Inverse Decision Modeling: Learning Interpretable Representations of Behavior

Daniel Jarrett^{1*} Alihan Hüyük^{1*} Mihaela van der Schaar¹²

Abstract

Decision analysis deals with modeling and enhancing decision processes. A principal challenge in improving behavior is in obtaining a transparent description of existing behavior in the first place. In this paper, we develop an expressive, unifying perspective on inverse decision modeling: a framework for learning parameterized representations of sequential decision behavior. First, we formalize the *forward* problem (as a normative standard), subsuming common classes of control behavior. Second, we use this to formalize the inverse problem (as a descriptive model), generalizing existing work on imitation/reward learning-while opening up a much broader class of research problems in behavior representation. Finally, we instantiate this approach with an example (inverse bounded rational control), illustrating how this structure enables learning (interpretable) representations of (bounded) rationality-while naturally capturing intuitive notions of suboptimal actions, biased beliefs, and imperfect knowledge of environments.

1. Introduction

Modeling and enhancing decision-making behavior is a fundamental concern in computational and behavioral science, with real-world applications to healthcare [1], economics [2], and cognition [3]. A principal challenge in improving decision processes is in obtaining a transparent *understanding* of existing behavior to begin with. In this pursuit, a key complication is that agents are often *boundedly rational* due to biological, psychological, and computational factors [4–8], the precise mechanics of which are seldom known. As such, how can we intelligibly characterize imperfect behavior? Consider the "lifecycle" of decision analysis [9] in the real world. First, *normative analysis* deals with modeling rational decision-making. It asks the question: What constitutes ideal behavior? To this end, a prevailing approach is given by von Neumann-Morgenstern's expected utility theory, and the study of optimal control is its incarnation in sequential decision-making [10]. But judgment rendered by real-world agents is often imperfect, so *prescriptive analysis* deals with improving existing decision behavior. It asks the question: How can we move closer toward the ideal? To this end, the study of decision engineering seeks to design "human-in-theloop" techniques that nudge or assist decision-makers, such as medical guidelines and best practices [11]. Importantly, however, this first requires a quantitative account of current practices and the imperfections that necessitate correcting.

To take this crucial first step, we must therefore start with *descriptive analysis*—that is, with understanding observed decision-making from demonstration. We ask the question: What does existing behavior look like—relative to the ideal? Most existing work on imitation learning (i.e. to replicate expert actions) [12] and apprenticeship learning (i.e. to match expert returns) [13] offers limited help, as our objective is instead in understanding (i.e. to interpret imperfect behavior). In particular, beyond the utility-driven nature of rationality for agent behaviors, we wish to quantify intuitive notions of *boundedness*—such as the apparent flexibility of decisions, tolerance for surprise, or optimism in beliefs. At the same time, we wish that such representations be *interpretable*—that is, that they be projections of observed behaviors onto parameterized spaces that are meaningful and parsimonious.

Contributions In this paper, our mission is to explicitly relax normative assumptions of optimality when modeling decision behavior from observations.³ First, we develop an expressive, unifying perspective on *inverse decision modeling*: a general framework for learning parameterized representations of sequential decision-making behavior. Specifically, we begin by formalizing the *forward problem* F (as

¹Department of Applied Mathematics and Theoretical Physics, University of Cambridge, UK; ²Department of Electrical Engineering, University of California, Los Angeles, USA. *Authors contributed equally. Correspondence to: <daniel.jarrett@maths.cam.ac.uk>.

Proceedings of the 38th International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

³Our terminology is borrowed from economics: By "descriptive" models, we refer to those that capture *observable* decision-making behavior as-is (e.g. an imitator policy in behavioral cloning), and by "normative" models, we refer to those that specify *optimal* decision-making behavior (e.g. with respect to some utility function).

Table 1. Inverse Decision Modeling. Comparison of primary classes of imitation/reward learning (IL/IRL) versus our prototypical example (i.e. inverse bounded rational control) as instantiations of inverse decision modeling. Constraints on agent behavior include: [†]environment dynamics (extrinsic), and [‡]bounded rationality (intrinsic). Legend: deterministic (Det.), stochastic (Stoc.), subjective dynamics (Subj.), behavioral cloning (BC), distribution matching (DM), risk-sensitive (RS), partially-observable (PO), maximum entropy (ME). All terms/notation are developed over Sections 3–4.

	$\textbf{Extrinsic}^{\dagger}$		Intrinsic [‡]							
Inverse Decision Model	Partially Controllable	Partially Observable	Purposeful Behavior	Subjective Dynamics	Action Stochasticity	Knowledge Uncertainty	Decision Complexity	Specification Complexity	Recognition Complexity	Examples
	$\tau_{\rm env}$	ω_{env}	v	$_{ au,\omega}$	π	$ ho,\sigma$	α	β	η	
BC-IL Subj. BC-IL	1	\ \	X X	×	\ \	X X	X X	X X	X X	[14–21] [22]
Det. DM-IL Stoc. DM-IL	\ \	X X	X X	X X	×	X X	X X	X X	X X	[23,24] [25–39]
Det. IRL Stoc. IRL Subj. IRL RS-IRL	\ \ \ \	X X X X	\$ \$ \$ \$	× ×	×	× × × ×	X X X X	X X X X	X X X X	[40–46] [47–66] [67] [68,69]
Det. PO-IRL Stoc. PO-IRL Subj. PO-IRL	\ \ \	\ \ \	\ \ \	× × ✓	× \$ \$	X X X	X X X	X X X	X X X	[70–73] [74–76] [77–80]
ME-IRL Subj. ME-IRL	\ \	X X	\ \	×	\ \	X X	<i>\</i> <i>\</i>	X X	X X	[81–92] [93,94]
Inverse Bounded Rational Control	l ✓	1	1	1	1	1	1	1	1	Section 4

a normative standard), showing that this subsumes common classes of control behavior in literature. Second, we use this to formalize the *inverse problem* G (as a descriptive model), showing that it generalizes existing work on imitation and reward learning. Importantly, this opens up a much broader variety of research problems in behavior representation learning—beyond simply learning optimal utility functions. Finally, we instantiate this approach with an example that we term *inverse bounded rational control*, illustrating how this structure enables learning (interpretable) representations of (bounded) rationality—capturing familiar notions of decision complexity, subjectivity, and uncertainty.

2. Related Work

As specific forms of descriptive modeling, imitation learning and apprenticeship learning are popular paradigms for learning policies that mimic the behavior of a demonstrator. *Imitation learning* focuses on replicating an expert's actions. Classically, "behavioral cloning" methods directly seek to learn a mapping from input states to output actions [14–16], using assistance from interactive experts or auxiliary regularization to improve generalization [17–21]. More recently, "distribution-matching" methods have been proposed for learning an imitator policy whose induced state-action occupancy measure is close to that of the demonstrator [23–39]. *Apprenticeship learning* focuses on matching the cumulative returns of the expert—on the basis of some ground-truth reward function not known to the imitator policy. This is most popularly approached by inverse reinforcement learning (IRL), which seeks to infer the reward function for which the demonstrated behavior appears most optimal, and using which an apprentice policy may itself be optimized via reinforcement learning. This includes maximum-margin methods based on feature expectations [13, 40–45], maximum likelihood soft policy matching [51, 52], maximum entropy policies [50, 89–92], and Bayesian maximum a posteriori inference [59–63], as well as methods that leverage preference models and additional annotations for assistance [95–99]. We defer to surveys of [12, 100] for more detailed overviews of imitation learning and inverse reinforcement learning.

Inverse decision modeling subsumes most of the standard approaches to imitation and apprenticeship learning as specific instantiations, as we shall see (cf. Table 1). Yet-with very few exceptions [78–80]—the vast majority of these works are limited to cases where demonstrators are assumed to be ideal or close to ideal. Inference is therefore limited to that of a single utility function; after all, its primary purpose is less for introspection than simply as a mathematical intermediary for mimicking the demonstrator's exhibited behavior. To the contrary, we seek to inspect and understand the demonstrator's behavior, rather than simply producing a faithful copy of it. In this sense, the novelty of our work is two-fold. First, we shall formally define "inverse decision models" much more generally as projections in the space of behaviors. These projections depend on our conscious choices for forward and inverse planners, and the explicit structure we choose for their parameterizations allows asking new classes of targeted research questions based on normative factors (which we impose) and descriptive factors (which we learn). Second, we shall model an agent's behavior as induced by both a recognition policy (committing observations to internal states) and a *decision policy* (emitting actions from internal states). Importantly, not only may an agent's mapping from internal states into actions be suboptimal (viz. the latter), but that their mapping from observations into beliefs may also be subjective (viz. the former). This greatly generalizes the idea of "boundedness" in sequential decision-making-that is, instead of commonlyassumed forms of noisy optimality, we arrive at precise notions of subjective dynamics and biased belief-updates. Appendix A gives a more detailed treatment of related work.

3. Inverse Decision Modeling

First, we describe our formalism for *planners* (Section 3.1) and *inverse planners* (Section 3.2)—together constituting our framework for inverse decision modeling (Section 3.3). Next, we instantiate this with a prototypical example to spotlight the wider class of research questions that this unified perspective opens up (Section 4). Table 1 summarizes related work subsumed, and contextualizes our later example.

Planner (F)	Setting (ψ)	Parameter (θ)	Optimization (π^*, ρ^*)	Examples
Decision-Rule CMP Policy	$\mathcal{S}, \mathcal{U}, \mathcal{T}$	χ	$\operatorname{argmax}_{\pi} \delta(\pi - f_{\operatorname{decision}}(\chi))$	[14]
Model-Free MDP Learner	$\mathcal{S}, \mathcal{U}, \mathcal{T}$	v,γ	$\operatorname{argmax}_{\pi} \mathbb{E}_{\pi, \tau_{env}} [\sum_{t} \gamma^{t} \upsilon(s_{t}, u_{t})]$	(any RL agent)
Max. Entropy MDP Learner	$\mathcal{S}, \mathcal{U}, \mathcal{T}$	$\upsilon, \gamma, lpha$	$\operatorname{argmax}_{\pi} \mathbb{E}_{\pi, \tau_{env}} [\sum_{t} \gamma^{t} \upsilon(s_{t}, u_{t}) + \alpha \mathcal{H}(\pi(\cdot s_{t}))]$	[101-104]
Model-Based MDP Planner	$\mathcal{S},\mathcal{U},\mathcal{T}$	υ,γ, au	$\operatorname{argmax}_{\pi} \mathbb{E}_{\pi,\tau} [\sum_{t} \gamma^{t} \upsilon(s_{t}, u_{t})]$	(any MDP solver)
Differentiable MDP Planner	$\mathcal{S}, \mathcal{U}, \mathcal{T}$	υ, γ, au	$\operatorname{argmax}_{\pi} \delta(\pi - \operatorname{neural-network}(\psi, \upsilon, \gamma, \tau))$	[105, 106]
KL-Regularized MDP Planner	$\mathcal{S},\mathcal{U},\mathcal{T}$	$\upsilon, \gamma, \tau, \alpha, \tilde{\pi}$	$\operatorname{argmax}_{\pi} \mathbb{E}_{\pi,\tau} \left[\sum_{t} \gamma^{t} (\upsilon(s_{t}, u_{t}) - \alpha D_{\mathrm{KL}}(\pi(\cdot s_{t}) \ \tilde{\pi})) \right]$	[107–111]
Decision-Rule IOHMM Policy	$\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O}$	χ, au, ω	$\operatorname{argmax}_{\pi} \delta(\pi - f_{\operatorname{decision}}(\chi), \rho - f_{\operatorname{recognition}}(\tau, \omega))$	[22]
Model-Free POMDP Learner	$\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O}$	v,γ	$\operatorname{argmax}_{\pi,\rho \in \{\rho \text{ is black-box}\}} \mathbb{E}_{\pi,\tau_{env},\rho}[\sum_{t} \gamma^{t} \upsilon(s_{t}, u_{t})]$	[112–117]
Model-Based POMDP Planner	$\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O}$	$\upsilon, \gamma, \tau, \omega$	$\operatorname{argmax}_{\pi,\rho\in\{\rho \text{ is unbiased}\}} \mathbb{E}_{\pi,\tau,\rho}[\sum_{t} \gamma^{t} \upsilon(s_{t}, u_{t})]$	[118–121]
Belief-Aware $\upsilon\text{-}\text{POMDP}$ Planner	$\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O}$	$\upsilon_{\scriptscriptstyle \mathcal{Z}}, \gamma, \tau, \omega$	$\operatorname{argmax}_{\pi,\rho\in\{\rho \text{ is unbiased}\}} \mathbb{E}_{\pi,\tau,\rho}[\sum_{t} \gamma^{t} \upsilon_{\mathbb{Z}}(s_{t}, z_{t}, u_{t})]$	[122, 123]
Bounded Rational Control	$\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O}$	$v, \gamma, \overline{lpha}, eta, \ \eta, ilde{\pi}, ilde{\sigma}, ilde{arrho}$	$ \begin{aligned} \operatorname*{argmax}_{\pi,\rho\in\{\rho\text{ is possibly-biased}\}} \mathbb{E}_{\pi,\rho}[\sum_{t} \gamma^{t} \upsilon(s_{t}, u_{t})] \\ - \alpha \mathbb{I}_{\pi,\rho}[\pi; \tilde{\pi}] - \beta \mathbb{I}_{\pi,\rho}[\sigma; \tilde{\sigma}] - \eta \mathbb{I}_{\pi,\rho}[\varrho; \tilde{\varrho}] \end{aligned} $	Theorems 4–5
General Formulation	$\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O}$	(any)	$\mathrm{argmax}_{\pi, ho}\mathcal{F}_\psi(\pi, ho; heta)$	Section 3.1

Table 2. Planners. Formulation of primary classes of planner algorithms in terms of our (forward) formalism, incl. the boundedly rational planner in our example (Section 4). *Legend:* controlled Markov process (CMP); Markov decision process (MDP); input-output hidden Markov model (IOHMM); partially-observable (PO); Dirac delta (δ); any mapping into policies (f); decision-rule parameterization (χ).

3.1. Forward Problem

Consider the standard setup for sequential decision-making, where an agent interacts with a (potentially partially-observable) environment. First, let $\psi \doteq (S, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O})$ give the problem setting, where S denotes the space of (external) environment states, \mathcal{X} of environment observables, \mathcal{Z} of (internal) agent states, \mathcal{U} of agent actions, $\mathcal{T} \doteq \Delta(\mathcal{S})^{\mathcal{S} \times \mathcal{U}}$ of environment transitions, and $\mathcal{O} \doteq \Delta(\mathcal{X})^{\mathcal{U} \times S}$ of environment emissions. Second, denote with θ the *planning parameter*: the parameterization of (subjective) factors that a planning algorithm uses to produce behavior, e.g. utility functions $v \in \mathbb{R}^{S \times U}$, discount factors $\gamma \in [0, 1)$, or any other biases that an agent might be subject to, such as imperfect knowledge τ, ω of true environment dynamics $\tau_{env}, \omega_{env} \in \mathcal{T} \times \mathcal{O}$. Note that access to the true dynamics is only (indirectly) possible via such knowledge, or by sampling online/from batch data. Now, a planner is a mapping producing observable behavior:

Definition 1 (Behavior) Denote the space of (observationaction) trajectories with $\mathcal{H} \doteq \bigcup_{t=0}^{\infty} (\mathcal{X} \times \mathcal{U})^t \times \mathcal{X}$. Then a *behavior* ϕ manifests as a distribution over trajectories (induced by an agent's policies interacting with the environment):

$$\Phi \doteq \Delta(\mathcal{H}) \tag{1}$$

Consider behaviors induced by an agent operating under a *recognition policy* $\rho \in \Delta(\mathcal{Z})^{\mathcal{Z} \times \mathcal{U} \times \mathcal{X}}$ (i.e. committing observation-action trajectories to internal states), together with a *decision policy* $\pi \in \Delta(\mathcal{U})^{\mathcal{Z}}$ (i.e. emitting actions from internal states). We shall denote behaviors induced by π, ρ :

$$\phi_{\pi,\rho}((x_0, u_0, \ldots)) \doteq \mathbb{P}_{\substack{u \sim \pi(\cdot|z) \\ s' \sim \tau_{\text{env}}(\cdot|s, u) \\ x' \sim \omega_{\text{env}}(\cdot|u, s') \\ z' \sim \rho(\cdot|z, u, x')}} (h = (x_0, u_0, \ldots)) \quad (2)$$

(*Note*: Our notation may not be immediately familiar as we seek to unify terminology across multiple fields. For reference, a summary of notation is provided in Appendix E).

Definition 2 (Planner) Given problem setting ψ and planning parameter θ , a *planner* is a mapping into behaviors:

$$F:\Psi\times\Theta\to\Phi\tag{3}$$

where Ψ indicates the space of settings, and Θ the space of parameters. Often, behaviors of the form $\phi_{\pi,\rho}$ can be naturally expressed in terms of the solution to an optimization:

 $F(\psi, \theta) \doteq \phi_{\pi^*, \rho^*} : \pi^*, \rho^* \doteq \operatorname{argmax}_{\pi, \rho} \mathcal{F}_{\psi}(\pi, \rho; \theta) \quad (4)$ of some real-valued function \mathcal{F}_{ψ} (e.g. this includes all cases where a utility function v is an element of θ). So, we shall write $\phi^* \doteq \phi_{\pi^*, \rho^*}$ to indicate the behavior produced by F.

This definition is very general: It encapsulates a wide range of standard algorithms in the literature (see Table 2), including decision-rule policies and neural-network planners. Importantly, however, observe that in most contexts, a global optimizer for ρ is (trivially) either an identity function, or perfect Bayesian inference (with the practical caveat, of course, that in model-free contexts actually reaching such an optimum may be difficult, such as with a deep recurrent network). Therefore in addition to just π , what Definition 2 makes explicit is the potential for ρ to be *biased*—that is, to deviate from (perfect) Bayes updates; this will be one of the important developments made in our subsequent example.

Note that by equating a planner with such a mapping, we are implicitly assuming that the embedded optimization (Equation 4) is *well-defined*—that is, that there exists a single global optimum. In general if the optimization is non-trivial, this requires that the spaces of policies $\pi, \rho \in \mathcal{P} \times \mathcal{R}$ be suitably restricted: This is satisfied by the usual (hard-/soft-Q) Boltzmann-rationality for decision policies, and by uniquely fixing the semantics of internal states as (subjective) beliefs, i.e. probability distributions over states, with recognition policies being (possibly-biased) Bayes updates.

A more practical question is whether this optimum is reach-

Figure 1. Forward, Inverse, and Projection Mappings. In the forward direction (i.e. generation): Given planning parameters θ , a planner F generates observable behavior ϕ (Definition 2). In the opposite direction (i.e. inference): Given observed behavior ϕ , an inverse planner G infers the planning parameters θ that produced it—subject to normative specifications (Definition 3). Finally, given observed behavior ϕ , the composition of F and G gives its projection onto the space of behaviors that are parameterizable by θ (Definition 4): This is the inverse decision model (Definition 5).



able. While this may seem more difficult (at least in the most general case), for our *interpretative* purposes it is rarely a problem, because (simple) human-understandable models are what we desire to be working with in the first instance. In healthcare, for example, diseases are often modeled in terms of discrete states, and subjective beliefs over those states are eminently transparent factors that medical practitioners can readily comprehend and reason about [124, 125]. This is prevalent in research and practice, e.g. two-to-four states in progressive dementia [126-128], cancer screening [129, 130], cystic fibrosis [131], as well as pulmonary disease [132]. Of course, this is not to say our exposition is incompatible with model-free, online settings with complex spaces and black-box approximators. But our focus here is to establish an interpretative paradigm-for which simple state-based models are most amenable to human reasoning.

3.2. Inverse Problem

Given any setting and appropriate planner, θ gives a complete account of $\phi^* = F(\psi, \theta)$: This deals with generation —that is, of behavior from its parameterization. In the opposite, given observed behavior ϕ_{demo} produced by some planner, we can ask what its θ appears to be: This now deals with *inference*—that is, of parameterizations from behavior.

First, note that absent any restrictions, this endeavor immediately falls prey to the celebrated "no free lunch" result: It is in general *impossible* to infer anything of use from ϕ_{demo} alone, if we posit nothing about θ (or F) to begin with [136, 137]. The only close attempt has recruited inductive biases requiring multiple environments, and is *not* interpretable due to the use of differentiable planners [105, 106]. On the other extreme, the vast literature on IRL has largely restricted attention to perfectly optimal agents—that is, with full visibility of states, certain knowledge of dynamics, and perfect ability to optimize v. While this indeed fends off the impossibility result, it is *overly restrictive* for understanding behavior: Summarizing ϕ_{demo} using v alone is not informative as to specific types of biases we may be interested in. How aggressive does this clinician seem? How flexible do their actions appear? It is difficult to tease out such nuances from just v—let alone comparing between agents [138, 139].

We take a generalized approach to allow any middle ground of choice. While some normative specifications are required to fend off the impossibility result [106, 136], they need not be so strong as to restrict us to perfect optimality. Formally:

Definition 3 (Inverse Planner) Let $\Theta \doteq \Theta_{norm} \times \Theta_{desc}$ decompose the parameter space into a *normative* component (i.e. whose values $\theta_{norm} \in \Theta_{norm}$ we wish to clamp), and a *descriptive* component (i.e. whose values $\theta_{desc} \in \Theta_{desc}$ we wish to infer). Then an *inverse planner* is given as follows:

$$G: \Phi \times \Theta_{\text{norm}} \to \Theta_{\text{desc}} \tag{5}$$

Often, the descriptive parameter can be naturally expressed as the solution to an optimization (of some real-valued \mathcal{G}_{ψ}):

$$G(\phi_{\text{demo}}, \theta_{\text{norm}}) \doteq \operatorname{argmin}_{\theta_{\text{decc}}} \mathcal{G}_{\psi}(\phi_{\text{demo}}, \phi_{\text{imit}}) \qquad (6)$$

where we denote by $\phi_{\text{imit}} \doteq F(\psi, (\theta_{\text{norm}}, \theta_{\text{desc}}))$ the *imitation* behavior generated on the basis of θ_{desc} . So, we shall write θ_{desc}^* for the (minimizing) descriptive parameter output by *G*.

As with the forward case, this definition is broad: It encapsulates a wide range of inverse optimization techniques in the literature (see Table 3). Although not all techniques entail learning imitating policies in the process, by far the most dominant paradigms do (i.e. maximum margin, soft policy matching, and distribution matching). Moreover, it is normatively flexible in the sense of the middle ground we wanted: $\theta_{\rm norm}$ can encode precisely the information we desire.⁴ This opens up new possibilities for interpretative research. For instance, contrary to IRL for imitation or apprenticeship, we may often *not* wish to recover v at all. Suppose—as an investigator—we believe that a certain v we defined is the "ought-to-be" ideal. By allowing v to be encoded in θ_{norm} (instead of θ_{desc}), we may now ask questions of the form: How "consistently" does ϕ_{demo} appear to be in pursuing v? Does it seem "optimistic" or "pessimistic" relative to neutral beliefs about the world? All that is required is for appropriate measures of such notions (and any others) to be represented in θ_{desc} . (Section 4 shall provide one such exemplar).

