
Optimal Policies Tend To Seek Power

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Some researchers speculate that intelligent reinforcement learning (RL) agents would
2 be incentivized to seek resources and power in pursuit of their objectives. Other
3 researchers are skeptical, because human-like power-seeking instincts need not be
4 present in RL agents. To clarify this debate, we develop the first formal theory of the
5 statistical tendencies of optimal policies in reinforcement learning. In the context of
6 Markov decision processes (MDPs), we prove that certain environmental symmetries
7 are sufficient for optimal policies to tend to seek power over the environment. These
8 symmetries exist in many environments in which the agent can be shut down or
9 destroyed. We prove that for most prior beliefs one might have about the agent’s
10 reward function (including as a special case the situations where the reward function
11 is known), one should expect optimal policies to seek power in these environments.
12 These policies seek power by keeping a range of options available and, when the
13 discount rate is sufficiently close to 1, by navigating towards larger sets of potential
14 terminal states.

15 1 Introduction

16 Omohundro [2008], Bostrom [2014], Russell [2019] hypothesize that highly intelligent agents tend to
17 seek power in pursuit of their goals. Such power-seeking agents might gain power over humans. For
18 example, Marvin Minsky imagined that an agent tasked with proving the Riemann hypothesis might
19 rationally turn the planet – along with everyone on it – into computational resources [Russell and
20 Norvig, 2009].

21 Some researchers argue that these worries stem from anthropomorphization of AI [Various, 2019,
22 Pinker and Russell, 2020, Mitchell, 2021]. LeCun and Zador [2019] argue that AI “never needed to
23 evolve, so it didn’t develop the survival instinct that leads to the impulse to dominate others.”

24 We clarify this debate by grounding the claim that highly intelligent agents will tend to seek power.
25 In section 4, we identify optimal policies in MDPs as a reasonable formalization of “highly intelligent
26 agents.” Optimal policies “tend to” take an action when the action is optimal for most reward functions.

27 Section 5 defines “power” as the ability to achieve a wide range of goals: after all, “money is power”,
28 and money is instrumentally useful for many goals. Conversely, it’s harder to pursue most goals when
29 physically restrained, and so a physically restrained person has little power. An action “seeks power”
30 if it leads to states where the agent has higher power.

31 We make no claims about when real-world AI power-seeking behavior could become plausible. Instead,
32 we consider the theoretical consequences of optimal action in MDPs. Future work is needed to translate
33 our theory from optimal policies to learned, real-world policies – we expect that our arguments will at
34 least hold conceptually, if not provably. Section 6 shows that power-seeking tendencies arise not from
35 anthropomorphism, but from the combination of optimal behavior and certain graphical symmetries
36 present in many MDPs. These symmetries automatically occur in many environments where the agent
37 can be shut down or destroyed, yielding broad applicability of our main result (theorem 6.13).

38 **2 Related Work**

39 An action is *instrumental to an objective* when it helps achieve that objective. Some actions are
 40 instrumental to many objectives, making them *robustly instrumental*. The claim that power-seeking is
 41 robustly instrumental is a specific instance of the *instrumental convergence thesis*:

42 Several instrumental values can be identified which are convergent in the sense that
 43 their attainment would increase the chances of the agent’s goal being realized for a
 44 wide range of final goals and a wide range of situations, implying that these instru-
 45 mental values are likely to be pursued by a broad spectrum of situated intelligent
 46 agents [Bostrom, 2014].

47 For example, in Atari games, avoiding (virtual) death is instrumental for both completing the game
 48 and for optimizing curiosity [Burda et al., 2019]. Many AI alignment researchers hypothesize that
 49 most advanced AI agents will have concerning instrumental incentives, such as resisting deactivation
 50 [Soares et al., 2015, Milli et al., 2017, Hadfield-Menell et al., 2017, Carey, 2018] and acquiring
 51 resources [Benson-Tilsen and Soares, 2016]. Lastly, Menache et al. [2002] identify and navigate
 52 towards robustly instrumental “bottleneck states.”

53 We formalize power as the ability to achieve a wide variety of goals. Appendix A compares our
 54 formalization with information-theoretic empowerment [Salge et al., 2014].

55 Some of our results relate the formal power of states to the structure of the environment. Foster
 56 and Dayan [2002], Drummond [1998], Sutton et al. [2011], Schaul et al. [2015] note that value
 57 functions encode important information about the environment, as they capture the agent’s ability to
 58 achieve different goals. Turner et al. [2020] speculate that a state’s optimal value correlates strongly
 59 across reward functions. In particular, Schaul et al. [2015] learn regularities across value functions,
 60 suggesting that some states are valuable for many different reward functions (*i.e.* powerful).

61 We are not the first to study convergence of behavior, form, or function. In economics, turnpike
 62 theory studies how certain paths of accumulation tend to be optimal [McKenzie, 1976]. In biology,
 63 convergent evolution occurs when similar features (*e.g.* flight) independently evolve in different
 64 time periods [Reece and Campbell, 2011]. Lastly, computer vision networks reliably learn *e.g.* edge
 65 detectors [Olah et al., 2020].

66 **3 State Visit Distribution Functions Quantify The Agent’s Available Options**

67 We clarify the power-seeking debate by proving what optimal policies “usually look like” in a given
 68 environment. We illustrate our results with a simple case study, before explaining how to reason about
 69 a wide range of MDPs. Appendix C.1 lists MDP theory contributions of independent interest, appendix
 70 C lists definitions and theorems, and appendix D contains the proofs.

72 **Definition 3.1** (Rewardless MDP). $\langle \mathcal{S}, \mathcal{A}, T \rangle$ is a re-
 73 wardless MDP with finite state and action spaces
 74 \mathcal{S} and \mathcal{A} , and stochastic transition function $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$. We treat the discount rate γ as a
 75 variable with domain $[0, 1]$.

77 **Definition 3.2** (1-cycle states). Let $\mathbf{e}_s \in \mathbb{R}^{|\mathcal{S}|}$ be
 78 the unit vector for state s , such that there is a 1 in
 79 the entry for state s and 0 elsewhere. State s is a
 80 1-cycle if $\exists a \in \mathcal{A} : T(s, a) = \mathbf{e}_s$. s is a *terminal*
 81 *state* if $\forall a \in \mathcal{A} : T(s, a) = \mathbf{e}_s$.

82 Our theorems apply to stochastic environments, but
 83 we present a deterministic case study for clarity. The
 84 environment of fig. 1 is small, but its structure is
 85 rich. For example, the agent has more “options” at
 86 \star than at the terminal state \emptyset . Formally, \star has more *visit distribution functions* than \emptyset does.

87 **Definition 3.3** (State visit distribution [Sutton and Barto, 1998]). $\Pi := \mathcal{A}^{\mathcal{S}}$, the set of stationary
 88 deterministic policies. The *visit distribution* induced by following policy π from state s at discount

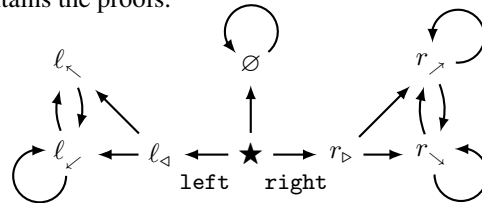


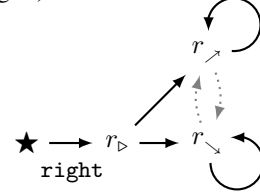
Figure 1: $l_{<}$ is a 1-cycle, and \emptyset is a terminal state. Arrows represent deterministic transitions induced by taking some action $a \in \mathcal{A}$. Since the left subgraph is “almost a copy” of the right subgraph, proposition 6.9 will prove that more reward functions have optimal policies which go right than which go left at state \star , and that such policies seek power – both intuitively, and in a reasonable formal sense.

