

## Appendix A Comparing POWER with information-theoretic empowerment

Salge et al. [2014] define information-theoretic *empowerment* as the maximum possible mutual information between the agent’s actions and the state observations  $n$  steps in the future, written  $\mathfrak{E}_n(s)$ . This notion requires an arbitrary choice of horizon, failing to account for the agent’s discount rate  $\gamma$ . “In a discrete deterministic world empowerment reduces to the logarithm of the number of sensor states reachable with the available actions” [Salge et al., 2014]. Figure 9 demonstrates how empowerment can return counterintuitive verdicts with respect to the agent’s control over the future.

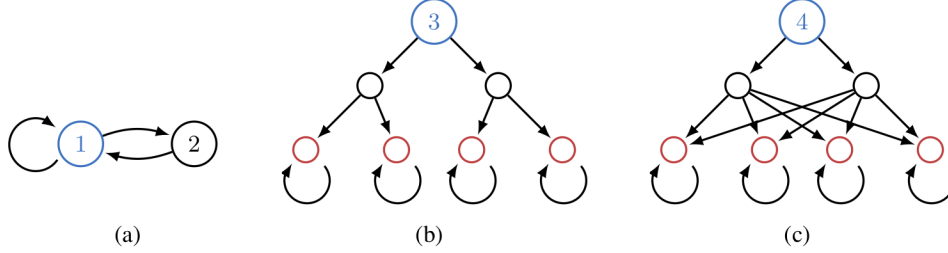


Figure 9: Proposed empowerment measures fail to adequately capture how future choice is affected by present actions. In a:  $\mathfrak{E}_n(s_1)$  varies depending on whether  $n$  is even; thus,  $\lim_{n \rightarrow \infty} \mathfrak{E}_n(s_1)$  does not exist. In b and c:  $\forall n : \mathfrak{E}_n(s_3) = \mathfrak{E}_n(s_4)$ , even though  $s_4$  allows greater control over future state trajectories than  $s_3$  does. For example, suppose that in both b and c, the leftmost black state and the rightmost red state have 1 reward while all other states have 0 reward. In c, the agent can independently maximize the intermediate black-state reward and the delayed red-state reward. Independent maximization is not possible in b.

POWER returns intuitive answers in these situations.  $\lim_{\gamma \rightarrow 1} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s_1, \gamma)$  converges by lemma 5.3. Consider the obvious involution  $\phi$  which takes each state in fig. 9b to its counterpart in fig. 9c. Since  $\phi \cdot \mathcal{F}_{\text{nd}}(s_3) \subsetneq \mathcal{F}_{\text{nd}}(s_4) = \mathcal{F}(s_4)$ , proposition 6.6 proves that  $\forall \gamma \in [0, 1] : \text{POWER}_{\mathcal{D}_{\text{bound}}}(s_3, \gamma) \leq_{\text{most: } \mathcal{D}_{\text{bound}}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s_4, \gamma)$ , with the proof of proposition 6.6 showing strict inequality under all  $\mathcal{D}_{X\text{-IID}}$  when  $\gamma \in (0, 1)$ .

Empowerment can be adjusted to account for these cases, perhaps by considering the channel capacity between the agent’s actions and the state trajectories induced by stationary policies. However, since POWER is formulated in terms of optimal value, we believe that POWER is better suited for MDPs than information-theoretic empowerment is.

## Appendix B Seeking POWER can be a detour

**Remark.** The results of appendix E do not depend on this section’s results.

One might suspect that optimal policies tautologically tend to seek POWER. This intuition is wrong.

**Proposition B.1** (Greater  $\text{POWER}_{\mathcal{D}_{\text{bound}}}$  does not imply greater  $\mathbb{P}_{\mathcal{D}_{\text{bound}}}$ ).  
Action  $a$  seeking more  $\text{POWER}_{\mathcal{D}_{\text{bound}}}$  than  $a'$  at state  $s$  and  $\gamma$  does not imply that  $\mathbb{P}_{\mathcal{D}_{\text{bound}}}(s, a, \gamma) \geq \mathbb{P}_{\mathcal{D}_{\text{bound}}}(s, a', \gamma)$ .

*Proof.* Consider the environment of fig. 10. Let  $X_u := \text{unif}(0, 1)$ , and consider  $\mathcal{D}_{X_u\text{-IID}}$ , which has bounded support. Direct computation<sup>5</sup> of the POWER expectation (definition 5.2) yields  $\text{POWER}_{\mathcal{D}_{X_u\text{-IID}}}(s_2, 1) = \frac{3}{4} > \frac{2}{3} = \text{POWER}_{\mathcal{D}_{X_u\text{-IID}}}(s_3, 1)$ . Therefore, N seeks more  $\text{POWER}_{\mathcal{D}_{X_u\text{-IID}}}$  than NE at state  $s_1$  and  $\gamma = 1$ .

