)\| \right. \\ &\quad \left. + \beta \left\| \sum_{z'} \sigma_{\hat{\theta}_2}(z', z, x, a) \right. \right. \\ &\quad \left. \left. \times \sum_{a'} \nabla_{\theta_1} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1} Q_{\theta_1}(z', x', a') \right\| \right).\end{aligned}\quad (\text{VIII.2})$$

Note that from the softmax structure of conditional choice probabilities,

$$\begin{aligned}\nabla_{\theta_1} \pi_{\theta_1}(a|z, x) &= \pi_{\theta_1}(a|z, x) \left[ \nabla_{\theta_1} Q_{\theta_1}(z, x, a) \right. \\ &\quad \left. - \sum_{a'} \pi_{\theta_1}(a'|z, x) \nabla_{\theta_1} Q_{\theta_1}(z, x, a') \right].\end{aligned}$$

Hence, the second term on the right hand side of (VIII.2) can be bounded by

$$\begin{aligned}\left\| \sum_{z'} \sigma_{\hat{\theta}_2}(z', z, x, a) \sum_{a'} \nabla_{\theta_1} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1} Q_{\theta_1}(z', x', a') \right\| \\ \leq \frac{2}{(1-\beta)^2} (L_{r,1})^2.\end{aligned}\quad (\text{VIII.3})$$

Combining (VIII.2) and (VIII.3) we obtain:

$$\|\nabla_{\theta_1}^2 Q_{\theta_1}(z, x, a)\| \leq \frac{1}{1-\beta} L_{r,2} + \frac{2\beta}{(1-\beta)^3} (L_{r,1})^2 = L_Q,$$

Similarly, from (VIII.1) we obtain:

$$\begin{aligned}\|\nabla_{\theta_1}^2 \bar{V}_{\theta_1}\| \\ \leq \left\| \sum_{a'} \nabla_{\theta_1} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1} Q_{\theta_1}(z', x', a') \right\| \\ \quad + \left\| \sum_{a'} \pi_{\theta_1}(a'|z', x') \nabla_{\theta_1}^2 Q_{\theta_1}(z', x', a') \right\| \\ \leq \left\| \sum_{a'} \pi_{\theta_1}(a'|z', x') [\nabla_{\theta_1} Q_{\theta_1}(z', x', a') - \sum_{a''} \pi_{\theta_1}(a''|z', x') \right. \\ \quad \left. \times \nabla_{\theta_1} Q_{\theta_1}(z', x', a'')] \nabla_{\theta_1} Q_{\theta_1}(z', x', a') \right\| + L_Q \\ \leq \frac{2}{(1-\beta)^3} (L_{r,1})^2 + \frac{1}{1-\beta} L_{r,2} \\ = L_{\bar{V}}.\end{aligned}$$

□

*Proof.* Proof of Theorem 4. Recall that

$$\begin{aligned}\nabla_{\theta_1}^2 \hat{\ell}(\theta_1) &= \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta_1}^2 \log \pi_{\theta_1}(a_{t,i}|z_{t,i}, x_{t,i}) \\ &= \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta_1}^2 Q_{\theta_1}(z_{t,i}, x_{t,i}, a_{t,i}) - \nabla_{\theta_1}^2 \bar{V}_{\theta_1}(z_{t,i}, x_{t,i})\end{aligned}$$

By Lemma 1, it follows that  $\|\nabla_{\theta_1}^2 \hat{\ell}(\theta_1)\| \leq L$  with

$$L := NT(L_Q + L_{\bar{V}})$$

Or equivalently,  $\nabla_{\theta_1} \hat{\ell}(\theta_1)$  is Lipschitz continuous in  $\theta_1$  with constant  $L$ .