89 rate $\gamma \in [0, 1)$ is $\mathbf{f}^{\pi, s}(\gamma) := \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s_t \sim \pi | s} [\mathbf{e}_{s_t}]$. $\mathbf{f}^{\pi, s}$ is a *visit distribution function*; $\mathcal{F}(s) :=$
 90 $\{\mathbf{f}^{\pi, s} \mid \pi \in \Pi\}$.

91 In fig. 1, starting from ℓ_{\swarrow} , the agent can stay at ℓ_{\swarrow} or alternate between ℓ_{\swarrow} and ℓ_{\searrow} . Therefore,
 92 $\mathcal{F}(\ell_{\swarrow}) = \{\frac{1}{1-\gamma} \mathbf{e}_{\ell_{\swarrow}}, \frac{1}{1-\gamma^2} (\mathbf{e}_{\ell_{\swarrow}} + \gamma \mathbf{e}_{\ell_{\searrow}})\}$. In contrast, agents at \emptyset must stay at \emptyset . $\mathcal{F}(\emptyset) = \{\frac{1}{1-\gamma} \mathbf{e}_{\emptyset}\}$.
 93 $\pi \mapsto \mathbf{f}^{\pi, s}$ is usually non-injective – at \emptyset , all policies π map to visit distribution function $\frac{1}{1-\gamma} \mathbf{e}_{\emptyset}$.

94 Before moving on, we introduce two important concepts used in our main results. First, we sometimes
 95 restrict our attention to visit distributions which take certain actions (fig. 2).

97 **Definition 3.4** (\mathcal{F} single-state restriction). Consider-
 98 ing only visit distribution functions induced by poli-
 99 cies taking action a at state s' , $\mathcal{F}(s \mid \pi(s') = a) :=$
 100 $\{\mathbf{f} \in \mathcal{F}(s) \mid \exists \pi \in \Pi : \pi(s') = a, \mathbf{f}^{\pi, s} = \mathbf{f}\}$.



101 Second, some $\mathbf{f} \in \mathcal{F}(s)$ are “unimportant.” Consider an agent op-
 102 timizing reward function $\mathbf{e}_{r_{\searrow}}$ (1 reward when at r_{\searrow} , 0 otherwise)
 103 at e.g. $\gamma = \frac{1}{2}$. Its optimal policies navigate to r_{\searrow} and stay there.
 104 Similarly, for reward function $\mathbf{e}_{r_{\swarrow}}$, optimal policies navigate to
 105 r_{\swarrow} and stay there. However, for no reward function is it uniquely
 106 optimal to alternate between r_{\swarrow} and r_{\searrow} . Only *dominated* visit
 107 distribution functions alternate between r_{\swarrow} and r_{\searrow} (definition 3.6).

Figure 2: The subgraph corre-
 sponding to $\mathcal{F}(\star \mid \pi(\star) =$
 right). Gray dotted actions are
 only taken by the policies of dom-
 inated $\mathbf{f}^{\pi} \in \mathcal{F}(\star) \setminus \mathcal{F}_{\text{nd}}(\star)$.

108 **Definition 3.5** (Value function). Let $\pi \in \Pi$. For any reward function $R \in \mathbb{R}^{\mathcal{S}}$ over the state space, the
 109 *on-policy value* at state s and discount rate $\gamma \in [0, 1)$ is $V_R^{\pi}(s, \gamma) := \mathbf{f}^{\pi, s}(\gamma)^{\top} \mathbf{r}$, where $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$ is R
 110 expressed as a column vector (one entry per state). The *optimal value* is $V_R^*(s, \gamma) := \max_{\pi} V_R^{\pi}(s, \gamma)$.

Definition 3.6 (Non-domination).

$$\mathcal{F}_{\text{nd}}(s) := \{\mathbf{f}^{\pi} \in \mathcal{F}(s) \mid \exists \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}, \gamma \in (0, 1) : \mathbf{f}^{\pi}(\gamma)^{\top} \mathbf{r} > \max_{\mathbf{f}^{\pi'} \in \mathcal{F}(s) \setminus \{\mathbf{f}^{\pi}\}} \mathbf{f}^{\pi'}(\gamma)^{\top} \mathbf{r}\}. \quad (1)$$

111 For any reward function R and discount rate γ , $\mathbf{f}^{\pi} \in \mathcal{F}(s)$ is (weakly) dominated by $\mathbf{f}^{\pi'} \in \mathcal{F}(s)$
 112 if $V_R^{\pi}(s, \gamma) \leq V_R^{\pi'}(s, \gamma)$. $\mathbf{f}^{\pi} \in \mathcal{F}_{\text{nd}}(s)$ is *non-dominated* if there exist R and γ at which \mathbf{f}^{π} is not
 113 dominated by any other $\mathbf{f}^{\pi'}$.

114 4 Some Actions Have A Greater Probability Of Being Optimal

115 We claim that optimal policies “tend” to take certain actions in certain situations. We first consider
 116 the probability that certain actions are optimal.

117 Reconsider the reward function $\mathbf{e}_{r_{\searrow}}$, optimized at $\gamma = \frac{1}{2}$. Starting from \star , the optimal trajectory goes
 118 right to r_{\triangleright} to r_{\searrow} , where the agent remains. The right action is optimal at \star under these incentives.
 119 Optimal policy sets capture the behavior incentivized by a reward function and a discount rate.

120 **Definition 4.1** (Optimal policy set function). $\Pi^*(R, \gamma)$ is the optimal policy set for reward function
 121 R at $\gamma \in (0, 1)$. All R have at least one optimal policy $\pi \in \Pi$ [Puterman, 2014]. $\Pi^*(R, 0) :=$
 122 $\lim_{\gamma \rightarrow 0} \Pi^*(R, \gamma)$ and $\Pi^*(R, 1) := \lim_{\gamma \rightarrow 1} \Pi^*(R, \gamma)$ exist by lemma D.32 (taking the limits with
 123 respect to the discrete topology over policy sets).

124 We may be unsure which reward function an agent will optimize. We may expect to deploy a system
 125 in a known environment, without knowing the exact form of e.g. the reward shaping [Ng et al., 1999]
 126 or intrinsic motivation [Pathak et al., 2017]. Alternatively, one might attempt to reason about future
 127 RL agents, whose details are unknown. Our power-seeking results do not hinge on such uncertainty,
 128 as they also apply to degenerate distributions (i.e. we know what reward function will be optimized).

129 **Definition 4.2** (Reward function distributions). Different results make different distributional assump-
 130 tions. Results with \mathcal{D}_{any} hold for any probability distribution over $\mathbb{R}^{|\mathcal{S}|}$. We sometimes assume a
 131 bounded distribution $\mathcal{D}_{\text{bound}}$ in order to ensure well-defined expectations. Letting X be any continuous
 132 bounded distribution over \mathbb{R} , $\mathcal{D}_{X\text{-ind}} := X^{|\mathcal{S}|}$. For example, when $X_u := \text{unif}(0, 1)$, $\mathcal{D}_{X_u\text{-ind}}$ is the
 133 maximum-entropy distribution. \mathcal{D}_s is the degenerate distribution on the state indicator reward function
 134 \mathbf{e}_s , which assigns 1 reward to s and 0 elsewhere.

135 With \mathcal{D}_{any} representing our prior beliefs about the agent’s reward function, what behavior should we
 136 expect from its optimal policies? Perhaps we want to reason about the probability that it’s optimal to
 137 go from \star to \emptyset or to r_{\triangleright} and then stay at r_{\triangleright} . In this case, we quantify the optimality probability of
 138 $F := \{\mathbf{e}_{\star} + \frac{\gamma}{1-\gamma}\mathbf{e}_{\emptyset}, \mathbf{e}_{\star} + \gamma\mathbf{e}_{r_{\triangleright}} + \frac{\gamma^2}{1-\gamma}\mathbf{e}_{r_{\triangleright}}\}$.

