
Models of human preference for learning reward functions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The utility of reinforcement learning is limited by the alignment of reward functions
2 with the interests of human stakeholders. One promising method for alignment is
3 to learn the reward function from human-generated preferences between pairs of
4 trajectory segments. These human preferences are typically assumed to be informed
5 solely by partial return, the sum of rewards along each segment. We find this
6 assumption to be flawed and propose modeling preferences instead as arising from
7 a different statistic: each segment’s regret, a measure of a segment’s deviation from
8 optimal decision-making. Given infinitely many preferences generated according
9 to regret, we prove that we can identify a reward function equivalent to the reward
10 function that generated those preferences. We also prove that the previous partial
11 return model lacks this identifiability property without preference noise that reveals
12 rewards’ relative proportions, and we empirically show that our proposed regret
13 preference model outperforms it with finite training data in otherwise the same
14 setting. Additionally, our proposed regret preference model better predicts real
15 *human* preferences and also learns reward functions from these preferences that
16 lead to policies that are better human-aligned. Overall, this work establishes that
17 the choice of preference model is impactful, and our proposed regret preference
18 model provides an improvement upon a core assumption of recent research.

19 1 Introduction

20 Improvements in reinforcement learning (RL) have led to notable recent achievements [1-6],
21 increasing its applicability to real-world problems. Yet, like all optimization algorithms, even *perfect*
22 RL optimization is limited by the objective it optimizes. For RL, this objective is created in large
23 part by the reward function. Poor alignment between reward functions and the interests of human
24 stakeholders limits the utility of RL and may even pose catastrophic risks [7, 8].

25 Influential recent research has focused on reward learning from preferences over pairs of fixed-length
26 trajectory segments. Nearly all of this recent work assumes that human preferences arise probabilis-
27 tically from *only* the sum of rewards over a segment, i.e., the segment’s **partial return** [9-16]. That is,
28 these works assume that people tend to prefer trajectory segments that yield greater rewards *during the*
29 *segment*. However, this preference model ignores seemingly important information about the segment’s
30 desirability, including the state values of the segment’s start and end states. Separately, this partial return
31 preference model can prefer suboptimal actions with lucky outcomes, like buying a lottery ticket.

32 This paper proposes an alternative preference model based on the **regret** of each segment, which is equiv-
33 alent to the negated sum of an optimal policy’s advantage of each transition in the segment (Section 2.2).

34 Figure 1 shows an intuitive example of when these two models disagree. Other classes of domains that
 35 the models will differ on are those with constant reward until the end, including competitive games like
 36 chess, go, and soccer as well as tasks for which the objective is to minimize time until reaching a goal.

37 For these two preference models, we first focus the-
 38oretically on a normative analysis (Section 3)—i.e.,
 39 what preference model would we *want* humans
 40 to use if we could choose—proving that reward
 41 learning on infinite, exhaustive preferences with
 42 our proposed regret preference model identifies a
 43 reward function with the same set of optimal poli-
 44 cies as the reward function with which the prefer-
 45 ences are generated. We also prove that the par-
 46 tial return preference model is not guaranteed to
 47 identify such a reward function without preference
 48 noise. We follow up with a descriptive analysis of
 49 how well each of these proposed models align with
 50 *actual* human preferences by collecting a human-
 51 labeled dataset of preferences in a rich grid world
 52 domain (Section 4) and showing that the regret pref-
 53 erence model better predicts these human prefer-
 54 ences (Section 5). Finally, we find that the poli-
 55 cies ultimately created through the regret preference
 56 model tend to outperform those from the partial
 57 return model learning—both when assessed with
 58 collected human preferences or when assessed with
 59 synthetic preferences (Section 6).

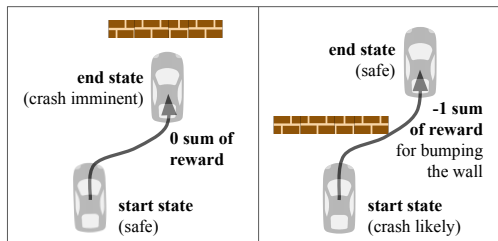


Figure 1: Two segments of a car moving at high speed near a brick wall. Assume the right segment is optimal and the left segment is suboptimal (as defined in Sec. 2.1). The left segment has a higher sum of reward, so the partial return preference model tends to prefer it. The regret preference model instead tends to prefer the right segment because optimal segments have minimal regret. If we also assume deterministic transitions, then the regret model includes the difference in values between the start state and the end state (Eq. 3), and the right segment would tend to be preferred because it greatly improves its state values from start to end, whereas the left segment’s state values greatly worsen. We suspect our human readers will also tend to prefer the right segment.

60 2 Preference models for learning reward functions

61 We assume that the task environment is a Markov decision process (MDP) specified by the tuple $(S, A,$
 62 $T, \gamma, D_0, r)$. S and A are the sets of possible states and actions, respectively. T is a transition function,
 63 $T: S \times A \rightarrow S$. γ is the discount factor and D_0 is the distribution of start states. Unless otherwise
 64 stated, we assume undiscounted tasks (i.e., $\gamma = 1$) that have terminal states, after which only 0 reward
 65 can be received. r is a reward function, $r: S \times A \times S \rightarrow \mathbb{R}$, where the reward r_t at time t is a function of
 66 s_t, a_t , and s_{t+1} . An $\text{MDP} \setminus r$ is an MDP without a reward function.

67 Throughout this paper, r refers to the ground-truth reward function for some MDP; \hat{r} refers to a learned
 68 approximation of r ; and \tilde{r} refers to any reward function (including r or \hat{r}). A policy $(\pi: S \times A \rightarrow [0,1])$
 69 specifies the probability of an action given a state. $Q_{\tilde{r}}^*$ and $V_{\tilde{r}}^*$ refer respectively to the state-action value
 70 function and state value function for an optimal policy, π^* , under \tilde{r} . The optimal advantage function is
 71 defined as $A_{\tilde{r}}^*(s, a) \triangleq Q_{\tilde{r}}^*(s, a) - V_{\tilde{r}}^*(s)$. Throughout this paper, the ground-truth reward function r
 72 is used to algorithmically generate preferences when they are not human-generated, is hidden during
 73 reward learning, and is used to evaluate the performance of optimal policies under a learned \hat{r} .

74 2.1 Reward learning from pairwise preferences

75 A reward function can be learned by minimizing the cross-entropy loss—i.e., maximizing the
 76 likelihood—of observed human preferences, a common approach in recent literature [9–11, 14, 16].