Note that parameter identifiability depends on the degrees of freedom in the target θ_{desc} and the nature of the identifi-

⁴We can verify that $\theta_{desc} = v$ alone recovers the usual IRL paradigm.

Table 3. Inverse Planners. Formulation of primary classes of identification strategies in terms of our (inverse) formalism. Legend: value functions for ϕ under θ ($V^{\phi}_{\theta}, Q^{\phi}_{\theta}$); regularizer (ζ); shaped-reward error (Δv); p-norm ($\|\cdot\|_p$); preference relation (\prec); f-divergence (D_f). Note that while our notation is general, virtually all original works here have $\theta_{desc} = v$ and assume full observability (whence $S = \mathcal{X} = \mathcal{Z}$).

Inverse Planner (G)	Demonstrator (ϕ_{demo})	Helper	Optimization (θ^*_{desc})	Examples
Minimum Perturbation	Deterministic, Optimal	Default $\tilde{\theta}_{\text{desc}}$	$\operatorname{argmin}_{\theta_{\text{desc}}} \ \theta_{\text{desc}} - \tilde{\theta}_{\text{desc}}\ _{p} : \phi_{\text{demo}} = F(\psi, \theta)$	[133]
Maximum Margin	Deterministic, Optimal	-	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} \mathbb{E}_{z \sim \rho_0}[V_{\theta}^{\phi_{\operatorname{imit}}}(z) - V_{\theta}^{\phi_{\operatorname{demo}}}(z)]$	[40-45, 53, 70-73]
Regularized Max. Margin	Stochastic, Optimal	-	$\operatorname{argmin}_{\theta_{\text{desc}}} \mathbb{E}_{z \sim \rho_0} \left[V_{\operatorname{soft}, \theta}^{\phi_{\text{imit}}}(z) - V_{\theta}^{\phi_{\text{demo}}}(z) \right] + \zeta(\theta)$	[25]
Multiple Experimentation	Deterministic, Optimal	Environments \mathcal{V}	$\operatorname{argmin}_{\theta_{\text{desc}}} \int \max_{\mathcal{V}, u} (Q_{\mathcal{V}, \theta}^{\phi_{\text{demo}}}(z, u) - V_{\mathcal{V}, \theta}^{\phi_{\text{demo}}}(z)) dx$	[134,135]
Distance Minimization	Individually-Scored	Scores $\tilde{v}(h) \in \mathbb{R}$	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} \mathbb{E}_{h \sim \phi_{\operatorname{demo}}} \ \tilde{v}(h) - \sum_{s, u \in h} v(s, u) \ _p$	[95,96]
Soft Policy Inversion	Stoc., Batch-Ordered	$\{\phi_{\rm demo}^{(1)},,\phi_{\rm demo}^{(K)}\}$	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} \sum_{k} \mathbb{E}_{s,u,s' \sim \phi^{(k)}} \ \Delta v^{(k)}(s,u,s')\ _{p}$	[97]
Preference Extrapolation	Stoc., Pairwise-Ranked	$\{(i,j) h_i \prec h_j\}$	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} \mathbb{E}_{(h_i \prec h_j) \sim \phi_{\operatorname{demo}}} \log \mathbb{P}_{\upsilon}(h_i \prec h_j)$	[98,99]
Soft Policy Matching	Stochastic, Optimal	-	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} D_{\operatorname{KL}}(\mathbb{P}_{\phi_{\operatorname{demo}}}(u_{0:T} \ x_{0:T}) \ \mathbb{P}_{\phi_{\operatorname{imit}}}(u_{0:T} \ x_{0:T}))$	[47-52, 76, 89-94]
Distribution Matching	Stochastic, Optimal	-	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} D_f(\phi_{\operatorname{demo}} \ \phi_{\operatorname{imit}})$	[23-39, 54, 81-88]
General Formulation	(any)	(any)	$\operatorname{argmin}_{\theta_{\operatorname{desc}}} \mathcal{G}_{\psi}(\phi_{\operatorname{demo}},\phi_{\operatorname{imit}})$	Section 3.2

cation strategy G. From our generalized standpoint, we simply note that—beyond the usual restrictions (e.g. on scaling, shifting, reward shaping) in conjunction with G—Bayesian inference remains a valid option to address ambiguities, as in [26] for distribution matching, [59–63, 74, 75] for soft policy matching, and [140, 141] for preference extrapolation.

3.3. Behavior Projection

Now we have the ingredients to formally define the business of inverse decision modeling. Compacting notation, denote $F_{\theta_{\text{norm}}}(\cdot) \doteq F(\psi, (\theta_{\text{norm}}, \cdot))$, and $G_{\theta_{\text{norm}}}(\cdot) \doteq G(\cdot, \theta_{\text{norm}})$. First, we require a projection operator that maps onto the space of behaviors that are *parameterizable* by θ given $F_{\theta_{\text{norm}}}$.

Definition 4 (Behavior Projection) Denote the image of Θ_{desc} under $F_{\theta_{\text{norm}}}$ by the following: $\Phi_{\theta_{\text{norm}}} \doteq F_{\theta_{\text{norm}}}[\Theta_{\text{desc}}] \leq \Phi$. Then the projection map onto this subspace is given by:

$$\operatorname{proj}_{\Phi_{\theta_{\operatorname{norm}}}} \doteq F_{\theta_{\operatorname{norm}}} \circ G_{\theta_{\operatorname{norm}}}$$
(7)

Definition 5 (Inverse Decision Model) Given a specified method of parameterization Θ , normative standards θ_{norm} , (and appropriate planner *F* and identification strategy *G*), the resulting *inverse decision model* of ϕ_{demo} is given by:

$$\phi_{\text{imit}}^* \doteq \operatorname{proj}_{\Phi_{\theta_{\text{norm}}}}(\phi_{\text{demo}}) \tag{8}$$

In other words, the model $\phi^*_{\rm imit}$ serves as a complete (generative) account of $\phi_{\rm demo}$ as its *behavior projection* onto $\Phi_{\theta_{\rm norm}}$.

Interpretability What dictates our choices? For pure imitation (i.e. replicating expert actions), a black-box decisionrule fitted by soft policy matching may do well. For apprenticeship (i.e. matching expert returns), a perfectly optimal planner inversed by distribution matching may do well. But for *understanding*, however, we wish to place appropriate structure on Θ depending on the question of interest: Precisely, the mission here is to choose some (interpretable) $F_{\theta_{norm}}$, $G_{\theta_{norm}}$ such that ϕ_{imit}^* is amenable to human reasoning.

Note that these are not passive *assumptions*: We are not making the (factual) claim that θ gives a scientific explanation of

the complex neurobiological processes in a clinician's head. Instead, these are active *specifications*: We are making the (effective) claim that the learned θ is a parameterized "as-if" interpretation of the observed behavior. For instance, while there exist a multitude of commonly studied human biases in psychology, it is difficult to measure their magnitudes much less compare them among agents. Section 4 shows an example of how inverse decision modeling can tackle this. (Figure 1 visualizes inverse decision modeling in a nutshell).

4. Bounded Rationality

We wish to understand observed behavior through the lens of *bounded rationality*. Specifically, let us account for the following facts: that (1) an agent's *knowledge* of the environment is uncertain and possibly biased; that (2) the agent's *capacity* for information processing is limited, both for decisions and recognition; and—as a result—that (3) the agent's (subjective) beliefs and (suboptimal) actions *deviate* from those expected of a perfectly rational agent. We shall see, this naturally allows quantifying such notions as flexibility of decisions, tolerance for surprise, and optimism in beliefs.

First, Section 4.1 describes inference and control under environment uncertainty (cf. 1). Then, 4.2 develops the forward model (F) for agents bounded by information constraints (cf. 2–3). Finally, 4.3 learns parameterizations of such boundedness from behavior by inverse decision modeling (G).

4.1. Inference and Control

Consider that an agent has *uncertain* knowledge of the environment, captured by a prior over dynamics $\tilde{\sigma} \in \Delta(\mathcal{T} \times \mathcal{O})$. As a normative baseline, let this be given by some (unbiased) posterior $\tilde{\sigma} \doteq p(\tau, \omega | \mathcal{E})$, where \mathcal{E} refers to any manner of experience (e.g. observed data about environment dynamics) with which we may come to form such a neutral belief.

Now, an agent may *deviate* from $\tilde{\sigma}$ depending on the situation, relying instead on $\tau, \omega \sim \sigma(\cdot | z, u)$ —where z, u allows the (biased) $\sigma \in \Delta(\mathcal{T} \times \mathcal{O})^{\mathcal{Z} \times \mathcal{U}}$ to be context-dependent. Consider recognition policies thereby parameterized by σ :

$$\rho(z'|z, u, x') \doteq \mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u)} \rho_{\tau, \omega}(z'|z, u, x') \qquad (9)$$

where $\rho_{\tau,\omega}$ denotes the policy for adapting z to x' given (a point value for) τ, ω . For interpretability, we let $\rho_{\tau,\omega}$ be the usual Bayes belief-update. Importantly, however, ρ can now effectively be biased (i.e. by σ) even while $\rho_{\tau,\omega}$ is Bayesian.

Forward Process The forward ("inference") process yields the occupancy measure. First, the *stepwise conditional* is:

$$p(z'|z) = \mathbb{E}_{\substack{u \sim \pi(\cdot|z) \\ \tau, \omega \sim \sigma(\cdot|z,u) \\ s' \sim \tau(\cdot|s,u) \\ x' \sim \omega(\cdot|u,s')}} \rho_{\tau,\omega}(z'|z,u,x')$$
(10)

Define Markov operator $\mathbb{M}_{\pi,\rho} \in \Delta(\mathcal{Z})^{\Delta(\mathcal{Z})}$ such that for any distribution $\mu \in \Delta(\mathcal{Z}) : (\mathbb{M}_{\pi,\rho}\mu)(z') \doteq \mathbb{E}_{z \sim \mu} p(z'|z)$. Then

$$\mu_{\pi,\rho}(z) \doteq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t p(z_t = z | z_0 \sim \rho_0)$$
 (11)

defines the *occupancy measure* $\mu_{\pi,\rho} \in \Delta(\mathcal{Z})$ for any initial (internal-state) distribution ρ_0 , and discount rate $\gamma \in [0, 1)$.

Lemma 1 (Forward Recursion) Define the forward operator $\mathbb{F}_{\pi,\rho}$: $\Delta(\mathcal{Z})^{\Delta(\mathcal{Z})}$ such that for any given $\mu \in \Delta(\mathcal{Z})$:

$$(\mathbb{F}_{\pi,\rho}\mu)(z) \doteq (1-\gamma)\rho_0(z) + \gamma(\mathbb{M}_{\pi,\rho}\mu)(z)$$
(12)

Then the occupancy $\mu_{\pi,\rho}$ is the (unique) fixed point of $\mathbb{F}_{\pi,\rho}$.

Backward Process The backward ("control") process yields the value function. We want that $\mu_{\pi,\rho}$ maximize utility:

$$\operatorname{maximize}_{\substack{\mu_{\pi,\rho}}} J_{\pi,\rho} \doteq \mathbb{E} \sum_{\substack{s \sim \mu_{\pi,\rho} \\ s \sim p(\cdot|z) \\ u \circ \pi(\cdot|z)}} v(s,u)$$
(13)

Using $V \in \mathbb{R}^{\mathcal{Z}}$ to denote the multiplier, the Lagrangian is given by $\mathcal{L}_{\pi,\rho}(\mu, V) \doteq J_{\pi,\rho} - \langle V, \mu - \gamma \mathbb{M}_{\pi,\rho} \mu - (1 - \gamma) \rho_0 \rangle$.

Lemma 2 (Backward Recursion) Define the backward operator $\mathbb{B}_{\pi,\rho} : \mathbb{R}^{\mathbb{Z}} \to \mathbb{R}^{\mathbb{Z}}$ such that for any given $V \in \mathbb{R}^{\mathbb{Z}}$:

$$(\mathbb{B}_{\pi,\rho}V)(z) \doteq \mathbb{E}_{\substack{s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z)}} [\upsilon(s,u) + \mathbb{E}_{\substack{\tau,\omega \sim \sigma(\cdot|z,u) \\ s' \sim \tau(\cdot|s,u) \\ x' \sim \omega(\cdot|u,s') \\ z' \sim \rho_{\tau,\omega}(\cdot|z,u,x')}} (14)$$

Then the (dual) optimal V is the (unique) fixed point of $\mathbb{B}_{\pi,\rho}$; this is the *value function* considering knowledge uncertainty:

$$V^{\phi_{\pi,\rho}}(z) \doteq \sum_{t=0}^{\infty} \gamma^{t} \mathbb{E}_{\substack{s_t \sim p(\cdot|z_t) \\ u_t \sim \pi(\cdot|z_t) \\ \tau, \omega \sim \sigma(\cdot|z_t, u_t) \\ s_{t+1} \sim \tau(\cdot|s_t, u_t) \\ x_{t+1} \sim \omega(\cdot|u_t, s_{t+1}) \\ z_{t+1} \sim \rho_{\tau,\omega}(\cdot|z_t, u_t, x_{t+1})}$$
(15)

so we can equivalently write targets $J_{\pi,\rho} = \mathbb{E}_{z \sim \rho_0} V^{\phi_{\pi,\rho}}(z)$. Likewise, we can also define the (state-action) value function $Q^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathcal{Z} \times \mathcal{U}}$ —that is, $Q^{\phi_{\pi,\rho}}(z,u) \doteq \mathbb{E}_{s \sim p(\cdot|z)}[\upsilon(s,u) + \mathbb{E}_{\tau,\omega \sim \sigma(\cdot|z,u),...,z' \sim \rho_{\tau,\omega}(\cdot|z,u,x')} \gamma V^{\phi_{\pi,\rho}}(z')]$ given an action.

4.2. Bounded Rational Control

For perfectly rational agents, the best *decision policy* given any z simply maximizes $V^{\phi_{\pi,\rho}}(z)$, thus it selects actions according to $\operatorname{argmax}_u Q^{\phi_{\pi,\rho}}(z, u)$. And the best *recognition policy* simply corresponds to their unbiased knowledge of the world, thus it sets $\sigma(\cdot|z, u) = \tilde{\sigma}, \forall z, u$ (in Equation 9).

Information Constraints But control is resource-intensive. We formalize an agent's boundedness in terms of capacities for processing information. First, *decision complexity* captures the informational effort in determining actions $\pi(\cdot|z)$, relative to some prior $\tilde{\pi}$ (e.g. baseline clinical guidelines):

$$\mathbb{I}_{\pi,\rho}[\pi;\tilde{\pi}] \doteq \mathbb{E}_{z \sim \mu_{\pi,\rho}} D_{\mathrm{KL}}(\pi(\cdot|z) \| \tilde{\pi}) \tag{16}$$

Second, *specification complexity* captures the average regret of their internal model $\sigma(\cdot|z, u)$ deviating from their prior (i.e. unbiased knowledge $\tilde{\sigma}$) about environment dynamics:

$$\mathbb{I}_{\pi,\rho}[\sigma;\tilde{\sigma}] \doteq \mathbb{E}_{\substack{z \sim \mu_{\pi,\rho} \\ u \sim \pi(\cdot|z)}} D_{\mathrm{KL}}(\sigma(\cdot|z,u) \| \tilde{\sigma}) \tag{17}$$

Finally, *recognition complexity* captures the statistical surprise in adapting to successive beliefs about the partially-observable states of the world (again, relative to some prior $\tilde{\varrho}$):

$$\mathbb{I}_{\pi,\rho}[\varrho;\tilde{\varrho}] \doteq \mathbb{E}_{\substack{z \sim \mu_{\pi,\rho} \\ u \sim \pi(\cdot|z) \\ \tau, \omega \sim \sigma(\cdot|z,u)}} D_{\mathrm{KL}}(\varrho_{\tau,\omega}(\cdot|z,u) \|\tilde{\varrho}) \qquad (18)$$

where $\rho_{\tau,\omega}(\cdot|z,u) \doteq \mathbb{E}_{s\sim p(\cdot|z),s'\sim \tau(\cdot|s,u),x'\sim \omega(\cdot|u,s')}\rho_{\tau,\omega}(\cdot|z,u,x')$ gives the internal-state update. We shall see, these measures generalize information-theoretic ideas in control.

Backward Process With capacity constraints, the maximization in Equation 13 now becomes subject to $\mathbb{I}_{\pi,\rho}[\pi; \tilde{\pi}] \leq A$, $\mathbb{I}_{\pi,\rho}[\sigma; \tilde{\sigma}] \leq B$, and $\mathbb{I}_{\pi,\rho}[\varrho; \tilde{\varrho}] \leq C$. So the Lagrangian (now with the additional multipliers $\alpha, \beta, \eta \in \mathbb{R}$) is given by $\mathcal{L}_{\pi,\rho}(\mu, \alpha, \beta, \eta, V) \doteq J_{\pi,\rho} - \langle V, \mu - \gamma \mathbb{M}_{\pi,\rho} \mu - (1-\gamma)\rho_0 \rangle - \alpha \cdot (\mathbb{I}_{\pi,\rho}[\pi; \tilde{\pi}] - A) - \beta \cdot (\mathbb{I}_{\pi,\rho}[\sigma; \tilde{\sigma}] - B) - \eta \cdot (\mathbb{I}_{\pi,\rho}[\varrho; \tilde{\varrho}] - C)$.

Proposition 3 (Backward Recursion) Define the backward operator $\mathbb{B}_{\pi,\rho} : \mathbb{R}^{\mathbb{Z}} \to \mathbb{R}^{\mathbb{Z}}$ such that for any given function $V \in \mathbb{R}^{\mathbb{Z}}$ and for any given coefficient values $\alpha, \beta, \eta \in \mathbb{R}$:

$$\begin{aligned}
(\mathbb{B}_{\pi,\rho}V)(z) &\doteq \mathbb{E}_{s \sim p(\cdot|z)} \left[-\alpha \log \frac{\pi(u|z)}{\tilde{\pi}(u)} + \upsilon(s, u) + \\
\mathbb{E}_{\tau,\omega \sim \sigma(\cdot|z,u)} \left[-\beta \log \frac{\sigma(\tau,\omega|z,u)}{\tilde{\sigma}(\tau,\omega)} + \\
\mathbb{E}_{\tau,\omega \sim \sigma(\cdot|z,u)} \left[-\eta \log \frac{\varrho_{\tau,\omega}(z'|z,u)}{\tilde{\varrho}(z')} + \gamma V(z') \right] \right]
\end{aligned}$$

Then the (dual) optimal V is the (unique) fixed point of $\mathbb{B}_{\pi,\rho}$; as before, this is the value function $V^{\phi_{\pi,\rho}}$ —which now includes the complexity terms. Likewise, we can also define the (state-action) $Q^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathbb{Z} \times \mathcal{U}}$ as the ¹/₃-step-ahead expectation, and the (state-action-model) $K^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathbb{Z} \times \mathcal{U} \times \mathcal{T} \times \mathcal{O}}$ as the ²/₃-steps-ahead expectation (which is new in this setup).

Policies and Values The (dis-)/utility-seeking decision policy (min-)/maximizes $V^{\phi_{\pi,\rho}}(z)$, and a pessimistic/optimis-

Table 4. Boundedly Rational Agents. Formulation of common decision agent	ts as instantiations of our (boundedly rational) formalism. Note
that either $\beta^{-1} \rightarrow 0$ or $\tilde{\sigma} = \delta$ is sufficient to guarantee $\forall z, u : \sigma(\cdot z, u) = \tilde{\sigma}$.	[†] Softmax added on top of deterministic, optimal Q-functions

Poundadly Dational Agant	Flexibility	Optimism	Adaptivity	(Action Prior) (Model Prior) (Belief Prior)		Obcowyobility	Evennles	
boundedry Kational Agent	α^{-1}	β^{-1}	η^{-1}	$\tilde{\pi}$	$\tilde{\sigma}$	õ	Observability	Examples
Uniformly Random Agent	$\rightarrow 0$	$\rightarrow 0$	$\rightarrow_{\pm\infty}$	Uniform	Dirac δ	-	Full/Partial	-
Deterministic, Optimal Agent	$ ightarrow\infty$	$\rightarrow 0$	$\rightarrow \pm \infty$	-	Dirac δ	-	Full/Partial	(any)
Boltzmann-Exploratory Agent [†]	$ ightarrow\infty$	$\rightarrow 0$	$ ightarrow \pm \infty$	-	Dirac δ	-	Full/Partial	[142–144]
Minimum-Information Agent	= 1	$\rightarrow 0$	= 1	(any)	Dirac δ	(any)	Full	[145–147]
Maximum Entropy Agent	$(0,\infty)$	$\rightarrow 0$	$\rightarrow \pm \infty$	Uniform	Dirac δ	-	Full	[101-104]
(Action) KL-Regularized Agent	$(0,\infty)$	$\rightarrow 0$	$\rightarrow \pm \infty$	(any)	Dirac δ	-	Full	[107–111]
KL-Penalized Robust Agent	$\rightarrow \infty$	$(-\infty, 0)$	$\rightarrow \pm \infty$	-	(any)	-	Full	[148–151]
General Formulation	$\mathbb{R} \setminus \{0\}$	$\mathbb{R} \setminus \{0\}$	$\mathbb{R} \setminus \{0\}$	(any)	(any)	(any)	Full/Partial	Section 4

tic recognition policy min-/maximizes $Q^{\phi_{\pi,\rho}}(z, u)$ via σ .⁵ These optimal policies depend on optimal value functions:

Theorem 4 (Boundedly Rational Values) Define the backward operator $\mathbb{B}^* : \mathbb{R}^{\mathcal{Z}} \to \mathbb{R}^{\mathcal{Z}}$ such that for any $V \in \mathbb{R}^{\mathcal{Z}}$:

$$(\mathbb{B}^*V)(z) \doteq \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \exp(\frac{1}{\alpha}Q(z,u))$$
(20)

$$Q(z,u) \doteq \beta \log \mathbb{E}_{\tau,\omega \sim \tilde{\sigma}} \exp(\frac{1}{\beta}K(z,u,\tau,\omega))$$

$$K(z,u,\tau,\omega) \doteq + \mathbb{E}_{s \sim p(\cdot|z)}v(s,u)$$

$$\mathbb{E}_{\substack{s \sim p(\cdot|z) \\ s' \sim \tau(\cdot|s,u) \\ x' \sim \omega(\cdot|u,s') \\ z' \sim \rho_{\tau,\omega}(\cdot|z,u,x')}} \left[-\eta \log \frac{\varrho_{\tau,\omega}(z'|z,u)}{\tilde{\varrho}(z')} + \gamma V(z')\right]$$

Then the *boundedly rational value function* V^* for the (primal) optimal π^* , ρ^* is the (unique) fixed point of $\mathbb{B}^*_{\pi,\rho}$. (Note that both Q^* and K^* are immediately obtainable from this).