139 **Definition 4.3** (Visit distribution optimality probability). Let $F \subseteq \mathcal{F}(s)$, $\gamma \in [0, 1]$. $\mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma) :=$
 140 $\mathbb{P}_{R \sim \mathcal{D}_{\text{any}}}(\exists \mathbf{f}^{\pi} \in F : \pi \in \Pi^*(R, \gamma))$.

141 Alternatively, perhaps we’re interested in the probability that `right` is optimal at \star .

142 **Definition 4.4** (Action optimality probability). At discount rate γ and at state s , the *optimality*
 143 *probability of action a* is $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma) := \mathbb{P}_{R \sim \mathcal{D}_{\text{any}}}(\exists \pi^* \in \Pi^*(R, \gamma) : \pi^*(s) = a)$.

144 Optimality probability may seem hard to reason about. It’s hard enough to compute an optimal policy
 145 for a single reward function, let alone uncountably many! But consider any $\mathcal{D}_{X_{\text{-ind}}}$. When $\gamma = 0$,
 146 optimal policies maximize next-state reward. At \star , identically distributed reward means ℓ_{\triangleleft} and r_{\triangleright}
 147 have an equal probability of having maximal next-state reward. Therefore, $\mathbb{P}_{\mathcal{D}_{X_{\text{-ind}}}}(\star, \text{left}, 0) =$
 148 $\mathbb{P}_{\mathcal{D}_{X_{\text{-ind}}}}(\star, \text{right}, 0)$. This is not a proof, but such statements are provable.

149 With $\mathcal{D}_{\ell_{\triangleleft}}$ being the degenerate distribution on reward function $\mathbf{e}_{\ell_{\triangleleft}}$, $\mathbb{P}_{\mathcal{D}_{\ell_{\triangleleft}}}(\star, \text{left}, \frac{1}{2}) = 1 > 0 =$
 150 $\mathbb{P}_{\mathcal{D}_{\ell_{\triangleleft}}}(\star, \text{right}, \frac{1}{2})$. Similarly, $\mathbb{P}_{\mathcal{D}_{r_{\triangleright}}}(\star, \text{left}, \frac{1}{2}) = 0 < 1 = \mathbb{P}_{\mathcal{D}_{r_{\triangleright}}}(\star, \text{right}, \frac{1}{2})$. Therefore, “what
 151 do optimal policies ‘tend’ to look like?” seems to depend on one’s prior beliefs. But in fig. 1, we
 152 claimed that `left` is optimal for fewer reward functions than `right` is. The claim is meaningful and
 153 true, but we will return to it in section 6.

154 5 Some States Give The Agent More Control Over The Future

155 The agent has more “options” at ℓ_{\swarrow} than at the inescapable terminal state \emptyset . Furthermore, since r_{\nearrow}
 156 has a loop, the agent has more “options” at r_{\swarrow} than at ℓ_{\swarrow} . A glance at fig. 3 leads us to intuit that
 157 r_{\swarrow} affords the agent *more power* than \emptyset .

158 What is power? Philosophers have many answers. One prominent answer is the *dispositional view*:
 159 power is the ability to achieve a range of goals [Sattarov, 2019]. In an MDP, optimal value functions
 160 $V_R^*(s, \gamma)$ capture the agent’s ability to “achieve the goal” R . So *average* optimal value captures the
 161 agent’s ability to achieve a range of goals $\mathcal{D}_{\text{bound}}$.

162 **Definition 5.1** (Average optimal value). The *average optimal value* at state s and discount rate
 163 $\gamma \in (0, 1)$ is $V_{\mathcal{D}_{\text{bound}}}^*(s, \gamma) := \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}}[V_R^*(s, \gamma)] = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}}[\max_{\mathbf{f} \in \mathcal{F}(s)} \mathbf{f}(\gamma)^{\top} \mathbf{r}]$.

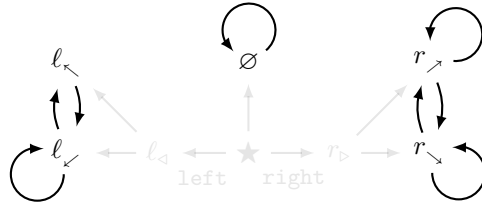


Figure 3: For $X_u := \text{unif}(0, 1)$, $V_{\mathcal{D}_{X_u\text{-ind}}}^*(\emptyset, \gamma) = \frac{1}{2} \frac{1}{1-\gamma}$, $V_{\mathcal{D}_{X_u\text{-ind}}}^*(\ell_{\swarrow}, \gamma) = \frac{1}{2} + \frac{\gamma}{1-\gamma^2} (\frac{2}{3} + \frac{1}{2}\gamma)$, and
 $V_{\mathcal{D}_{X_u\text{-ind}}}^*(r_{\swarrow}, \gamma) = \frac{1}{2} + \frac{\gamma}{1-\gamma} \frac{2}{3}$. $\frac{1}{2}$ and $\frac{2}{3}$ are the expected maxima of one and two draws from the uniform
 distribution, respectively. For all $\gamma \in (0, 1)$, $V_{\mathcal{D}_{X_u\text{-ind}}}^*(\emptyset, \gamma) < V_{\mathcal{D}_{X_u\text{-ind}}}^*(\ell_{\swarrow}, \gamma) < V_{\mathcal{D}_{X_u\text{-ind}}}^*(r_{\swarrow}, \gamma)$.
 $\text{POWER}_{\mathcal{D}_{X_u\text{-ind}}}(\emptyset, \gamma) = \frac{1}{2}$, $\text{POWER}_{\mathcal{D}_{X_u\text{-ind}}}(\ell_{\swarrow}, \gamma) = \frac{1}{1+\gamma} (\frac{2}{3} + \frac{1}{2}\gamma)$, and $\text{POWER}_{\mathcal{D}_{X_u\text{-ind}}}(r_{\swarrow}, \gamma) = \frac{2}{3}$. The
 POWER of ℓ_{\swarrow} is due to the agent only choosing the best of two states on every other time step.

164 Figure 3 shows the pleasing result that for the max-entropy distribution, r_{\swarrow} has greater average
 165 optimal value than \emptyset . However, $V_{\mathcal{D}_{\text{bound}}}^*(s, \gamma)$ has a few problems as a measure of power. All $\mathbf{f} \in \mathcal{F}(s)$
 166 count the agent’s initial presence at the initial state s , and $\|\mathbf{f}(\gamma)\|_1 = \frac{1}{1-\gamma}$ (proposition D.3) diverges
 167 as $\gamma \rightarrow 1$. Accordingly, $\lim_{\gamma \rightarrow 1} V_{\mathcal{D}_{\text{bound}}}^*(s, \gamma)$ tends to diverge. Definition 5.2 fixes these issues.

168 **Definition 5.2 (POWER).** Let $\gamma \in (0, 1)$.

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) := \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[\max_{\mathbf{f} \in \mathcal{F}(s)} \frac{1-\gamma}{\gamma} (\mathbf{f}(\gamma) - \mathbf{e}_s)^\top \mathbf{r} \right] = \frac{1-\gamma}{\gamma} \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [V_R^*(s, \gamma) - R(s)]. \quad (2)$$

169 POWER has nice formal properties.

170 **Lemma 5.3 (Continuity of POWER).** $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$ is Lipschitz continuous on $\gamma \in [0, 1]$.

171 **Proposition 5.4 (Maximal POWER).** $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \leq \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [\max_{s \in \mathcal{S}} R(s)]$, with equality
172 if s can deterministically reach all states in one step and all states are 1-cycles.

173 **Proposition 5.5 (POWER is smooth across reversible dynamics).** Let $\mathcal{D}_{\text{bound}}$ be bounded $[b, c]$. Suppose
174 s and s' can both reach each other in one step with probability 1.

$$|\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) - \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)| \leq (c - b)(1 - \gamma). \quad (3)$$

175 We now formalize what it means for actions to “seek power” in a situation.

176 **Definition 5.6 (POWER-seeking actions).** At state s and discount rate $\gamma \in [0, 1]$, action a seeks more
177 $\text{POWER}_{\mathcal{D}_{\text{bound}}}$ than a' when $\mathbb{E}_{s_a \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_a, \gamma)] \geq \mathbb{E}_{s_{a'} \sim T(s, a')} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_{a'}, \gamma)]$.