77 **Segments** Let σ denote a segment starting at state $s_{\sigma,0}$. Its length $|\sigma|$ is the number of transitions within
 78 the segment. A segment includes $|\sigma| + 1$ states and $|\sigma|$ actions: $(s_{\sigma,0}, a_{\sigma,0}, s_{\sigma,1}, a_{\sigma,1}, \dots, s_{\sigma,|\sigma|})$. In this
 79 problem setting, segments lack any reward information. As shorthand, we define $\sigma_t \triangleq (s_{\sigma,t}, a_{\sigma,t}, s_{\sigma,t+1})$.
 80 A segment σ is **optimal** with respect to \tilde{r} if, for every $i \in \{1, \dots, |\sigma| - 1\}$, $Q_{\tilde{r}}^*(s_{\sigma,i}, a_{\sigma,i}) = V_{\tilde{r}}^*(s_{\sigma,i})$. A
 81 segment that is not optimal is **suboptimal**. Given some \tilde{r} and a segment σ , $\tilde{r}_t \triangleq \tilde{r}(s_{\sigma,t}, a_{\sigma,t}, s_{\sigma,t+1})$,
 82 and the **partial return** of a segment σ is $\sum_{t=0}^{|\sigma|-1} \gamma^t \tilde{r}_t$, denoted in shorthand as $\Sigma_{\sigma} r$.

83 **Preference datasets** Each preference over a pair of segments creates a sample $(\sigma_1, \sigma_2, \mu)$ in a
 84 preference dataset D_{\succ} . Vector $\mu = \langle \mu_1, \mu_2 \rangle$ represents the preference; specifically, if σ_1 is preferred
 85 over σ_2 , denoted $\sigma_1 \succ \sigma_2$, $\mu = \langle 1, 0 \rangle$. μ is $\langle 0, 1 \rangle$ if $\sigma_1 \prec \sigma_2$ and is $\langle 0.5, 0.5 \rangle$ for $\sigma_1 \sim \sigma_2$ (no preference).

86 **Loss function** To learn a reward function from a preference dataset, D_{\succ} , a common assumption
 87 is that these preferences were generated by a preference model P that arises from an unobservable
 88 *ground-truth* reward function r . We approximate r by minimizing cross-entropy loss to learn \hat{r} :

$$\text{loss}(\hat{r}, D_{\succ}) = - \sum_{(\sigma_1, \sigma_2, \mu) \in D_{\succ}} \mu_1 \log P(\sigma_1 \succ \sigma_2 | \hat{r}) + \mu_2 \log P(\sigma_1 \prec \sigma_2 | \hat{r}) \quad (1)$$

89 This loss is under-specified until $P(\sigma_1 \succ \sigma_2 | \hat{r})$ is defined, which is the focus of this paper. We show that
 90 the common model of preference probabilities is flawed and introduce an improved preference model.

91 **Preference models** A preference model determines the probability of one trajectory segment being
 92 preferred over another, $P(\sigma_1 \succ \sigma_2 | \hat{r})$. Preference models could be applied to model preferences
 93 provided by humans or other systems. Preference models can also directly generate preferences, and in
 94 such cases we refer to them as **preference generators**.

95 2.2 Choice of preference model: partial return and regret

96 **Partial return** Recent work assumes human preferences are generated by a Boltzmann distribution
 97 over the two segments' partial returns [9-16], expressed here as a logistic function [1]:

$$P_{\Sigma_r}(\sigma_1 \succ \sigma_2 | \tilde{r}) = \text{logistic}(\Sigma_{\sigma_1} \tilde{r} - \Sigma_{\sigma_2} \tilde{r}). \quad (2)$$

98 **Regret** We introduce an alternative preference model based on the regret of each transition in a
 99 segment. We first focus on segments with deterministic transitions. For a transition (s_t, a_t, s_{t+1}) in a
 100 deterministic segment, $\text{regret}_d(\sigma_t | \tilde{r}) \triangleq V_{\tilde{r}}^*(s_{\sigma,t}) - [\tilde{r}_t + V_{\tilde{r}}^*(s_{\sigma,t+1})]$. For a full deterministic segment,
 101

$$\text{regret}_d(\sigma | \tilde{r}) \triangleq \sum_{t=0}^{|\sigma|-1} \text{regret}_d(\sigma_t | \tilde{r}) = V_{\tilde{r}}^*(s_{\sigma,0}) - (\Sigma_{\sigma} \tilde{r} + V_{\tilde{r}}^*(s_{\sigma,|\sigma|})), \quad (3)$$

102 with the right-hand expression arising from cancelling out intermediate state values. Therefore,
 103 deterministic regret measures how much the segment reduces expected return from $V_{\tilde{r}}^*(s_{\sigma,0})$. An
 104 optimal segment, σ^* , always has 0 regret, and a suboptimal segment, $\sigma^{\neg*}$, will always have positive
 105 regret, a intuitively appealing property that also plays a role in the identifiability proof of Theorem 3.1

106 Stochastic transitions, however, can result in $\text{regret}_d(\sigma^* | \tilde{r}) > \text{regret}_d(\sigma^{\neg*} | \tilde{r})$, losing the property
 107 above. To retain it, we note that the effect on expected return of transition stochasticity from a
 108 transition (s_t, a_t, s_{t+1}) is $[\tilde{r}_t + V_{\tilde{r}}^*(s_{t+1})] - Q_{\tilde{r}}^*(s_t, a_t)$ and add this expression once per transition to
 109 get $\text{regret}(\sigma)$, removing the subscript d that refers to determinism. The regret for a single transition
 110 becomes $\text{regret}(\sigma_t | \tilde{r}) = [V_{\tilde{r}}^*(s_{\sigma,t}) - [\tilde{r}_t + V_{\tilde{r}}^*(s_{\sigma,t+1})]] + [[\tilde{r}_t + V_{\tilde{r}}^*(s_{\sigma,t+1})] - Q_{\tilde{r}}^*(s_{\sigma,t}, a_{\sigma,t})] =$
 111 $V_{\tilde{r}}^*(s_{\sigma,t}) - Q_{\tilde{r}}^*(s_{\sigma,t}, a_{\sigma,t}) = -A_{\tilde{r}}^*(s_{\sigma,t}, a_{\sigma,t})$. Regret for a full segment is

$$\text{regret}(\sigma | \tilde{r}) = \sum_{t=0}^{|\sigma|-1} \text{regret}(\sigma_t | \tilde{r}) = \sum_{t=0}^{|\sigma|-1} [V_{\tilde{r}}^*(s_{\sigma,t}) - Q_{\tilde{r}}^*(s_{\sigma,t}, a_{\sigma,t})] = \sum_{t=0}^{|\sigma|-1} -A_{\tilde{r}}^*(s_{\sigma,t}, a_{\sigma,t}). \quad (4)$$

112 The regret preference model is the Boltzmann distribution over negated regret:

$$P_{\text{regret}}(\sigma_1 \succ \sigma_2 | \tilde{r}) \triangleq \text{logistic}(\text{regret}(\sigma_2 | \tilde{r}) - \text{regret}(\sigma_1 | \tilde{r})). \quad (5)$$

113 Lastly, we note that if two segments have deterministic transitions, end in terminal states, and have the
 114 same starting state, this regret model reduces to the partial return model: $P_{\text{regret}}(\cdot | \tilde{r}) = P_{\Sigma_r}(\cdot | \tilde{r})$.

