
Inverse Decision Modeling: Learning Interpretable Representations of Behavior

Daniel Jarrett^{1*} Alihan Hüyük^{1*} Mihaela van der Schaar^{1,2}

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-the-loop” 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>.

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.

Inverse Decision Model	Extrinsic [†]		Intrinsic [‡]							Examples
	Partially Controllable	Partially Observable	Purposeful Behavior	Subjective Dynamics	Action Stochasticity	Knowledge Uncertainty	Decision Complexity	Specification Complexity	Recognition Complexity	
	τ_{env}	ω_{env}	v	τ, ω	π	ρ, σ	α	β	η	
BC-IL	✓	✓	×	×	✓	×	×	×	×	[14–21]
Subj. BC-IL	✓	✓	×	✓	✓	×	×	×	×	[22]
Det. DM-IL	✓	×	×	×	×	×	×	×	×	[23, 24]
Stoc. DM-IL	✓	×	×	×	✓	×	×	×	×	[25–39]
Det. IRL	✓	×	✓	×	×	×	×	×	×	[40–46]
Stoc. IRL	✓	×	✓	×	✓	×	×	×	×	[47–66]
Subj. IRL	✓	×	✓	✓	✓	×	×	×	×	[67]
RS-IRL	✓	×	✓	✓	×	✓	×	×	×	[68, 69]
Det. PO-IRL	✓	✓	✓	×	×	×	×	×	×	[70–73]
Stoc. PO-IRL	✓	✓	✓	×	✓	×	×	×	×	[74–76]
Subj. PO-IRL	✓	✓	✓	✓	✓	×	×	×	×	[77–80]
ME-IRL	✓	×	✓	×	✓	×	✓	×	×	[81–92]
Subj. ME-IRL	✓	×	✓	✓	✓	×	✓	×	×	[93, 94]
Inverse Bounded Rational Control	✓	✓	✓	✓	✓	✓	✓	✓	✓	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 re-

ward 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 commonly-assumed 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.

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).

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

3.1. Forward Problem

Consider the standard setup for sequential decision-making, where an agent interacts with a (potentially partially-observable) environment. First, let $\psi \doteq (\mathcal{S}, \mathcal{X}, \mathcal{Z}, \mathcal{U}, \mathcal{T}, \mathcal{O})$ give the *problem setting*, where \mathcal{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 $\mathcal{O} \doteq \Delta(\mathcal{X})^{\mathcal{U} \times \mathcal{S}}$ 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 $v \in \mathbb{R}^{\mathcal{S} \times \mathcal{U}}$, 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 $\tau_{\text{env}}, \omega_{\text{env}} \in \mathcal{T} \times \mathcal{O}$. 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 (observation-action) trajectories with $\mathcal{H} \doteq \cup_{t=0}^{\infty} (\mathcal{X} \times \mathcal{U})^t \times \mathcal{X}$. Then a *behavior* ϕ manifests as a distribution over trajectories (induced by an agent’s policies interacting with the environment):

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

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 :

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

(*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:

$$F : \Psi \times \Theta \rightarrow \Phi \quad (3)$$

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:

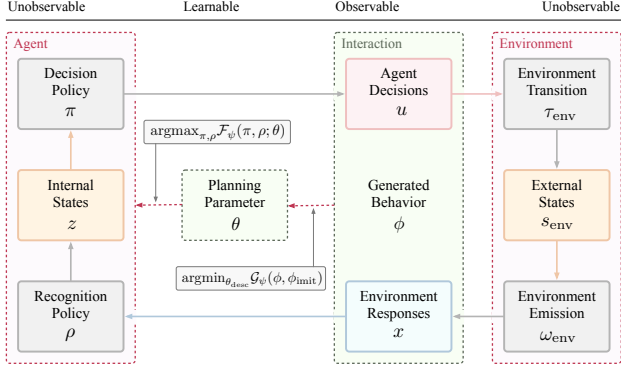
$$F(\psi, \theta) \doteq \phi_{\pi^*, \rho^*} : \pi^*, \rho^* \doteq \operatorname{argmax}_{\pi, \rho} \mathcal{F}_{\psi}(\pi, \rho; \theta) \quad (4)$$

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.

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).