By Lipschitz continuous gradients,

$$\begin{aligned}\hat{\ell}(\theta_1^{k+1}) &\geq \hat{\ell}(\theta_1^k) + \nabla \hat{\ell}(\theta_1^k)^\top (\theta_1^{k+1} - \theta_1^k) - \frac{L}{2} \|\theta_1^{k+1} - \theta_1^k\|^2 \\ &= \hat{\ell}(\theta_1^k) + \rho \left(1 - \frac{\rho L}{2}\right) \left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\|^2.\end{aligned}$$

Hence,

$$\rho \left(1 - \frac{\rho L}{2}\right) \left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\|^2 \leq \hat{\ell}(\theta_1^{k+1}) - \hat{\ell}(\theta_1^k).$$

Adding over  $k = 1, \dots, K$  we obtain

$$\rho \left(1 - \frac{\rho L}{2}\right) \sum_{k=1}^K \left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\|^2 \leq \hat{\ell}(\theta_1^K) - \hat{\ell}(\theta_1^0) \leq \hat{\ell}(\theta_1^*) - \hat{\ell}(\theta_1^0)$$

where  $\theta_1^*$  is a maximizer of log-likelihood. It follows that

$$\begin{aligned}\min_{k \in \{1, \dots, K\}} \left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\|^2 &\leq \frac{1}{K} \sum_{k=1}^K \left\| \nabla_{\theta_1} \hat{\ell}(\theta_1^k) \right\|^2 \\ &\leq \frac{1}{\rho \left(1 - \frac{\rho L}{2}\right)} \frac{\hat{\ell}(\theta_1^*) - \hat{\ell}(\theta_1^0)}{K}.\end{aligned}$$

□

### C. Identification Results

*Proof.* Proof of Theorem 5. For any given two system dynamics  $P_{\theta_2}(z', z, a)$ ,  $P_{\theta_2'}(z', z, a)$ , where  $P(z', z, a) = \{P_{ij}(z', z, a)\}$ ,  $P_{ij}(z', z, a) = P(z_{t+1} = z', s_{t+1} = j | z_t = z, s_t = i, a_t = a)$ , we show that they can be distinguished by the data. Since the dataset contains  $x_0, z_0, a_0$  and  $|S|$  is known, we can obtain

$$\begin{aligned}\sigma_{\theta_2}^0(z_1, z_0, x_0, a_0) &= \sum_{s_1} \sum_s x_0(s) P_{\theta_2}(z_1, s_1 | z_0, s, a_0) \\ &= \sum_s x_0(s) P_{\theta_2}(z_1 | z_0, s, a_0),\end{aligned}$$

and

$$\begin{aligned}\sigma_{\theta_2'}^0(z_1, z_0, x_0, a_0) &= \sum_{s_1} \sum_s x_0(s) P_{\theta_2'}(z_1, s_1 | z_0, s, a_0) \\ &= \sum_s x_0(s) P_{\theta_2'}(z_1 | z_0, s, a_0).\end{aligned}$$

Note that

$$\sigma^0(z_1, z_0, x_0, a_0) = \hat{P}(z_1 | \zeta_0, a_0),$$

where  $\hat{P}(z_1 | \zeta_0, a_0)$  is a function of the first period data. Thus,  $\sigma$  can be obtained from the data and  $\sigma_{\theta_2}^0 = \sigma_{\theta_2'}^0$  if and only if  $P_{\theta_2}(z' | z, s, a) = P_{\theta_2'}(z' | z, s, a)$ . If  $\sigma_{\theta_2}^0 \neq \sigma_{\theta_2'}^0$ , we are done. However, it is possible that there exists  $s'$  such that  $P_{\theta_2}(z', s' | z, s, a) \neq P_{\theta_2'}(z', s' | z, s, a)$  but  $P_{\theta_2}(z' | z, s, a) = P_{\theta_2'}(z' | z, s, a)$ . In this case, update belief by Eq. (III.4),

$$x_{1, \theta_2} = \lambda_{\theta_2}(z_1, z_0, x_0, a_0) = \frac{x_0 P_{\theta_2}(z_1, z_0, a_0)}{\sigma_{\theta_2}^0(z_1, z_0, x_0, a_0)},$$

$$x_{1, \theta_2'} = \lambda_{\theta_2'}(z_1, z_0, x_0, a_0) = \frac{x_0 P_{\theta_2'}(z_1, z_0, a_0)}{\sigma_{\theta_2'}^0(z_1, z_0, x_0, a_0)}.$$