115 **Algorithms in this paper** All algorithms in the body of this paper are defined simply as “minimize
 116 Equation [1]”. They differ only in how the preference probabilities are calculated. All reward function
 117 learning via partial return uses Equation [2]. We use two algorithms for reward function learning

¹See Appendix B for a derivation of this logistic expression from a Boltzmann distribution with a temperature of 1. Unless otherwise stated, we ignore the temperature because scaling reward has the same effect.

118 via regret. The theory in Section 3 assumes exact measurement of regret, using Equation 5. Our
 119 experimental results in Section 6 use Equation 6 to approximate regret. Appendix B introduces other
 120 algorithms that use Equation 1 as well as one in Appendix B.2 that generalizes Equation 1.

121 **Regret as a model for human preference** P_{regret} makes at least three assumptions worth noting.
 122 First, it keeps the assumption that human preferences follow a Boltzmann distribution over some
 123 statistic, which is a common model of choice behavior in economics and psychology, where it is
 124 called the Luce-Shepard choice rule [17, 18]. Second, P_{regret} implicitly assumes humans can identify
 125 optimal and suboptimal segments when they see them, which will less true in domains where the human
 126 has less expertise. Lastly, P_{regret} assumes that in stochastic settings where the best *outcome* may only
 127 result from suboptimal decisions (e.g., buying a lottery ticket), humans instead prefer optimal *decisions*.
 128 We suspect humans are capable of expressing either type of preference—based on decision quality
 129 or desirability of outcomes—and can be influenced by training or the preference elicitation interface.
 130 In practice we determine that the regret model produces improvements over the partial-return model
 131 (Section 6), and its assumptions represent an opportunity for follow-up research.

132 **Alternative methods for learning reward functions** Other methods for learning reward functions
 133 include inverse reinforcement learning from demonstrations [19, 20] (discussed in Appendix B.5) and
 134 inverse reward design from trial-and-error reward design in multiple instances of a task domain [21].

135 3 Theoretical comparisons

136 In this section, we consider how different ways of generating preferences affect reward inference, setting
 137 aside whether humans can be influenced to give preferences in accordance with a specific preference
 138 method. In economic terms, this analysis—and all of our analyses with synthetic preferences—could
 139 be considered a normative analysis. In artificial intelligence, this analysis might be cast as a step
 140 towards defining criteria for a rational preference model.

141 **Definition 3.1** (An identifiable preference model). *For a preference model P , assume an infinite*
 142 *dataset D_{\succ} of n -length pairs of segments is constructed by repeatedly choosing (σ_1, σ_2) and sampling*
 143 *a label $\mu \sim P(\sigma_1 \succ \sigma_2 | r)$, using P as a preference generator. Further assume that in this dataset, all*
 144 *possible n -length segment pairs appear infinitely many times. For some MDP $\setminus r$ M , let $M_{\tilde{r}}$ be M with*
 145 *the reward function \tilde{r} . Let $\Pi_{\tilde{r}}^*$ be the set of optimal policies in $M_{\tilde{r}}$. Let reward-equivalence class \mathfrak{R} be*
 146 *the set of all reward functions such that if $r_1, r_2 \in \mathfrak{R}$ then $\Pi_{r_1}^* = \Pi_{r_2}^*$. Preference model P is **identifiable***
 147 *if, for any choice of n and M_r , any $\hat{r} = \operatorname{argmin}_{\tilde{r} \in D_{\succ}} [\operatorname{loss}(\tilde{r})]$ —for the cross-entropy loss (Eqn. 1)—*
 148 *with P as the preference model—is in the same reward equivalence class as r . I.e., $\Pi_{\hat{r}}^* = \Pi_r^*$.*

149 **Theorem 3.1** (P_{regret} is identifiable). *Let P_{regret} be any function such that if $\operatorname{regret}(\sigma_1 | \tilde{r}) <$*
 150 *$\operatorname{regret}(\sigma_2 | \tilde{r})$, $P_{regret}(\sigma_1 \succ \sigma_2 | \tilde{r}) > 0.5$, and if $\operatorname{regret}(\sigma_1 | \tilde{r}) = \operatorname{regret}(\sigma_2 | \tilde{r})$, $P_{regret}(\sigma_1 \succ \sigma_2 | \tilde{r}) =$*
 151 *0.5. P_{regret} is identifiable.*

152 This class of regret preference models includes but is not limited to the Boltzmann distribution of Eqn. 5
 153 and the narrower class that Theorem 3.1 focuses upon.

154 **Theorem 3.2** (Noiseless P_{Σ_r} is not identifiable). *Let P_{Σ_r} be any function such that if $\Sigma_{\sigma_1} \tilde{r} > \Sigma_{\sigma_2} \tilde{r}$,*
 155 *$P_{\Sigma_r}(\sigma_1 \succ \sigma_2 | \tilde{r}) = 1$, and if $\Sigma_{\sigma_1} \tilde{r} = \Sigma_{\sigma_2} \tilde{r}$, $P_{\Sigma_r}(\sigma_1 \succ \sigma_2 | \tilde{r}) = 0.5$. There exists an MDP in which P_{Σ_r} is*
 156 *not identifiable.*

157 Appendix C contains a proof of Theorem 3.1 and two proofs by example for Theorem 3.2 each
 158 focusing on a different weakness of P_{Σ_r} . The first proof by example reveals issues when learning
 159 reward functions with stochastic transitions with either P_{Σ_r} or deterministic P_{regret} . These issues
 160 directly correspond to the need for preferences over distributions over outcomes (i.e., lotteries) to
 161 construct a cardinal utility function (see Russell and Norvig [22, Ch. 16]). Note that the noiseless
 162 version of P_{Σ_r} in Theorem 3.2 is achieved in the limit as reward values are scaled higher; equivalently,
 163 one could include a Boltzmann temperature parameter in Equation 2 and scale it towards 0. Intuitively,
 164 Theorem 3.2 says that P_{Σ_r} is not identifiable without the distribution over preferences providing
 165 information about the proportions of rewards with respect to each other. In contrast, to be identifiable,
 166 the regret preference model does not require this preference error (though it can presumably benefit
 167 from it in certain contexts).