Theorem 5 (Boundedly Rational Policies) The *bounded-ly rational decision policy* (i.e. primal optimal) is given by:

$$\pi^*(u|z) = \frac{\tilde{\pi}(u)}{Z_{Q^*}(z)} \exp\left(\frac{1}{\alpha}Q^*(z,u)\right)$$
(21)

and the boundedly rational recognition policy is given by:

$$\rho^*(z'|z, u, x') = \mathbb{E}_{\tau, \omega \sim \sigma^*(\cdot|z, u)} \rho_{\tau, \omega}(z'|z, u, x') \text{, where}$$

$$\sigma^*(\tau, \omega|z, u) \doteq \frac{\tilde{\sigma}(\tau, \omega)}{Z_{K^*}(z, u)} \exp\left(\frac{1}{\beta}K^*(z, u, \tau, \omega)\right) \tag{22}$$

where $Z_{OP}(z) = \mathbb{E}_{\tau, \omega} \exp\left(\frac{1}{2}O^*(z, u)\right)$ and $Z_{VP}(z, u) = \mathbb{E}_{\tau, \omega} \exp\left(\frac{1}{2}O^*(z, u)\right)$

where $Z_{Q^*}(z) = \mathbb{E}_{u \sim \tilde{\pi}} \exp(\frac{1}{\alpha}Q^*(z, u))$ and $Z_{K^*}(z, u) = \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \exp(\frac{1}{\beta}K^*(z, u, \tau, \omega))$ give the partition functions.

Interpretation of Parameters This articulation of bounded rationality reflects the fact that imperfect behavior results from two sources of "boundedness": Firstly, that (1) given a mental model ρ for comprehending the world, an agent's information-processing capacities distort their decision-making π (cf. suboptimal actions); and secondly, that (2) the agent's mental model ρ itself is an imperfect characterization of the world—because prior knowledge $\tilde{\sigma}$ is uncertain, and internal states can be biased by σ (cf. subjective beliefs).

Concretely, the parameters in Theorems 4–5 admit intuitive interpretations. First, α^{-1} captures *flexibility* of decision-making, from a completely inflexible agent ($\alpha^{-1} \rightarrow 0$) to an

infinitely flexible, utility-seeking $(\alpha^{-1} \rightarrow \infty)$ or disutilityseeking $(\alpha^{-1} \rightarrow -\infty)$ one. Second, β^{-1} captures *optimism* in internal models, from a completely neutral agent $(\beta^{-1} \rightarrow 0)$ to an infinitely optimistic $(\beta^{-1} \rightarrow \infty)$ or pessimistic $(\beta^{-1} \rightarrow -\infty)$ one. Lastly, η^{-1} captures *adaptivity* of beliefs, from a perfectly adaptive agent $(\eta^{-1} \rightarrow \pm \infty)$ to one with infinite intolerance $(\eta^{-1} \rightarrow 0^+)$ or affinity $(\eta^{-1} \rightarrow 0^-)$ for surprise. Table 4 underscores the generality of this parameterization.

4.3. Inverse Bounded Rational Control

We hark back to our framework of Section 3: In bounded rational control ("BRC"), the *planning parameter* θ^{BRC} represents $\{v, \gamma, \alpha, \beta, \eta, \tilde{\pi}, \tilde{\sigma}, \tilde{\varrho}\}$, and the space Θ^{BRC} is again decomposable as $\Theta^{BRC}_{norm} \times \Theta^{BRC}_{desc}$. The *forward problem* is encapsulated by Theorems 4–5 (which also yield a straightforward algorithm, i.e. iterate 4 until convergence, then plug into 5). Therefore the *forward planner* is given as follows:

$$F_{\theta_{\text{norm}}^{\text{BRC}}}(\theta_{\text{desc}}^{\text{BRC}}) \doteq \phi_{\pi^*,\rho^*} : \pi^*, \rho^* \leftarrow \text{Theorems 4-5}$$
(23)

In the opposite direction, the problem is of *inverse bounded* rational control. Consider a minimal setting where we are given access to logged data $\mathcal{D} \doteq \{h_n \sim \phi_{\text{demo}}\}_{n=1}^N$ with no additional annotations. While several options from Table 3 are available, for simplicity we select soft policy matching for illustration. Thus the *inverse planner* is given as follows:

$$G_{\theta_{\text{norm}}^{\text{BRC}}}(\phi) \doteq \operatorname{argmin}_{\theta_{\text{dec}}^{\text{BRC}}} \mathbb{E}_{h \sim \phi} \log \mathbb{P}_{\phi_{\text{imit}}}(u_{0:T} \| x_{0:T})$$
(24)

where $\mathbb{P}_{\phi_{\pi,\rho}}(u_{0:T} || x_{0:T})$ is the causally-conditioned probability [152–155] $\prod_{t=0}^{T} \mathbb{P}_{\phi_{\pi,\rho}}(u_t | x_{1:t}, u_{1:t-1})$ —with the conditioning as induced by π, ρ . In the most general case where $\rho_{\tau,\omega}$ may be stochastic, $G_{\theta_{\text{norm}}^{\text{BRC}}}$ would require an EM approach; however, since we selected $\rho_{\tau,\omega}$ to be the (deterministic) Bayes update for interpretability, the likelihood is:

$$\log \mathbb{P}_{\phi_{\pi,\rho}}(u_{0:T} \| x_{0:T}) \propto \sum_{t=0}^{T} \log \pi(u_t | z_t)$$
(25)

where the z_t terms are computed recursively by ρ (see Appendix C). Finally, here the *inverse decision model* of any ϕ_{demo} is given by its projection $\phi^*_{\text{imit}} = F_{\theta^{\text{BRC}}_{\text{norm}}} \circ G_{\theta^{\text{BRC}}_{\text{norm}}}(\phi_{\text{demo}})$ onto the space $\Phi_{\theta^{\text{BRC}}_{\text{norm}}}$ of behaviors thereby *interpretably* parameterized—i.e. by the structure we designed for Θ^{BRC} , and by the normative standards $\theta^{\text{BRC}}_{\text{norm}}$ we may choose to specify.

⁵In general, flipping the direction of optimization for π or ρ corresponds to the *signs* of α or β , but does not change Theorems 4–5.

Inverse Decision Modeling



Figure 2. Bounded Rational Control. Decision agents in DIAG: In each panel, the boundedly rational decision policy π is shown in terms of action probabilities (y-axis) for different subjective beliefs (x-axis). To visualize the boundedly rational recognition policy ρ , each panel shows an example trajectory of beliefs (z_0 , z_1 , z_2 , z_3) for the case where three consecutive positive outcomes are observed (\blacktriangle markers).

5. Illustrative Use Case

So far, we have argued for a systematic, unifying perspective on inverse decision modeling ("IDM") for behavior representation learning, and presented inverse bounded rational control ("IBRC") as a concrete example of the formalism. Three aspects of this approach deserve empirical illustration:

- <u>Interpretability</u>: IBRC gives a *transparent* parameterization of behavior that can be successfully learned from data.
- <u>Expressivity</u>: IBRC more finely differentiates between imperfect behaviors, while standard reward learning cannot.
- <u>Applicability</u>: IDM can be used in real-world settings, as an investigative device for understanding *human* decisions.

Normative-Descriptive Questions Consider medical diagnosis, where there is often remarkable regional, institutional, and subgroup-level variability in practice [156-158], rendering detection and quantification of biases crucial [159–161]. Now in modeling an agent's behavior, reward learning asks: (1) "What does this (perfectly rational) agent appear to be optimizing?" And the answer takes the form of a function v. However, while v alone is often sufficient as an *intermediary* for imitation/apprenticeship, it is seldom what we actually want as an end by itself-for introspective understanding. Importantly, we often can articulate some version of what our preferences v are. In medical diagnosis, for instance, from the view of an investigator, the average relative healthcare cost/benefit of in-/correct diagnoses is certainly specifiable as a normative standard. So instead, we wish to ask: (2) "Given that this (boundedly rational) agent should optimize this v, how suboptimally do they appear to behave?" Clearly, such normative-descriptive questions are only possible with the generalized perspective of IDM (and IBRC): Here, v is specified (in θ_{norm}), whereas one or more behavioral parameters α , β , η are what we wish to recover (in θ_{desc}).

Decision Environments For our simulated setting (**DIAG**), we consider a POMDP where patients are diseased (s_+) or healthy (s_-) , and vital-signs measurements taken at each step

are noisily indicative of being disease-positive (x_+) or negative (x_-) . Actions consist of the decision to continue monitoring the patient $(u_=)$ —which yields evidence, but is also costly; or stopping and declaring a final diagnosis—and if so, a diseased (u_+) or healthy (u_-) call. Importantly, note that since we simulate $\tau, \omega \sim \sigma(\cdot|z, u)$, DIAG is a strict generalization of the diagnostic environment from [22] with a pointvalued, subjective $\tau, \omega \neq \tau_{env}, \omega_{env}$, and of the classic Tiger Problem in POMDP literature where $\tau, \omega = \tau_{env}, \omega_{env}$ [162].

For our real-world setting, we consider 6-monthly clinical data for 1,737 patients in the Alzheimer's Disease Neuroimaging Initiative [163] study (**ADNI**). The state space consists of normal function (s_{norm}), mild cognitive impairment (s_{mild}), and dementia (s_{dem}). For the action space, we consider ordering/not ordering an MRI—which yields evidence, but is costly. Results are classified per hippocampal volume: average (x_{avg}^{MRI}), high (x_{high}^{MRI}), low (x_{low}^{MRI}), not ordered (x_{none}^{MRI}); separately, the cognitive dementia rating test result—which is always measured—is classified as normal (x_{norm}^{CDR}), questionable impairment (x_{ques}^{CDR}), and suspected dementia (x_{susp}^{CDR}). So the observation space consists of such 12 combinations.

In DIAG, our normative specification (for v) is that diagnostic tests cost -1, correct diagnoses award 10, incorrect -36, and $\gamma = 0.95$. Accuracies are 70% in both directions (ω_{env}), and patients arrive in equal proportions (τ_{env}). But this is unknown to the agent: We simulate $\tilde{\sigma}$ by discretizing the space of models such that probabilities vary in $\pm 10\%$ increments from the (highest-likelihood) truth. In ADNI, the configuration is similar—except each MRI costs -1, while 2.5 is awarded once beliefs reach >90% certainty in any direction; also, $\tilde{\sigma}$ is centered at the IOHMM learned from the data. For simplicity, for $\tilde{\pi}$, $\tilde{\rho}$ we use uniform priors in both settings.

Computationally, inference is performed via MCMC in log-parameter space (i.e. $\log \alpha$, $\log \beta$, $\log \eta$) using standard methods, similar to e.g. Bayesian IRL [59,61,74]. In DIAG, we use 1,000 generated trajectories as basis for learning. Appendix B provides further details on experimental setup.



(a) Learned α for Various Flexibility Levels (b) Learned β , η for Non-adaptive Behavior (c) Learned β , η for Optimistic Behavior Figure 3. Inverse Bounded Rational Control. (a) Posteriors of α learned from extremely flexible ($\alpha^{true} = 10^{-5}$), flexible ($\alpha^{true} = 0.5$), and inflexible ($\alpha^{true} = 10$) behaviors (with β , η fixed as neutral and adaptive; similar plots can be obtained for those as well). (b) Joint posterior of β , η for neutral but non-adaptive behavior ($\beta^{true} = 10^3$, $\eta^{true} = 75$), and for (c) optimistic but adaptive behavior ($\beta^{true} = 1.25$, $\eta^{true} = 10^{-3}$).

5.1. Interpretability Figure 2 verifies (for DIAG) that different BRC behaviors accord with our intuitions. First, cet*eris paribus*, the flexibility (α) dimension manifests in how deterministically/stochastically optimal actions are selected (cf. willingness to deviate from action prior $\tilde{\pi}$): This is the notion of behavioral consistency [164] in psychology. Second, the optimism (β) dimension manifests in the illusion that diagnostic tests are more/less informative for subjective beliefs (cf. willingness to deviate from knowledge prior $\tilde{\sigma}$): This is the phenomenon of *over-/underreaction* [165]. Third, the adaptivity (η) dimension manifests in how much/little evidence is required for declaring a final diagnosis: This corresponds to base-rate neglect/confirmation bias [166]. Hence by *learning* the parameters α , β , η from data, IBRC provides an eminently interpretable example of behavior representation learning—one that exercises the IDM perspective (much more than just reward learning). Taking a Bayesian approach to the likelihood (Equation 25), Figure 3(a) verifies that—as expected—IBRC is capable of recovering different parameter values from their generated behaviors.

5.2. Expressivity Consider (i.) an agent who is biased towards optimism, but otherwise flexible and adaptive (Figure 2(b), top), and (ii.) an agent who is non-adaptive, but otherwise flexible and neutral (2(c), bottom). Now, to an external observer, both types of boundedness lead to similar styles of behavior: They both tend to declare final diagnoses earlier than a neutral and adaptive agent would (2(c), top)—that is, $\pi(u_+|z) \approx 1$ after only 2 (not 3) positive tests. Of course, the former does so due to overreaction (evaluating the evidence incorrectly), whereas the latter does so due to a lower threshold for stopping (despite correctly evaluating the evidence). As shown by Figures 3(b)–(c), IBRC does differentiate between these two different types of biased behaviors: This is revealing, if not necessarily surprising. Crucially, however, this distinction is not possible with conventional IRL. All else equal, let us perform Bayesian IRL on the very same behaviors—that is, to learn an effectively skewed v (while implicitly setting α, β, η to their perfectly rational limits). As it turns out, the recovered v for (i.) gives a cost-benefit ratio (of incorrect/correct diagnoses) of -2.70 ± 0.31 , and the recovered v for (ii.) gives a ratio of -2.60 ± 0.29 . Both

agents appear to penalize incorrect diagnoses much less than the normative specification of -3.60, which is consistent with them tending to commit to final diagnoses earlier than they should. However, this fails to *differentiate* between the two distinct underlying reasons for behaving in this manner.

5.3. Applicability Lastly, we highlight the potential utility of IDM in real-world settings as an investigative device for auditing and understanding human decision-making. Consider diagnostic patterns for identifying dementia in ADNI, for patients from different risk groups. For instance, we discover that while $\beta = 3.86$ for all patients, clinicians appear to be significantly less optimistic when diagnosing patients with the ApoE4 genetic risk factor ($\beta = 601.74$), for female patients ($\beta = 920.70$), and even more so for patients aged >75 ($\beta = 2,265.30$). Note that such attitudes toward risk factors align with prevailing medical knowledge [167-169]. Moreover, in addition to obtaining such agent-level interpretations of biases (i.e. using the learned parameters), we can also obtain trajectory-level interpretations of decisions (i.e. using the evolution of beliefs). Appendix D gives examples of ADNI patients using diagrams of trajectories in the belief simplex, to contextualize actions the taken by clinical decision-makers and identify potentially belated diagnoses.

6. Conclusion

In this paper, we motivated the importance of descriptive models of behavior as the bridge between normative and prescriptive decision analysis, and formalized a unifying perspective on inverse decision modeling for behavior representation learning. For future work, an important question lies in exploring differently structured parameterizations Θ that are *interpretable* for different purposes. After all, IBRC is only one prototype that exercises the IDM formalism more fully. Another question is to what extent different forms of the inverse problem is *identifiable* to begin with. For instance, it is well-known that even with perfect knowledge of a demonstrator's policy, in single environments we can only infer utility functions up to reward shaping. Thus balancing complexity, interpretability, and identifiability of decision models would be a challenging direction of work.

Acknowledgments

We would like to thank the reviewers for their generous feedback. This work was supported by Alzheimer's Research UK, The Alan Turing Institute under the EPSRC grant EP/N510129/1, the US Office of Naval Research, as well as the National Science Foundation under grant numbers 1407712, 1462245, 1524417, 1533983, and 1722516.

References

- Aiping Li, Songchang Jin, Lumin Zhang, and Yan Jia. A sequential decision-theoretic model for medical diagnostic system. *Technology and Healthcare*, 2015.
- [2] John A Clithero. Response times in economics: Looking through the lens of sequential sampling models. *Journal of Economic Psychology*, 2018.
- [3] Jan Drugowitsch, Rubén Moreno-Bote, and Alexandre Pouget. Relation between belief and performance in perceptual decision making. *PloS one*, 2014.
- [4] Gregory Wheeler. Bounded rationality. *SEP: Stanford Center for the Study of Language and Information*, 2018.
- [5] Thomas L Griffiths, Falk Lieder, and Noah D Goodman. Rational use of cognitive resources: Levels of analysis between the computational and the algorithmic. *Topics in cognitive science*, 2015.
- [6] Tim Genewein, Felix Leibfried, Jordi Grau-Moya, and Daniel Alexander Braun. Bounded rationality, abstraction, and hierarchical decision-making: An information-theoretic optimality principle. *Frontiers in Robotics and AI*, 2015.
- [7] Ned Augenblick and Matthew Rabin. Belief movement, uncertainty reduction, and rational updating. UC Berkeley-Haas and Harvard University Mimeo, 2018.
- [8] Pedro A Ortega, Daniel A Braun, Justin Dyer, Kee-Eung Kim, and Naftali Tishby. Information-theoretic bounded rationality. *arXiv preprint*, 2015.
- [9] L Robin Keller. The role of generalized utility theories in descriptive, prescriptive, and normative decision analysis. *Information and Decision Technologies*, 1989.
- [10] Ludwig Johann Neumann, Oskar Morgenstern, et al. *Theory of games and economic behavior*. Princeton university press Princeton, 1947.
- [11] Barbara A Mellers, Alan Schwartz, and Alan DJ Cooke. Judgment and decision making. *Annual review of psychology*, 1998.

- [12] Yisong Yue and Hoang M Le. Imitation learning (presentation). *International Conference on Machine Learning (ICML)*, 2018.
- [13] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. *International conference on Machine learning (ICML)*, 2004.
- [14] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation (NC)*, 1991.
- [15] Michael Bain and Claude Sammut. A framework for behavioural cloning. *Machine Intelligence (MI)*, 1999.
- [16] Umar Syed and Robert E Schapire. Imitation learning with a value-based prior. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2007.
- [17] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. *International conference* on artificial intelligence and statistics (AISTATS), 2010.
- [18] Umar Syed and Robert E Schapire. A reduction from apprenticeship learning to classification. Advances in neural information processing systems (NeurIPS), 2010.
- [19] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. *International conference on artificial intelligence and statistics* (AISTATS), 2011.
- [20] Bilal Piot, Matthieu Geist, and Olivier Pietquin. Boosted and reward-regularized classification for apprenticeship learning. *International conference on Autonomous agents and multi-agent systems (AA-MAS)*, 2014.
- [21] Daniel Jarrett, Ioana Bica, and Mihaela van der Schaar. Strictly batch imitation learning by energybased distribution matching. *Advances in neural information processing systems (NeurIPS)*, 2020.
- [22] Alihan Hüyük, Daniel Jarrett, Cem Tekin, and Mihaela van der Schaar. Explaining by imitating: Understanding decisions by interpretable policy learning. *International Conference on Learning Representations (ICLR)*, 2021.
- [23] Lionel Blondé and Alexandros Kalousis. Sampleefficient imitation learning via gans. *International* conference on artificial intelligence and statistics (AISTATS), 2019.

- [24] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation. *International Conference on Learning Representations* (*ICLR*), 2019.
- [25] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems (NeurIPS)*, 2016.
- [26] Wonseok Jeon, Seokin Seo, and Kee-Eung Kim. A bayesian approach to generative adversarial imitation learning. Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [27] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. Understanding the relation of bc and irl through divergence minimization. *ICML Workshop on Deep Generative Models for Highly Structured Data*, 2019.
- [28] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. *Conference on Robot Learning (CoRL)*, 2019.
- [29] Liyiming Ke, Matt Barnes, Wen Sun, Gilwoo Lee, Sanjiban Choudhury, and Siddhartha Srinivasa. Imitation learning as *f*-divergence minimization. *arXiv* preprint, 2019.
- [30] Liyiming Ke, Matt Barnes, Wen Sun, Gilwoo Lee, Sanjiban Choudhury, and Siddhartha Srinivasa. Imitation learning as *f*-divergence minimization. *International Workshop on the Algorithmic Foundations of Robotics (WAFR)*, 2020.
- [31] Kee-Eung Kim and Hyun Soo Park. Imitation learning via kernel mean embedding. AAAI Conference on Artificial Intelligence (AAAI), 2018.
- [32] Huang Xiao, Michael Herman, Joerg Wagner, Sebastian Ziesche, Jalal Etesami, and Thai Hong Linh. Wasserstein adversarial imitation learning. *arXiv* preprint, 2019.
- [33] Robert Dadashi, Leonard Hussenot, Matthieu Geist, and Olivier Pietquin. Primal wasserstein imitation learning. *International Conference on Learning Representations (ICLR)*, 2021.
- [34] Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribution matching. *International Conference on Learning Representations (ICLR)*, 2020.
- [35] Oleg Arenz and Gerhard Neumann. Non-adversarial imitation learning and its connections to adversarial methods. *arXiv preprint*, 2020.

- [36] Srivatsan Srinivasan and Finale Doshi-Velez. Interpretable batch irl to extract clinician goals in icu hypotension management. *AMIA Summits on Translational Science Proceedings*, 2020.
- [37] Xin Zhang, Yanhua Li, Ziming Zhang, and Zhi-Li Zhang. f-gail: Learning f-divergence for generative adversarial imitation learning. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [38] Nir Baram, Oron Anschel, and Shie Mannor. Modelbased adversarial imitation learning. *arXiv preprint*, 2016.
- [39] Nir Baram, Oron Anschel, and Shie Mannor. Modelbased adversarial imitation learning. *International Conference on Machine Learning (ICML)*, 2017.
- [40] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. *International conference on Machine learning (ICML)*, 2000.
- [41] Umar Syed and Robert E Schapire. A game-theoretic approach to apprenticeship learning. Advances in neural information processing systems (NeurIPS), 2008.
- [42] Umar Syed, Michael Bowling, and Robert E Schapire. Apprenticeship learning using linear programming. *International conference on Machine learning (ICML)*, 2008.
- [43] Edouard Klein, Matthieu Geist, and Olivier Pietquin. Batch, off-policy and model-free apprenticeship learning. *European Workshop on Reinforcement Learning (EWRL)*, 2011.
- [44] Takeshi Mori, Matthew Howard, and Sethu Vijayakumar. Model-free apprenticeship learning for transfer of human impedance behaviour. *IEEE-RAS International Conference on Humanoid Robots*, 2011.
- [45] Donghun Lee, Srivatsan Srinivasan, and Finale Doshi-Velez. Truly batch apprenticeship learning with deep successor features. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [46] Bilal Piot, Matthieu Geist, and Olivier Pietquin. Bridging the gap between imitation learning and irl. *IEEE transactions on neural networks and learning* systems, 2017.
- [47] Edouard Klein, Matthieu Geist, Bilal Piot, and Olivier Pietquin. Irl through structured classification. *Advances in neural information processing systems* (*NeurIPS*), 2012.