A more practical question is whether this optimum is reachable. 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 $\Theta \doteq \Theta_{\text{norm}} \times \Theta_{\text{desc}}$ decompose the parameter space into a *normative* component (i.e. whose values $\theta_{\text{norm}} \in \Theta_{\text{norm}}$ we wish to clamp), and a *descriptive* component (i.e. whose values $\theta_{\text{desc}} \in \Theta_{\text{desc}}$ we wish to infer). Then an *inverse planner* is given as follows:

$$G : \Phi \times \Theta_{\text{norm}} \rightarrow \Theta_{\text{desc}} \quad (5)$$

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

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

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: θ_{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 $\theta_{\text{desc}} = v$ 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 ($V_\theta^\phi, Q_\theta^\phi$); 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 $\theta_{\text{desc}} = v$ and assume full observability (whence $\mathcal{S} = \mathcal{X} = \mathcal{Z}$).

Inverse Planner (G)	Demonstrator (ϕ_{demo})	Helper	Optimization (θ_{desc}^*)	Examples
Minimum Perturbation	Deterministic, Optimal	Default $\tilde{\theta}_{\text{desc}}$	$\text{argmin}_{\theta_{\text{desc}}} \ \theta_{\text{desc}} - \tilde{\theta}_{\text{desc}}\ _p : \phi_{\text{demo}} = F(\psi, \theta)$	[133]
Maximum Margin	Deterministic, Optimal	-	$\text{argmin}_{\theta_{\text{desc}}} \mathbb{E}_{z \sim \rho_0} [V_\theta^{\phi_{\text{init}}}(z) - V_\theta^{\phi_{\text{demo}}}(z)]$	[40–45, 53, 70–73]
Regularized Max. Margin	Stochastic, Optimal	-	$\text{argmin}_{\theta_{\text{desc}}} \mathbb{E}_{z \sim \rho_0} [V_{\text{soft}, \theta}^{\phi_{\text{init}}}(z) - V_\theta^{\phi_{\text{demo}}}(z)] + \zeta(\theta)$	[25]
Multiple Experimentation	Deterministic, Optimal	Environments \mathcal{V}	$\text{argmin}_{\theta_{\text{desc}}} \int \max_{\mathcal{V}, u} (Q_{\mathcal{V}, \theta}^{\phi_{\text{demo}}}(z, u) - V_{\mathcal{V}, \theta}^{\phi_{\text{demo}}}(z)) dx$	[134, 135]
Distance Minimization	Individually-Scored	Scores $\tilde{v}(h) \in \mathbb{R}$	$\text{argmin}_{\theta_{\text{desc}}} \mathbb{E}_{h \sim \phi_{\text{demo}}} \ \tilde{v}(h) - \sum_{s, u \in h} v(s, u)\ _p$	[95, 96]
Soft Policy Inversion	Stoc., Batch-Ordered	$\{\phi_{\text{demo}}^{(1)}, \dots, \phi_{\text{demo}}^{(K)}\}$	$\text{argmin}_{\theta_{\text{desc}}} \sum_k \mathbb{E}_{s, u, s' \sim \phi_{\text{demo}}^{(k)}} \ \Delta v^{(k)}(s, u, s')\ _p$	[97]
Preference Extrapolation	Stoc., Pairwise-Ranked	$\{(i, j) h_i \prec h_j\}$	$\text{argmin}_{\theta_{\text{desc}}} \mathbb{E}_{(h_i \prec h_j) \sim \phi_{\text{demo}}} \log \mathbb{P}_v(h_i \prec h_j)$	[98, 99]
Soft Policy Matching	Stochastic, Optimal	-	$\text{argmin}_{\theta_{\text{desc}}} D_{\text{KL}}(\mathbb{P}_{\phi_{\text{demo}}}(u_{0:T} x_{0:T}) \mathbb{P}_{\phi_{\text{init}}}(u_{0:T} x_{0:T}))$	[47–52, 76, 89–94]
Distribution Matching	Stochastic, Optimal	-	$\text{argmin}_{\theta_{\text{desc}}} D_f(\phi_{\text{demo}} \phi_{\text{init}})$	[23–39, 54, 81–88]
General Formulation	(any)	(any)	$\text{argmin}_{\theta_{\text{desc}}} \mathcal{G}_\psi(\phi_{\text{demo}}, \phi_{\text{init}})$	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: $\Phi_{\theta_{\text{norm}}} \doteq F_{\theta_{\text{norm}}}[\Theta_{\text{desc}}] \subseteq \Phi$. Then the projection map onto this subspace is given by:

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

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:

$$\phi_{\text{init}}^* \doteq \text{proj}_{\Phi_{\theta_{\text{norm}}}}(\phi_{\text{demo}}) \quad (8)$$

In other words, the model ϕ_{init}^* serves as a complete (generative) account of ϕ_{demo} as its *behavior projection* onto $\Phi_{\theta_{\text{norm}}}$.