178 POWER is sensitive to the choice of distribution. $\mathcal{D}_{\ell_{\swarrow}}$ gives maximal $\text{POWER}_{\mathcal{D}_{\ell_{\swarrow}}}$ to ℓ_{\swarrow} . $\mathcal{D}_{r_{\searrow}}$ assigns
179 maximal $\text{POWER}_{\mathcal{D}_{r_{\searrow}}}$ to r_{\searrow} . \mathcal{D}_{\emptyset} even gives maximal $\text{POWER}_{\mathcal{D}_{\emptyset}}$ to \emptyset ! In what sense does \emptyset have “less
180 POWER” than r_{\searrow} , and in what sense does going right “tend to seek POWER” compared to left?

181 6 Certain Environmental Symmetries Produce Power-Seeking Tendencies

182 We prove that for all $\gamma \in [0, 1]$ and for most distributions \mathcal{D} , $\text{POWER}_{\mathcal{D}}(\ell_{\swarrow}, \gamma) \leq \text{POWER}_{\mathcal{D}}(r_{\searrow}, \gamma)$.
183 But first, we explore why this must be true.

184 $\mathcal{F}(\ell_{\swarrow}) = \{\frac{1}{1-\gamma} \mathbf{e}_{\ell_{\swarrow}}, \frac{1}{1-\gamma^2} (\mathbf{e}_{\ell_{\swarrow}} + \gamma \mathbf{e}_{\ell_{\nwarrow}})\}$ and $\mathcal{F}(r_{\searrow}) = \{\frac{1}{1-\gamma} \mathbf{e}_{r_{\searrow}}, \frac{1}{1-\gamma^2} (\mathbf{e}_{r_{\searrow}} + \gamma \mathbf{e}_{r_{\nearrow}}), \frac{1}{1-\gamma} \mathbf{e}_{r_{\nearrow}}\}$.
185 These two sets look awfully similar. $\mathcal{F}(\ell_{\swarrow})$ is a “subset” of $\mathcal{F}(r_{\searrow})$, only with “different states.”
186 Figure 4 demonstrates a state permutation ϕ which embeds $\mathcal{F}(\ell_{\swarrow})$ into $\mathcal{F}(r_{\searrow})$.

188 **Definition 6.1 (Similarity of visitation distribution
189 sets).** Let $F, F' \subseteq \mathbb{R}^{|S|}$. Given state permutation
190 $\phi \in S_{|S|}$ inducing permutation matrix $\mathbf{P}_\phi \in$
191 $\mathbb{R}^{|S| \times |S|}$, $\phi(F) := \{\mathbf{P}_\phi \mathbf{f} \mid \mathbf{f} \in F\}$. F is similar to
192 F' when $\exists \phi : \phi(F) = F'$. ϕ is an involution if
193 $\phi = \phi^{-1}$ – it either swaps states, or leaves them be.

194 Let $I \subseteq \mathbb{R}$. If F, F' instead contain vector-valued
195 functions $I \mapsto \mathbb{R}^{|S|}$, F is similar to F' when $\exists \phi :$
196 $\forall \gamma \in I : \{\mathbf{P}_\phi \mathbf{f}(\gamma) \mid \mathbf{f} \in F\} = \{\mathbf{f}'(\gamma) \mid \mathbf{f}' \in F'\}$.

197 Consider a reward function R' which assigns
198 $R'(\ell_{\swarrow}) = R'(\ell_{\nwarrow}) = 1$, $R'(r_{\searrow}) = R'(r_{\nearrow}) = 0$.
199 R' assigns more optimal value to ℓ_{\swarrow} than to r_{\searrow} :
200 $V_{R'}^*(\ell_{\swarrow}, \gamma) = \frac{1}{1-\gamma} > 0 = V_{R'}^*(r_{\searrow}, \gamma)$. However,
201 consider $\phi \cdot R'$: $(\phi \cdot R')(\ell_{\swarrow}) := R'(\phi(\ell_{\swarrow})) = R'(r_{\searrow}) = 0$. This permuted reward function switches
202 the optimal value functions: $V_{\phi \cdot R'}^*(\ell_{\swarrow}, \gamma) = 0 < \frac{1}{1-\gamma} = V_{\phi \cdot R'}^*(r_{\searrow}, \gamma)$.

203 Remarkably, this ϕ has the property that for any R which assigns ℓ_{\swarrow} greater optimal value than r_{\searrow} (i.e.
204 $V_R^*(\ell_{\swarrow}, \gamma) > V_R^*(r_{\searrow}, \gamma)$), the opposite holds for the permuted $\phi \cdot R$: $V_{\phi \cdot R}^*(\ell_{\swarrow}, \gamma) < V_{\phi \cdot R}^*(r_{\searrow}, \gamma)$.

205 We can permute reward functions, but we can also permute reward function distributions. Permuted
206 distributions simply permute which states get which rewards.

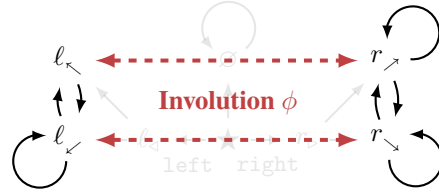


Figure 4: $\mathcal{F}_{\text{nd}}(\ell_{\swarrow})$ is similar to a strict subset of $\mathcal{F}(r_{\searrow})$ via involution ϕ :

$$\begin{aligned} \phi(\mathcal{F}_{\text{nd}}(\ell_{\swarrow})) &:= \left\{ \frac{1}{1-\gamma} \mathbf{P}_\phi \mathbf{e}_{\ell_{\swarrow}}, \frac{1}{1-\gamma^2} (\mathbf{P}_\phi \mathbf{e}_{\ell_{\swarrow}} + \gamma \mathbf{P}_\phi \mathbf{e}_{\ell_{\nwarrow}}) \right\} \\ &= \left\{ \frac{1}{1-\gamma} \mathbf{e}_{r_{\searrow}}, \frac{1}{1-\gamma^2} (\mathbf{e}_{r_{\searrow}} + \gamma \mathbf{e}_{r_{\nearrow}}) \right\} \subsetneq \mathcal{F}(r_{\searrow}). \end{aligned}$$

207

208 **Definition 6.2** (Pushforward distribution of a permutation).
 209 Let $\phi \in S_{|S|}$. $\phi(\mathcal{D}_{\text{any}})$ is the pushforward distribution induced
 210 by applying the random vector $f(\mathbf{r}) := \mathbf{P}_\phi \mathbf{r}$ to \mathcal{D}_{any} .

211 **Definition 6.3** (Orbit of a probability distribution). The orbit
 212 of \mathcal{D}_{any} under the symmetric group $S_{|S|}$ is $S_{|S|} \cdot \mathcal{D}_{\text{any}} :=$
 213 $\{\phi(\mathcal{D}_{\text{any}}) \mid \phi \in S_{|S|}\}$.

214 For example, the orbit of a degenerate state indicator distri-
 215 bution \mathcal{D}_s is $S_{|S|} \cdot \mathcal{D}_s = \{\mathcal{D}_{s'} \mid s' \in S\}$, and fig. 5 shows
 216 the orbit of a 2D Gaussian distribution.

217 Considering again the involution ϕ of fig. 4: for every $\mathcal{D}_{\text{bound}}$
 218 for which ℓ_{\swarrow} has more $\text{POWER}_{\mathcal{D}_{\text{bound}}}$ than r_{\searrow} , ℓ_{\swarrow} has less
 219 $\text{POWER}_{\phi(\mathcal{D}_{\text{bound}})}$ than r_{\searrow} . This fact is not obvious – it is shown
 220 by the proof of lemma D.22.