- [48] Edouard Klein, Bilal Piot, Matthieu Geist, and Olivier Pietquin. A cascaded supervised learning approach to inverse reinforcement learning. *Joint European conference on machine learning and knowledge discovery in databases (ECML)*, 2013.
- [49] Aristide CY Tossou and Christos Dimitrakakis. Probabilistic inverse reinforcement learning in unknown environments. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2013.
- [50] Vinamra Jain, Prashant Doshi, and Bikramjit Banerjee. Model-free irl using maximum likelihood estimation. *AAAI Conference on Artificial Intelligence* (*AAAI*), 2019.
- [51] Gergely Neu and Csaba Szepesvári. Apprenticeship learning using irl and gradient methods. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2007.
- [52] Monica Babes, Vukosi Marivate, and Michael L Littman. Apprenticeship learning about multiple intentions. *International conference on Machine learning (ICML)*, 2011.
- [53] Jonathan Ho, Jayesh Gupta, and Stefano Ermon. Model-free imitation learning with policy optimization. *International Conference on Machine Learning* (*ICML*), 2016.
- [54] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. *International conference on machine learning (ICML)*, 2016.
- [55] Matteo Pirotta and Marcello Restelli. Inverse reinforcement learning through policy gradient minimization. *AAAI Conference on Artificial Intelligence* (*AAAI*), 2016.
- [56] Alberto Maria Metelli, Matteo Pirotta, and Marcello Restelli. Compatible reward inverse reinforcement learning. Advances in Neural Information Processing Systems (NeurIPS), 2017.
- [57] Davide Tateo, Matteo Pirotta, Marcello Restelli, and Andrea Bonarini. Gradient-based minimization for multi-expert inverse reinforcement learning. *IEEE Symposium Series on Computational Intelligence* (SSCI), 2017.
- [58] Gergely Neu and Csaba Szepesvári. Training parsers by inverse reinforcement learning. *Machine learning* (*ML*), 2009.
- [59] Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2007.

- [60] Jaedeug Choi and Kee-Eung Kim. Map inference for bayesian irl. *Advances in Neural Information Processing Systems (NeurIPS)*, 2011.
- [61] Christos Dimitrakakis and Constantin A Rothkopf. Bayesian multitask irl. European workshop on reinforcement learning (EWRL), 2011.
- [62] Constantin A Rothkopf and Christos Dimitrakakis. Preference elicitation and inverse reinforcement learning. Joint European conference on machine learning and knowledge discovery in databases (ECML), 2011.
- [63] Sreejith Balakrishnan, Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Harold Soh. Efficient exploration of reward functions in inverse reinforcement learning via bayesian optimization. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [64] Ajay Kumar Tanwani and Aude Billard. Inverse reinforcement learning for compliant manipulation in letter handwriting. *National Center of Competence in Robotics (NCCR)*, 2013.
- [65] McKane Andrus. Inverse reinforcement learning for dynamics. *Dissertation, University of California at Berkeley*, 2019.
- [66] Stav Belogolovsky, Philip Korsunsky, Shie Mannor, Chen Tessler, and Tom Zahavy. Learning personalized treatments via irl. *arXiv preprint*, 2019.
- [67] Sid Reddy, Anca Dragan, and Sergey Levine. Where do you think you're going?: Inferring beliefs about dynamics from behavior. *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [68] Anirudha Majumdar, Sumeet Singh, Ajay Mandlekar, and Marco Pavone. Risk-sensitive inverse reinforcement learning via coherent risk models. *Robotics: Science and Systems*, 2017.
- [69] Sumeet Singh, Jonathan Lacotte, Anirudha Majumdar, and Marco Pavone. Risk-sensitive inverse reinforcement learning via semi-and non-parametric methods. *International Journal of Robotics Research*, 2018.
- [70] Jaedeug Choi and Kee-Eung Kim. Inverse reinforcement learning in partially observable environments. *International Joint Conference on Artificial Intelli*gence (IJCAI), 2009.
- [71] Jaedeug Choi and Kee-Eung Kim. Inverse reinforcement learning in partially observable environments. *Journal of Machine Learning Research* (*JMLR*), 2011.

- [72] Hamid R Chinaei and Brahim Chaib-Draa. An inverse reinforcement learning algorithm for partially observable domains with application on healthcare dialogue management. *International Conference on Machine Learning and Applications*, 2012.
- [73] Ioana Bica, Daniel Jarrett, Alihan Hüyük, and Mihaela van der Schaar. Learning what-if explanations for sequential decision-making. *International Conference on Learning Representations (ICLR)*, 2021.
- [74] Takaki Makino and Johane Takeuchi. Apprenticeship learning for model parameters of partially observable environments. *International Conference on Machine Learning (ICML)*, 2012.
- [75] Daniel Jarrett and Mihaela van der Schaar. Inverse active sensing: Modeling and understanding timely decision-making. *International Conference on Machine Learning*, 2020.
- [76] Kunal Pattanayak and Vikram Krishnamurthy. Inverse reinforcement learning for sequential hypothesis testing and search. *International Conference on Information Fusion (FUSION)*, 2020.
- [77] Matthew Golub, Steven Chase, and Byron Yu. Learning an internal dynamics model from control demonstration. *International Conference on Machine Learning (ICML)*, 2013.
- [78] Zhengwei Wu, Paul Schrater, and Xaq Pitkow. Inverse pomdp: Inferring what you think from what you do. *arXiv preprint*, 2018.
- [79] Saurabh Daptardar, Paul Schrater, and Xaq Pitkow. Inverse rational control with partially observable continuous nonlinear dynamics. *arXiv preprint*, 2019.
- [80] Minhae Kwon, Saurabh Daptardar, Paul Schrater, and Xaq Pitkow. Inverse rational control with partially observable continuous nonlinear dynamics. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [81] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. AAAI Conference on Artificial Intelligence (AAAI), 2008.
- [82] Abdeslam Boularias, Jens Kober, and Jan Peters. Relative entropy inverse reinforcement learning. *International conference on artificial intelligence and statistics (AISTATS)*, 2011.
- [83] Mrinal Kalakrishnan, Peter Pastor, Ludovic Righetti, and Stefan Schaal. Learning objective functions for

manipulation. International Conference on Robotics and Automation (ICRA), 2013.

- [84] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse reinforcement learning. *arXiv preprint*, 2015.
- [85] Chelsea Finn, Paul Christiano, Pieter Abbeel, and Sergey Levine. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *NeurIPS Workshop on Adversarial Training*, 2016.
- [86] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2018.
- [87] Ahmed H Qureshi, Byron Boots, and Michael C Yip. Adversarial imitation via variational inverse reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2019.
- [88] Paul Barde, Julien Roy, Wonseok Jeon, Joelle Pineau, Christopher Pal, and Derek Nowrouzezahrai. Adversarial soft advantage fitting: Imitation learning without policy optimization. Advances in neural information processing systems (NeurIPS), 2020.
- [89] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of maximum causal entropy. *International conference on Machine learning (ICML)*, 2010.
- [90] Zhengyuan Zhou, Michael Bloem, and Nicholas Bambos. Infinite time horizon maximum causal entropy inverse reinforcement learning. *IEEE Transactions on Automatic Control (TACON)*, 2017.
- [91] Kyungjae Lee, Sungjoon Choi, and Songhwai Oh. Maximum causal tsallis entropy imitation learning. Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [92] Tien Mai, Kennard Chan, and Patrick Jaillet. Generalized maximum causal entropy for inverse reinforcement learning. AAAI Conference on Artificial Intelligence (AAAI), 2019.
- [93] Michael Herman, Tobias Gindele, Jörg Wagner, Felix Schmitt, and Wolfram Burgard. Inverse reinforcement learning with simultaneous estimation of rewards and dynamics. *International conference on artificial intelligence and statistics (AISTATS)*, 2016.
- [94] Michael Herman. Simultaneous estimation of rewards and dynamics in irl. *Dissertation, Albert-Ludwigs-Universitat Freiburg*, 2016.

- [95] Layla El Asri, Bilal Piot, Matthieu Geist, Romain Laroche, and Olivier Pietquin. Score-based inverse reinforcement learning. *International conference* on Autonomous agents and multi-agent systems (AA-MAS), 2016.
- [96] Benjamin Burchfiel, Carlo Tomasi, and Ronald Parr. Distance minimization for reward learning from scored trajectories. AAAI Conference on Artificial Intelligence (AAAI), 2016.
- [97] Alexis Jacq, Matthieu Geist, Ana Paiva, and Olivier Pietquin. Learning from a learner. *International Conference on Machine Learning (ICML)*, 2019.
- [98] Daniel S Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations. *International Conference on Machine Learning (ICML)*, 2019.
- [99] Daniel S Brown, Wonjoon Goo, and Scott Niekum. Better-than-demonstrator imitation learning via automatically-ranked demonstrations. *Conference on Robot Learning (CoRL)*, 2020.
- [100] Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J Andrew Bagnell, Pieter Abbeel, and Jan Peters. An algorithmic perspective on imitation learning. *Foundations and Trends in Robotics*, 2018.
- [101] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. *International Conference on Machine Learning (ICML)*, 2017.
- [102] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *International Conference on Machine Learning (ICML)*, 2018.
- [103] Benjamin Eysenbach and Sergey Levine. If maxent rl is the answer, what is the question? *arXiv preprint*, 2019.
- [104] Wenjie Shi, Shiji Song, and Cheng Wu. Soft policy gradient method for maximum entropy deep reinforcement learning. *International Joint Conference* on Artificial Intelligence (IJCAI), 2019.
- [105] Rohin Shah, Noah Gundotra, Pieter Abbeel, and Anca Dragan. Inferring reward functions from demonstrators with unknown biases. *OpenReview*, 2018.
- [106] Rohin Shah, Noah Gundotra, Pieter Abbeel, and Anca D Dragan. On the feasibility of learning, rather than assuming, human biases for reward inference. *International Conference on Machine Learning* (*ICML*), 2019.

- [107] Jonathan Rubin, Ohad Shamir, and Naftali Tishby. Trading value and information in mdps. *Decision Making with Imperfect Decision Makers (Springer)*, 2012.
- [108] Alexandre Galashov, Siddhant M Jayakumar, Leonard Hasenclever, Dhruva Tirumala, Jonathan Schwarz, Guillaume Desjardins, Wojciech M Czarnecki, Yee Whye Teh, Razvan Pascanu, and Nicolas Heess. Information asymmetry in kl-regularized rl. International Conference on Learning Representations (ICLR), 2019.
- [109] Mark K Ho, David Abel, Jonathan D Cohen, Michael L Littman, and Thomas L Griffiths. The efficiency of human cognition reflects planned information processing. AAAI Conference on Artificial Intelligence (AAAI), 2020.
- [110] Stas Tiomkin and Naftali Tishby. A unified bellman equation for causal information and value in markov decision processes. *arXiv preprint arXiv:1703.01585*, 2017.
- [111] Felix Leibfried, Jordi Grau-Moya, and Haitham Bou-Ammar. An information-theoretic optimality principle for deep reinforcement learning. *NeurIPS Workshop on Deep Reinforcement Learning*, 2017.
- [112] Matthew Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable mdps. AAAI Conference on Artificial Intelligence (AAAI), 2015.
- [113] Pengfei Zhu, Xin Li, Pascal Poupart, and Guanghui Miao. On improving deep reinforcement learning for pomdps. *arXiv preprint*, 2017.
- [114] Maximilian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, and Shimon Whiteson. Deep variational reinforcement learning for pomdps. *International Conference on Machine Learning (ICML)*, 2018.
- [115] Amy Zhang, Zachary C Lipton, Luis Pineda, Kamyar Azizzadenesheli, Anima Anandkumar, Laurent Itti, Joelle Pineau, and Tommaso Furlanello. Learning causal state representations of partially observable environments. arXiv preprint, 2019.
- [116] Dongqi Han, Kenji Doya, and Jun Tani. Variational recurrent models for solving partially observable control tasks. arXiv preprint arXiv:1912.10703, 2019.
- [117] Joseph Futoma, Michael C Hughes, and Finale Doshi-Velez. Popcorn: Partially observed prediction constrained reinforcement learning. *International conference on artificial intelligence and statistics (AIS-TATS)*, 2020.

- [118] Richard D Smallwood and Edward J Sondik. The optimal control of partially observable markov processes over a finite horizon. *Operations research*, 1973.
- [119] Milos Hauskrecht. Value-function approximations for partially observable markov decision processes. *Journal of Artificial Intelligence Research (JAIR)*, 2000.
- [120] Joelle Pineau, Geoff Gordon, Sebastian Thrun, et al. Point-based value iteration: An anytime algorithm for pomdps. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2003.
- [121] Hanna Kurniawati, David Hsu, and Wee Sun Lee. Sarsop: Efficient point-based pomdp planning by approximating optimally reachable belief spaces. *Robotics: Science and systems*, 2008.
- [122] Mauricio Araya, Olivier Buffet, Vincent Thomas, and Françcois Charpillet. A pomdp extension with beliefdependent rewards. Advances in Neural Information Processing Systems (NeurIPS), 2010.
- [123] Mathieu Fehr, Olivier Buffet, Vincent Thomas, and Jilles Dibangoye. rho-pomdps have lipschitzcontinuous epsilon-optimal value functions. Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [124] F A Sonnenberg and J R Beck. Markov models in medical decision making: a practical guide. *Health Econ.*, 1983.
- [125] C H Jackson, L D Sharples, S G Thompson, S W Duffy, and E Couto. Multistate Markov models for disease progression with classification error. *Statistician*, 2003.
- [126] S E O'Bryant, S C Waring, C M Cullum, J Hall, L Lacritz, P J Massman, P J Lupo, J S Reisch, and R Doody. Staging dementia using Clinical Dementia Rating Scale Sum of Boxes scores: a Texas Alzheimer's research consortium study. *Arch. of Neurology*, 2008.
- [127] D Jarrett, J Yoon, and M van der Schaar. Matchnet: Dynamic prediction in survival analysis using convolutional neural networks. *NeurIPS Workshop* on Machine Learning for Health, 2018.
- [128] Daniel Jarrett, Jinsung Yoon, and Mihaela van der Schaar. Dynamic prediction in clinical survival analysis using temporal convolutional networks. *IEEE Journal of Biomedical and Health Informatics*, 2019.

- [129] P Petousis, A Winter, W Speier, D R Aberle, W Hsu, and A A T Bui. Using sequential decision making to improve lung cancer screening performance. *IEEE Access*, 2019.
- [130] F Cardoso, S Kyriakides, S Ohno, F Penault-Llorca, P Poortmans, I T Rubio, S Zackrisson, and E Senkus. Early breast cancer: ESMO Clinical Practice Guidelines for diagnosis, treatment and follow-up. *Anna. Oncology*, 2019.
- [131] A M Alaa and M van der Schaar. Attentive statespace modeling of disease progression. Advances in neural information processing systems (NeurIPS), 2019.
- [132] X Wang, D Sontag, and F Wang. Unsupervised learning of disease progression models. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2014.
- [133] Clemens Heuberger. Inverse combinatorial optimization. *Journal of combinatorial optimization*, 2004.
- [134] Kareem Amin and Satinder Singh. Towards resolving unidentifiability in inverse reinforcement learning. *arXiv preprint*, 2016.
- [135] Kareem Amin, Nan Jiang, and Satinder Singh. Repeated inverse reinforcement learning. Advances in neural information processing systems (NeurIPS), 2017.
- [136] Stuart Armstrong and Sören Mindermann. Occam's razor is insufficient to infer the preferences of irrational agents. Advances in neural information processing systems (NeurIPS), 2018.
- [137] Paul Christiano. The easy goal inference problem is still hard. *AI Alignment*, 2015.
- [138] Eric J Michaud, Adam Gleave, and Stuart Russell. Understanding learned reward functions. *NeurIPS Workshop on Deep Reinforcement Learning*, 2020.
- [139] Adam Gleave, Michael Dennis, Shane Legg, Stuart Russell, and Jan Leike. Quantifying differences in reward functions. *International Conference on Learning Representations (ICLR)*, 2021.
- [140] Daniel S Brown and Scott Niekum. Deep bayesian reward learning from preferences. *NeurIPS Workshop* on Safety and Robustness in Decision-Making, 2019.
- [141] Daniel S Brown, Russell Coleman, Ravi Srinivasan, and Scott Niekum. Safe imitation learning via fast bayesian reward inference from preferences. *International Conference on Machine Learning (ICML)*, 2020.

- [142] Nicolas Heess, David Silver, and Yee Whye Teh. Actor-critic reinforcement learning with energybased policies. *European Workshop on Reinforcement Learning (EWRL)*, 2013.
- [143] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. *International Conference on Machine Learning (ICML)*, 2016.
- [144] Nicolò Cesa-Bianchi, Claudio Gentile, Gábor Lugosi, and Gergely Neu. Boltzmann exploration done right. Advances in neural information processing systems (NeurIPS), 2017.
- [145] Amir Globerson, Eran Stark, Eilon Vaadia, and Naftali Tishby. The minimum information principle and its application to neural code analysis. *Proceedings* of the National Academy of Sciences, 2009.
- [146] Naftali Tishby and Daniel Polani. Information theory of decisions and actions. *Perception-action cycle* (*Springer*), 2011.
- [147] Pedro A Ortega and Daniel A Braun. Thermodynamics as a theory of decision-making with informationprocessing costs. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 2013.
- [148] Ian R Petersen, Matthew R James, and Paul Dupuis. Minimax optimal control of stochastic uncertain systems with relative entropy constraints. *IEEE Transactions on Automatic Control*, 2000.
- [149] Charalambos D Charalambous, Farzad Rezaei, and Andreas Kyprianou. Relations between information theory, robustness, and statistical mechanics of stochastic systems. *IEEE Conference on Decision* and Control (CDC), 2004.
- [150] Takayuki Osogami. Robustness and risk-sensitivity in markov decision processes. *Advances in Neural Information Processing Systems (NeurIPS)*, 2012.
- [151] Jordi Grau-Moya, Felix Leibfried, Tim Genewein, and Daniel A Braun. Planning with informationprocessing constraints and model uncertainty in markov decision processes. *Joint European conference on machine learning and knowledge discovery in databases (ECML)*, 2016.
- [152] Brian D Ziebart. Modeling purposeful adaptive behavior with the principle of maximum causal entropy. *Dissertation, Carnegie Mellon University*, 2010.
- [153] Gerhard Kramer. Directed information for channels with feedback. *Dissertation, ETH Zurich*, 1998.

- [154] James Massey. Causality, feedback and directed information. *International Symposium on Information Theory and Its Applications*, 1990.
- [155] Hans Marko. The bidirectional communication theory-a generalization of information theory. *IEEE Transactions on Communications*, 1973.
- [156] John B McKinlay, Carol L Link, et al. Sources of variation in physician adherence with clinical guidelines. *Journal of general internal medicine*, 2007.
- [157] Matthias Bock, Gerhard Fritsch, and David L Hepner. Preoperative laboratory testing. *Anesthesiology clinics*, 2016.
- [158] Jack W O'Sullivan, Carl Heneghan, Rafael Perera, Jason Oke, Jeffrey K Aronson, Brian Shine, and Ben Goldacre. Variation in diagnostic test requests and outcomes: a preliminary metric for openpathology. net. *Nature Scientific Reports*, 2018.
- [159] Yunjie Song, Jonathan Skinner, Julie Bynum, Jason Sutherland, John E Wennberg, and Elliott S Fisher. Regional variations in diagnostic practices. *New England Journal of Medicine*, (1), 2010.
- [160] Shannon K Martin and Adam S Cifu. Routine preoperative laboratory tests for elective surgery. *Journal* of the American Medical Association (JAMA), 2017.
- [161] M. Allen. Unnecessary tests and treatment explain why health care costs so much. *Scientific American*, 2017.
- [162] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 1998.
- [163] Razvan V Marinescu, Neil P Oxtoby, Alexandra L Young, Esther E Bron, Arthur W Toga, Michael W Weiner, Frederik Barkhof, Nick C Fox, Stefan Klein, Daniel C Alexander, et al. Tadpole challenge: Prediction of longitudinal evolution in alzheimer's disease. arXiv preprint, 2018.
- [164] Edi Karni and Zvi Safra. Behavioral consistency in sequential decisions. *Progress in Decision, Utility* and Risk Theory, 1991.
- [165] Kent Daniel, David Hirshleifer, and Avanidhar Subrahmanyam. Investor psychology and security market under-and overreactions. *The Journal of Finance*, 1998.
- [166] Amos Tversky and Daniel Kahneman. Evidential impact of base rates. *Stanford University Department Of Psychology*, 1981.

- [167] Charlotte L Allan and Klaus P Ebmeier. The influence of apoe4 on clinical progression of dementia: a meta-analysis. *International journal of geriatric psychiatry*, 2011.
- [168] Sylvaine Artero, Marie-Laure Ancelin, Florence Portet, A Dupuy, Claudine Berr, Jean-François Dartigues, Christophe Tzourio, Olivier Rouaud, Michel Poncet, Florence Pasquier, et al. Risk profiles for mild cognitive impairment and progression to dementia are gender specific. *Journal of Neurology, Neurosurgery & Psychiatry*, 2008.
- [169] Xue Hua, Derrek P Hibar, Suh Lee, Arthur W Toga, Clifford R Jack Jr, Michael W Weiner, Paul M Thompson, Alzheimer's Disease Neuroimaging Initiative, et al. Sex and age differences in atrophic rates: an adni study with n= 1368 mri scans. *Neurobiology* of aging, 2010.
- [170] Brendan O'Donoghue, Remi Munos, Koray Kavukcuoglu, and Volodymyr Mnih. Combining policy gradient and q-learning. *International Conference on Learning Representations (ICLR)*, 2017.
- [171] Momchil Tomov. Structure learning and uncertaintyguided exploration in the human brain. *Dissertation, Harvard University*, 2020.
- [172] Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Benjamin Eysenbach. F-irl: Inverse reinforcement learning via state marginal matching. *Conference on Robot Learning (CoRL)*, 2020.
- [173] Hong Jun Jeon, Smitha Milli, and Anca D Dragan. Reward-rational (implicit) choice: A unifying formalism for reward learning. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [174] Jeffrey Ely, Alexander Frankel, and Emir Kamenica. Suspense and surprise. *Journal of Political Economy*, 2015.
- [175] Ahmed M Alaa and Mihaela van der Schaar. Balancing suspense and surprise: Timely decision making with endogenous information acquisition. Advances in neural information processing systems (NeurIPS), 2016.
- [176] Owain Evans and Noah D Goodman. Learning the preferences of bounded agents. *NeurIPS Workshop* on Bounded Optimality, 2015.
- [177] Tan Zhi-Xuan, Jordyn L Mann, Tom Silver, Joshua B Tenenbaum, and Vikash K Mansinghka. Online

bayesian goal inference for boundedly-rational planning agents. Advances in neural information processing systems (NeurIPS), 2020.