Then  $x_{1, \theta_2} = x_{1, \theta_2'}$  if and only if  $P_{\theta_2}(z', z, a) = P_{\theta_2'}(z', z, a)$ . Now,  $\sigma_{\theta_2}^1(z_2, z_1, x_{1, \theta_2}, a_1) = \sum_s x_{1, \theta_2}(s) P_{\theta_2}(z_2 | z_1, s, a_1)$  and  $\sigma_{\theta_2'}^1(z_2, z_1, x_{1, \theta_2'}, a_1) = \sum_s x_{1, \theta_2'}(s) P_{\theta_2'}(z_2 | z_1, s, a_1)$ , and again  $\sigma^1$  is obtainable from the two periods of data as

$$\sigma^1(z_2, z_1, x_1, a_1) = \hat{P}(z_2 | \zeta_1, a_1).$$

Now,  $\sigma_{\theta_2}^1 = \sigma_{\theta_2'}^1$  if and only if  $x_{1, \theta_2} = x_{1, \theta_2'}$ , indicating  $P_{\theta_2}(z', z, a) = P_{\theta_2'}(z', z, a)$  assuming  $P(z', z, a)$  is not rank-1. □

*Proof.* Proof of Theorem 6. By Theorem 1 and Theorem 5, both  $\pi(a | z_t, x_t) = \hat{P}(a | \zeta_t)$  and hidden dynamics  $\{P_{\theta_2}(z', z, a)\}$  can be identified from the data. Treating belief as the state and since  $|\zeta_t| < \infty$ , Proposition 1 in [12] and [22] show that there is a one-to-one mapping  $q(Z, X) : R^{|A|} \rightarrow R^{|A|}$ , only depending on  $\mu$ , which maps the choice probability set  $\{\pi_\theta(a | z, x)\}$  to the set of the difference in action-specific value function  $\{Q_\theta(z, x, a) - Q_\theta(z, x, a^0)\}$ , namely,

$$Q_\theta(z, x, a) - Q_\theta(z, x, a^0) = q_a(\{\pi_\theta(a' | z, x)\}; \mu), \quad (\text{VIII.4})$$

$q_0(\cdot) = 0$ , and  $q = (q_0, \dots, q_{|A|-1})$ . Thus, if we know  $Q_\theta(z, x, a^0)$ , we can recover  $Q_\theta(z, x, a), \forall a \in A \setminus a^0$ . Note that

$$\begin{aligned}Q_\theta(z, x, a) &= r_{\theta_1}(z, x, a) + \beta E_{z' | z, x, a} \left[ E_{\epsilon' | z'} \max_{a' \in A} \{ \dots \right. \\ &\quad \left. Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a') + \epsilon'(a') \} \right] \\ &= r_{\theta_1}(z, x, a) \\ &\quad + \beta E_{z' | z, x, a} \left[ E_{\epsilon' | z'} \left[ \max_{a' \in A} \{ Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a') \right. \right. \\ &\quad \left. \left. - Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a^0) + \epsilon'(a') \} \right. \right. \\ &\quad \left. \left. + Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a^0) \right] \right] \quad (\text{VIII.5})\end{aligned}$$

$$\begin{aligned}&= r_{\theta_1}(z, x, a) \\ &\quad + \beta E_{z' | z, x, a} \left[ E_{\epsilon' | z'} \max_{a' \in A} \{ Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a') \right. \\ &\quad \left. - Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a^0) + \epsilon'(a') \} \right] \\ &\quad + \beta E_{z' | z, x, a} [Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a^0)] \quad (\text{VIII.6})\end{aligned}$$

Because of the mapping  $q$  in Eq. (VIII.4),  $\pi_\theta(a | z_t, x_t) = \hat{P}(a | \zeta_t)$ , and that both the hidden dynamic  $\{P_{\theta_2}(z', z, a)\}$  and  $\mu$  are known, the quantity

$$\begin{aligned}C &= \beta E_{z' | z, x, a} \left[ E_{\epsilon' | z'} \max_{a' \in A} \{ Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a') \right. \\ &\quad \left. - Q_\theta(z', \lambda_{\theta_2}(z', z, x, a), a^0) + \epsilon'(a') \} \right]\end{aligned}$$