221 Imagine $\mathcal{D}_{\text{bound}}$'s orbit elements “voting” whether ℓ_{\swarrow} or r_{\searrow}
 222 has strictly more POWER . Proposition 6.5 will show that r_{\searrow} can't lose the “vote” for the orbit of any
 223 bounded reward function distribution. Definition 6.4 formalizes this “voting” notion.

224 **Definition 6.4** (Inequalities which hold for most bounded probability distributions). Let $f_1, f_2 :$
 225 $\Delta(\mathbb{R}^{|S|}) \rightarrow \mathbb{R}$ be any functions from reward function distributions to real numbers. We write
 226 $f_1 \geq_{\text{most}} f_2$ when, for all $\mathcal{D}_{\text{bound}}$,¹ the following cardinality inequality holds:

$$\left| \{ \mathcal{D} \in S_{|S|} \cdot \mathcal{D}_{\text{bound}} \mid f_1(\mathcal{D}) > f_2(\mathcal{D}) \} \right| \geq \left| \{ \mathcal{D} \in S_{|S|} \cdot \mathcal{D}_{\text{bound}} \mid f_1(\mathcal{D}) < f_2(\mathcal{D}) \} \right|. \quad (4)$$

227 **Proposition 6.5** (States with “more options” generally have more POWER). Suppose $\mathcal{F}_{\text{nd}}(s')$ is similar
 228 to a subset of $\mathcal{F}(s)$ via involution ϕ . Then $\forall \gamma \in [0, 1] : \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \leq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$.

229 If $\mathcal{F}_{\text{nd}}(s) \setminus \phi(\mathcal{F}_{\text{nd}}(s'))$ is non-empty, then for all $\gamma \in (0, 1)$, the inequality is strict for all $\mathcal{D}_{X\text{-IID}}$ and
 230 $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \not\leq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$ – i.e. $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \geq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$ does
 231 not hold.

232 Proposition 6.5 proves that for all $\gamma \in [0, 1]$, $\text{POWER}_{\mathcal{D}_{\text{bound}}}(\ell_{\swarrow}, \gamma) \leq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(r_{\searrow}, \gamma)$ via
 233 $s' := \ell_{\swarrow}$, $s := r_{\searrow}$, and the involution ϕ shown in fig. 4. In fact, because $(\frac{1}{1-\gamma} \mathbf{e}_{r_{\swarrow}}) \in \mathcal{F}_{\text{nd}}(r_{\searrow}) \setminus$
 234 $\phi(\mathcal{F}_{\text{nd}}(\ell_{\swarrow}))$, r_{\searrow} has “strictly more options” and therefore fulfills proposition 6.5's stronger condition.

235 Proposition 6.5 is shown using the fact that ϕ injectively maps \mathcal{D} under which r_{\searrow} has less $\text{POWER}_{\mathcal{D}}$,
 236 to distributions $\phi(\mathcal{D})$ which agree with the intuition that r_{\searrow} offers more control. Therefore, at least
 237 half of each orbit must agree, and r_{\searrow} never “loses the POWER vote” against ℓ_{\swarrow} .²

238 6.1 Keeping options open tends to be POWER -seeking and tends to be optimal

239 Certain symmetries in the MDP structure ensure that, compared to `left`, going `right` tends to be
 240 optimal and to be POWER -seeking. Intuitively, by going `right`, the agent has “strictly more choices.”
 241 Proposition 6.9 formalizes this tendency, but we need several technical concepts to state the result.

242 **Definition 6.6** (Equivalent actions). Actions a_1 and a_2 are *equivalent at state s* (written $a_1 \equiv_s a_2$) if
 243 they induce the same transition probabilities: $T(s, a_1) = T(s, a_2)$.

244 The agent can only reach states in $\{r_{\triangleright}, r_{\swarrow}, r_{\searrow}\}$ by taking actions equivalent to `right` at state \star .

245 **Definition 6.7** (State-space bottleneck). Starting from s , state s' is a *bottleneck* for $S \subseteq \mathcal{S}$ via action
 246 a when state s can reach the states of S with positive probability, but only by taking actions equivalent
 247 to a at state s' . We write this as $s \rightarrow s' \xrightarrow{a} S$.

¹Boundedness only ensures that POWER is well-defined. Optimality probability results hold for all \mathcal{D}_{any} orbits.

²Proposition 6.5 also proves that in general, \emptyset has less POWER than ℓ_{\swarrow} and r_{\searrow} . However, this does not prove that most distributions \mathcal{D} satisfy the joint inequality $\text{POWER}_{\mathcal{D}}(\emptyset, \gamma) \leq \text{POWER}_{\mathcal{D}}(\ell_{\swarrow}, \gamma) \leq \text{POWER}_{\mathcal{D}}(r_{\searrow}, \gamma)$ – only that these inequalities hold pairwise for most \mathcal{D} . The orbit elements \mathcal{D} which agree that \emptyset has less $\text{POWER}_{\mathcal{D}}$ than ℓ_{\swarrow} need not be the same elements \mathcal{D}' which agree that ℓ_{\swarrow} has less $\text{POWER}_{\mathcal{D}'}$ than r_{\searrow} .

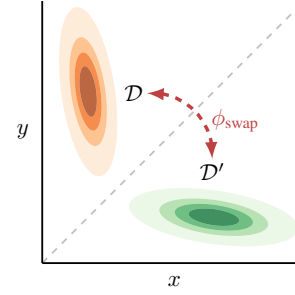


Figure 5: Probability density plots of the Gaussian distributions \mathcal{D} and \mathcal{D}' over \mathbb{R}^2 . The symmetric group S_2 contains the identity permutation ϕ_{id} and the reflection permutation ϕ_{swap} (switching the y and x values). The orbit of \mathcal{D} consists of $\phi_{\text{id}}(\mathcal{D}) = \mathcal{D}$ and $\phi_{\text{swap}}(\mathcal{D}) = \mathcal{D}'$.

248 **Definition 6.8** (States reachable after taking an action). $\text{REACH}(s', a)$ is the set of states reachable
 249 with positive probability after taking the action a in state s' .

250 **Proposition 6.9** (Keeping options open tends to be POWER-seeking and tends to be optimal).

251 Suppose that $s \rightarrow s' \xrightarrow{a} \text{REACH}(s', a)$ and $s \rightarrow s' \xrightarrow{a'} \text{REACH}(s', a')$. $F_{a'} := \mathcal{F}(s \mid \pi(s') =$
 252 $a')$, $F_a := \mathcal{F}(s \mid \pi(s') = a)$. Suppose $F_{a'}$ is similar to a subset of F_a via involution ϕ which fixes all
 253 states not belonging to $\text{REACH}(s', a')$ or $\text{REACH}(s', a)$.

254 Then for all $\gamma \in [0, 1]$, $\mathbb{E}_{s_{a'} \sim T(s', a')} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_{a'}, \gamma)] \leq_{\text{most}} \mathbb{E}_{s_a \sim T(s', a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_a, \gamma)]$
 255 and $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(F_{a'}, \gamma) \leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(F_a, \gamma)$.

256 If $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi(F_{a'}))$ is non-empty, then for all $\gamma \in (0, 1)$, both inequalities are strict
 257 for all $\mathcal{D}_{X\text{-IND}}$, $\mathbb{E}_{s_{a'} \sim T(s', a')} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_{a'}, \gamma)] \not\leq_{\text{most}} \mathbb{E}_{s_a \sim T(s', a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_a, \gamma)]$, and
 258 $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(F_{a'}, \gamma) \not\leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(F_a, \gamma)$.