168 4 Creating a human-labeled preference dataset

169 To empirically investigate the consequences of each preference model when learning reward from
170 *human* preferences, we created a preference dataset labeled by human subjects via Amazon Mechanical
171 Turk. This data collection was IRB-approved. Appendix D adds detail to the content below.

172 4.1 The general delivery domain

173 The delivery domain consists of a grid of cells, each of a specific road surface type. The delivery agent’s
174 state is its location. The agent’s action space is moving one cell in one of the four cardinal directions.
175 The episode can terminate either at the destination for +50 reward or in failure at a sheep for −50
176 reward. The reward for a non-terminal transition is the sum of any reward components. Cells with a
177 white road surface have a −1 reward component, and cells with brick surface have a −2 component.
178 Additionally, each cell may contain a coin (+1) or a roadblock (−1). Coins do not disappear and at
179 best cancel out the road surface cost. Actions that would move the agent into a house or beyond the
180 grid’s perimeter result in no motion and receive reward that includes the current cell’s surface reward
181 component but not any coin or roadblock components. In this work, the start state distribution, D_0 , is
182 always uniformly random over non-terminal states. This domain was designed to permit subjects to
183 easily identify bad behavior yet also to be difficult for them to determine *optimal* behavior from most
184 states, which is representative of many common tasks.

185 4.1.1 The delivery task

186 We chose one instantiation of the delivery domain for gathering our dataset of human preferences. This specific MDP
187 has a 10×10 grid. From every state, the highest return possible involves reaching the goal, rather than hitting a sheep or
188 perpetually avoiding termination. Figure 2 shows this task.
189

191 4.2 The user interface and survey

192 This subsection describes the three main stages of the experimental session. A video showing the full experimental
193 protocol can be seen at bit.ly/humanprefs
194

195 **Teaching subjects about the task** Subjects first view instructions describing the general domain. To avoid the jargon
196 of “return” and “reward,” these terms are mapped to equivalent values in US dollars, and the instructions describe the
197 goal of the task as maximizing the delivery vehicle’s financial outcome, where the reward components
198 are specific financial impacts. This information is shared amongst interspersed interactive episodes,
199 in which the subject controls the agent in domain maps that are each designed to teach one or two
200 concepts. Our intention during this stage is to inform the later preferences of the subject by teaching
201 them about the domain’s dynamics and its reward function, as well as to develop the subject’s sense of
202 how desirable various behaviors are. At the end of this stage, the subject controls the agent for two
203 episodes in the specific delivery task shown in Figure 2.
204
205

206 **Preference elicitation** After each subject is trained to understand the task, they indicate their
207 preferences between 40–50 randomly-ordered pairs of segments, using the interface shown in Figure 3.
208 The users select a preference, no preference (“same”), or “can’t tell”. In this work, we exclude responses
209 labeled “can’t tell”, though one might alternatively try to extract information from these responses.

210 **Users’ task comprehension** Subjects then answered questions testing their understanding of the task,
211 and we removed their data if they scored poorly. We also removed a subject’s data if they preferred
212 colliding the vehicle into a sheep over not doing so, which we interpreted as poor task understanding or
213 inattentiveness. This filtered dataset contains 1812 preferences from 50 subjects.

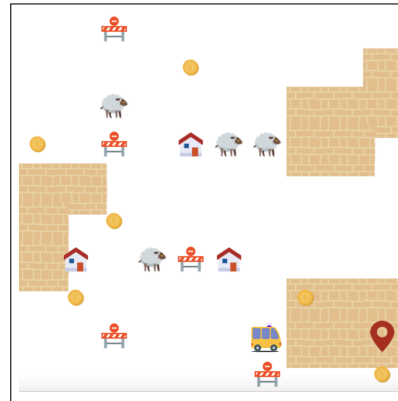


Figure 2: The delivery task used to gather human preferences. The yellow van is the agent and the red inverted teardrop is the destination.

214 4.3 Selection of segment pairs for labeling

215 We collected human preferences
 216 in two stages, each
 217 with different methods for
 218 selecting which segment
 219 pairs to present for labeling.
 220 The second stage’s
 221 sole purpose was to improve
 222 the reward-learning
 223 performance of P_{Σ_r} . Without
 224 second-stage data, P_{Σ_r}
 225 compared even worse to
 226 P_{regret} than in the results
 227 described in Section 6 (see
 228 Appendix ??). Both stages’
 229 data are combined and used as
 230 a single dataset. These methods
 and their justification are described in Appendix D.3

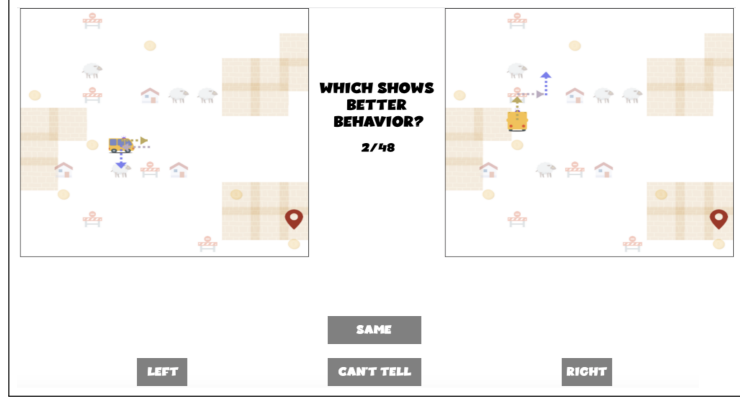


Figure 3: Interface shown to subjects during preference elicitation.

231 5 Descriptive results

232 This section considers how well different preference
 233 models explain our dataset of human preferences.
 234

235 5.1 Correlations 236 between preferences and segment statistics

237 We hypothesize that the values of segments’ start
 238 and end states—which are included in P_{regret}
 239 but not in P_{Σ} —affect human preferences, independent
 240 of partial return. To simplify analysis, we combine the
 241 two parts of $regret_d(\sigma|r)$ that are additional to
 242 $\Sigma_{\sigma}\tilde{r}$ and introduce the following shorthand:
 243 $\Delta_{\sigma}V_{\tilde{r}} \triangleq V_{\tilde{r}}^*(s_{\sigma,|\sigma|}) - V_{\tilde{r}}^*(s_{\sigma,0})$. Note that
 244 with an algebraic manipulation (see Appendix E.1),
 245 $regret_d(\sigma_2|\tilde{r}) - regret_d(\sigma_1|\tilde{r}) = (\Delta_{\sigma_1}V_{\tilde{r}} - \Delta_{\sigma_2}V_{\tilde{r}}) + (\Sigma_{\sigma_1}\tilde{r} - \Sigma_{\sigma_2}\tilde{r})$.
 246 Therefore, on the diagonal line in Figure 4, $regret_d(\sigma_2|r) =$
 247 $regret_d(\sigma_1|r)$, making the P_{regret_d} preference model indifferent.