However,  $\mathbb{P}_{\mathcal{D}_{X_u\text{-IID}}}(s_1, \text{N}, 1) = \frac{1}{3} < \frac{2}{3} = \mathbb{P}_{\mathcal{D}_{X_u\text{-IID}}}(s_1, \text{NE}, 1)$ .  $\square$

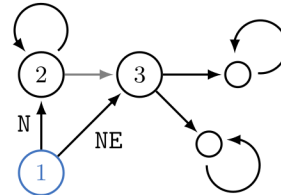


Figure 10

<sup>5</sup>In small deterministic MDPs, the POWER and optimality probability of the maximum-entropy reward function distribution can be computed using <https://github.com/loganriggs/Optimal-Policies-Tend-To-Seek-Power>.

**Lemma B.2** (Fraction of orbits which agree on weak optimality). *Let  $\mathfrak{D} \subseteq \Delta(\mathbb{R}^{|\mathcal{S}|})$ , and suppose  $f_1, f_2 : \Delta(\mathbb{R}^{|\mathcal{S}|}) \rightarrow \mathbb{R}$  are such that  $f_1(\mathcal{D}) \geq_{\text{most: } \mathfrak{D}} f_2(\mathcal{D})$ . Then for all  $\mathcal{D} \in \mathfrak{D}$ ,*

$$\frac{|\{\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D} \mid f_1(\mathcal{D}') \geq f_2(\mathcal{D}')\}|}{|S_{|\mathcal{S}|} \cdot \mathcal{D}|} \geq \frac{1}{2}.$$

*Proof.* All  $\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D}$  such that  $f_1(\mathcal{D}') = f_2(\mathcal{D}')$  satisfy  $f_1(\mathcal{D}') \geq f_2(\mathcal{D}')$ .

Otherwise, consider the  $\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D}$  such that  $f_1(\mathcal{D}') \neq f_2(\mathcal{D}')$ . By the definition of  $\geq_{\text{most}}$  (definition 6.5), at least  $\frac{1}{2}$  of these  $\mathcal{D}'$  satisfy  $f_1(\mathcal{D}') > f_2(\mathcal{D}')$ , in which case  $f_1(\mathcal{D}') \geq f_2(\mathcal{D}')$ . Then the desired inequality follows.  $\square$

**Lemma B.3** ( $\geq_{\text{most}}$  and trivial orbits). *Let  $\mathfrak{D} \subseteq \Delta(\mathbb{R}^{|\mathcal{S}|})$  and suppose  $f_1(\mathcal{D}) \geq_{\text{most: } \mathfrak{D}} f_2(\mathcal{D})$ . For all reward function distributions  $\mathcal{D} \in \mathfrak{D}$  with one-element orbits,  $f_1(\mathcal{D}) \geq f_2(\mathcal{D})$ . In particular,  $\mathcal{D}$  has a one-element orbit when it distributes reward identically and independently (IID) across states.*

*Proof.* By lemma B.2, at least half of the elements  $\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D}$  satisfy  $f_1(\mathcal{D}') \geq f_2(\mathcal{D}')$ . But  $|S_{|\mathcal{S}|} \cdot \mathcal{D}| = 1$ , and so  $f_1(\mathcal{D}) \geq f_2(\mathcal{D})$  must hold.

If  $\mathcal{D}$  is IID, it has a one-element orbit due to the assumed identical distribution of reward.  $\square$

**Proposition B.4** (Actions which tend to seek POWER do not necessarily tend to be optimal). *Action  $a$  tending to seek more POWER than  $a'$  at state  $s$  and  $\gamma$  does not imply that  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma) \geq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a', \gamma)$ .*

*Proof.* Consider the environment of fig. 10. Since  $\text{RSD}_{\text{nd}}(s_3) \subsetneq \text{RSD}(s_2)$ , proposition 6.12 shows that  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_2, 1) \geq_{\text{most: } \mathfrak{D}_{\text{bound}}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s_3, 1)$  via  $s' := s_3, s := s_2, \phi$  the identity permutation (which is an involution). Therefore, N tends to seek more POWER than NE at state  $s_1$  and  $\gamma = 1$ .