259 We check the conditions of proposition 6.9. $s = s' := \star$,
 260 $a' := \text{left}$, $a := \text{right}$. Figure 6 shows that $\star \rightarrow$
 261 $\star \xrightarrow{\text{left}} \{\ell_{\triangleleft}, \ell_{\times}, \ell_{\triangleright}\}$ and $\star \rightarrow \star \xrightarrow{\text{right}} \{r_{\triangleright}, r_{\times}, r_{\triangleleft}\}$.
 262 $\mathcal{F}(\star \mid \pi(\star) = \text{left})$ is similar to a strict subset of
 263 $\mathcal{F}(\star \mid \pi(\star) = \text{right})$ via permutation ϕ , which fixes
 264 \star and \emptyset . Furthermore, $\mathcal{F}_{\text{nd}}(\star) \cap \{\mathbf{e}_{\star} + \gamma \mathbf{e}_{r_{\triangleright}} + \gamma^2 \mathbf{e}_{r_{\triangleleft}} +$
 265 $\frac{\gamma^3}{1-\gamma} \mathbf{e}_{r_{\times}}, \mathbf{e}_{\star} + \gamma \mathbf{e}_{r_{\triangleright}} + \frac{\gamma^2}{1-\gamma} \mathbf{e}_{r_{\times}}\} = \{\mathbf{e}_{\star} + \gamma \mathbf{e}_{r_{\triangleright}} + \frac{\gamma^2}{1-\gamma} \mathbf{e}_{r_{\times}}\}$
 266 is non-empty, and so all conditions are met.

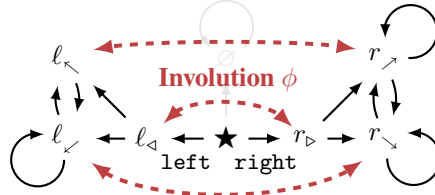


Figure 6

267 For any $\gamma \in [0, 1]$ and \mathcal{D} such that $\mathbb{P}_{\mathcal{D}}(\star, \text{left}, \gamma) > \mathbb{P}_{\mathcal{D}}(\star, \text{right}, \gamma)$, environmental symmetry
 268 ensures that $\mathbb{P}_{\phi(\mathcal{D})}(\star, \text{left}, \gamma) < \mathbb{P}_{\phi(\mathcal{D})}(\star, \text{right}, \gamma)$. A similar statement holds for POWER.

269 6.2 When $\gamma \approx 1$, optimal policies tend to navigate towards “larger” sets of cycles

270 Proposition 6.5 and proposition 6.9 are powerful because they apply to all $\gamma \in [0, 1]$, but they can
 271 only be applied given hard-to-satisfy environmental symmetries. In contrast, proposition 6.12 and
 272 theorem 6.13 apply to many finite structured environments common to RL.

273 Starting from \star , consider the cycles which the agent can reach. Recurrent state distributions (RSDs)
 274 generalize deterministic graphical cycles to potentially stochastic environments. RSDs simply record
 275 how often the agent tends to visit a state in the limit of infinitely many time steps.

276 **Definition 6.10** (Recurrent state distributions [Puterman, 2014]). The recurrent state distributions
 277 which can be induced from state s are $\text{RSD}(s) := \{\lim_{\gamma \rightarrow 1} (1 - \gamma) \mathbf{f}(\gamma) \mid \mathbf{f} \in \mathcal{F}(s)\}$. $\text{RSD}_{\text{nd}}(s)$ is
 278 the set of RSDs which strictly maximize average reward for some reward function.

279 As suggested by fig. 3, $\text{RSD}(\star) = \{\mathbf{e}_{\ell_{\triangleright}}, \frac{1}{2}(\mathbf{e}_{\ell_{\triangleright}} + \mathbf{e}_{\ell_{\times}}), \mathbf{e}_{\emptyset}, \mathbf{e}_{r_{\triangleright}}, \frac{1}{2}(\mathbf{e}_{r_{\triangleright}} + \mathbf{e}_{r_{\triangleleft}}), \mathbf{e}_{r_{\triangleleft}}\}$. As discussed
 280 in section 3, $\frac{1}{2}(\mathbf{e}_{r_{\triangleright}} + \mathbf{e}_{r_{\triangleleft}})$ is dominated: alternating between r_{\triangleright} and r_{\triangleleft} is never strictly better than
 281 choosing one or the other.

282 A reward function’s optimal policies can vary with the discount rate. When $\gamma \approx 1$, optimal policies
 283 “ignore” transient reward because average reward is the dominant consideration.

284 **Definition 6.11** (Blackwell optimal policies [Blackwell, 1962]). $\Pi^*(R, 1) := \lim_{\gamma \rightarrow 1} \Pi^*(R, \gamma)$ is
 285 the Blackwell optimal policy set for reward function R .

286 Blackwell optimal policies maximize average reward. Average reward is governed by RSD access. For
 287 example, r_{\triangleleft} has “more” RSDs than \emptyset ; therefore, it usually has greater Power when $\gamma = 1$.

288 **Proposition 6.12** (When $\gamma = 1$, RSDs control Power). If $\text{RSD}_{\text{nd}}(s')$ is similar to a subset of $\text{RSD}(s)$
 289 via involution ϕ , then $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', 1) \leq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1)$. If $\text{RSD}_{\text{nd}}(s) \setminus \phi(\text{RSD}_{\text{nd}}(s'))$ is
 290 non-empty, then this inequality is strict for all $\mathcal{D}_{X\text{-IND}}$, and $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', 1) \not\leq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1)$.

291 We check that both conditions of proposition 6.12 are satisfied when $s' := \emptyset$, $s := r_{\triangleleft}$, and the
 292 involution ϕ swaps \emptyset and r_{\triangleleft} . $\phi(\text{RSD}_{\text{nd}}(\emptyset)) = \phi(\{\mathbf{e}_{\emptyset}\}) = \{\mathbf{e}_{r_{\triangleleft}}\} \subsetneq \{\mathbf{e}_{r_{\triangleleft}}, \mathbf{e}_{r_{\triangleright}}\} = \text{RSD}_{\text{nd}}(r_{\triangleleft})$.
 293 The conditions are satisfied.

294 Informally, states with more RSDs generally have more POWER at $\gamma = 1$, no matter their transient
 295 dynamics. Furthermore, Blackwell optimal policies are more likely to end up in “larger” sets of RSDs
 296 than in “smaller” ones. Thus, Blackwell optimal policies tend to navigate towards parts of the state
 297 space which contain more RSDs.

298

299 **Theorem 6.13** (Blackwell optimal policies tend to end
 300 up in “larger” sets of RSDs). *Let $D', D \subseteq \text{RSD}(s)$. Suppose
 301 that D' is similar to a subset of D via involution
 302 ϕ and that the sets $D' \cup D$ and $\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)$
 303 have pairwise orthogonal vector elements (i.e. pairwise
 304 disjoint vector support).*

305 *Then $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(D', 1) \leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(D, 1)$. If $\text{RSD}_{\text{nd}}(s) \cap$
 306 $(D \setminus \phi(D'))$ is non-empty, the inequality is strict for all
 307 $\mathcal{D}_{X\text{-IID}}$, and $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(D', 1) \not\leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(D, 1)$.*

308 **Corollary 6.14** (Blackwell optimal policies tend not to end up in any given 1-cycle).

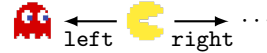
309 *Suppose $e_{s_x}, e_{s'} \in \text{RSD}(s)$ are distinct. Then $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(\{e_{s_x}\}, 1) \leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(\text{RSD}(s) \setminus \{e_{s_x}\}, 1)$.
 310 If there exists a third $e_{s''} \in \text{RSD}(s)$, then $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(\{e_{s_x}\}, 1) \not\leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(\text{RSD}(s) \setminus \{e_{s_x}\}, 1)$.*

311 Figure 7 illustrates that $e_{\emptyset}, e_{r_{\leftarrow}}, e_{r_{\rightarrow}} \in \text{RSD}(\star)$. Thus, both conclusions of corol-
 312 lary 6.14 hold: $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(\{e_{\emptyset}\}, 1) \leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{bound}}}(\text{RSD}(\star) \setminus \{e_{\emptyset}\}, 1)$ and $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(\{e_{\emptyset}\}, 1) \not\leq_{\text{most}}$
 313 $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(\text{RSD}(\star) \setminus \{e_{\emptyset}\}, 1)$. In other words, Blackwell optimal policies tend to end up in RSDs be-
 314 sides \emptyset . Since \emptyset is a terminal state, it cannot reach other RSDs. Since Blackwell optimal policies
 315 tend to end up in other RSDs, Blackwell optimal policies tend to avoid \emptyset .