249 The dataset of preferences is visualized in Figure 4. This
 250 plot shows how $\Delta_{\sigma}V_r$ has influence independent of
 251 partial return by focusing only on points at a chosen
 252 y -axis value; if the colors along the corresponding
 253 horizontal line reddens as the x -axis value increases,
 254 then $\Delta_{\sigma}V_r$ appears to have independent influence. To
 255 statistically test for independent influence of $\Delta_{\sigma}V_r$
 256 on preferences, we consider subsets of data where
 257 $\Sigma_{\sigma_1}r - \Sigma_{\sigma_2}r$ is constant. For $\Sigma_{\sigma_1}r - \Sigma_{\sigma_2}r = -1$
 258 and $\Sigma_{\sigma_1}r - \Sigma_{\sigma_2}r = -2$, the only values with more
 259 than 30 samples that also include informative samples
 260 with both negative and positive values of $regret(\sigma_1|r) -$
 261 $regret(\sigma_2|r)$, the Spearman’s rank correlations between
 $\Delta_{\sigma}V_r$ and the preferences are significant ($r > 0.3$,
 $p < 0.0001$). This result indicates that $\Delta_{\sigma}V_r$
 influences human preferences independent of partial
 return, validating our hypothesis that humans form
 preferences based on information about segments’ start
 states and end states, not only partial returns.

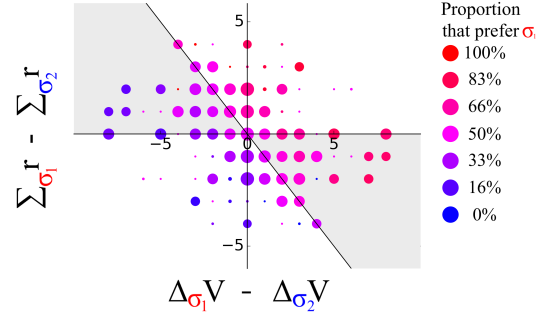


Figure 4: Proportions at which subjects preferred each segment in a pair, plotted by the difference in the segments’ changes in state values (x-axis) and partial returns (y-axis). The diagonal line shows points of preference indifference for P_{regret} . Points of indifference for P_{Σ} lie on the x-axis. The shaded gray area indicates where the two models disagree, each giving a different segment a preference probability greater than 0.5. Each circle’s area is proportional to the number of samples it describes.

| Preference model | Loss |
|------------------------------------|-------------|
| $P(\cdot) \equiv 0.5$ (uninformed) | 0.69 |
| P_{Σ_r} (partial return) | 0.62 |
| P_{regret} | 0.57 |

Table 1: Mean cross-entropy test loss over 10-fold cross validation (n=1812) from predicting human preferences. Lower is better.

262 5.2 Likelihood of human preferences under different preference models

263 To examine how well each preference model predicts human preferences, we calculate the cross-
 264 entropy loss for each model (Eqn. 1)—i.e., the negative log likelihood—of the preferences in our
 265 dataset. Scaling reward by a constant factor does not affect the set of optimal policies. Therefore,
 266 throughout this work we ensure that our analyses of preference models are insensitive to reward scaling.
 267 To do so for this specific analysis, we conduct 10-fold cross validation to learn a reward scaling factor
 268 for each of P_{regret} and P_{Σ_r} . Table 1 shows that the loss of P_{regret} is lower than that of P_{Σ_r} , indicating
 269 that it is more reflective of how people actually express preferences.

270 6 Results from learning reward functions

271 Analysis of a preference model’s predictions of human preferences is informative, but such predictions
 272 are a means to the ends of learning human-aligned reward functions and policies. We now examine each
 273 preference model’s performance on these ends. In all cases, we learn a reward function \hat{r} according
 274 to Eqn. 1 and apply value iteration [23] to find the approximately optimal $Q_{\hat{r}}^*$ function. For this $Q_{\hat{r}}^*$,
 275 we then evaluate the mean return of the maximum-entropy optimal policy—which chooses uniformly
 276 randomly among all *optimal* actions—with respect to the ground-truth reward function r , over D_0 .
 277 To compare performance across different MDPs, the mean return of a policy π , V_r^π , is normalized
 278 to $(V_r^\pi - V_r^U)/V_r^*$, where V_r^* is the optimal expected return and V_r^U is the expected return of the
 279 uniformly random policy (both given D_0). Normalized mean return above 0 is better than V_r^U . Optimal
 280 policies have a normalized mean return of 1, and we consider above 0.9 to be *near optimal*.

281 6.1 An algorithm to learn reward functions with $regret(\sigma_\sigma|\hat{r})$

282 Algorithm 1 is a general algorithm for learning a *linear* reward function according to P_{regret} . This
 283 regret-specific algorithm only changes the regret-based algorithm from Section 2.2 by replacing
 284 Equation 5 with a tractable approximation of regret, avoiding expensive repeated evaluation of $V_{\hat{r}}^*(\cdot)$
 285 and $Q_{\hat{r}}^*(\cdot, \cdot)$ to compute $P_{regret}(\cdot|\hat{r})$ during reward learning. Specifically, successor features for a set
 286 of policies are used to approximate the optimal state values and state-action values for *any* reward
 287 function.

288 **Approximating P_{regret} with successor features** Following the notation of Barreto et al. [24], assume
 289 the ground-truth reward is linear with respect to a feature vector extracted by $\phi: S \times A \times S \rightarrow \mathbb{R}^d$ and
 290 a weight vector $\mathbf{w}_r \in \mathbb{R}^d$: $r(s, a, s') = \phi(s, a, s')^\top \mathbf{w}_r$. During learning, $\mathbf{w}_{\hat{r}}$ similarly expresses \hat{r} as
 291 $\hat{r}(s, a, s') = \phi(s, a, s')^\top \mathbf{w}_{\hat{r}}$.

292 Given a policy π , the successor features for (s, a) are the expectation of discounted reward features
 293 from that state-action pair when following π : $\psi_Q^\pi(s, a) = E^\pi[\sum_{i=t}^{\infty} \gamma^{i-t} \phi(s_t, a_t, s_{t+1}) | s_t = s, a_t = a]$.
 294 Therefore, $Q_{\hat{r}}^\pi(s, a) = \psi_Q^\pi(s, a)^\top \mathbf{w}_{\hat{r}}$. Additionally, state-based successor features can be calculated
 295 from the ψ_Q^π above as $\psi_V^\pi(s) = \sum_{a \in A} \pi(a|s) \psi_Q^\pi(s, a)$, making $V_{\hat{r}}^\pi(s) = \psi_V^\pi(s)^\top \mathbf{w}_{\hat{r}}$.