- [178] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy extraction. Advances in neural information processing systems (NeurIPS), 2018.
- [179] Herman Yau, Chris Russell, and Simon Hadfield. What did you think would happen? explaining agent behaviour through intended outcomes. Advances in neural information processing systems (NeurIPS), 2020.
- [180] Tom Bewley, Jonathan Lawry, and Arthur Richards. Modelling agent policies with interpretable imitation learning. *TAILOR Workshop at ECAI*, 2020.
- [181] Tathagata Chakraborti, Sarath Sreedharan, Yu Zhang, and Subbarao Kambhampati. Plan explanations as model reconciliation: Moving beyond explanation as soliloquy. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.
- [182] Tathagata Chakraborti, Sarath Sreedharan, Sachin Grover, and Subbarao Kambhampati. Plan explanations as model reconciliation: an empirical study. *International Conference on Human-Robot Interaction (HRI)*, 2019.
- [183] Bradley Hayes and Julie A Shah. Improving robot controller transparency through autonomous policy explanation. *International Conference on Human-Robot Interaction (HRI)*, 2017.
- [184] Sarath Sreedharan, Utkash Soni, Mudit Verma, Siddharth Srivastava, and Subbarao Kambhampati. Bridging the gap: Providing post-hoc symbolic explanations for sequential decision-making problems with black box simulators. *ICML Workshop on Human-inthe-Loop Learning*, 2020.
- [185] Roy Fox and Naftali Tishby. Minimum-information lqg control part i: Memoryless controllers. *IEEE Conference on Decision and Control (CDC)*, 2016.
- [186] Roy Fox and Naftali Tishby. Minimum-information lqg control part ii: Retentive controllers. *IEEE Conference on Decision and Control (CDC)*, 2016.
- [187] Robert Babuska. Model-based imitation learning. Springer Encyclopedia of the Sciences of Learning, 2012.
- [188] Yoshua Bengio and Paolo Frasconi. An input output hmm architecture. *Advances in neural information processing systems (NeurIPS)*, 1995.

Appendices

Appendix A gives a longer discussion of merits and caveats; Appendix B gives further experiment details; Appendix C gives derivations of propositions; Appendix D shows illustrative trajectories; Appendix E gives a summary of notation.

A. Discussion

In this paper, we motivated the importance of *descriptive* models of behavior as the bridge between normative and prescriptive decision analysis [9–11] (Figure 4). On account of this, we formalized a unifying perspective on inverse decision modeling for behavior representation learning. Precisely, the *inverse decision model* of any observed behavior ϕ_{demo} is given by its projection $\phi_{\text{imit}}^* = F_{\theta_{\text{norm}}} \circ G_{\theta_{\text{norm}}}(\phi_{\text{demo}})$ onto the space $\Phi_{\theta_{\text{norm}}}$ of behaviors parameterizable by the structure designed for Θ and normative standards θ_{norm} specified. This formulation is general. For instance, it is agnostic as to the nature of agent and environment state spaces (whichamong other properties—are encoded in ψ); it is also agnostic as to whether the underlying forward problem is modelfree or model-based (which-among other properties-is encoded in θ). Per the priorities of the investigator (cf. imitation, apprenticeship, understanding, and other objectives), different choices can and should be made to balance the expressivity, interpretability, and tractability of learned models.

Partial Observability At first glance, our choice to accommodate partial observability may have appeared inconsequential. However, its significance becomes immediately apparent once we view an agent's behavior as induced by both a *decision* policy π as well as a *recognition* policy ρ , and—importantly—that not only may an agent's mapping from internal states into actions be suboptimal (viz. the former), but that their mapping from observations into beliefs may also be subjective (viz. the latter). Therefore in addition to the oft-studied, purely utility-centric nature of (perfectly rational) behavior, this generalized formalism immediately invites consideration of (boundedly rational) behaviors—that is, agents acting under knowledge uncertainty, biased by optimism/robustness, with policies distorted by the complexities of information processing required for decision-making.

Bounded Rationality While the IDM formalism subsumes most standard approaches to imitation learning, apprenticeship learning, and reward learning (cf. Table 1 and Table 3), we emphasize that—with very few exceptions [78–80]—the vast majority of original studies in these areas are limited to cases where $\theta_{desc} = v$ alone, or assume fully-observable environments (whence $S = \mathcal{X} = \mathcal{Z}$, and ρ simply being the identity function). Therefore our concrete example of *inverse bounded rational control* was presented as a prototypical instantiation of IDM that much more fully exercises the flexibility afforded by this generalized perspective. Importantly, while our notion of bounded rationality has (implicitly) been Figure 4. Normative, Prescriptive, and Descriptive Modeling. Recall the "lifecycle" of decision analysis (Section 1). As a paradigm of optimal behavior, normative standards serve as a theoretical benchmark. To guide imperfect agents toward this ideal, prescriptive advice serves to engineer behavior from humans in the loop. Importantly, however, this first requires an understanding of the imperfections—relative to the normative ideal—that require correcting. This is the goal of descriptive modeling—that is, to obtain an empirical account of existing behavior from observed data. Precisely, inverse decision modeling (middle) leverages a normative standard (left) to obtain an interpretable account of demonstrated behavior, thereby enabling the introspection of existing practices, which may inform construction of prescriptive guidelines (right).



present to varying degrees in (forward) control and reinforcement learning (cf. Table 2 and Table 4), "boundedness" has largely been limited to mean "noisy actions". To be precise, we may differentiate between three "levels" of boundedness:

- Imperfect Response: This is the shallowest form of boundedness, and includes Boltzmann-exploratory [142–144] and (locally) entropy-regularized [170] behaviors: It considers first that agents are perfect in their ability to compute the optimal values/policies; however, their actions are ultimately executed with an artificial layer of stochasticity.
- Capacity Constraints: Given an agent's model (e.g. τ, Q-network, etc.), the information processing needed in computing actions on the go is costly. We may view soft-opt-imal [101–104] and KL-regularized [107–111] planning and learning as examples. However, these do not model subjectivity of beliefs, adaptivity, or optimism/robustness.
- Model Imperfection: The agent's mental model itself is systematically flawed, due to uncertainty in knowledge, and to biases from optimism or pessimism. We may view certain robust MDPs (with penalties for deviating from priors) [148–151] as examples. However, these still do not account for partial observability (and biased recognition).

Now in the inverse direction, imitation/apprenticeship learning has typically viewed reward learning as but an intermediary, so classical methods have worked with perfectly rational planners [13, 40–45, 70–73]. Approaches that leverage probabilistic methods have usually simply used Boltzmannexploratory policies on top of optimal action-value functions (viz. imperfect response) [49–52, 59–63, 74, 75], or worked within maximum entropy planning/learning frameworks (viz. capacity constraints) [81–92]. Crucially, however, the corresponding parameters (i.e. inverse temperatures) have largely been treated as *pre-specified* parameters for learning v alone—not *learnable* parameters of interest by themselves. In contrast, what IDM allows (and what IBRC illustrates) is the "fullest" extent of boundedness—that is, where stochastic actions and subjective beliefs are endogenously the result of knowledge uncertainty and information processing constraints. Importantly, while recent work in imitation/apprenticeship have studied aspects of subjective dynamics that can be jointly learnable [67–69,93,94], they are limited to environments that are fully-observable and/or agents that have point-valued knowledge of environments—substantial simplifications that ignore how humans can and do make imperfect inferences from recognizing environment signals.

A.1. Important Distinctions

Our goal of *understanding* in IDM departs from the standard objectives of imitation and apprenticeship learning. As a result, some caveats and distinctions warrant special attention as pertains assumptions, subjectivity, and model accuracy.

Decision-maker vs. Investigator As noted in Section 3.3, the design of Θ (and specification of θ_{norm}) are not *assumptions*: We are not making "factual" claims concerning the underlying psychological processes that govern human behavior; these are hugely complex, and are the preserve of neuroscience and biology [171]. Instead, such specifications are active *design choices*: We seek to make the "effective" claim that an agent is behaving *as if* their generative mechanism were parameterized by the (interpretable) structure we designed for Θ . Therefore when we speak of "assumptions", it is important to distinguish between assumptions about the *agent* (of which we make none), versus assumptions about the *investigator* performing IDM (of which, by construction, we assume they have the ability to specify values for θ_{norm}).

In IBRC, for example, in learning β we are asking the question: "How much (optimistic/pessimistic) deviation from neutral knowledge does the agent appear to tolerate?" For this question to be meaningfully answered, we—as the investigator—must be able to produce a meaningful value for $\tilde{\sigma}$ to specify as part of θ_{norm} . In most cases, we are interested in deviations from some notion of "current medical knowledge", or what knowledge an "ideal" clinician may be expected to possess; thus we may—for instance—supply a value for $\tilde{\sigma}$ via models learned a priori from data. Of course, coming up such values for θ_{norm} is not trivial (not to mention entirely dependent on the problem and the investigator's objectives regarding interpretability); however, we emphasize that this does not involve assumptions regarding the *agent*.

Subjective vs. Objective Dynamics In imitation and apprenticeship learning, parameterizations of utilities and dynamics models are simply *intermediaries* for the downstream *Figure 5. Graphical Model.* In general, the environment's states (top) are only accessible via its emissions in response to actions (middle), which the agent incorporates by way of internal states (bottom). However, note that—unlike classic POMDP/IOHMM settings, here the agent's knowledge of the dynamics is subjective.



Figure 6. Backup Diagram. In IBRC, the backup operation (Theorem 4) transfers value information across three recursive "layers"—that is, of successor values for agent states (V), state-action pairs (Q), and state-action-model tuples (K). Indicated below are the utility and penalty terms collected along these backup operations.



task (of replicating expert actions or matching expert returns). As a result, no distinction needs be made between the "external" environment (with *objective* dynamics τ_{env} , ω_{env}) and the "internal" environment model that an agent works with (with *subjective* dynamics τ , ω). Indeed, if the learned model were to be evaluated based on live deployment in the real environment (as is the case in IL/IRL), it only makes sense that we stipulate τ , $\omega = \tau_{env}$, ω_{env} for the best results.

However, in IDM (and IBRC) we are precisely accounting for how an agent may appear to deviate from such perfect, point-valued knowledge of the environment. Disentangling subjective and objective dynamics is now critical: Both the forward recursion (Lemma 1) for occupancy measures and the backward recursion (Theorem 4) for value functions are computations *internal* to the agent's mind—and need not correspond to any notion of true environment dynamics. The *external* dynamics only comes into play when considering the distribution of trajectories $h \sim \phi_{\pi,\rho}$ induced by an agent's policies, which—by definition—manifests through (actual or potential) interaction with the real environment.

Demonstrated vs. Projected Behavior As advanced throughout, a primary benefit of the generalized perspective we develop is that we may ask *normative-descriptive questions* taking the form: "Given that this (boundedly rational) agent should optimize this v, how suboptimally do they appear to behave?" Precisely, as pertains IBRC we noted thatas the investigator—we are free to specify (what we deem) "meaningful" values for v within θ_{norm} , while recovering one or more behavioral parameters α, β, γ from θ_{desc} . Clearly, however, we are not at liberty to specify completely random values for v (or, more generally, that we are not at liberty to design Θ and θ_{norm} in an entirely arbitrary fashion). For one, the resulting inverse decision model may simply be a poor reflection the original behavior (i.e. the projection ϕ_{imit}^* onto $\Phi_{\theta_{norm}}$ may simply lose too much information from ϕ_{demo} .⁶

Without doubt, the usefulness of the inverse decision model (i.e. in providing valid interpretations of observed behavior) depends entirely on the design and specification of Θ and θ_{norm} , which requires care in practice. Most importantly, it should be verified that—under our designed parameterization—the *projected* behavior ϕ_{imit}^* is still a faithful model of the *demonstrated* behavior ϕ_{demo}^* . In particular, compared with fitting a black-box model for imitating behavior—or any standard method for imitation/apprenticeship learning, for that matter—it should be verified that our (interpretably parameterized) model does not suffer inordinately in terms of accuracy measures (i.e. in predicting *u* from *h*); otherwise the model (and its interpretation) would not be meaningful. In Appendix B, we perform precisely such a sanity check for IBRC, using a variety of standard benchmarks (Table 5).

A.2. Further Related Work

While relevant works have been noted throughout the manuscript, here we provide additional context for IDM and IBRC, and how notable techniques/frameworks relate to our work.

Inverse Decision Modeling Pertinent methods subsumed by our forward and inverse formalisms have been noted in Tables 2–3. In particular, techniques that can be formalized as instantiations of IDM are enumerated in Table 1. Broadly, for *imitation learning* these include behavioral cloning-like methods [14–21], as well as distribution-matching methods that directly match occupancy measures [23–39]; we defer to [12, 100] for more thorough surveys. For *apprenticeship learning* by inverse reinforcement learning, these include classic maximum-margin methods based on feature expectations [13, 40–45], maximum likelihood soft policy matching using Boltzmann-rational policies [51, 52], maximum entropy policies [50, 89–92], and Bayesian maximum a posteriori inference [59–63], as well as methods that leverage preference models and annotations for learning [95–99].

In this context, the novelty of the IDM formalism is two-fold. First, in defining a unifying framework that generalizes all prior techniques, IDM simultaneously opens up a new class of problems in *behavior representation learning* with consciously designed parameterizations. Specifically, in defining inverse decision models as projections in Φ -space induced by F, G, and Θ , the structure and decomposition chosen for $\Theta_{norm} \times \Theta_{desc}$ allows asking normative-descriptive questions that seek to *understand* observed decision-making behavior. Second, in elevating *recognition policies* to firstclass citizenship in partially-observable environments, IDM greatly generalizes the notion of "boundedness" in decisionmaking—that is, from the existing focus on noisy optimality in π , to the ideas of subjective dynamics σ and biased beliefupdates ρ (viz. discussion in the beginning of this section).

Orthogonal Frameworks Multiple studies have proposed frameworks that provide generalized treatments of different aspects of inverse reinforcement learning [28, 30, 35, 58, 60, 172, 173]. However, these are *orthogonal* to our purposes in the sense that they are primarily concerned with establishing connections between different aspects/subsets of the imitation/apprenticeship learning literature. These include loss-function perspectives [58] and Bayesian MAP perspectives [60] on inverse reinforcement learning, *f*-divergence minimization perspectives [28, 30] on distribution matching, connections between adversarial and non-adversarial methods for distribution matching [35], as well as different problem settings for learning reward functions [173]. But relative to the IDM formalism, all such frameworks operate within the special case of $\theta_{desc} = v$ (and full observability).

Case Study: GAIL Beyond aforementioned distinctions, another implication is that IDM defines a single language for understanding key results in such prior works. For example, we revisit the well-known result in [25] that gives rise to generative adversarial imitation learning ("GAIL"): It is instructive to recast it in more general—but simpler—terms. First, consider a *maximum entropy* learner in the MDP setting (cf. Table 2), paired with a *maximum margin* identification strategy with a parameter regularizer ζ (cf. Table 3):

$$F_{\theta_{\text{norm}}}^{\text{ME}}(\theta_{\text{desc}}) \doteq \phi_{\pi^*} \text{ where } \pi^* \doteq \operatorname{argmax}_{\pi} \mathbb{E}_{z \sim \rho_0} V_{\text{soft},\theta}^{\phi_{\pi}}(z) (26)$$

$$G_{\theta_{\text{norm}}}^{\text{MM}}(\phi) \doteq \operatorname{argmin}_{\theta_{\text{desc}}} \mathbb{E}_{z \sim \rho_0} [V_{\text{soft},\theta}^{\phi_{\text{imit}}}(z) - V_{\theta}^{\phi}(z)] + \zeta(\theta_{\text{desc}}) (27)$$

Second, consider a black-box *decision-rule* policy (cf. Table 2), where neural-network weights χ directly parameterize a policy network f_{decision} (and $\theta_{\text{desc}} = \chi$); this is paired with a *distribution matching* identification strategy (cf. Table 3):

$$F_{\theta_{\text{norm}}}^{\text{DR}}(\theta_{\text{desc}}) \doteq \operatorname{argmax}_{\pi} \delta(\pi - f_{\text{decision}}(\chi))$$
(28)

$$G_{\theta_{\text{norm}}}^{\text{DM}}(\phi) \doteq \operatorname{argmin}_{\theta,\zeta}^{*}(\phi_{\text{demo}} - \phi_{\text{imit}}) - \mathcal{H}_{\text{imit}}$$
 (29)

where distance measures are given by the convex conjugate ζ^* , and \mathcal{H}_{imit} gives the causal entropy of the imitating policy. Now, the primary motivation behind generative adversarial imitation learning is the observation that ζ -regularized maximum-margin soft IRL implicitly seeks a policy whose occupancy is close to the demonstrator's as measured by ζ^* . In IDM, this corresponds to a remarkably simple statement:

⁶Abstractly, this is not dissimilar to any type of model fitting problem: If the mismatch between the (unknown) data generating process and the (imposed) structure of the model is too great, then the quality of the model—by any reasonable measure—would suffer.

Proposition 6 (Ho and Ermon, Recast) Define the *beha-vior projections* induced by the composition of each pairing:

$$\operatorname{proj}_{\Phi_{\theta_{\operatorname{norm}}}}^{\operatorname{ME},\operatorname{MM}} \doteq F_{\theta_{\operatorname{norm}}}^{\operatorname{ME}} \circ G_{\theta_{\operatorname{norm}}}^{\operatorname{MM}}$$
(30)

$$\operatorname{proj}_{\Phi_{\theta_{\operatorname{norm}}}}^{\operatorname{DR},\operatorname{DM}} \doteq F_{\theta_{\operatorname{norm}}}^{\operatorname{DR}} \circ G_{\theta_{\operatorname{norm}}}^{\operatorname{DM}}$$
(31)

Then these projections are identical: $\text{proj}_{\Phi_{\theta_{norm}}}^{\text{ME,MM}} = \text{proj}_{\Phi_{\theta_{norm}}}^{\text{DR,DM}}$ (and inverse decision models thereby obtained are identical).

In their original context, the significance of this lies in the fact that the first pairing explicitly requires parameterizations via reward functions (which—in classic apprenticeship methods—is restricted to be linear/convex), whereas the second pairing allows arbitrary parameterization by neural networks (which—while black-box—are more flexible). In our language, this simply means that the first projection requires $\theta_{desc} = v$, while the second projection allows $\theta_{desc} = \chi$.

Inverse Bounded Rational Control Pertaining to IBRC, methods that are comparable and/or subsumed have been noted in Tables 1 and 4. In addition, the context of IBRC within existing notions of bounded rationality have been discussed in detail in the beginning of this section. Now, more broadly, we note that the study of imperfect behaviors [4] spans multiple disciplines: in cognitive science [5], biological systems [6], behavioral economics [7], and information theory [8]. Specifically, IBRC generalizes this latter class of information-theoretic approaches to bounded rationality. First, the notion of *flexibility* in terms of the informational effort in determining successive actions (cf. decision complexity) is present in maximum entropy [101–104] and KLregularized [107-111] agents. Second, the notion of tolerance in terms of the statistical surprise in adapting to successive beliefs (cf. recognition complexity) is present in behavioral economics [7, 174] and decision theory [75, 146, 175]. Third, the notions of optimism and pessimism in terms of the average regret in deviating from prior knowledge (cf. specification complexity) are present in robust planners [148–151].

On account of this, the novelty of the IBRC example is threefold. First, it is the first to present generalized recursions incorporating all three notions of complexity-that is, in the mappings into internal states, models, and actions. Second, IBRC does so in the partially-observable setting, whichas noted in above discussions-crucially generalizes the idea of subjective dynamics into subjective beliefs, thereby accounting for boundedness in the recognition process itself. Third (perhaps most importantly), IBRC is the first to consider the inverse problem-that is, of turning the entire formalism on its head to *learn* the parameterizations of such boundedness, instead of simply assuming known parameters as required by the forward problem. Finally, it is important to note that IBRC is simply one example: There are of course many possibilities for formulating boundedness, including such aspects as myopia and temporal inconsistency [176, 176]; we leave such applications for future work. Interpretable Behavior Representations Lastly, a variety of works have approached the task of representing behaviors in an interpretable manner. In inverse reinforcement learning, multiple works have focused on the reward function itself, specifying interpretable structures that explicitly express a decision-maker's preferences [62], behavior under time pressure [75], consideration of counterfactual outcomes [73], as well as intended goals [177]. Separately, another strand of research has focused on imposing interpretable structures onto *policy functions* themselves, such as representing policies in terms of decision trees [178] and intended outcomes [179] in the forward problem, or—in the inverse case-learning imitating policies based on decision trees [180] or decision boundaries [22]. In the context of IDM, both of these approaches can naturally be viewed as instantiations of our more general approach of learning representations of behavior through interpretably parameterized planners and inverse planners (as noted throughout Tables 1-3). Finally, for completeness also note that an orthogonal branch of research is dedicated to generating autonomous explanations of artificial behavior, as suggested updates to human models [181, 182], and also as responses to human queries in a shared [183] or user-specified vocabulary [184].

A.3. Future Work

A clear source of potential research lies in exploring differently structured parameterizations Θ to allow interpretable representation learning of behaviors. After all, beyond the black-box and reward-centric approaches in Table 1 and the handful of works that have sought to account for subjective dynamics [22, 67, 80, 93], our example of IBRC is only one such prototype that exercises the IDM formalism more fully. In developing more complex and/or expressive forward models, an important question to bear in mind is to what extent the inverse problem is identifiable. In most existing cases we have seen, the usual strategies-such as constraining scaling, shifting, reward shaping, as well as the use of Bayesian inference-is sufficient to recover meaningful values. However, we have also seen that in the extreme case of an arbitrary differentiable planner, any inverse problem immediately falls prey to the "no free lunch" result [105, 106, 136, 137]. Thus balancing aspects of complexity, interpretability, and identifiability of decision models would be an interesting direction of work. Finally, in this work we primarily focused on the idea of limited intentionality-that is, in the goalseeking nature of an agent and how they may be constrained in this respect. But the flip side is also interesting: One can explore the idea of limited attentionality-that is, in how an agent may be constrained in their ability to focus on sequences of past events. This idea is explored in [185,186] by analogy with information bottlenecks in sensors and memory capacities; however, there is much room for developing more human-interpretable parameterizations of how an agent may pay selective attention to observations over time.