is known. Under the assumption of  $r_{\theta_1}(Z, S, a^0) = 0$ , we have

$$Q_\theta(z, x, a^0) = C + \beta E_{z' | z, x, a^0} [Q_\theta(z', \lambda_{\theta_2}(z', z, x, a^0), a^0)] \quad (\text{VIII.7})$$

with only unknown  $Q_\theta(Z, X, a^0)$ . It is easy to see there is a unique solution to Eq. (VIII.7) due to the contraction

mapping theorem. Consequently, all  $Q_\theta(Z, X, a)$ ,  $a \in A$  can be recovered by Eq. (VIII.4). Lastly,

$$r_{\theta_1}(z, x, a) = Q_\theta(z, x, a) - C \\ - \beta E_{z'|z, x, a^0} [Q_\theta(z', \lambda_{\theta_2}(z', z, x, a^0), a^0)].$$

As

$$r_{\theta_1}(z, x, a) = \sum_{s \in S} x(s) r_{\theta_1}(z, s, a), \{r_{\theta_1}(s, a) : s \in S, a \in A\}$$

can be uniquely determined.  $\square$

*Proof.* Proof of Corollary 1. When  $Q_{T, \theta}(Z, X, a^0)$  and  $r_{\theta_1}(Z, S, a^0)$  are known, we can obtain  $Q_{t, \theta}(Z, X, a^0)$  via

$$Q_{t, \theta}(z, x, a^0) = r_{\theta_1}(z, x, a^0) \\ + \beta E_{z'|z, x, a^0} \left[ E_{\epsilon'|z'} \max_{a' \in A} \{Q_{t+1, \theta}(z', \lambda_{\theta_2}(z', z, x, a^0), a') \\ + \epsilon'(a') - Q_{t+1, \theta}(z', \lambda_{\theta_2}(z', z, x, a^0), a^0)\} \right] \\ + \beta E_{z'|z, x, a^0} [Q_{t+1, \theta}(z', \lambda_{\theta_2}(z', z, x, a^0), a^0)] \quad (\text{VIII.8})$$

The rest follows exactly as in the proof of Theorem 6.  $\square$

*Proof.* Proof of Theorem 7. The proof of Theorem 5 shows that  $\theta_2$  can be uniquely determined by

$$\sigma_{\theta_2}(z_{t+1}, z_t, \lambda_{\theta_2}(z_t, z_{t-1}, x_{t-1}^*, a_{t-1}), a_t) = \hat{P}(z_{t+1} | \zeta_t, a_t),$$

given the true value  $x_{t-1}^*$  is known for  $\zeta_t$ . Namely, there exists a unique  $\theta_2^*$  such that

$$\sigma_{\theta_2^*}(z_{t+1}, z_t, \lambda_{\theta_2^*}(z_t, z_{t-1}, x_{t-1}^*, a_{t-1}), a_t) = \hat{P}(z_{t+1} | \zeta_t, a_t),$$

and  $\forall \theta_2 \neq \theta_2^*$ , there is an  $\epsilon > 0$  such that

$$D(\sigma_{\theta_2}(\cdot, z_t, \lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}^*), a_t), \hat{P}(\cdot | \zeta_t, a_t)) \geq \epsilon,$$

where  $x_{t-M}^*$  is the true belief at  $t-M$ . Due to the contraction coefficient  $\eta(P_{\theta_2}(z', z, a)) < 1$ , we have

$$D(\lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}^*), \lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}^*)) \leq \eta^M,$$

where  $\eta = \max_{z', z', a} \eta(P_{\theta_2}(z', z, a)) < 1$  since  $|Z| < \infty, |A| < \infty$ . Thus,  $\forall \theta_2 \neq \theta_2^*$ , we have

$$\lim_{M \rightarrow \infty} D(\sigma_{\theta_2}(\cdot, z_t, \lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}^*), a_t), \hat{P}(\cdot | \zeta_t, a_t)) \\ = \lim_{M \rightarrow \infty} D(\sigma_{\theta_2}(\cdot, z_t, \lambda_{\theta_2}^M(z_{t-M}^t, a_{t-M}^{t-1}, x_{t-M}^*), a_t), \hat{P}(\cdot | \zeta_t, a_t)) \\ \geq \epsilon, \forall x_{t-M} \in X,$$

indicating  $\theta_2$  can be distinguished by the data given  $M$  is sufficiently large. Once  $\theta_2$  is determined,  $\theta_1$  can be determined by Theorem 6, which completes the proof.  $\square$