296 Given a set Ψ_Q of state-action successor feature functions and a set Ψ_V of state successor feature func-
 297 tions for various policies and given a reward function via $\mathbf{w}_{\hat{r}}$, $Q_{\hat{r}}^*(s, a) \geq \max_{\psi_Q \in \Psi_Q} [\psi_Q^\pi(s, a)^\top \mathbf{w}_{\hat{r}}]$
 298 and $V_{\hat{r}}^*(s) \geq \max_{\psi_V \in \Psi_V} [\psi_V^\pi(s)^\top \mathbf{w}_{\hat{r}}]$ [24], so we use these two maximizations as approximations of
 299 $Q_{\hat{r}}^*(s, a)$ and $V_{\hat{r}}^*(s)$, respectively. In practice, to enable gradient-based optimization with current tools,
 300 the maximization in this expression is replaced with the softmax-weighted average, making the loss
 301 function linear. Focusing first on the approximation of $V_{\hat{r}}^*(s)$, for each $\psi_V \in \Psi_V$, a softmax weight is
 302 calculated for $\psi_V^\pi(s)$: $\text{softmax}_{\Psi_V}(\psi_V^\pi(s)^\top \mathbf{w}_{\hat{r}}) \triangleq [(\psi_V^\pi(s)^\top \mathbf{w}_{\hat{r}})^{1/T}] / [(\sum_{\psi_V' \in \Psi_V} \psi_V'^\pi(s)^\top \mathbf{w}_{\hat{r}})^{1/T}]$,
 303 where temperature T is a constant hyperparameter. The resulting approximation of $V_{\hat{r}}^*(s)$ is there-
 304 fore defined as $\tilde{V}_{\hat{r}}^*(s) \triangleq \sum_{\psi_V \in \Psi_V} \text{softmax}_{\Psi_V}(\psi_V^\pi(s)^\top \mathbf{w}_{\hat{r}}) [\psi_V^\pi(s)^\top \mathbf{w}_{\hat{r}}]$. Similarly, to approxi-
 305 mate $Q_{\hat{r}}^*(s, a)$, $\text{softmax}_{\Psi_Q}(\psi_Q^\pi(s, a)^\top \mathbf{w}_{\hat{r}}) \triangleq [(\psi_Q^\pi(s, a)^\top \mathbf{w}_{\hat{r}})^{1/T}] / [(\sum_{\psi_Q' \in \Psi_Q} \psi_Q'^\pi(s, a)^\top \mathbf{w}_{\hat{r}})^{1/T}]$
 306 and $\tilde{Q}_{\hat{r}}^*(s, a) \triangleq \sum_{\psi_Q \in \Psi_Q} \text{softmax}_{\Psi_Q}(\psi_Q^\pi(s, a)^\top \mathbf{w}_{\hat{r}}) [\psi_Q^\pi(s, a)^\top \mathbf{w}_{\hat{r}}]$. Consequently, from Eqns. 4

Algorithm 1 Linear reward learning with regret preference model (P_{regret}), using successor features

- 1: Input: a set of reward functions and a set of policies (where one set can be \emptyset)
 - 2: $\Psi \leftarrow \emptyset$
 - 3: **for** each reward function r_{SF} or policy π_{SF} in the input sets **do**
 - 4: **if** r_{SF} **then** $\pi_{SF} \leftarrow$ estimate of optimal maximum-entropy policy for r_{SF}
 - 5: estimate $\psi_Q^{\pi_{SF}}$ and $\psi_V^{\pi_{SF}}$ (if not estimated already during step 4)
 - 6: add $\psi_Q^{\pi_{SF}}$ to Ψ_Q
 - 7: add $\psi_V^{\pi_{SF}}$ to Ψ_V
 - 8: **end for**
 - 9: **repeat**
 - 10: optimize $w_{\hat{r}}$ by loss of Eqn. [1](#), calculating $\tilde{P}_{regret}(\sigma_1 \succ \sigma_2 | \hat{r})$ via Eqn. [6](#) using Ψ_Q and Ψ_V
 - 11: **until** stopping criteria are met
 - 12: **return** $w_{\hat{r}}$
-

307 and [5](#) the corresponding approximation \tilde{P}_{regret} of the regret preference model is:

$$\tilde{P}_{regret}(\sigma_1 \succ \sigma_2 | \hat{r}) = \text{logistic} \left(\sum_{t=0}^{|\sigma_2|-1} [\tilde{V}_{\hat{r}}^*(s_{\sigma_2,t}) - \tilde{Q}_{\hat{r}}^*(s_{\sigma_2,t}, a_{\sigma_2,t})] - \sum_{t=0}^{|\sigma_1|-1} [\tilde{V}_{\hat{r}}^*(s_{\sigma_1,t}) - \tilde{Q}_{\hat{r}}^*(s_{\sigma_1,t}, a_{\sigma_1,t})] \right) \quad (6)$$

308 **The algorithm** In Algorithm [1](#), lines 9–12 describe the supervised-learning optimization using
309 the approximation \tilde{P}_{regret} , and the prior lines create Ψ_Q and Ψ_V . Specifically, given a set of reward
310 functions, a corresponding set of policies is created (line 4), where each policy is an estimate of the
311 maximum entropy policy for a reward function. Standard policy improvement methods can be used to
312 create each such policy. Alternatively, some or all of the set of policies can be given as input directly,
313 not derived from input reward functions. For each such policy π_{SF} , successor feature functions $\Psi_Q^{\pi_{SF}}$
314 and $\Psi_V^{\pi_{SF}}$ are estimated (line 5), which by default would be performed by a minor extension of a
315 standard policy evaluation algorithm as detailed by Barreto et al. [\[24\]](#). Note that the reward function
316 that is ultimately learned is not restricted to be in the input set of reward functions, which is used only
317 to create an approximation of regret.

318 The details of our instantiation of Algorithm [1](#) for the delivery domain can be found in Appendix [F.1](#)
319 along with guidance for extending it to reward functions that might be non-linear.

320 6.2 Results from synthetic preferences

321 Before considering human preferences, we first ask how each preference model performs when it is
322 correct. In other words, we investigate empirically how well the preference model could perform if
323 humans perfectly adhered to it. Recall that the ground-truth reward function, r , is used to create these
324 preferences but is inaccessible to the reward-learning algorithms.