B. Experiment Details

Computation In IBRC, we define the space of agent states (i.e. subjective beliefs) as $\mathcal{Z} \doteq \mathbb{R}^k$, where k is the number of world states (k=3 for ADNI, and k=2 for DIAG). To implement the backward recursion (Theorem 4), each dimension of Z is discretized with a resolution of 100, and the values V(z) in the resulting lattice are updated iteratively exactly according to the backup operator \mathbb{B}^* —until convergence (which is guaranteed by the fact that \mathbb{B}^* is contractive, therefore the fixed point is unique; see Appendix C). For evaluation at any point z, we (linearly) interpolate between the closest neighboring grid points. In terms of implementing the inverse problem in a Bayesian manner (i.e. to recover posterior distributions over Θ_{desc}), we perform MCMC in log-parameter space (i.e. $\log \alpha, \log \beta, \log \eta$). Specifically, the proposal distribution is zero-mean Gaussian with standard deviation 0.1, with every 10th step collected as a sample. In each instance, the initial 1,000 burn-in samples are discarded, and a total of 10,000 steps are taken after burn-in.

Recognition In the manuscript, we make multiple references to the *Bayes update*, in particular within the context of our (possibly-biased) belief-update (Equation 9). For completeness, we state this explicitly: Given point-valued knowledge of τ , ω , update $\rho_{\tau,\omega}(z'|z, u, x')$ is the Dirac delta centered at

$$p(s'|z, u, x', \tau, \omega) \doteq \mathbb{E}_{s \sim p(\cdot|z)} \left[\frac{\tau(s'|s, u)\omega(x'|u, s')}{\mathbb{E}_{s' \sim \tau(\cdot|s, u)}\omega(x'|u, s')} \right]$$
(32)

and the overall recognition policy is the expectation over such values of τ, ω (Equation 9). As noted in Section 4.1, in general $\tilde{\sigma}$ represents any prior distribution the agent is specified to have, and in particular can be some Bayesian posterior $p(\tau, \omega | \mathcal{E})$ given any form of experience \mathcal{E} . This can be modeled in any manner, and is not the focus of our work; what matters here is simply that the agent may *deviate* optimistically/pessimistically from such a prior. As noted in Section 5, for our purposes we simulate $\tilde{\sigma}$ by discretizing the space of models such that probabilities vary in $\pm 10\%$ increments from the (highest-likelihood) truth. In ADNI, this means $\tilde{\sigma}$ is centered at the IOHMM learned from the data.

Model Accuracy In Appendix A.1 we discussed the caveat: In order for an inverse decision model to provide valid *interpretations* of observed behavior, it should be verified that—under the designed parameterization—the projected behavior ϕ_{imit}^* is still an *accurate* model of the demonstrated behavior ϕ_{demo}^* . Here we perform such a sanity check for our IBRC example using the ADNI environment. We consider the following standard *benchmark algorithms*. First, in terms of black-box models for imitation learning, we consider behavioral cloning [15] with a recurrent neural network for observation-action histories (**RNN-Based BC-IL**); an adaptation of model-based imitation learning [187] to partially-observable settings, using the learned IOHMM as Table 5. Comparison of Model Accuracies. IBRC performs similarly to all benchmark algorithms in matching demonstrated actions. Results are computed using held-out samples based on 5-fold cross-validation. IBRC is slightly better-calibrated, and similar in precision-recall scores (differences are statistically insignificant).

Inverse Decision Model	Calibration (Low is Better)	PRC Score (High is Better)	
Black-Box Model:			
RNN-Based BC-IL	0.18 ± 0.05	0.81 ± 0.08	
IOHMM-Based BC-IL	0.19 ± 0.07	0.79 ± 0.11	
Joint IOHMM-Based BC-IL	0.17 ± 0.05	0.81 ± 0.09	
Reward-Centric Model:			
Bayesian PO-IRL	0.23 ± 0.01	0.78 ± 0.09	
Joint Bayesian PO-IRL	0.24 ± 0.01	0.79 ± 0.09	
Boundedly Rational Model:			
IBRC (with learned α, β, η)	0.16 ± 0.00	0.77 ± 0.01	

model (IOHMM-Based BC-IL); and a recently-proposed model-based imitation learning that allows for subjective dynamics [22] by jointly learning the agent's possibly-biased internal model and their probabilistic decision boundaries (Joint IOHMM-Based BC-IL). Second, in terms of classic reward-centric methods for apprenticeship learning, we consider Bayesian inverse reinforcement learning in partiallyobservable environments [75] equipped with the learned IOHMM as model (Bayesian PO-IRL); and-analogous to the black-box case—the equivalent of this method that trains the dynamics model jointly along with the agent's apprenticeship policy [74] (Joint Bayesian PO-IRL). Algorithms requiring learned models are given IOHMMs estimated using conventional methods [188]—which is the same method by which the true model is estimated in IBRC (that is, as part of the space of candidate models in the support of $\tilde{\sigma}$).

<u>Results</u>. Table 5 shows results of this comparison on predicting actions, computed using held-out samples based on 5fold cross-validation. Crucially, while IBRC has the advantage in terms of interpretability of parameterization, its performance—purely in terms of predicting actions—does not degrade: IBRC is slightly better in terms of calibration, and similar in precision-recall (differences are statistically insignificant), which—for our ADNI example—affirms the validity of IBRC as an (interpretable) representation of ϕ_{demo} .

Data Selection From the ADNI data, we first selected out anomalous cases without a cognitive dementia rating test result, which is almost always taken at every visit by every patient. Second, we also truncated patient trajectories at points where a visit is skipped (that is, if the next visit of a patient does not occur immediately after the 6-monthly period following the previous visit). This selection process leaves 1,626 patients out of the original 1,737, and the median number of consecutive visits for each patient is three. In measuring MRI outcomes, the "average" is defined to be within half a standard deviation of the population mean. Note that this is the same pre-processing method employed for ADNI in [22].

Implementation Details of implementation for benchmark algorithms follow the setup in [22], and are reproduced here: RNN-Based BC-IL: We train an RNN whose inputs are the observed histories h and whose outputs are the predicted probabilities $\hat{\pi}(u|h)$ of taking action u given the observed history h. The network consists of an LSTM unit of size 64and a fully-connected hidden layer of size 64. The crossentropy $\mathcal{L} = -\sum_{n=1}^{N} \sum_{t=1}^{T} \sum_{u \in \mathcal{U}} \mathbb{I}\{u_t = u\} \log \hat{\pi}(u|h)$ is minimized using the Adam optimizer with a learning rate of 0.001 until convergence (that is, when the loss does not improve for 100 consecutive iterations). Bayesian PO-IRL: The IOHMM parameters are initialized by sampling uniformly at random. Then, they are estimated and fixed using conventional IOHMM methods. The utility v is initialized as $\hat{v}^0(s, u) = \varepsilon_{s, u}$, where $\varepsilon_{s, u} \sim \mathcal{N}(0, 0.001^2)$. Then, it is estimated via MCMC sampling, during which new candidate samples are generated by adding Gaussian noise with standard deviation 0.001 to the previous sample. To form the final estimate, we average every 10th sample among the second set of 500 samples, ignoring the first 500 samples. To compute optimal Q-values, we use an off-the-shelf POMDP solver https://www.pomdp.org/code/index.html. Joint Bayesian PO-IRL: All parameters are initialized exactly the same way as in Bayesian PO-IRL. Then, both the IOHMM parameters and the utility are estimated jointly via MCMC sampling. In order to generate new candidate samples, with equal probabilities we either sample new IOHMM parameters from the posterior (but without changing v) or obtain a new v the same way we do in Bayesian PO-IRL (but without changing the IOHMM parameters). A final estimate is formed the same way as in Bayesian PO-IRL. IOHMM-Based BC-IL: The IOHMM parameters are initialized by sampling them uniformly at random. Then, they are estimated and fixed using conventional IOHMM methods. Given the IOHMM parameters, we parameterize policies using the method of [22], with the policy parameters $\{\mu_u\}_{u \in \mathcal{U}}$ (not to be confused with the occupancy measure " μ " as defined in the present work) initialized as $\hat{\mu}_u^0(s) = (1/|S| + \varepsilon_{u,s}) / \sum_{s' \in S} (1/|S| + \varepsilon_{u,s'})$, where $\varepsilon_{u,s'} \sim \mathcal{N}(0, 0.001^2)$. Then, they are estimated according solely to the action likelihoods in using the EM algorithm. The expected log-posterior is maximized using the Adam optimizer with learning rate 0.001 until convergence (that is, when the expected log-posterior does not improve for 100 consecutive iterations). Joint IOHMM-Based BC-IL: This corresponds exactly to the proposed method of [22] itself, which is similar to IOHMM-Based BC-IL except parameters are trained jointly. All parameters are initialized exactly the same way as before; then, the IOHMM parameters and the policy parameters are estimated jointly according to both the action likelihoods and the observation likelihoods simultaneously. The expected log-posterior is again maximized using the Adam optimizer with a learning rate of 0.001 until convergence (non-improvement for 100 consecutive iterations).

C. Proofs of Propositions

Lemma 1 (Forward Recursion) Define the forward operator $\mathbb{F}_{\pi,\rho}$: $\Delta(\mathcal{Z})^{\Delta(\mathcal{Z})}$ such that for any given $\mu \in \Delta(\mathcal{Z})$:

$$(\mathbb{F}_{\pi,\rho}\mu)(z) \doteq (1-\gamma)\rho_0(z) + \gamma(\mathbb{M}_{\pi,\rho}\mu)(z)$$
(12)

Then the occupancy $\mu_{\pi,\rho}$ is the (unique) fixed point of $\mathbb{F}_{\pi,\rho}$.

Proof. Start from the definition of $\mathbb{M}_{\pi,\rho}$; episodes are restarted on completion ad infinitum, so we can write $\mu_{\pi,\rho}$ as:

$$\mu_{\pi,\rho}(z) \doteq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t p(z_t = z | z_0 \sim \rho_0) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t ((\mathbb{M}_{\pi,\rho})^t \rho_0)(z)$$
(33)

Then we obtain the result by simple algebraic manipulation:

$$(1 - \gamma)\rho_{0}(z) + \gamma(\mathbb{M}_{\pi,\rho}\mu_{\pi,\rho})(z)$$

= $(1 - \gamma)\rho_{0}(z) + \gamma(1 - \gamma)\sum_{t=0}^{\infty}\gamma^{t}((\mathbb{M}_{\pi,\rho})^{t+1}\rho_{0})(z)$
= $(1 - \gamma)(\rho_{0}(z) + \sum_{t=0}^{\infty}\gamma^{t+1}((\mathbb{M}_{\pi,\rho})^{t+1}\rho_{0})(z))$
= $(1 - \gamma)\sum_{t=0}^{\infty}\gamma^{t}((\mathbb{M}_{\pi,\rho})^{t}\rho_{0})(z)$
= $\mu_{\pi,\rho}(z)$ (34)

For uniqueness, we use the usual conditions—that is, that the process induced by the environment and the agent's policies is ergodic, with a single closed communicating class.

Lemma 2 (Backward Recursion) Define the backward operator $\mathbb{B}_{\pi,\rho} : \mathbb{R}^{\mathcal{Z}} \to \mathbb{R}^{\mathcal{Z}}$ such that for any given $V \in \mathbb{R}^{\mathcal{Z}}$:

$$(\mathbb{B}_{\pi,\rho}V)(z) \doteq \mathbb{E}_{\substack{s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z)}} [v(s,u) + \mathbb{E}_{\substack{\tau,\omega \sim \sigma(\cdot|z,u) \\ s' \sim \tau(\cdot|s,u) \\ x' \sim \omega(\cdot|u,s') \\ z' \sim \rho_{\tau,\omega}(\cdot|z,u,x')}} (14)$$

Then the (dual) optimal V is the (unique) fixed point of $\mathbb{B}_{\pi,\rho}$; this is the *value function* considering knowledge uncertainty:

$$V^{\phi_{\pi,\rho}}(z) \doteq \sum_{t=0}^{\infty} \gamma^{t} \mathbb{E}_{\substack{s_t \sim p(\cdot|z_t) \\ u_t \sim \pi(\cdot|z_t) \\ \tau, \omega \sim \sigma(\cdot|z_t, u_t) \\ s_{t+1} \sim \tau(\cdot|s_t, u_t) \\ x_{t+1} \sim \omega(\cdot|u_t, s_{t+1}) \\ z_{t+1} \sim \rho_{\tau,\omega}(\cdot|z_t, u_t, x_{t+1})}$$
(15)

so we can equivalently write targets $J_{\pi,\rho} = \mathbb{E}_{z \sim \rho_0} V^{\phi_{\pi,\rho}}(z)$. Likewise, we can also define the (state-action) value function $Q^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathcal{Z} \times \mathcal{U}}$ —that is, $Q^{\phi_{\pi,\rho}}(z,u) \doteq \mathbb{E}_{s \sim p(\cdot|z)}[v(s,u) + \mathbb{E}_{\tau,\omega \sim \sigma(\cdot|z,u),...,z' \sim \rho_{\tau,\omega}(\cdot|z,u,x')} \gamma V^{\phi_{\pi,\rho}}(z')]$ given an action.

Proof. Start with the Lagrangian, with $V \in \mathbb{R}^{\mathbb{Z}}$: $\mathcal{L}_{\pi,\rho}(\mu, V)$

$$\begin{split} & \doteq J_{\pi,\rho} - \langle V, \mu - \gamma \mathbb{M}_{\pi,\rho} \mu - (1-\gamma)\rho_0 \rangle \\ &= \mathbb{E} \sum_{\substack{s \sim \mu_{\pi,\rho} \\ s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z)}} \upsilon(s, u) - \langle V, \mu - \gamma \mathbb{M}_{\pi,\rho} \mu - (1-\gamma)\rho_0 \rangle \\ &= \mathbb{E} \sum_{\substack{z \sim \mu_{\pi,\rho} \\ s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z) \\ u \sim \pi(\cdot|z) \\ v \sim \pi(\cdot|z, u) \\ s' \sim \tau(\cdot|s, u) \\ s' \sim \omega(\cdot|u, s') \\ z' \sim \rho(\cdot|z, u, x') \\ - \mathbb{E}_{z \sim \mu_{\pi,\rho}} V(z) + \langle V, (1-\gamma)\rho_0 \rangle \end{split}$$
(35)

$$= \mathbb{E}_{\substack{z \sim \mu_{\pi,\rho} \\ s \sim p(\cdot|z)}} [\mathbb{E}_{u \sim \pi(\cdot|z)}[\upsilon(s,u) \qquad (36)$$

$$+ \mathbb{E}_{\substack{\tau,\omega \sim \sigma(\cdot|z,u) \\ s' \sim \tau(\cdot|s,u) \\ x' \sim \omega(\cdot|u,s') \\ z' \sim \rho(\cdot|z,u,x')}} (36)$$

Then taking the gradient w.r.t. μ and setting it to zero yields:

$$V(z) = \mathbb{E}_{\substack{s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z)}} [\upsilon(s, u) + \mathbb{E}_{\substack{\tau, \omega \sim \sigma(\cdot|z, u) \\ s' \sim \tau(\cdot|s, u) \\ x' \sim \omega(\cdot|u, s') \\ z' \sim \rho(\cdot|z, u, x')}} (37)$$

For uniqueness, observe as usual that $\mathbb{B}_{\pi,\rho}$ is γ -contracting:

$$\begin{split} \|\mathbb{B}_{\pi,\rho}V - \mathbb{B}_{\pi,\rho}V'\|_{\infty} \\ &= \max_{z} \|\mathbb{E}_{\substack{u \sim \pi(\cdot|z) \\ \tau, \omega \sim \sigma(\cdot|z,u) \\ z' \sim \varrho_{\tau,\omega}(\cdot|z,u)}} \left[\gamma V(z') - \gamma V'(z')\right] \\ &\leq \max_{z} \mathbb{E}_{\substack{u \sim \pi(\cdot|z) \\ \tau, \omega \sim \sigma(\cdot|z,u) \\ z' \sim \varrho_{\tau,\omega}(\cdot|z,u)}} \left[\left|\gamma V(z') - \gamma V'(z')\right|\right] \quad (38) \\ &\leq \max_{z'} |\gamma V(z') - \gamma V'(z')| \\ &= \gamma \|V - V'\|_{\infty} \end{split}$$

which allows appealing to the contraction mapping theorem.

Proposition 3 (Backward Recursion) Define the backward operator $\mathbb{B}_{\pi,\rho} : \mathbb{R}^{\mathbb{Z}} \to \mathbb{R}^{\mathbb{Z}}$ such that for any given function $V \in \mathbb{R}^{\mathbb{Z}}$ and for any given coefficient values $\alpha, \beta, \eta \in \mathbb{R}$:

$$(\mathbb{B}_{\pi,\rho}V)(z) \doteq \mathbb{E}_{s \sim p(\cdot|z)} \left[-\alpha \log \frac{\pi(u|z)}{\tilde{\pi}(u)} + \upsilon(s, u) + \mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u)} \left[-\beta \log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} + \right] \right]$$

$$\mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u)} \left[-\eta \log \frac{\rho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V(z') \right]$$

$$(19)$$

Then the (dual) optimal V is the (unique) fixed point of $\mathbb{B}_{\pi,\rho}$; as before, this is the value function $V^{\phi_{\pi,\rho}}$ —which now includes the complexity terms. Likewise, we can also define the (state-action) $Q^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathbb{Z} \times \mathcal{U}}$ as the ¹/₃-step-ahead expectation, and the (state-action-model) $K^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathbb{Z} \times \mathcal{U} \times \mathcal{T} \times \mathcal{O}}$ as the ²/₃-steps-ahead expectation (which is new in this setup).

Proof. Start with the Lagrangian, now with the new multipliers $\alpha, \beta, \eta \in \mathbb{R}$ in addition to $V \in \mathbb{R}^{\mathcal{Z}}$: $\mathcal{L}_{\pi,\rho}(\mu, \alpha, \beta, \eta, V)$

$$\begin{split} & \doteq J_{\pi,\rho} - \langle V, \mu - \gamma \mathbb{M}_{\pi,\rho} \mu - (1-\gamma)\rho_0 \rangle \\ & -\alpha \cdot (\mathbb{I}_{\pi,\rho}[\pi; \tilde{\pi}] - A) - \beta \cdot (\mathbb{I}_{\pi,\rho}[\sigma; \tilde{\sigma}] - B) \\ & -\eta \cdot (\mathbb{I}_{\pi,\rho}[\varrho; \tilde{\varrho}] - C) \\ & = \mathbb{E} \sum_{\substack{z \sim \mu_{\pi,\rho} \\ s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z)}} v(s, u) - \langle V, \mu - \gamma \mathbb{M}_{\pi,\rho} \mu - (1-\gamma)\rho_0 \rangle \\ & -\alpha \cdot (\mathbb{E}_{z \sim \mu_{\pi,\rho}} D_{\mathrm{KL}}(\pi(\cdot|z) \| \tilde{\pi}) - A) \\ & -\beta \cdot (\mathbb{E} \sum_{\substack{z \sim \mu_{\pi,\rho} \\ u \sim \pi(\cdot|z)}} D_{\mathrm{KL}}(\sigma(\cdot|z, u) \| \tilde{\sigma}) - B) \\ & -\eta \cdot (\mathbb{E} \sum_{\substack{z \sim \mu_{\pi,\rho} \\ u \sim \pi(\cdot|z)}} D_{\mathrm{KL}}(\varrho_{\tau,\omega}(\cdot|z, u) \| \tilde{\varrho}) - C) \\ & \sum_{\substack{z \sim \mu(\cdot|z) \\ \tau, \omega \sim \sigma(\cdot|z, u)}} (39) \end{split}$$

$$= \mathbb{E}_{z \sim \mu_{\pi,\rho}} \upsilon(s, u) + \mathbb{E}_{z \sim \mu_{\pi,\rho}} \gamma V(z')$$

$$s \sim p(\cdot|z) \qquad u \sim \pi(\cdot|z) \qquad u \sim \pi(\cdot|z) \qquad \tau_{\infty} \sim \sigma(\cdot|z, u) \qquad s' \sim \tau(\cdot|s, u) \qquad x' \sim \omega(\cdot|u, s') \qquad z' \sim \rho(\cdot|z, u, x') \qquad - \mathbb{E}_{z \sim \mu_{\pi,\rho}} V(z) + \langle V, (1 - \gamma)\rho_0 \rangle \\ - \alpha \cdot (\mathbb{E}_{z \sim \mu_{\pi,\rho}} \log \frac{\pi(u|z)}{\tilde{\pi}(u)} - A) \\ - \beta \cdot (\mathbb{E}_{z \sim \mu_{\pi,\rho}} \log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} - B) \\ u \sim \pi(\cdot|z) \qquad \tau, \omega \sim \sigma(\cdot|z, u) \qquad s' \sim \tau(\cdot|s, u) \qquad s' \sim \tau(\cdot|z, u, s') \qquad z' \sim \rho(\cdot|z, u, x') \qquad = \mathbb{E}_{z \sim \mu_{\pi,\rho}} \left[\mathbb{E}_{u \sim \pi(\cdot|z)} \left[\upsilon(s, u) - \alpha \cdot (\log \frac{\pi(u|z)}{\tilde{\pi}(u)} - A) \right] \\ + \mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u)} \left[-\beta \cdot (\log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} - B) \right] \\ + \mathbb{E}_{s' \sim \tau(\cdot|s, u)} \left[-\eta \cdot (\log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} - B) \right] \\ + \mathbb{E}_{s' \sim \tau(\cdot|s, u)} \left[-\eta \cdot (\log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\varrho}(z')} - C) \right] \\ + \gamma V(z') \right] \left] - V(z) \right] + \langle V, (1 - \gamma)\rho_0 \rangle \qquad (40)$$

Then taking the gradient w.r.t. μ and setting it to zero yields:

$$V(z) = \mathbb{E}_{s \sim p(\cdot|z)} \left[-\alpha \log \frac{\pi(u|z)}{\tilde{\pi}(u)} + \upsilon(s, u) + \mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u)} \left[-\beta \log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} + (41) \right] \right]$$
$$\mathbb{E}_{\tau, \omega \sim \sigma(\cdot|s, u)} \left[-\eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V(z') \right]$$

For uniqueness, observe as before that $\mathbb{B}_{\pi,\rho}$ is γ -contracting: $\|\mathbb{B}_{\pi,\rho}V - \mathbb{B}_{\pi,\rho}V'\|_{\infty} \leq \gamma \|V - V'\|_{\infty}$; then appeal to the contraction mapping theorem for uniqueness of fixed point. The only change from before is the additional log terms, which—like the utility term—cancel out of the differences.