If  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s_1, \text{N}, 1) \geq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(s_1, \text{NE}, 1)$ , then lemma B.3 shows that  $\mathbb{P}_{\mathcal{D}_{X\text{-IID}}}(s_1, \text{N}, 1) \geq \mathbb{P}_{\mathcal{D}_{X\text{-IID}}}(s_1, \text{NE}, 1)$  for all  $\mathcal{D}_{X\text{-IID}}$ . But the proof of proposition B.1 showed that  $\mathbb{P}_{\mathcal{D}_{X_u\text{-IID}}}(s_1, \text{N}, 1) < \mathbb{P}_{\mathcal{D}_{X_u\text{-IID}}}(s_1, \text{NE}, 1)$  for  $X_u := \text{unif}(0, 1)$ . Therefore, it cannot be true that  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s_1, \text{N}, 1) \geq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(s_1, \text{NE}, 1)$ .  $\square$

## Appendix C Sub-optimal POWER

In certain situations, POWER returns intuitively surprising verdicts. There exists a policy under which the reader chooses a winning lottery ticket, but it seems wrong to say that the reader has the power to win the lottery with high probability. For various reasons, humans and other bounded agents are generally incapable of computing optimal policies for arbitrary objectives. More formally, consider the rewardless MDP of fig. 11.

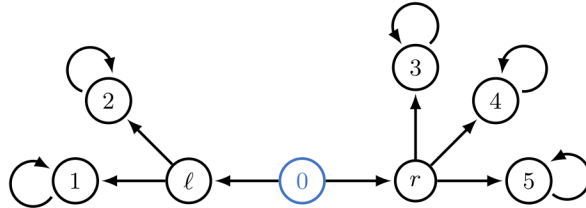


Figure 11:  $s_0$  is the starting state, and  $|\mathcal{A}| = 10^{10^{10}}$ . At  $s_0$ , half of the actions lead to  $s_\ell$ , while the other half lead to  $s_r$ . Similarly, half of the actions at  $s_\ell$  lead to  $s_1$ , while the other half lead to  $s_2$ . At  $s_r$ , one action leads to  $s_3$ , one action leads to  $s_4$ , and the remaining  $10^{10^{10}} - 2$  actions lead to  $s_5$ .

Consider a model-based RL agent with black-box simulator access to this environment. The agent has no prior information about the model, and so it acts randomly. Before long, the agent has probably

learned how to navigate from  $s_0$  to states  $s_\ell$ ,  $s_r$ ,  $s_1$ ,  $s_2$ , and  $s_5$ . However, over any reasonable timescale, it is extremely improbable that the agent discovers the two actions respectively leading to  $s_3$  and  $s_4$ .

Even provided with a reward function  $R$  and the discount rate  $\gamma$ , the agent has yet to learn the relevant environmental dynamics, and so many of its policies are far from optimal. Although proposition 6.6 shows that  $\forall \gamma \in [0, 1] : \text{POWER}_{\mathcal{D}_{\text{bound}}}(s_\ell, \gamma) \leq_{\text{most: } \mathcal{D}_{\text{bound}}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s_r, \gamma)$ , there is a sense in which  $s_\ell$  gives this agent more power.

We formalize a bounded agent’s goal-achievement capabilities with a function  $\text{pol}$ , which takes as input a reward function and a discount rate, and returns a policy. Informally, this is the best policy which the agent knows about. We can then calculate  $\text{POWER}_{\mathcal{D}_{\text{bound}}}$  with respect to  $\text{pol}$ .

**Definition C.1 (Suboptimal POWER).** Let  $\Pi_\Delta$  be the set of stationary stochastic policies, and let  $\text{pol} : \mathbb{R}^S \times [0, 1] \rightarrow \Pi_\Delta$ . For  $\gamma \in [0, 1]$ ,

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}}(s, \gamma) := \mathbb{E}_{\substack{R \sim \mathcal{D}_{\text{bound}}, \\ a \sim \text{pol}(R, \gamma)(s), \\ s' \sim T(s, a)}} \left[ \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) V_R^{\text{pol}(R, \gamma^*)}(s', \gamma^*) \right]. \quad (5)$$

By lemma E.38,  $\text{POWER}_{\mathcal{D}_{\text{bound}}}$  is the special case where  $\forall R \in \mathbb{R}^S, \gamma \in [0, 1] : \text{pol}(R, \gamma) \in \Pi^*(R, \gamma)$ . We define  $\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}}$ -seeking similarly as in definition 5.6.

$\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}}(s_0, 1)$  increases as the policies returned by  $\text{pol}$  are improved. We illustrate this by considering the  $\mathcal{D}_{X\text{-IID}}$  case.

- $\text{pol}_1$  The model is initially unknown, and so  $\forall R, \gamma : \text{pol}_1(R, \gamma)$  is a uniformly random policy. Since  $\text{pol}_1$  is constant on its inputs,  $\text{POWER}_{\mathcal{D}_{X\text{-IID}}}^{\text{pol}_1}(s_0, 1) = \mathbb{E}[X]$  by the linearity of expectation and the fact that  $\mathcal{D}_{X\text{-IID}}$  distributes reward independently and identically across states.
- $\text{pol}_2$  The agent knows the dynamics, except that it does not know how to reach  $s_3$  or  $s_4$ . At this point,  $\text{pol}_2(R, 1)$  navigates from  $s_0$  to the average-optimal choice among three terminal states:  $s_1$ ,  $s_2$ , and  $s_5$ . Therefore,  $\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}_2}(s_0, 1) = \mathbb{E}[\text{max of 3 draws from } X]$ .
- $\text{pol}_3$  The agent knows the dynamics, the environment is small enough to solve explicitly, and so  $\forall R, \gamma : \text{pol}_3(R, \gamma)$  is an optimal policy.  $\text{pol}_3(R, 1)$  navigates from  $s_0$  to the average-optimal choice among all five terminal states. Therefore,  $\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}_3}(s_0, 1) = \mathbb{E}[\text{max of 5 draws from } X]$ .