325 For these evaluations, either a stochastic or
326 noiseless preference model acts a preference
327 generator to create a preference dataset, and
328 then the stochastic version of the same model
329 is used for reward learning. For the noiseless
330 case, the deterministic preference generator compares a segment pair's $\Sigma_{\sigma} r$ values for P_{Σ_r} or
331 their $regret(\sigma|r)$ values for P_{regret} . Note that
332 through reward scaling the preference generators
333 approach determinism in the limit, so this noiseless
334 analysis examines minimal-entropy versions
335 of the two preference-generating models. (The opposite extreme, uniformly random preferences,
336 would remove all information from preferences and therefore is not examined.) In the stochastic case,
337 for each preference model, each segment pair is labeled by sampling from that preference generator's
338 output distribution (Eqs [2](#) or [5](#)), using the unscaled ground-truth reward function.
339

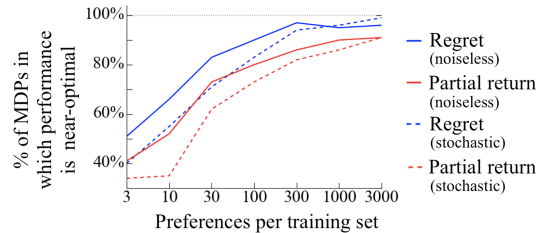


Figure 5: Performance comparison over 100 randomly generated deterministic MDPs

340 We created 100 deterministic MDPs that instantiate variants of our delivery domain (see Section 4.1).
 341 To create each MDP, we sampled from sets of possible widths, heights, and reward component values,
 342 and the resultant grid cells were randomly populated with a destination, objects, and road surface types
 343 (see Appendix F.2 for details). Each segment in the preference datasets for each MDP was generated
 344 by choosing a start state and three actions, all uniformly randomly. For a set number of preferences,
 345 each method had the same set of segment pairs in its preference dataset. Figure 5 shows the percentage
 346 of MDPs in which each preference model results in near-optimal performance. The regret preference
 347 model outperforms the partial return model at every dataset size, both with and without noise. By a
 348 Wilcoxon paired signed-rank test on normalized mean returns, $p < 0.05$ for 86% of these comparisons
 349 and $p < 0.01$ for 57% of them, as reported in Appendix F.2.

350 Further analyses can be found in Appendix F.2 including with stochastic transitions, with different
 351 segment lengths, and while artificially lowering the discount factor (as is common in deep RL and
 352 recent work on deep reward learning from preferences).

353 6.3 Results from human preferences

354 We randomly assign human preferences from our gathered dataset to different numbers of same-sized partitions,
 355 resulting in different training set sizes, and test each preference model on each partition. Figure 6
 356 shows the results. With smaller training sets (20–100 partitions), the regret preference model results in near-
 357 optimal performance more often. With larger training sets (1–10 partitions), both preference models always
 358 reach near-optimal return, but the mean return from the regret preference model is higher for all of these
 359 partitions except for 3 partitions in the 10-partition test. Applying a Wilcoxon paired signed-rank test on normalized mean return to each group with 5 or
 360 more partitions, $p < 0.05$ for all numbers of partitions except 100 and $p < 0.01$ for 20 and 50 partitions.
 361
 362
 363
 364
 365
 366

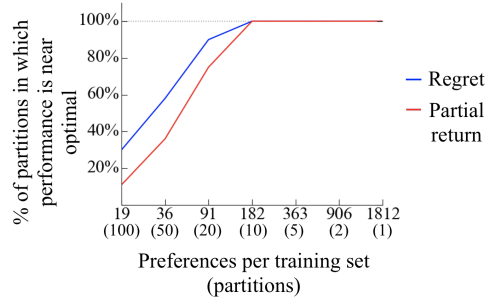


Figure 6: Performance comparison over various amounts of human preferences. Each partition has the number of preferences shown or one less.

367 7 Conclusion

368 Over numerous evaluations with human preferences, our proposed regret preference model (P_{regret})
 369 shows improvements summarized below over the previous partial return preference model (P_{Σ_r}).
 370 When each preference model generates the preferences for its own infinite and exhaustive training set,
 371 we prove that P_{regret} identifies the set of optimal policies, whereas P_{Σ_r} is not guaranteed to do so
 372 without preference noise that reveals the proportions of rewards with respect to each other. With finite
 373 training data of synthetic preferences, P_{regret} also empirically results in learned policies that tend to
 374 outperform those resulting from P_{Σ_r} . This superior performance of P_{regret} is also seen with human
 375 preferences. In summary, our analyses suggest that regret preference models are more effective both
 376 descriptively with respect to human preferences and also normatively, as the model we want humans to
 377 follow if we had the choice.

378 Independent of P_{regret} , this paper also reveals that segments’ changes in state values provide informa-
 379 tion about human preferences that is not fully provided by partial return. More generally, we show that
 380 the choice of preference model impacts the performance of learned reward functions.

381 This study motivates several new directions for research. Future work could address any of the
 382 limitations detailed in Appendix A.1. Specifically, future work could further test the general superiority
 383 of P_{regret} or apply it to deep learning settings. Additionally, *prescriptive* methods could be developed
 384 via the user interface or elsewhere to nudge humans to conform more to P_{regret} or to other normatively
 385 appealing preference models. Lastly, subsequent efforts could seek preference models that are even
 386 more effective with preferences from actual humans, now that this work has provided conclusive
 387 evidence that the choice of preference model is impactful.