For Theorems 4 and 5, we give a single derivation for both:

Theorem 4 (Boundedly Rational Values) Define the backward operator $\mathbb{B}^* : \mathbb{R}^{\mathbb{Z}} \to \mathbb{R}^{\mathbb{Z}}$ such that for any $V \in \mathbb{R}^{\mathbb{Z}}$:

$$(\mathbb{B}^*V)(z) \doteq \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \exp(\frac{1}{\alpha}Q(z,u))$$
(20)

$$Q(z,u) \doteq \beta \log \mathbb{E}_{\tau,\omega \sim \tilde{\sigma}} \exp(\frac{1}{\beta}K(z,u,\tau,\omega))$$

$$K(z,u,\tau,\omega) \doteq + \mathbb{E}_{s \sim p(\cdot|z)}v(s,u)$$

$$\mathbb{E}_{\substack{s \sim p(\cdot|z) \\ s' \sim \tau(\cdot|s,u) \\ x' \sim \omega(\cdot|u,s') \\ z' \sim \rho_{\tau,\omega}(\cdot|z,u,x')}} \left[-\eta \log \frac{\varrho_{\tau,\omega}(z'|z,u)}{\tilde{\varrho}(z')} + \gamma V(z')\right]$$

Then the *boundedly rational value function* V^* for the (primal) optimal π^* , ρ^* is the (unique) fixed point of $\mathbb{B}^*_{\pi,\rho}$. (Note that both Q^* and K^* are immediately obtainable from this).

Theorem 5 (Boundedly Rational Policies) The *bounded-ly rational decision policy* (i.e. primal optimal) is given by:

$$\pi^*(u|z) = \frac{\tilde{\pi}(u)}{Z_Q^*(z)} \exp\left(\frac{1}{\alpha}Q^*(z,u)\right)$$
(21)

and the boundedly rational recognition policy is given by:

$$\rho^*(z'|z, u, x') = \mathbb{E}_{\tau, \omega \sim \sigma^*(\cdot|z, u)} \rho_{\tau, \omega}(z'|z, u, x') \text{, where}$$

$$\sigma^*(\tau, \omega|z, u) \doteq \frac{\tilde{\sigma}(\tau, \omega)}{Z_{K^*}(z, u)} \exp\left(\frac{1}{\beta} K^*(z, u, \tau, \omega)\right) (22)$$

where $Z_{Q^*}(z) = \mathbb{E}_{u \sim \tilde{\pi}} \exp(\frac{1}{\alpha}Q^*(z, u))$ and $Z_{K^*}(z, u) = \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \exp(\frac{1}{\beta}K^*(z, u, \tau, \omega))$ give the partition functions.

Proof. From Proposition 3, the (state) value $V^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathcal{Z}}$ is:

$$V^{\phi_{\pi,\rho}}(z) = \mathbb{E}_{s \sim p(\cdot|z)} \left[-\alpha \log \frac{\pi(u|z)}{\tilde{\pi}(u)} + v(s,u) + \mathbb{E}_{\tau,\omega \sim \sigma(\cdot|z,u)} \left[-\beta \log \frac{\sigma(\tau,\omega|z,u)}{\tilde{\sigma}(\tau,\omega)} + (42) \right] \right]$$
$$\mathbb{E}_{\tau,\omega \sim \sigma(\cdot|s,u)} \left[-\eta \log \frac{\rho_{\tau,\omega}(z'|z,u)}{\tilde{\varrho}(z')} + \gamma V^{\phi_{\pi,\rho}}(z') \right]$$

Define (state-action) $Q^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathcal{Z} \times \mathcal{U}}$ to be ahead by 1/3 steps:

$$Q^{\phi_{\pi,\rho}}(z,u) \doteq \mathbb{E}_{s\sim p(\cdot|z)} \left[\upsilon(s,u) + \\ \mathbb{E}_{\tau,\omega\sim\sigma(\cdot|z,u)} \left[-\beta \log \frac{\sigma(\tau,\omega|z,u)}{\tilde{\sigma}(\tau,\omega)} + \\ \mathbb{E}_{\substack{\tau,\omega\sim\sigma(\cdot|z,u)\\ x'\sim\omega(\cdot|u,s')\\ z'\sim\rho_{\tau,\omega}(\cdot|z,u,x')}} \left[-\eta \log \frac{\varrho_{\tau,\omega}(z'|z,u)}{\tilde{\varrho}(z')} + \gamma V^{\phi_{\pi,\rho}}(z') \right] \right]$$
(43)

and (state-action-model) $K^{\phi_{\pi,\rho}} \in \mathbb{R}^{\mathcal{Z} \times \mathcal{U} \times \mathcal{T} \times \mathcal{O}}$ by 2/3 steps:

$$K^{\phi_{\pi,\rho}}(z, u, \tau, \omega) \doteq \tag{44}$$

$$\mathbb{E}_{\substack{s \sim p(\cdot|z) \\ s' \sim \tau(\cdot|s, u) \\ x' \sim \omega(\cdot|u, s') \\ z' \sim \rho_{\tau,\omega}(\cdot|z, u, x')}} \left[-\eta \log \frac{\varrho_{\tau,\omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V^{\phi_{\pi,\rho}}(z') \right] \right]$$

The decision and recognition policies seek the optimizations:

extremize_{$$\pi$$} $V^{\phi_{\pi,\rho}}(z)$
s.t. $\mathbb{E}_{u \sim \pi(\cdot|z)} 1 = 1$ (45)

extremize_{$$\sigma$$} $Q^{\phi_{\pi,\rho}}(z,u)$
s.t. $\mathbb{E}_{\tau,\omega\sim\sigma(\cdot|z,u)}1 = 1$ (46)

Equations 42-44 are true in particular for optimal values, so

$$V^*(z) = \mathbb{E}_{u \sim \pi^*(\cdot|z)} \left[-\alpha \log \frac{\pi^*(u|z)}{\tilde{\pi}(u)} + Q^*(z,u) \right]$$
(47)

$$Q^{*}(z, u) = \mathbb{E}_{s \sim p(\cdot|z)} [\upsilon(s, u)] + \mathbb{E}_{\tau, \omega \sim \sigma^{*}(\cdot|z, u)} [-\beta \log \frac{\sigma^{*}(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} + K^{*}(z, u, \tau, \omega)]$$
(48)

Therefore for the extremizations we write the Lagrangians

$$\mathcal{L}(\pi^*, \lambda) \doteq V^*(z) + \lambda \cdot (\mathbb{E}_{u \sim \pi^*(\cdot|z)} 1 - 1)$$
(49)

$$\mathcal{L}(\sigma^*,\nu) \doteq Q^*(z,u) + \nu \cdot (\mathbb{E}_{\tau,\omega \sim \sigma^*(\cdot|z,u)}1 - 1) \quad (50)$$

Straightforward algebraic manipulation yields the policies:

$$\pi^*(u|z) = \frac{\tilde{\pi}(u_t)}{Z_{Q^*}(z)} \exp\left(\frac{1}{\alpha}Q^*(z,u)\right) \tag{51}$$

$$\sigma^*(\tau,\omega|z,u) = \frac{\sigma(\tau,\omega)}{Z_{K^*}(z,u)} \exp\left(\frac{1}{\beta}K^*(z,u,\tau,\omega)\right)$$
(52)

where partition functions $Z_{Q^*}(z)$ and $Z_{K^*}(z)$ are given by:

$$Z_{Q^*}(z) = \mathbb{E}_{u \sim \tilde{\pi}} \exp(\frac{1}{\alpha} Q^*(z, u))$$
(53)

$$Z_{K^*}(z, u) = \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \exp(\frac{1}{\beta} K^*(z, u, \tau, \omega))$$
(54)

which proves Theorem 5. Then Theorem 4 is obtained by plugging back into the backward recursion (Proposition 3).

For uniqueness, we want $\|\mathbb{B}V - \mathbb{B}V'\|_{\infty} \leq \gamma \|V - V'\|_{\infty}$. Let $\|V - V'\|_{\infty} = \varepsilon (\max_{z'}|V(z') - V'(z')| = \varepsilon)$. Now, $(\mathbb{B}^*V)(z)$

$$\begin{split} &\doteq \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp \left(\frac{1}{\alpha} \left(\mathbb{E}_{s \sim p(\cdot|z)} \upsilon(s, u) \right. \right. \\ &+ \beta \log \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \left[\exp \left(\frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \right[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V(z') \right] \right) \right] \right) \right] \\ &\leq \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp \left(\frac{1}{\alpha} \left(\mathbb{E}_{s \sim p(\cdot|z)} \upsilon(s, u) \right. \\ &+ \beta \log \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \left[\exp \left(\frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \right[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma (V'(z') + \varepsilon) \right] \right) \right]) \right] \right] \\ &= \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp \left(\frac{1}{\alpha} \left(\mathbb{E}_{s \sim p(\cdot|z)} \upsilon(s, u) \right. \\ &+ \beta \log \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \left[\exp \left(\frac{1}{\beta} \gamma \varepsilon + \frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \left[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V'(z') \right] \right) \right]) \right] \right] \\ &= \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp \left(\frac{1}{\alpha} \left(\mathbb{E}_{s \sim p(\cdot|z)} \upsilon(s, u) \right. \\ &+ \beta \log \left(\exp(\frac{1}{\beta} \gamma \varepsilon) \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \left[\exp \left(\frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \left[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V'(z') \right] \right] \right]) \right] \right] \\ &= \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp \left(\frac{1}{\alpha} \gamma \varepsilon + \frac{1}{\alpha} \left(\mathbb{E}_{s \sim p(\cdot|z)} \upsilon(s, u) \right. \\ &+ \beta \log \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \left[\exp \left(\frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \left[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V'(z') \right] \right] \right]) \right] \right] \\ &= \alpha \log \left(\exp(\frac{1}{\alpha} \gamma \varepsilon) \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp\left(\frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \left[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V'(z') \right] \right] \right]) \right] \right] \\ &= \gamma \varepsilon + \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \left[\exp\left(\frac{1}{\beta} \mathbb{E}_{z' \sim \varrho_{\tau, \omega}(\cdot|z, u)} \left[\right. \\ &- \eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V'(z') \right] \right] \right]) \right] \\ &= \gamma \varepsilon + (\mathbb{B}^* V')(z) \tag{55}$$

Likewise, we can show that $(\mathbb{B}^*V)(z) \ge (\mathbb{B}^*V')(z) - \gamma \varepsilon$. Hence $\max_z |(\mathbb{B}V)(z) - (\mathbb{B}V')(z)| = ||\mathbb{B}V - \mathbb{B}V'||_{\infty} \le \gamma \epsilon$. Note on Equation 24: Note that we originally formulated "soft policy matching" in Table 3 as a forward Kullback-Leibler divergence expression. However, analogously to maximum likelihood in supervised learning, the entropy terms drop out of the optimization, which yields Equation 24. To see this, note that the causally-conditioned probability is simply the product of conditional probabilities at each time step, and each conditional is "Markovianized" using beliefs z_t (i.e. Equation 25).

D. Illustrative Trajectories

Here we direct attention to the potential utility of IBRC (and—more generally—instantiations of the IDM paradigm) as an "investigative device" for auditing and quantifying individual decisions. In Figure 7, we see that modeling the evolution of a decision-maker's subjective beliefs provides a concrete basis for analyzing the corresponding sequence of actions chosen. Each vertex of the belief simplex corresponds to one of the three stable Alzheimer's diagnoses, and each point within the simplex corresponds to a unique belief (i.e. probability distribution). The closer the point is to a vertex (i.e. disease state), the higher the probability assigned to that state. For instance, if the belief is located exactly in the middle of the simplex (i.e. equidistant from all vertices), then all states are believed to be equally likely. Note that this is visual presentation is done similarly to [22], where decision trajectories within belief simplices are first visualized in this manner—with the core difference here being that the decision policies (hence decision boundaries thereby induced) are computed using a different technique.



Figure 7. Decision Trajectories. Examples of apparent beliefs and actions of a clinical decision-maker regarding real patients, including cases where: (a) the clinician's decisions coincide with those that would have been dictated by a "perfectly-rational" policy—despite their bounded rationality; (b) the clinician fails to make "perfectly-rational" decisions (in this context, the "boundedness" of the clinician could be due to any number of issues encountered during the diagnostic process); and (c) a patient who—apparently—could have been diagnosed much earlier than they actually were, but for the clinician not having followed the decisions prescribed by the "perfectly-rational" policy.

E. Summary of Notation

Notation	Meaning	(first defined in)	Notation	Meaning	(first defined in)
ψ	problem setting	Section 3.1	s	environment state	Section 3.1
x	environment emission	Section 3.1	z	agent state, i.e. belief	Section 3.1
u	agent emission, i.e. action	Section 3.1	$ au_{ m env}$	environment transition	Section 3.1
au	subjective transition	Section 3.1	$\omega_{ m env}$	environment emission	Section 3.1
ω	subjective emission	Section 3.1	v	utility (i.e. reward) function	Section 3.1
γ	discount factor	Section 3.1	ϕ	behavior	Section 3.1
$\phi_{ m demo}$	demonstrated behavior	Section 3.2	$\phi_{ m imit}$	imitation behavior	Section 3.2
θ	planning parameter	Section 3.1	$\dot{\theta}_{\rm norm}$	normative parameter	Section 3.2
$\theta_{ m desc}$	descriptive parameter	Section 3.2	π	decision policy	Section 3.1
ρ	recognition policy	Section 3.1	σ	specification policy	Section 4.1
F	forward planner	Section 3.1	G	inverse planner	Section 3.2
α^{-1}	flexibility coefficient	Section 4.2	β^{-1}	optimism coefficient	Section 4.2
η^{-1}	adaptivity coefficient	Section 4.2	$ ilde{\pi}$	action prior	Section 4.2
$\tilde{\sigma}$	model prior	Section 4.2	$\tilde{\varrho}$	belief prior	Section 4.2

References

- Aiping Li, Songchang Jin, Lumin Zhang, and Yan Jia. A sequential decision-theoretic model for medical diagnostic system. *Technology and Healthcare*, 2015.
- [2] John A Clithero. Response times in economics: Looking through the lens of sequential sampling models. *Journal of Economic Psychology*, 2018.
- [3] Jan Drugowitsch, Rubén Moreno-Bote, and Alexandre Pouget. Relation between belief and performance in perceptual decision making. *PloS one*, 2014.
- [4] Gregory Wheeler. Bounded rationality. SEP: Stanford Center for the Study of Language and Information, 2018.
- [5] Thomas L Griffiths, Falk Lieder, and Noah D Goodman. Rational use of cognitive resources: Levels of analysis between the computational and the algorithmic. *Topics in cognitive science*, 2015.
- [6] Tim Genewein, Felix Leibfried, Jordi Grau-Moya, and Daniel Alexander Braun. Bounded rationality, abstraction, and hierarchical decision-making: An information-theoretic optimality principle. *Frontiers in Robotics and AI*, 2015.
- [7] Ned Augenblick and Matthew Rabin. Belief movement, uncertainty reduction, and rational updating. UC Berkeley-Haas and Harvard University Mimeo, 2018.
- [8] Pedro A Ortega, Daniel A Braun, Justin Dyer, Kee-Eung Kim, and Naftali Tishby. Information-theoretic bounded rationality. *arXiv preprint*, 2015.
- [9] L Robin Keller. The role of generalized utility theories in descriptive, prescriptive, and normative decision analysis. *Information and Decision Technologies*, 1989.
- [10] Ludwig Johann Neumann, Oskar Morgenstern, et al. *Theory of games and economic behavior*. Princeton university press Princeton, 1947.
- [11] Barbara A Mellers, Alan Schwartz, and Alan DJ Cooke. Judgment and decision making. *Annual review of psychology*, 1998.
- [12] Yisong Yue and Hoang M Le. Imitation learning (presentation). *International Conference on Machine Learning (ICML)*, 2018.
- [13] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. *International conference on Machine learning (ICML)*, 2004.

- [14] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation (NC)*, 1991.
- [15] Michael Bain and Claude Sammut. A framework for behavioural cloning. *Machine Intelligence (MI)*, 1999.
- [16] Umar Syed and Robert E Schapire. Imitation learning with a value-based prior. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2007.
- [17] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. *International conference* on artificial intelligence and statistics (AISTATS), 2010.
- [18] Umar Syed and Robert E Schapire. A reduction from apprenticeship learning to classification. Advances in neural information processing systems (NeurIPS), 2010.
- [19] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. *International conference on artificial intelligence and statistics* (AISTATS), 2011.
- [20] Bilal Piot, Matthieu Geist, and Olivier Pietquin. Boosted and reward-regularized classification for apprenticeship learning. *International conference on Autonomous agents and multi-agent systems (AA-MAS)*, 2014.
- [21] Daniel Jarrett, Ioana Bica, and Mihaela van der Schaar. Strictly batch imitation learning by energybased distribution matching. Advances in neural information processing systems (NeurIPS), 2020.
- [22] Alihan Hüyük, Daniel Jarrett, Cem Tekin, and Mihaela van der Schaar. Explaining by imitating: Understanding decisions by interpretable policy learning. *International Conference on Learning Representations (ICLR)*, 2021.
- [23] Lionel Blondé and Alexandros Kalousis. Sampleefficient imitation learning via gans. *International conference on artificial intelligence and statistics* (AISTATS), 2019.
- [24] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation. *International Conference on Learning Representations* (*ICLR*), 2019.

- [25] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. Advances in neural information processing systems (NeurIPS), 2016.
- [26] Wonseok Jeon, Seokin Seo, and Kee-Eung Kim. A bayesian approach to generative adversarial imitation learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [27] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. Understanding the relation of bc and irl through divergence minimization. *ICML Workshop on Deep Generative Models for Highly Structured Data*, 2019.
- [28] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. *Conference on Robot Learning (CoRL)*, 2019.
- [29] Liyiming Ke, Matt Barnes, Wen Sun, Gilwoo Lee, Sanjiban Choudhury, and Siddhartha Srinivasa. Imitation learning as *f*-divergence minimization. *arXiv* preprint, 2019.
- [30] Liyiming Ke, Matt Barnes, Wen Sun, Gilwoo Lee, Sanjiban Choudhury, and Siddhartha Srinivasa. Imitation learning as *f*-divergence minimization. *International Workshop on the Algorithmic Foundations of Robotics (WAFR)*, 2020.
- [31] Kee-Eung Kim and Hyun Soo Park. Imitation learning via kernel mean embedding. AAAI Conference on Artificial Intelligence (AAAI), 2018.
- [32] Huang Xiao, Michael Herman, Joerg Wagner, Sebastian Ziesche, Jalal Etesami, and Thai Hong Linh. Wasserstein adversarial imitation learning. *arXiv* preprint, 2019.
- [33] Robert Dadashi, Leonard Hussenot, Matthieu Geist, and Olivier Pietquin. Primal wasserstein imitation learning. *International Conference on Learning Representations (ICLR)*, 2021.
- [34] Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribution matching. *International Conference on Learning Representations (ICLR)*, 2020.
- [35] Oleg Arenz and Gerhard Neumann. Non-adversarial imitation learning and its connections to adversarial methods. *arXiv preprint*, 2020.
- [36] Srivatsan Srinivasan and Finale Doshi-Velez. Interpretable batch irl to extract clinician goals in icu hypotension management. *AMIA Summits on Translational Science Proceedings*, 2020.

- [37] Xin Zhang, Yanhua Li, Ziming Zhang, and Zhi-Li Zhang. f-gail: Learning f-divergence for generative adversarial imitation learning. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [38] Nir Baram, Oron Anschel, and Shie Mannor. Modelbased adversarial imitation learning. *arXiv preprint*, 2016.
- [39] Nir Baram, Oron Anschel, and Shie Mannor. Modelbased adversarial imitation learning. *International Conference on Machine Learning (ICML)*, 2017.
- [40] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. *International conference on Machine learning (ICML)*, 2000.
- [41] Umar Syed and Robert E Schapire. A game-theoretic approach to apprenticeship learning. Advances in neural information processing systems (NeurIPS), 2008.
- [42] Umar Syed, Michael Bowling, and Robert E Schapire. Apprenticeship learning using linear programming. *International conference on Machine learning (ICML)*, 2008.
- [43] Edouard Klein, Matthieu Geist, and Olivier Pietquin. Batch, off-policy and model-free apprenticeship learning. *European Workshop on Reinforcement Learning (EWRL)*, 2011.
- [44] Takeshi Mori, Matthew Howard, and Sethu Vijayakumar. Model-free apprenticeship learning for transfer of human impedance behaviour. *IEEE-RAS International Conference on Humanoid Robots*, 2011.
- [45] Donghun Lee, Srivatsan Srinivasan, and Finale Doshi-Velez. Truly batch apprenticeship learning with deep successor features. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [46] Bilal Piot, Matthieu Geist, and Olivier Pietquin. Bridging the gap between imitation learning and irl. *IEEE transactions on neural networks and learning* systems, 2017.
- [47] Edouard Klein, Matthieu Geist, Bilal Piot, and Olivier Pietquin. Irl through structured classification. Advances in neural information processing systems (NeurIPS), 2012.
- [48] Edouard Klein, Bilal Piot, Matthieu Geist, and Olivier Pietquin. A cascaded supervised learning approach to inverse reinforcement learning. *Joint Eu*ropean conference on machine learning and knowledge discovery in databases (ECML), 2013.

- [49] Aristide CY Tossou and Christos Dimitrakakis. Probabilistic inverse reinforcement learning in unknown environments. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2013.
- [50] Vinamra Jain, Prashant Doshi, and Bikramjit Banerjee. Model-free irl using maximum likelihood estimation. *AAAI Conference on Artificial Intelligence* (*AAAI*), 2019.
- [51] Gergely Neu and Csaba Szepesvári. Apprenticeship learning using irl and gradient methods. *Conference* on Uncertainty in Artificial Intelligence (UAI), 2007.
- [52] Monica Babes, Vukosi Marivate, and Michael L Littman. Apprenticeship learning about multiple intentions. *International conference on Machine learning (ICML)*, 2011.
- [53] Jonathan Ho, Jayesh Gupta, and Stefano Ermon. Model-free imitation learning with policy optimization. *International Conference on Machine Learning* (*ICML*), 2016.
- [54] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. *International conference on machine learning (ICML)*, 2016.
- [55] Matteo Pirotta and Marcello Restelli. Inverse reinforcement learning through policy gradient minimization. *AAAI Conference on Artificial Intelligence* (*AAAI*), 2016.
- [56] Alberto Maria Metelli, Matteo Pirotta, and Marcello Restelli. Compatible reward inverse reinforcement learning. Advances in Neural Information Processing Systems (NeurIPS), 2017.
- [57] Davide Tateo, Matteo Pirotta, Marcello Restelli, and Andrea Bonarini. Gradient-based minimization for multi-expert inverse reinforcement learning. *IEEE Symposium Series on Computational Intelligence* (SSCI), 2017.
- [58] Gergely Neu and Csaba Szepesvári. Training parsers by inverse reinforcement learning. *Machine learning* (*ML*), 2009.
- [59] Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2007.
- [60] Jaedeug Choi and Kee-Eung Kim. Map inference for bayesian irl. *Advances in Neural Information Processing Systems (NeurIPS)*, 2011.
- [61] Christos Dimitrakakis and Constantin A Rothkopf. Bayesian multitask irl. *European workshop on reinforcement learning (EWRL)*, 2011.