## REFERENCES

- [1] R. Smallwood and E. Sondik, "The Optimal Control of Partially Observable Markov Processes over a Finite Horizon," *Operations Research*, vol. 21, no. 5, pp. 1071–1088, 1973.
- [2] G. E. Monahan, "State of the Art—A Survey of Partially Observable Markov Decision Processes: Theory, Models, and Algorithms," *Management Science*, vol. 28, no. 1, pp. 1–16, 1982.
- [3] C. C. White, "A survey of solution techniques for the partially observed Markov decision process," *Annals of Operations Research*, vol. 32, pp. 215–230, 1991.
- [4] V. Krishnamurthy, *Partially observed Markov decision processes*. Cambridge University Press, 2016.
- [5] L. Winterer, S. Junges, R. Wimmer, N. Jansen, U. Topcu, J. P. Katoen, and N. Becker, "Strategy synthesis for POMDPs in robot planning via game-based abstractions," *IEEE Transactions on Automatic Control*, vol. 66, pp. 1040–1054, 2021.
- [6] D. Silver and J. Veness, "Monte-Carlo planning in large POMDPs," *Advanced in Neural Information Process Systems*, vol. 23, pp. 2164–2172, 2010.
- [7] V. Krishnamurthy and D. Djonin, "Structured threshold policies for dynamic sensor scheduling - a POMDP approach," *IEEE Trans Signal Processing*, vol. 55, pp. 4938–4957, 2007.
- [8] A. Foka and P. Trahanias, "Real-time hierarchical POMDPs for autonomous robot navigation," *Robotics and Autonomous Systems*, vol. 55, pp. 561–571, 2007.
- [9] S. Young, M. Gasic, S. Keizer, F. Mairesse, J. Schatzmann, B. Thomson, and K. Yu, "The hidden information state model: A practical framework for POMDP-based spoken dialogue management," *Computer Speech & Language*, vol. 24, pp. 150–174, 2010.
- [10] J. Rust, "Optimal replacement of GMC bus engines: An empirical model of Harold Zurcher," *Econometrica*, pp. 999–1033, 1987.
- [11] J. Rust, "Maximum likelihood estimation of discrete choice processes," *SIAM Journal on Control and Optimization*, vol. 26, no. 5, pp. 1006–1024, 1988.
- [12] V. J. Hotz and R. A. Miller, "Conditional choice probabilities and the estimation of dynamic models," *Review of Economic Studies*, vol. 60, no. 3, pp. 497–529, 1993.
- [13] V. J. Hotz, R. A. Miller, S. Sanders, and J. Smith, "A simulation estimator for dynamic models of discrete choice," *The Review of Economic Studies*, vol. 61, pp. 265–289, 1994.
- [14] V. Aguirregabiria and P. Mira, "Swapping the nested fixed point algorithm: A class of estimators for discrete Markov decision models," *Econometrica*, vol. 70, pp. 1519–1543, 2002.
- [15] V. Aguirregabiria and P. Mira, "Dynamic discrete choice structural models: A survey," *Journal of Econometrics*, vol. 156, pp. 38–67, 2010.
- [16] C. L. Su and K. L. Judd, "Constrained optimization approaches to estimation of structural models," *Econometrica*, vol. 80, no. 5, pp. 2213–2230, 2012.
- [17] H. Kasahara and K. Shimotsu, "Estimation of discrete choice dynamic programming models," *Journal of Applied Econometrics*, vol. 69, no. 1, pp. 28–58, 2018.
- [18] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence*, pp. 1433–1438, 2008.
- [19] A. Boularias, J. Kober, and J. Peters, "Relative entropy inverse reinforcement learning," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, vol. 15, pp. 182–189, PMLR, 2011.
- [20] C. Finn, S. Levine, and P. Abbeel, "Guided cost learning: deep inverse optimal control via policy optimization," in *International conference on Machine learning*, 2016.
- [21] J. Choi and K. Kim, "Inverse reinforcement learning in partially observable environments," *Journal of Machine Learning Research*, vol. 12, pp. 691–730, 2011.
- [22] T. Magnac and D. Thesmar, "Identifying dynamic discrete choice processes," *Econometrica*, vol. 70, no. 2, p. 801–816, 2002.
- [23] L. K. Platzman, "Optimal infinite-horizon undiscounted control of finite probabilistic systems," *SIAM Journal on Control and Optimization*, vol. 18, pp. 362–380, 1980.
- [24] C. C. White and W. T. Scherer, "Finite-memory suboptimal design for partially observed Markov decision processes," *Operations Research*, vol. 42, no. 3, pp. 439–455, 1994.
- [25] J. E. Eckles, "Optimum maintenance with incomplete information," *Operations Research*, vol. 16, no. 5, pp. 1058–1067, 1968.
- [26] Z. Fang, X. Wang, and L. Wang, "Operational state evaluation and maintenance decision-making method for multi-state cnc machine tools based on partially observable Markov decision process," in *2020 International Conference on Sensing, Diagnostics, Prognostics, and Control (SDPC)*, pp. 120–124, IEEE, 2020.
- [27] W. Xu and L. Cao, "Optimal maintenance control of machine tools for energy efficient manufacturing," *The International Journal of Advanced Manufacturing Technology*, vol. 104, no. 9, pp. 3303–3311, 2019.
- [28] C. Shi, R. Wan, R. Song, W. Lu, and L. Leng, "Does the Markov decision process fit the data: Testing for Markov property in sequential decision making," *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [29] E. Zhou, M. C. Fu, and S. Marcus, "Solving continuous state POMDPs via density projection," *IEEE Transactions on Automatic Control*, vol. 55, pp. 1101–1116, 2010.