388 **References**

- 389 [1] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian
390 Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go
391 with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 392 [2] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin,
393 Augustin Žídek, Alexander WR Nelson, Alex Bridgland, et al. Improved protein structure prediction using
394 potentials from deep learning. *Nature*, 577(7792):706–710, 2020.
- 395 [3] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung
396 Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in StarCraft
397 II using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 398 [4] Marc G Bellemare, Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C Machado, Subhodeep
399 Moitra, Sameera S Ponda, and Ziyu Wang. Autonomous navigation of stratospheric balloons using rein-
400 forcement learning. *Nature*, 588(7836):77–82, 2020.
- 401 [5] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison,
402 David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement
403 learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 404 [6] Jonas Degraeve, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese,
405 Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de Las Casas, et al. Magnetic control of tokamak
406 plasmas through deep reinforcement learning. *Nature*, 602(7897):414–419, 2022.
- 407 [7] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete
408 problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 409 [8] W Bradley Knox, Alessandro Allievi, Holger Banzhaf, Felix Schmitt, and Peter Stone. Reward (mis)design
410 for autonomous driving. *arXiv preprint arXiv:2104.13906*, 2021.
- 411 [9] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforce-
412 ment learning from human preferences. In *Advances in Neural Information Processing Systems (NIPS)*,
413 pages 4299–4307, 2017.
- 414 [10] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward learning
415 from human preferences and demonstrations in atari. *arXiv preprint arXiv:1811.06521*, 2018.
- 416 [11] Xiaofei Wang, Kimin Lee, Kourosh Hakhamaneshi, Pieter Abbeel, and Michael Laskin. Skill preferences:
417 Learning to extract and execute robotic skills from human feedback. In *Conference on Robot Learning*,
418 pages 1259–1268. PMLR, 2022.
- 419 [12] Erdem Bıyık, Dylan P Losey, Malayandi Palan, Nicholas C Landolfi, Gleb Shevchuk, and Dorsa Sadigh.
420 Learning reward functions from diverse sources of human feedback: Optimally integrating demonstrations
421 and preferences. *The International Journal of Robotics Research*, page 02783649211041652, 2021.
- 422 [13] Dorsa Sadigh, Anca D Dragan, Shankar Sastry, and Sanjit A Seshia. Active preference-based learning of
423 reward functions. *Robotics: Science and Systems*, 2017.
- 424 [14] Kimin Lee, Laura Smith, and Pieter Abbeel. Pebble: Feedback-efficient interactive reinforcement learning
425 via relabeling experience and unsupervised pre-training. *arXiv preprint arXiv:2106.05091*, 2021.
- 426 [15] Kimin Lee, Laura Smith, Anca Dragan, and Pieter Abbeel. B-pref: Benchmarking preference-based
427 reinforcement learning. *arXiv preprint arXiv:2111.03026*, 2021.
- 428 [16] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang,
429 Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with
430 human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- 431 [17] R Duncan Luce. *Individual choice behavior: A theoretical analysis*. John Wiley, 1959.
- 432 [18] Roger N Shepard. Stimulus and response generalization: A stochastic model relating generalization to
433 distance in psychological space. *Psychometrika*, 22(4):325–345, 1957.
- 434 [19] A.Y. Ng and S. Russell. Algorithms for inverse reinforcement learning. In *Seventeenth International
435 Conference on Machine Learning (ICML)*, 2000.
- 436 [20] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse
437 reinforcement learning. In *Twenty-third AAAI Conference on Artificial Intelligence*, volume 8, pages
438 1433–1438, 2008.
- 439 [21] Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J Russell, and Anca Dragan. Inverse reward
440 design. In *Advances in Neural Information Processing Systems (NIPS)*, pages 6765–6774, 2017.
- 441 [22] Stuart Russell and Peter Norvig. *Artificial intelligence: a modern approach*. 2020.
- 442 [23] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

- 443 [24] André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado Van Hasselt, and David
444 Silver. Successor features for transfer in reinforcement learning. *arXiv preprint arXiv:1606.05312*, 2016.
- 445 [25] Riad Akrou, Marc Schoenauer, and Michele Sebag. Preference-based policy learning. In *Joint European*
446 *Conference on Machine Learning and Knowledge Discovery in Databases*, pages 12–27. Springer, 2011.
- 447 [26] Daniel Brown, Russell Coleman, Ravi Srinivasan, and Scott Niekum. Safe imitation learning via fast
448 bayesian reward inference from preferences. In *International Conference on Machine Learning*, pages
449 1165–1177. PMLR, 2020.
- 450 [27] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings*
451 *of the twenty-first international conference on Machine learning*, page 1, 2004.
- 452 [28] Kuno Kim, Shivam Garg, Kirankumar Shiragur, and Stefano Ermon. Reward identification in inverse
453 reinforcement learning. In *International Conference on Machine Learning*, pages 5496–5505. PMLR, 2021.
- 454 [29] Saurabh Arora and Prashant Doshi. A survey of inverse reinforcement learning: Challenges, methods and
455 progress. *Artificial Intelligence*, 297:103500, 2021.
- 456 [30] John Von Neumann and Oskar Morgenstern. *Theory of games and economic behavior*. Princeton university
457 press, 1944.
- 458 [31] Yuchen Cui, Qiping Zhang, Alessandro Allievi, Peter Stone, Scott Niekum, and W Bradley Knox. The
459 empathic framework for task learning from implicit human feedback. *arXiv preprint arXiv:2009.13649*,
460 2020.
- 461 [32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
462 Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary
463 DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and
464 Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach,
465 H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural*
466 *Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- 467 [33] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,
468 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.
469 Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

470 Checklist

- 471 1. For all authors...
- 472 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
473 contributions and scope? [Yes]
- 474 (b) Did you describe the limitations of your work? [Yes] See Appendix [A.1](#)
- 475 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
476 Appendix [A.2](#)
- 477 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
478 them? [Yes]
- 479 2. If you are including theoretical results...
- 480 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Sections [3](#) and [C](#)
481 include all assumptions.
- 482 (b) Did you include complete proofs of all theoretical results? [Yes] See Section [3](#) and
483 Appendix [C](#)
- 484 3. If you ran experiments...
- 485 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
486 imental results (either in the supplemental material or as a URL)? [No] However, the
487 learning code, the code for running experiments, the code and UI elements for gathering
488 human preferences on Mechanical Turk, and the anonymized human preferences data
489 will be opened. We are particularly excited to provide the first open dataset of human
490 preferences over pairs of trajectory segments.
- 491 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
492 chosen)? [Yes] Appendix [F.1](#)

- 493 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
494 multiple times)? [Yes] Error bars do not seem applicable to our plots, which do not show
495 the exact data that we do statistical testing on. However, statistical significance testing
496 was reported, in Sections 5.1 and 6.2 (with a pointer to the appendix for details).
- 497 (d) Did you include the total amount of compute and the type of resources used (e.g., type of
498 GPUs, internal cluster, or cloud provider)? [Yes] See Appendix F.1.
- 499 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 500 (a) If your work uses existing assets, did you cite the creators? [Yes] Appendix D does so
501 for visual assets used to visualize the delivery task.
- 502 (b) Did you mention the license of the assets? [Yes] Appendix D mentions the license for
503 visual assets used to visualize the delivery task.
- 504 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 505 (d) Did you discuss whether and how consent was obtained from people whose data you're
506 using/curating? [Yes] See Appendix D.
- 507 (e) Did you discuss whether the data you are using/curating contains personally identifiable
508 information or offensive content? [Yes] See Appendix D.
- 509 5. If you used crowdsourcing or conducted research with human subjects...
- 510 (a) Did you include the full text of instructions given to participants and screenshots, if
511 applicable? [Yes] Section 4.1.1 includes a link to a video of a full experimental session
512 (with an author acting as the subject).
- 513 (b) Did you describe any potential participant risks, with links to Institutional Review Board
514 (IRB) approvals, if applicable? [Yes] We discuss participant risks from our crowdsourced
515 study and provide a link to the IRB approval in Appendix D.
- 516 (c) Did you include the estimated hourly wage paid to participants and the total amount
517 spent on participant compensation? [Yes] See Appendix D.