Interpretability What dictates our choices? For pure imitation (i.e. replicating expert actions), a black-box decision-rule fitted by soft policy matching may do well. For apprenticeship (i.e. matching expert returns), a perfectly optimal planner inverted 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) $F_{\theta_{\text{norm}}}, G_{\theta_{\text{norm}}}$ such that ϕ_{init}^* 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 :

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

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:

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

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

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

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 $\mu \in \Delta(\mathcal{Z})$:

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

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*:

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

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 $\mathbb{B}_{\pi, \rho} : \mathbb{R}^{\mathcal{Z}} \rightarrow \mathbb{R}^{\mathcal{Z}}$ such that for any given $V \in \mathbb{R}^{\mathcal{Z}}$:

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

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

$$V^{\phi_{\pi, \rho}}(z) \doteq \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\substack{s_t \sim p(\cdot|z_t) \\ u_t \sim \pi(\cdot|z_t) \\ \tau, \omega \sim \sigma(\cdot|z_t, u_t) \\ s_{t+1} \sim \tau(\cdot|s_t, u_t) \\ x_{t+1} \sim \omega(\cdot|u_t, s_{t+1}) \\ z_{t+1} \sim \rho_{\tau, \omega}(\cdot|z_t, u_t, x_{t+1})}} [v(s_t, u_t) | z_0 = z] \quad (15)$$

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, $Q^{\phi_{\pi, \rho}}(z, u) \doteq \mathbb{E}_{s \sim p(\cdot|z)} [v(s, u) + \mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u), \dots, z' \sim \rho_{\tau, \omega}(\cdot|z, u, x')} \gamma V^{\phi_{\pi, \rho}}(z')]$ 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 $\text{argmax}_u Q^{\phi_{\pi, \rho}}(z, u)$. 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):

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

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:

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

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 $\tilde{\rho}$):

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

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 $\mathbb{I}_{\pi, \rho}[\rho; \tilde{\rho}] \leq C$. So the Lagrangian (now with the additional multipliers $\alpha, \beta, \eta \in \mathbb{R}$) is given by $\mathcal{L}_{\pi, \rho}(\mu, \alpha, \beta, \eta, V) \doteq J_{\pi, \rho} - \langle V, \mu - \gamma \mathbb{M}_{\pi, \rho} \mu - (1 - \gamma) \rho_0 \rangle - \alpha \cdot (\mathbb{I}_{\pi, \rho}[\pi; \tilde{\pi}] - A) - \beta \cdot (\mathbb{I}_{\pi, \rho}[\sigma; \tilde{\sigma}] - B) - \eta \cdot (\mathbb{I}_{\pi, \rho}[\rho; \tilde{\rho}] - C)$.

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

$$(\mathbb{B}_{\pi, \rho} V)(z) \doteq \mathbb{E}_{\substack{s \sim p(\cdot|z) \\ u \sim \pi(\cdot|z)}} \left[-\alpha \log \frac{\pi(u|z)}{\tilde{\pi}(u)} + v(s, u) + \mathbb{E}_{\tau, \omega \sim \sigma(\cdot|z, u)} \left[-\beta \log \frac{\sigma(\tau, \omega|z, u)}{\tilde{\sigma}(\tau, \omega)} + \mathbb{E}_{\substack{s' \sim \tau(\cdot|s, u) \\ x' \sim \omega(\cdot|u, s') \\ z' \sim \rho_{\tau, \omega}(\cdot|z, u, x')}} \left[-\eta \log \frac{\rho_{\tau, \omega}(z'|z, u)}{\tilde{\rho}(z')} + \gamma V(z') \right] \right] \right] \quad (19)$$

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) $Q^{\phi_{\pi, \rho}} \in \mathbb{R}^{\mathcal{Z} \times \mathcal{U}}$ as the 1/3-step-ahead expectation, and the (state-action-model) $K^{\phi_{\pi, \rho}} \in \mathbb{R}^{\mathcal{Z} \times \mathcal{U} \times \mathcal{T} \times \mathcal{O}}$ 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 agents 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.

Boundedly Rational Agent	Flexibility	Optimism	Adaptivity	(Action Prior)	(Model Prior)	(Belief Prior)	Observability	Examples
	α^{-1}	β^{-1}	η^{-1}	$\tilde{\pi}$	$\tilde{\sigma}$	$\tilde{\varrho}$		
Uniformly Random Agent	$\rightarrow 0$	$\rightarrow 0$	$\rightarrow \pm\infty$	Uniform	Dirac δ	-	Full/Partial	-
Deterministic, Optimal Agent	$\rightarrow \infty$	$\rightarrow 0$	$\rightarrow \pm\infty$	-	Dirac δ	-	Full/Partial	(any)
Boltzmann-Exploratory Agent [†]	$\rightarrow \infty$	$\rightarrow 0$	$\rightarrow \pm\infty$	-	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 $\mathbb{B}^* : \mathbb{R}^{\mathcal{Z}} \rightarrow \mathbb{R}^{\mathcal{Z}}$ such that for any $V \in \mathbb{R}^{\mathcal{Z}}$:

$$\begin{aligned}
 (\mathbb{B}^*V)(z) &\doteq \alpha \log \mathbb{E}_{u \sim \tilde{\pi}} \exp\left(\frac{1}{\alpha} Q(z, u)\right) & (20) \\
 Q(z, u) &\doteq \beta \log \mathbb{E}_{\tau, \omega \sim \tilde{\sigma}} \exp\left(\frac{1}{\beta} K(z, u, \tau, \omega)\right) \\
 K(z, u, \tau, \omega) &\doteq \mathbb{E}_{s \sim p(\cdot|z)} v(s, u) \\
 &\quad + \mathbb{E}_{s \sim p(\cdot|z)} \left[-\eta \log \frac{\varrho_{\tau, \omega}(z'|z, u)}{\tilde{\varrho}(z')} + \gamma V(z') \right] \\
 &\quad \mathbb{E}_{\substack{s' \sim \tau(\cdot|s, u) \\ x' \sim \omega(\cdot|u, s') \\ z' \sim \rho_{\tau, \omega}(\cdot|z, u, x')}}
 \end{aligned}$$

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 *boundedly rational decision policy* (i.e. primal optimal) is given by:

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

and the *boundedly rational recognition policy* is given by:

$$\begin{aligned}
 \rho^*(z'|z, u, x') &= \mathbb{E}_{\tau, \omega \sim \sigma^*(\cdot|z, u)} \rho_{\tau, \omega}(z'|z, u, x'), \text{ where} \\
 \sigma^*(\tau, \omega|z, u) &\doteq \frac{\tilde{\sigma}(\tau, \omega)}{Z_{K^*}(z, u)} \exp\left(\frac{1}{\beta} K^*(z, u, \tau, \omega)\right) \quad (22)
 \end{aligned}$$

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

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.

infinitely flexible, utility-seeking ($\alpha^{-1} \rightarrow \infty$) or disutility-seeking ($\alpha^{-1} \rightarrow -\infty$) 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 $\Theta_{\text{norm}}^{\text{BRC}} \times \Theta_{\text{desc}}^{\text{BRC}}$. The *forward problem* is encapsulated by Theorems 4–5 (which also yield a straight-forward algorithm, i.e. iterate 4 until convergence, then plug into 5). Therefore the *forward planner* is given as follows:

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

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:

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

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:

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

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 $\phi_{\text{imit}}^* = F_{\theta_{\text{norm}}^{\text{BRC}}} \circ G_{\theta_{\text{norm}}^{\text{BRC}}}(\phi_{\text{demo}})$ onto the space $\Phi_{\theta_{\text{norm}}^{\text{BRC}}}$ of behaviors thereby *interpretablely* parameterized—i.e. by the structure we designed for Θ^{BRC} , and by the normative standards $\theta_{\text{norm}}^{\text{BRC}}$ we may choose to specify.

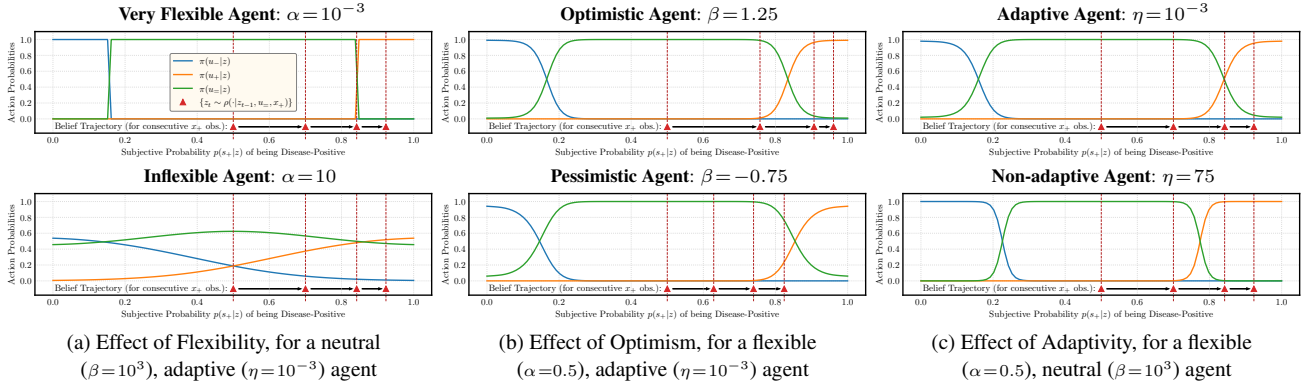


Figure 2. *Bounded Rational Control*. Decision agents in DIAG: In each panel, the boundedly rational decision policy π is shown in terms of action probabilities (y -axis) for different subjective beliefs (x -axis). To visualize the boundedly rational recognition policy ρ , each panel shows an example 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:

- *Interpretability*: IBRC gives a *transparent* parameterization of behavior that can be successfully learned from data.
- *Expressivity*: IBRC more finely *differentiates* between imperfect behaviors, while standard reward learning cannot.
- *Applicability*: 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 point-valued, subjective $\tau, \omega \neq \tau_{\text{env}}, \omega_{\text{env}}$, and of the classic Tiger Problem in POMDP literature where $\tau, \omega = \tau_{\text{env}}, \omega_{\text{env}}$ [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 ($x_{\text{avg}}^{\text{MRI}}$), high ($x_{\text{high}}^{\text{MRI}}$), low ($x_{\text{low}}^{\text{MRI}}$), not ordered ($x_{\text{none}}^{\text{MRI}}$); separately, the cognitive dementia rating test result—which is always measured—is classified as normal ($x_{\text{norm}}^{\text{CDR}}$), questionable impairment ($x_{\text{ques}}^{\text{CDR}}$), and suspected dementia ($x_{\text{susp}}^{\text{CDR}}$). 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 $\bar{\sigma}$ 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, $\bar{\sigma}$ is centered at the IOHMM learned from the data. For simplicity, for $\bar{\pi}, \bar{q}$ 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.

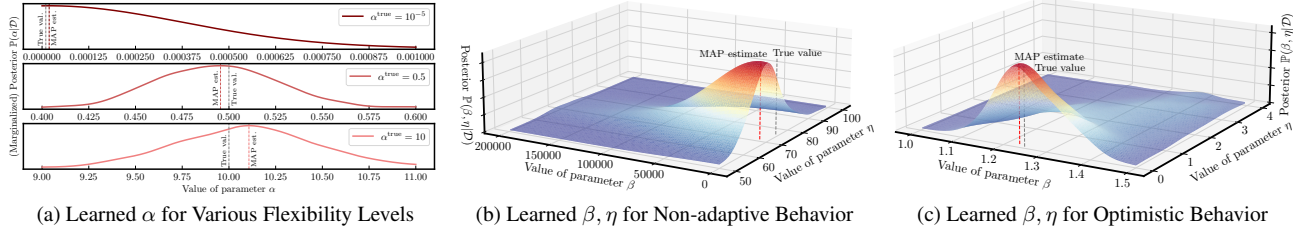


Figure 3. *Inverse Bounded Rational Control*. (a) Posteriors of α learned from extremely flexible ($\alpha^{\text{true}} = 10^{-5}$), flexible ($\alpha^{\text{true}} = 0.5$), and inflexible ($\alpha^{\text{true}} = 10$) behaviors (with β, η fixed as neutral and adaptive; similar plots can be obtained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

5.1. Interpretability Figure 2 verifies (for DIAG) that different BRC behaviors accord with our intuitions. First, *ceteris 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.

References

- [1] 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 energy-based 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. Sample-efficient 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 Ansel, and Shie Mannor. Model-based adversarial imitation learning. *arXiv preprint*, 2016.
- [39] Nir Baram, Oron Ansel, and Shie Mannor. Model-based 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 Pirota and Marcello Restelli. Inverse reinforcement learning through policy gradient minimization. *AAAI Conference on Artificial Intelligence (AAAI)*, 2016.
- [56] Alberto Maria Metelli, Matteo Pirota, and Marcello Restelli. Compatible reward inverse reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [57] Davide Tateo, Matteo Pirota, 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 Intelligence (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-Universität 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 (AAMAS)*, 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 (AISTATS)*, 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 François Charpillet. A pomdp extension with belief-dependent rewards. *Advances in Neural Information Processing Systems (NeurIPS)*, 2010.
- [123] Mathieu Fehr, Olivier Buffet, Vincent Thomas, and Jilles Dibangoye. rho-pomdps have lipschitz-continuous 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. *Statistical*, 2003.
- [126] S E O’Byrant, 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. *Annals of Oncology*, 2019.
- [131] A M Alaa and M van der Schaar. Attentive state-space 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 energy-based 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 information-processing 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 information-processing 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 Avaniidhar 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 uncertainty-guided 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-in-the-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.

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.