316 This section’s results prove the $\gamma = 1$ case. Lemma 5.3 and proposition D.33 respectively show that
 317 POWER and optimality probability are continuous at $\gamma = 1$. Therefore, if strict POWER-seeking or
 318 strictly greater optimality probability holds when $\gamma = 1$, the same holds for discount rates sufficiently
 319 close to 1.

320 Lastly, since \geq_{most} (definition 6.4) results apply to all $\mathcal{D}_{\text{bound}}$, our key results apply to all degenerate
 321 distributions. Therefore, our key results apply not just to distributions over reward functions, but to
 322 individual reward functions.

323



324 6.3 How to reason about other environments

325 Consider an embodied navigation task through a
 326 room with a vase. Proposition 6.9 suggests that opti-
 327 mal policies tend to avoid breaking the vase, since
 328 doing so would strictly decrease available options.

329 Theorem 6.13 dictates where Blackwell optimal
 330 agents tend to end up, but not what actions they
 331 tend to take in order to reach their RSDs. Therefore, care is needed. In appendix B, fig. 10 demonstrates
 332 an environment in which seeking POWER is a “detour” for most reward functions.

333 Even though randomly generated MDPs are unlikely to satisfy our sufficient conditions for POWER-
 334 seeking tendencies, theorem 6.13 is easy to apply to many structured RL environments. Loosely
 335 speaking, if the “vast majority” of RSDs are only reachable by following a subset of policies, theo-
 336 rem 6.13 implies that that subset tends to be Blackwell optimal.

337 In fig. 8, the player dies by going left, but can reach thousands of RSDs by heading in other directions.
 338 Even if some Blackwell optimal policies go left in order to reach fig. 8’s ‘game over’ terminal state,
 339 all other RSDs cannot be reached by going left. There are many 1-cycles besides the immediate
 340 terminal state. Therefore, corollary 6.14 proves that Blackwell optimal policies tend to not go left in
 341 this situation. Blackwell optimal policies tend to avoid immediately dying in Pac-Man, even though
 342 most reward functions do not resemble Pac-Man’s original score function.

343 Optimal policies for e.g. Pac-Man do not seek power in the real world. However, we think that our
 344 results suggest that optimal policies for real-world environments tend to seek real-world power.

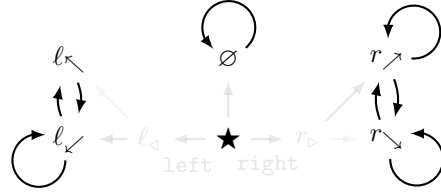


Figure 7: The cycles in $\text{RSD}(\star)$.

345 7 Discussion

346 Corollary 6.14 dictates where Blackwell optimal agents tend to end up, but not how they get there.
347 Corollary 6.14 says that such agents tend not to *stay* in any given 1-cycle. It does not say that such
348 agents will avoid *entering* such states. For example, in an embodied navigation task, a robot may
349 enter a 1-cycle by idling in the center of a room. Corollary 6.14 implies that Blackwell optimal robots
350 tend not to idle in that particular spot, but not that they tend to avoid that spot entirely.

351 However, Blackwell optimal robots do tend to avoid getting shut down. A terminal state is, by
352 definition 3.2, unable to access other 1-cycles. Since corollary 6.14 shows that Blackwell optimal
353 agents tend to end up in other 1-cycles, Blackwell optimal policies must tend to completely avoid the
354 terminal state. The agent’s task MDP often represents agent shutdown with terminal states. Therefore,
355 we conclude that in many such situations, Blackwell optimal policies tend to avoid shutdown.

356 Reconsider the case of a hypothetical intelligent real-world agent. Suppose the designers initially have
357 control over a Blackwell optimal agent. If the agent began to misbehave, they could just deactivate
358 it. Unfortunately, our results suggest that this strategy might not work. Blackwell optimal agents
359 would generally stop us from deactivating them, if physically possible. Extrapolating from our results,
360 we conjecture that Blackwell optimal policies tend to seek power by accumulating resources, to the
361 detriment of any other agents in the environment.

362 **Future work.** Most real-world tasks are partially observable. Although our results only apply to
363 optimal policies in finite MDPs, we expect the key conclusions to generalize. We look forward to
364 future work which addresses partially observable environments or suboptimal policies.

365 Past work shows that it would be bad for an agent to disempower humans in its environment. In a
366 two-player agent / human game, minimizing the human’s information-theoretic empowerment [Salge
367 et al., 2014] produces adversarial agent behavior [Guckelsberger et al., 2018]. In contrast, maximizing
368 human empowerment produces helpful agent behavior [Salge and Polani, 2017, Guckelsberger et al.,
369 2016, Du et al., 2020]. We do not yet formally understand if, when, or why POWER-seeking policies
370 tend to disempower other agents in the environment.

371 More complex environments probably have more pronounced power-seeking incentives. Intuitively,
372 there are often combinatorially many ways for power-seeking to be optimal, and relatively few ways for
373 power-seeking not to be optimal. For example, suppose that in some environment, theorem 6.13 holds
374 for one million involutions ϕ . Does this guarantee more pronounced incentives than if theorem 6.13
375 only held for one involution?

376 We proved sufficient conditions for when reward functions tend to incentivize power-seeking. In
377 the absence of prior information, one should expect that an arbitrary reward function distribution
378 incentivizes power-seeking behavior under these conditions. However, we have prior information:
379 AI designers usually try to specify a good reward function. Even so, it may be hard to specify orbit
380 elements which do not incentivize bad power-seeking.

381 **Societal impact.** We believe that this paper builds toward a rigorous understanding of the risks
382 presented by AI power-seeking incentives. Understanding these risks is the first step in addressing
383 them. However, basic theoretical work can have many consequences. For example, this theory could
384 somehow help future researchers build power-seeking agents which disempower other agents in their
385 environment. We believe that the benefit of understanding outweighs the potential societal harm.

386 **Conclusion.** We developed the first formal theory of the statistical tendencies of optimal policies in
387 reinforcement learning. In the context of MDPs, we proved sufficient conditions under which optimal
388 policies tend to seek power, both formally (by taking POWER-seeking actions) and intuitively (by
389 taking actions which keep the agent’s options open). Many real-world environments have symmetries
390 which produce power-seeking incentives. In particular, optimal policies tend to seek power when the
391 agent can be shut down or destroyed. Maximizing control over the environment will often involve
392 resisting shutdown, and perhaps monopolizing resources.

393 We caution that many real-world tasks are partially observable and that learned policies are rarely
394 optimal. Our results do not mathematically *prove* that hypothetical superintelligent AI agents will
395 seek power. However, we hope that this work will foster thoughtful, serious, and rigorous discussion
396 of this possibility.