- [30] M. Igl, Z. L., A. Tuan, F. Wood, and S. Whiteson, “Deep variational reinforcement learning for POMDPs,” in *Proceedings of the 35th International Conference on Machine Learning*, PMLR, 2018.
- [31] K. Train, *Discrete choice methods with simulation*. Cambridge University Press, 2002.



**Yanling Chang** is an assistant professor in the Department of Engineering Technology and Industrial Distribution and the Department of Industrial and Systems Engineering, Texas A&M University, College Station, TX, USA. She received her Bachelor’s degree in Electronic and Information Science and Technology from Peking University, Beijing, China; the Master’s degree in Mathematics and the Ph.D. degree in Operations Research both from Georgia Institute of Technology, Atlanta, GA, USA. Her research interests include partially observable Markov

decision processes, game theory, and optimization.



**Alfredo Garcia** received the Degree in electrical engineering from the Universidad de los Andes, Bogotá, Colombia, in 1991, the Diplôme d’Etudes Approfondies in automatic control from the Université Paul Sabatier, Toulouse, France, in 1992, and the Ph.D. degree in industrial and operations engineering from the University of Michigan, Ann Arbor, MI, USA, in 1997. From 1997 to 2001, he was a consultant to government agencies and private utilities in the electric power industry. From 2001 to 2015, he was a Faculty with the Department of Systems and

Information Engineering, University of Virginia, Charlottesville, VA, USA. From 2015 to 2017, he was a Professor with the Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL, USA. In 2018, he joined the Department of Industrial and System Engineering, Texas A&M University, College Station, TX, USA. His research interests include game theory and dynamic optimization with applications in communications and energy networks.



**Zhide Wang** is a PhD candidate in the Department of Industrial and System Engineering, Texas A&M University, College Station, TX, USA. He received his Bachelor’s degree in Industrial Engineering and Management from Shanghai Jiao Tong University, Shanghai, China. His research interests include partially observable Markov decision processes and structural estimation of dynamic discrete choice models with applications in cognitive psychology.



**Lu Sun** is a PhD student in the Department of Industrial and System Engineering, Texas A&M University, College Station, TX, USA. She received her Bachelor’s degree in Mathematics from Beihang University, Beijing, China. Her research interests include partially observable Markov decision processes and game theory.