- [62] Constantin A Rothkopf and Christos Dimitrakakis. Preference elicitation and inverse reinforcement learning. Joint European conference on machine learning and knowledge discovery in databases (ECML), 2011.
- [63] Sreejith Balakrishnan, Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Harold Soh. Efficient exploration of reward functions in inverse reinforcement learning via bayesian optimization. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [64] Ajay Kumar Tanwani and Aude Billard. Inverse reinforcement learning for compliant manipulation in letter handwriting. *National Center of Competence in Robotics (NCCR)*, 2013.
- [65] McKane Andrus. Inverse reinforcement learning for dynamics. *Dissertation, University of California at Berkeley*, 2019.
- [66] Stav Belogolovsky, Philip Korsunsky, Shie Mannor, Chen Tessler, and Tom Zahavy. Learning personalized treatments via irl. *arXiv preprint*, 2019.
- [67] Sid Reddy, Anca Dragan, and Sergey Levine. Where do you think you're going?: Inferring beliefs about dynamics from behavior. Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [68] Anirudha Majumdar, Sumeet Singh, Ajay Mandlekar, and Marco Pavone. Risk-sensitive inverse reinforcement learning via coherent risk models. *Robotics: Science and Systems*, 2017.
- [69] Sumeet Singh, Jonathan Lacotte, Anirudha Majumdar, and Marco Pavone. Risk-sensitive inverse reinforcement learning via semi-and non-parametric methods. *International Journal of Robotics Research*, 2018.
- [70] Jaedeug Choi and Kee-Eung Kim. Inverse reinforcement learning in partially observable environments. *International Joint Conference on Artificial Intelli*gence (IJCAI), 2009.
- [71] Jaedeug Choi and Kee-Eung Kim. Inverse reinforcement learning in partially observable environments. *Journal of Machine Learning Research* (*JMLR*), 2011.
- [72] Hamid R Chinaei and Brahim Chaib-Draa. An inverse reinforcement learning algorithm for partially observable domains with application on healthcare dialogue management. *International Conference on Machine Learning and Applications*, 2012.

- [73] Ioana Bica, Daniel Jarrett, Alihan Hüyük, and Mihaela van der Schaar. Learning what-if explanations for sequential decision-making. *International Conference on Learning Representations (ICLR)*, 2021.
- [74] Takaki Makino and Johane Takeuchi. Apprenticeship learning for model parameters of partially observable environments. *International Conference on Machine Learning (ICML)*, 2012.
- [75] Daniel Jarrett and Mihaela van der Schaar. Inverse active sensing: Modeling and understanding timely decision-making. *International Conference on Machine Learning*, 2020.
- [76] Kunal Pattanayak and Vikram Krishnamurthy. Inverse reinforcement learning for sequential hypothesis testing and search. *International Conference on Information Fusion (FUSION)*, 2020.
- [77] Matthew Golub, Steven Chase, and Byron Yu. Learning an internal dynamics model from control demonstration. *International Conference on Machine Learning (ICML)*, 2013.
- [78] Zhengwei Wu, Paul Schrater, and Xaq Pitkow. Inverse pomdp: Inferring what you think from what you do. *arXiv preprint*, 2018.
- [79] Saurabh Daptardar, Paul Schrater, and Xaq Pitkow. Inverse rational control with partially observable continuous nonlinear dynamics. *arXiv preprint*, 2019.
- [80] Minhae Kwon, Saurabh Daptardar, Paul Schrater, and Xaq Pitkow. Inverse rational control with partially observable continuous nonlinear dynamics. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [81] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. AAAI Conference on Artificial Intelligence (AAAI), 2008.
- [82] Abdeslam Boularias, Jens Kober, and Jan Peters. Relative entropy inverse reinforcement learning. *International conference on artificial intelligence and statistics (AISTATS)*, 2011.
- [83] Mrinal Kalakrishnan, Peter Pastor, Ludovic Righetti, and Stefan Schaal. Learning objective functions for manipulation. *International Conference on Robotics* and Automation (ICRA), 2013.
- [84] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse reinforcement learning. *arXiv preprint*, 2015.

- [85] Chelsea Finn, Paul Christiano, Pieter Abbeel, and Sergey Levine. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *NeurIPS Workshop on Adversarial Training*, 2016.
- [86] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2018.
- [87] Ahmed H Qureshi, Byron Boots, and Michael C Yip. Adversarial imitation via variational inverse reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2019.
- [88] Paul Barde, Julien Roy, Wonseok Jeon, Joelle Pineau, Christopher Pal, and Derek Nowrouzezahrai. Adversarial soft advantage fitting: Imitation learning without policy optimization. Advances in neural information processing systems (NeurIPS), 2020.
- [89] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of maximum causal entropy. *International conference on Machine learning (ICML)*, 2010.
- [90] Zhengyuan Zhou, Michael Bloem, and Nicholas Bambos. Infinite time horizon maximum causal entropy inverse reinforcement learning. *IEEE Transactions on Automatic Control (TACON)*, 2017.
- [91] Kyungjae Lee, Sungjoon Choi, and Songhwai Oh. Maximum causal tsallis entropy imitation learning. Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [92] Tien Mai, Kennard Chan, and Patrick Jaillet. Generalized maximum causal entropy for inverse reinforcement learning. *AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [93] Michael Herman, Tobias Gindele, Jörg Wagner, Felix Schmitt, and Wolfram Burgard. Inverse reinforcement learning with simultaneous estimation of rewards and dynamics. *International conference on artificial intelligence and statistics (AISTATS)*, 2016.
- [94] Michael Herman. Simultaneous estimation of rewards and dynamics in irl. *Dissertation, Albert-Ludwigs-Universitat Freiburg*, 2016.
- [95] Layla El Asri, Bilal Piot, Matthieu Geist, Romain Laroche, and Olivier Pietquin. Score-based inverse reinforcement learning. *International conference* on Autonomous agents and multi-agent systems (AA-MAS), 2016.

- [96] Benjamin Burchfiel, Carlo Tomasi, and Ronald Parr. Distance minimization for reward learning from scored trajectories. AAAI Conference on Artificial Intelligence (AAAI), 2016.
- [97] Alexis Jacq, Matthieu Geist, Ana Paiva, and Olivier Pietquin. Learning from a learner. *International Conference on Machine Learning (ICML)*, 2019.
- [98] Daniel S Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations. *International Conference on Machine Learning (ICML)*, 2019.
- [99] Daniel S Brown, Wonjoon Goo, and Scott Niekum. Better-than-demonstrator imitation learning via automatically-ranked demonstrations. *Conference on Robot Learning (CoRL)*, 2020.
- [100] Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J Andrew Bagnell, Pieter Abbeel, and Jan Peters. An algorithmic perspective on imitation learning. *Foundations and Trends in Robotics*, 2018.
- [101] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. *International Conference on Machine Learning (ICML)*, 2017.
- [102] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *International Conference on Machine Learning (ICML)*, 2018.
- [103] Benjamin Eysenbach and Sergey Levine. If maxent rl is the answer, what is the question? *arXiv preprint*, 2019.
- [104] Wenjie Shi, Shiji Song, and Cheng Wu. Soft policy gradient method for maximum entropy deep reinforcement learning. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [105] Rohin Shah, Noah Gundotra, Pieter Abbeel, and Anca Dragan. Inferring reward functions from demonstrators with unknown biases. *OpenReview*, 2018.
- [106] Rohin Shah, Noah Gundotra, Pieter Abbeel, and Anca D Dragan. On the feasibility of learning, rather than assuming, human biases for reward inference. *International Conference on Machine Learning* (*ICML*), 2019.
- [107] Jonathan Rubin, Ohad Shamir, and Naftali Tishby. Trading value and information in mdps. *Decision Making with Imperfect Decision Makers (Springer)*, 2012.

- [108] Alexandre Galashov, Siddhant M Jayakumar, Leonard Hasenclever, Dhruva Tirumala, Jonathan Schwarz, Guillaume Desjardins, Wojciech M Czarnecki, Yee Whye Teh, Razvan Pascanu, and Nicolas Heess. Information asymmetry in kl-regularized rl. International Conference on Learning Representations (ICLR), 2019.
- [109] Mark K Ho, David Abel, Jonathan D Cohen, Michael L Littman, and Thomas L Griffiths. The efficiency of human cognition reflects planned information processing. AAAI Conference on Artificial Intelligence (AAAI), 2020.
- [110] Stas Tiomkin and Naftali Tishby. A unified bellman equation for causal information and value in markov decision processes. *arXiv preprint arXiv:1703.01585*, 2017.
- [111] Felix Leibfried, Jordi Grau-Moya, and Haitham Bou-Ammar. An information-theoretic optimality principle for deep reinforcement learning. *NeurIPS Workshop on Deep Reinforcement Learning*, 2017.
- [112] Matthew Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable mdps. *AAAI Conference on Artificial Intelligence (AAAI)*, 2015.
- [113] Pengfei Zhu, Xin Li, Pascal Poupart, and Guanghui Miao. On improving deep reinforcement learning for pomdps. arXiv preprint, 2017.
- [114] Maximilian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, and Shimon Whiteson. Deep variational reinforcement learning for pomdps. *International Conference on Machine Learning (ICML)*, 2018.
- [115] Amy Zhang, Zachary C Lipton, Luis Pineda, Kamyar Azizzadenesheli, Anima Anandkumar, Laurent Itti, Joelle Pineau, and Tommaso Furlanello. Learning causal state representations of partially observable environments. arXiv preprint, 2019.
- [116] Dongqi Han, Kenji Doya, and Jun Tani. Variational recurrent models for solving partially observable control tasks. arXiv preprint arXiv:1912.10703, 2019.
- [117] Joseph Futoma, Michael C Hughes, and Finale Doshi-Velez. Popcorn: Partially observed prediction constrained reinforcement learning. *International conference on artificial intelligence and statistics (AIS-TATS)*, 2020.
- [118] Richard D Smallwood and Edward J Sondik. The optimal control of partially observable markov processes over a finite horizon. *Operations research*, 1973.

- [119] Milos Hauskrecht. Value-function approximations for partially observable markov decision processes. *Journal of Artificial Intelligence Research (JAIR)*, 2000.
- [120] Joelle Pineau, Geoff Gordon, Sebastian Thrun, et al. Point-based value iteration: An anytime algorithm for pomdps. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2003.
- [121] Hanna Kurniawati, David Hsu, and Wee Sun Lee. Sarsop: Efficient point-based pomdp planning by approximating optimally reachable belief spaces. *Robotics: Science and systems*, 2008.
- [122] Mauricio Araya, Olivier Buffet, Vincent Thomas, and Françcois Charpillet. A pomdp extension with beliefdependent rewards. Advances in Neural Information Processing Systems (NeurIPS), 2010.
- [123] Mathieu Fehr, Olivier Buffet, Vincent Thomas, and Jilles Dibangoye. rho-pomdps have lipschitzcontinuous epsilon-optimal value functions. Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [124] F A Sonnenberg and J R Beck. Markov models in medical decision making: a practical guide. *Health Econ.*, 1983.
- [125] C H Jackson, L D Sharples, S G Thompson, S W Duffy, and E Couto. Multistate Markov models for disease progression with classification error. *Statistician*, 2003.
- [126] S E O'Bryant, S C Waring, C M Cullum, J Hall, L Lacritz, P J Massman, P J Lupo, J S Reisch, and R Doody. Staging dementia using Clinical Dementia Rating Scale Sum of Boxes scores: a Texas Alzheimer's research consortium study. *Arch. of Neurology*, 2008.
- [127] D Jarrett, J Yoon, and M van der Schaar. Matchnet: Dynamic prediction in survival analysis using convolutional neural networks. *NeurIPS Workshop* on Machine Learning for Health, 2018.
- [128] Daniel Jarrett, Jinsung Yoon, and Mihaela van der Schaar. Dynamic prediction in clinical survival analysis using temporal convolutional networks. *IEEE Journal of Biomedical and Health Informatics*, 2019.
- [129] P Petousis, A Winter, W Speier, D R Aberle, W Hsu, and A A T Bui. Using sequential decision making to improve lung cancer screening performance. *IEEE Access*, 2019.

- [130] F Cardoso, S Kyriakides, S Ohno, F Penault-Llorca, P Poortmans, I T Rubio, S Zackrisson, and E Senkus. Early breast cancer: ESMO Clinical Practice Guidelines for diagnosis, treatment and follow-up. *Anna. Oncology*, 2019.
- [131] A M Alaa and M van der Schaar. Attentive statespace modeling of disease progression. Advances in neural information processing systems (NeurIPS), 2019.
- [132] X Wang, D Sontag, and F Wang. Unsupervised learning of disease progression models. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2014.
- [133] Clemens Heuberger. Inverse combinatorial optimization. *Journal of combinatorial optimization*, 2004.
- [134] Kareem Amin and Satinder Singh. Towards resolving unidentifiability in inverse reinforcement learning. arXiv preprint, 2016.
- [135] Kareem Amin, Nan Jiang, and Satinder Singh. Repeated inverse reinforcement learning. Advances in neural information processing systems (NeurIPS), 2017.
- [136] Stuart Armstrong and Sören Mindermann. Occam's razor is insufficient to infer the preferences of irrational agents. Advances in neural information processing systems (NeurIPS), 2018.
- [137] Paul Christiano. The easy goal inference problem is still hard. *AI Alignment*, 2015.
- [138] Eric J Michaud, Adam Gleave, and Stuart Russell. Understanding learned reward functions. *NeurIPS Workshop on Deep Reinforcement Learning*, 2020.
- [139] Adam Gleave, Michael Dennis, Shane Legg, Stuart Russell, and Jan Leike. Quantifying differences in reward functions. *International Conference on Learning Representations (ICLR)*, 2021.
- [140] Daniel S Brown and Scott Niekum. Deep bayesian reward learning from preferences. *NeurIPS Workshop* on Safety and Robustness in Decision-Making, 2019.
- [141] Daniel S Brown, Russell Coleman, Ravi Srinivasan, and Scott Niekum. Safe imitation learning via fast bayesian reward inference from preferences. *International Conference on Machine Learning (ICML)*, 2020.
- [142] Nicolas Heess, David Silver, and Yee Whye Teh. Actor-critic reinforcement learning with energybased policies. *European Workshop on Reinforcement Learning (EWRL)*, 2013.

- [143] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. *International Conference on Machine Learning (ICML)*, 2016.
- [144] Nicolò Cesa-Bianchi, Claudio Gentile, Gábor Lugosi, and Gergely Neu. Boltzmann exploration done right. *Advances in neural information processing systems* (*NeurIPS*), 2017.
- [145] Amir Globerson, Eran Stark, Eilon Vaadia, and Naftali Tishby. The minimum information principle and its application to neural code analysis. *Proceedings* of the National Academy of Sciences, 2009.
- [146] Naftali Tishby and Daniel Polani. Information theory of decisions and actions. *Perception-action cycle* (*Springer*), 2011.
- [147] Pedro A Ortega and Daniel A Braun. Thermodynamics as a theory of decision-making with informationprocessing costs. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 2013.
- [148] Ian R Petersen, Matthew R James, and Paul Dupuis. Minimax optimal control of stochastic uncertain systems with relative entropy constraints. *IEEE Transactions on Automatic Control*, 2000.
- [149] Charalambos D Charalambous, Farzad Rezaei, and Andreas Kyprianou. Relations between information theory, robustness, and statistical mechanics of stochastic systems. *IEEE Conference on Decision and Control (CDC)*, 2004.
- [150] Takayuki Osogami. Robustness and risk-sensitivity in markov decision processes. *Advances in Neural Information Processing Systems (NeurIPS)*, 2012.
- [151] Jordi Grau-Moya, Felix Leibfried, Tim Genewein, and Daniel A Braun. Planning with informationprocessing constraints and model uncertainty in markov decision processes. *Joint European conference on machine learning and knowledge discovery in databases (ECML)*, 2016.
- [152] Brian D Ziebart. Modeling purposeful adaptive behavior with the principle of maximum causal entropy. *Dissertation, Carnegie Mellon University*, 2010.
- [153] Gerhard Kramer. Directed information for channels with feedback. *Dissertation, ETH Zurich*, 1998.
- [154] James Massey. Causality, feedback and directed information. International Symposium on Information Theory and Its Applications, 1990.

- [155] Hans Marko. The bidirectional communication theory-a generalization of information theory. *IEEE Transactions on Communications*, 1973.
- [156] John B McKinlay, Carol L Link, et al. Sources of variation in physician adherence with clinical guidelines. *Journal of general internal medicine*, 2007.
- [157] Matthias Bock, Gerhard Fritsch, and David L Hepner. Preoperative laboratory testing. *Anesthesiology clinics*, 2016.
- [158] Jack W O'Sullivan, Carl Heneghan, Rafael Perera, Jason Oke, Jeffrey K Aronson, Brian Shine, and Ben Goldacre. Variation in diagnostic test requests and outcomes: a preliminary metric for openpathology. net. *Nature Scientific Reports*, 2018.
- [159] Yunjie Song, Jonathan Skinner, Julie Bynum, Jason Sutherland, John E Wennberg, and Elliott S Fisher. Regional variations in diagnostic practices. *New England Journal of Medicine*, (1), 2010.
- [160] Shannon K Martin and Adam S Cifu. Routine preoperative laboratory tests for elective surgery. *Journal* of the American Medical Association (JAMA), 2017.
- [161] M. Allen. Unnecessary tests and treatment explain why health care costs so much. *Scientific American*, 2017.
- [162] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 1998.
- [163] Razvan V Marinescu, Neil P Oxtoby, Alexandra L Young, Esther E Bron, Arthur W Toga, Michael W Weiner, Frederik Barkhof, Nick C Fox, Stefan Klein, Daniel C Alexander, et al. Tadpole challenge: Prediction of longitudinal evolution in alzheimer's disease. arXiv preprint, 2018.
- [164] Edi Karni and Zvi Safra. Behavioral consistency in sequential decisions. *Progress in Decision, Utility* and Risk Theory, 1991.
- [165] Kent Daniel, David Hirshleifer, and Avanidhar Subrahmanyam. Investor psychology and security market under-and overreactions. *The Journal of Finance*, 1998.
- [166] Amos Tversky and Daniel Kahneman. Evidential impact of base rates. *Stanford University Department Of Psychology*, 1981.
- [167] Charlotte L Allan and Klaus P Ebmeier. The influence of apoe4 on clinical progression of dementia: a meta-analysis. *International journal of geriatric psychiatry*, 2011.

- [168] Sylvaine Artero, Marie-Laure Ancelin, Florence Portet, A Dupuy, Claudine Berr, Jean-François Dartigues, Christophe Tzourio, Olivier Rouaud, Michel Poncet, Florence Pasquier, et al. Risk profiles for mild cognitive impairment and progression to dementia are gender specific. *Journal of Neurology, Neurosurgery & Psychiatry*, 2008.
- [169] Xue Hua, Derrek P Hibar, Suh Lee, Arthur W Toga, Clifford R Jack Jr, Michael W Weiner, Paul M Thompson, Alzheimer's Disease Neuroimaging Initiative, et al. Sex and age differences in atrophic rates: an adni study with n= 1368 mri scans. *Neurobiology* of aging, 2010.
- [170] Brendan O'Donoghue, Remi Munos, Koray Kavukcuoglu, and Volodymyr Mnih. Combining policy gradient and q-learning. *International Conference on Learning Representations (ICLR)*, 2017.
- [171] Momchil Tomov. Structure learning and uncertaintyguided exploration in the human brain. *Dissertation, Harvard University*, 2020.
- [172] Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Benjamin Eysenbach. F-irl: Inverse reinforcement learning via state marginal matching. *Conference on Robot Learning (CoRL)*, 2020.
- [173] Hong Jun Jeon, Smitha Milli, and Anca D Dragan. Reward-rational (implicit) choice: A unifying formalism for reward learning. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [174] Jeffrey Ely, Alexander Frankel, and Emir Kamenica. Suspense and surprise. *Journal of Political Economy*, 2015.
- [175] Ahmed M Alaa and Mihaela van der Schaar. Balancing suspense and surprise: Timely decision making with endogenous information acquisition. Advances in neural information processing systems (NeurIPS), 2016.
- [176] Owain Evans and Noah D Goodman. Learning the preferences of bounded agents. *NeurIPS Workshop on Bounded Optimality*, 2015.
- [177] Tan Zhi-Xuan, Jordyn L Mann, Tom Silver, Joshua B Tenenbaum, and Vikash K Mansinghka. Online bayesian goal inference for boundedly-rational planning agents. Advances in neural information processing systems (NeurIPS), 2020.

- [178] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy extraction. Advances in neural information processing systems (NeurIPS), 2018.
- [179] Herman Yau, Chris Russell, and Simon Hadfield. What did you think would happen? explaining agent behaviour through intended outcomes. Advances in neural information processing systems (NeurIPS), 2020.
- [180] Tom Bewley, Jonathan Lawry, and Arthur Richards. Modelling agent policies with interpretable imitation learning. *TAILOR Workshop at ECAI*, 2020.
- [181] Tathagata Chakraborti, Sarath Sreedharan, Yu Zhang, and Subbarao Kambhampati. Plan explanations as model reconciliation: Moving beyond explanation as soliloquy. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.
- [182] Tathagata Chakraborti, Sarath Sreedharan, Sachin Grover, and Subbarao Kambhampati. Plan explanations as model reconciliation: an empirical study. *International Conference on Human-Robot Interaction (HRI)*, 2019.
- [183] Bradley Hayes and Julie A Shah. Improving robot controller transparency through autonomous policy explanation. *International Conference on Human-Robot Interaction (HRI)*, 2017.
- [184] Sarath Sreedharan, Utkash Soni, Mudit Verma, Siddharth Srivastava, and Subbarao Kambhampati. Bridging the gap: Providing post-hoc symbolic explanations for sequential decision-making problems with black box simulators. *ICML Workshop on Human-inthe-Loop Learning*, 2020.
- [185] Roy Fox and Naftali Tishby. Minimum-information lqg control part i: Memoryless controllers. *IEEE Conference on Decision and Control (CDC)*, 2016.
- [186] Roy Fox and Naftali Tishby. Minimum-information lqg control part ii: Retentive controllers. *IEEE Conference on Decision and Control (CDC)*, 2016.
- [187] Robert Babuska. Model-based imitation learning. Springer Encyclopedia of the Sciences of Learning, 2012.
- [188] Yoshua Bengio and Paolo Frasconi. An input output hmm architecture. *Advances in neural information* processing systems (NeurIPS), 1995.