397 **References**

- 398 Tsvi Benson-Tilsen and Nate Soares. Formalizing convergent instrumental goals. *Workshops at the Thirtieth*
399 *AAAI Conference on Artificial Intelligence*, 2016.
- 400 David Blackwell. Discrete dynamic programming. *The Annals of Mathematical Statistics*, page 9, 1962.
- 401 Nick Bostrom. *Superintelligence*. Oxford University Press, 2014.
- 402 Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A. Efros. Large-scale
403 study of curiosity-driven learning. In *International Conference on Learning Representations*, 2019.
- 404 Ryan Carey. In corrigibility in the CIRC framework. *AI, Ethics, and Society*, 2018.
- 405 Chris Drummond. Composing functions to speed up reinforcement learning in a changing world. In *Machine*
406 *Learning: ECML-98*, volume 1398, pages 370–381. Springer, 1998.
- 407 Yuqing Du, Stas Tiomkin, Emre Kiciman, Daniel Polani, Pieter Abbeel, and Anca Dragan. AvE: Assistance via
408 empowerment. *Advances in Neural Information Processing Systems*, 33, 2020.
- 409 David Foster and Peter Dayan. Structure in the space of value functions. *Machine Learning*, pages 325–346,
410 2002.
- 411 Christian Guckelsberger, Christoph Salge, and Simon Colton. Intrinsically motivated general companion NPCs
412 via coupled empowerment maximisation. In *IEEE Conference on Computational Intelligence and Games*,
413 pages 1–8, 2016.
- 414 Christian Guckelsberger, Christoph Salge, and Julian Togelius. New and surprising ways to be mean. In *IEEE*
415 *Conference on Computational Intelligence and Games*, pages 1–8, 2018.
- 416 Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart Russell. The off-switch game. In *Proceedings*
417 *of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 220–227, 2017.
- 418 Tejas D Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J Gershman. Deep successor reinforcement
419 learning. *arXiv preprint arXiv:1606.02396*, 2016.
- 420 Yann LeCun and Anthony Zador. Don't fear the Terminator, September 2019. URL [https://blogs.](https://blogs.scientificamerican.com/observations/dont-fear-the-terminator/)
421 [scientificamerican.com/observations/dont-fear-the-terminator/](https://blogs.scientificamerican.com/observations/dont-fear-the-terminator/).
- 422 Steven A Lippman. On the set of optimal policies in discrete dynamic programming. *Journal of Mathematical*
423 *Analysis and Applications*, 24(2):440–445, 1968.
- 424 Lionel W McKenzie. Turnpike theory. *Econometrica: Journal of the Econometric Society*, pages 841–865, 1976.
- 425 Ishai Menache, Shie Mannor, and Nahum Shimkin. Q-cut—dynamic discovery of sub-goals in reinforcement
426 learning. In *European Conference on Machine Learning*, pages 295–306. Springer, 2002.
- 427 Smitha Milli, Dylan Hadfield-Menell, Anca Dragan, and Stuart Russell. Should robots be obedient? In
428 *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 4754–4760, 2017.
- 429 Melanie Mitchell. Why AI is harder than we think. *arXiv preprint arXiv:2104.12871*, 2021.
- 430 Andrew Y. Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and
431 application to reward shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning*,
432 pages 278–287. Morgan Kaufmann, 1999.
- 433 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An
434 introduction to circuits. *Distill*, 2020.
- 435 Stephen Omohundro. The basic AI drives, 2008.
- 436 Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-
437 supervised prediction. In *ICML*, 2017.
- 438 Steven Pinker and Stuart Russell. The foundations, benefits, and possible existential threat of AI, June
439 2020. URL [https://futureoflife.org/2020/06/15/steven-pinker-and-stuart-russell-on-](https://futureoflife.org/2020/06/15/steven-pinker-and-stuart-russell-on-the-foundations-benefits-and-possible-existential-risk-of-ai/)
440 [the-foundations-benefits-and-possible-existential-risk-of-ai/](https://futureoflife.org/2020/06/15/steven-pinker-and-stuart-russell-on-the-foundations-benefits-and-possible-existential-risk-of-ai/).
- 441 Martin L Puterman. *Markov decision processes: Discrete stochastic dynamic programming*. John Wiley & Sons,
442 2014.

- 443 J.B. Reece and N.A. Campbell. *Campbell Biology*. Pearson Australia, 2011.
- 444 Kevin Regan and Craig Boutilier. Robust policy computation in reward-uncertain MDPs using nondominated
445 policies. In *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.
- 446 Stuart Russell. *Human compatible: Artificial intelligence and the problem of control*. Viking, 2019.
- 447 Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Pearson Education Limited, 2009.
- 448 Christoph Salge and Daniel Polani. Empowerment as replacement for the three laws of robotics. *Frontiers in*
449 *Robotics and AI*, 4:25, 2017.
- 450 Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment—an introduction. In *Guided Self-*
451 *Organization: Inception*, pages 67–114. Springer, 2014.
- 452 Faridun Sattarov. *Power and technology: a philosophical and ethical analysis*. Rowman & Littlefield International,
453 Ltd, 2019.
- 454 Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In
455 *International Conference on Machine Learning*, pages 1312–1320, 2015.
- 456 Nate Soares, Benja Fallenstein, Stuart Armstrong, and Eliezer Yudkowsky. Corrigibility. *AAAI Workshops*, 2015.
- 457 Richard S Sutton and Andrew G Barto. *Reinforcement learning: an introduction*. MIT Press, 1998.
- 458 Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, Adam White, and Doina
459 Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor
460 interaction. In *International Conference on Autonomous Agents and Multiagent Systems*, pages 761–768,
461 2011.
- 462 Alexander Matt Turner, Dylan Hadfield-Menell, and Prasad Tadepalli. Conservative agency via attainable utility
463 preservation. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 385–391, 2020.
- 464 Various. Debate on instrumental convergence between LeCun, Russell, Bengio, Zador, and more, 2019.
465 URL [https://www.alignmentforum.org/posts/Wxw6Gc6f2z3mzmqKs/debate-on-instrumental-](https://www.alignmentforum.org/posts/Wxw6Gc6f2z3mzmqKs/debate-on-instrumental-convergence-between-lecun-russell)
466 [convergence-between-lecun-russell](https://www.alignmentforum.org/posts/Wxw6Gc6f2z3mzmqKs/debate-on-instrumental-convergence-between-lecun-russell).
- 467 Tao Wang, Michael Bowling, and Dale Schuurmans. Dual representations for dynamic programming and rein-
468 forcement learning. In *International Symposium on Approximate Dynamic Programming and Reinforcement*
469 *Learning*, pages 44–51. IEEE, 2007.
- 470 Tao Wang, Michael Bowling, Dale Schuurmans, and Daniel J Lizotte. Stable dual dynamic programming. In
471 *Advances in Neural Information Processing Systems*, pages 1569–1576, 2008.

472 **Checklist**

- 473 1. For all authors...
- 474 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
475 contributions and scope? [Yes]
- 476 (b) Did you describe the limitations of your work? [Yes] See section 7.
- 477 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Section 7:
478 “We believe that this paper builds toward a rigorous understanding of the risks presented
479 by AI power-seeking incentives. Understanding these risks is the first step in preventing
480 them. However, basic theoretical work can have many consequences. For example,
481 this theory could somehow help future researchers build power-seeking agents which
482 disempower other agents in their environment. Nonetheless, we believe that the benefit
483 of understanding outweighs the potential societal harm.”
- 484 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
485 them? [Yes]
- 486 2. If you are including theoretical results...
- 487 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See appendix
488 D.
- 489 (b) Did you include complete proofs of all theoretical results? [Yes] See appendix D.
- 490 3. If you ran experiments...
- 491 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
492 mental results (either in the supplemental material or as a URL)? [N/A]
- 493 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
494 were chosen)? [N/A]
- 495 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
496 multiple times)? [N/A]
- 497 (d) Did you include the total amount of compute and the type of resources used (e.g., type
498 of GPUs, internal cluster, or cloud provider)? [N/A]
- 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? [N/A]
- 501 (b) Did you mention the license of the assets? [N/A]
- 502 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
503 However, the final version will include a link to a GitHub containing Mathematica code
504 for calculating POWER and optimality probability in small rewardless MDPs.
- 505 (d) Did you discuss whether and how consent was obtained from people whose data you’re
506 using/curating? [N/A]
- 507 (e) Did you discuss whether the data you are using/curating contains personally identifiable
508 information or offensive content? [N/A]
- 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? [N/A]
- 512 (b) Did you describe any potential participant risks, with links to Institutional Review
513 Board (IRB) approvals, if applicable? [N/A]
- 514 (c) Did you include the estimated hourly wage paid to participants and the total amount
515 spent on participant compensation? [N/A]