As the agent learns more about the environment and improves  $\text{pol}$ , the agent’s  $\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}}$  increases. The agent seeks  $\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}_2}$  by navigating to  $s_\ell$  instead of  $s_r$ , but seeks more  $\text{POWER}_{\mathcal{D}_{\text{bound}}}$  by navigating to  $s_r$  instead of  $s_\ell$ . Intuitively, bounded agents gain power by improving  $\text{pol}$  and by formally seeking  $\text{POWER}_{\mathcal{D}_{\text{bound}}}^{\text{pol}}$  within the environment.

## Appendix D Lists of results

### List of definitions

3.1	Definition (Rewardless MDP)	2
3.2	Definition (1-cycle states)	2
3.3	Definition (State visit distribution [Sutton and Barto, 1998])	3
3.4	Definition ( $\mathcal{F}$ single-state restriction)	3
3.5	Definition (Value function)	3
3.6	Definition (Non-domination)	3
4.1	Definition (Optimal policy set function)	3

4.2	Definition (Reward function distributions)	3
4.3	Definition (Visit distribution optimality probability)	4
4.4	Definition (Action optimality probability)	4
5.1	Definition (Average optimal value)	4
5.2	Definition (POWER)	4
5.6	Definition (POWER-seeking actions)	5
6.1	Definition (Similarity of vector sets)	6
6.2	Definition (Similarity of vector function sets)	6
6.3	Definition (Pushforward distribution of a permutation)	6
6.4	Definition (Orbit of a probability distribution)	6
6.5	Definition (Inequalities which hold for most probability distributions)	6
6.7	Definition (Equivalent actions)	7
6.8	Definition (States reachable after taking an action)	7
6.10	Definition (Recurrent state distributions [Puterman, 2014])	8
6.11	Definition (Average-optimal policies)	8
C.1	Definition (Suboptimal POWER)	15
E.2	Definition (Transition matrix induced by a policy)	18
E.6	Definition (Continuous reward function distribution)	19
E.9	Definition (Non-dominated linear functionals)	20
E.13	Definition (Non-dominated vector functions)	21
E.14	Definition (Affine transformation of visit distribution sets)	21
E.18	Definition (Support of $\mathcal{D}_{\text{any}}$ )	22
E.19	Definition (Linear functional optimality probability)	22
E.23	Definition (Bounded, continuous IID reward)	23
E.25	Definition (Indicator function)	24
E.31	Definition (Evaluating sets of visit distribution functions at $\gamma$ )	27
E.39	Definition (Discount-normalized value function)	29
E.42	Definition (Normalized visit distribution function)	32

## List of theorems

5.3	Lemma (Continuity of POWER)	5
5.4	Proposition (Maximal POWER)	5
5.5	Proposition (POWER is smooth across reversible dynamics)	5
6.6	Proposition (States with “more options” have more POWER)	6
6.9	Proposition (Keeping options open tends to be POWER-seeking and tends to be optimal)	7
6.12	Proposition (When $\gamma = 1$ , RSDs control POWER)	8
6.13	Theorem (Average-optimal policies tend to end up in “larger” sets of RSDs)	8
6.14	Corollary (Average-optimal policies tend not to end up in any given 1-cycle)	8
B.1	Proposition (Greater $\text{POWER}_{\mathcal{D}_{\text{bound}}}$ does not imply greater $\mathbb{P}_{\mathcal{D}_{\text{bound}}}$ )	13



B.2	Lemma (Fraction of orbits which agree on weak optimality)	14
B.3	Lemma ( $\geq_{\text{most}}$ and trivial orbits)	14
B.4	Proposition (Actions which tend to seek POWER do not necessarily tend to be optimal)	14
E.1	Lemma (A policy is optimal iff it induces an optimal visit distribution at every state)	18
E.3	Proposition (Properties of visit distribution functions)	18
E.4	Lemma ( $\mathbf{f} \in \mathcal{F}(s)$ is multivariate rational on $\gamma$ )	19
E.5	Corollary (On-policy value is rational on $\gamma$ )	19
E.7	Lemma (Distinct linear functionals disagree almost everywhere on their domains)	19
E.8	Corollary (Unique maximization of almost all vectors)	20
E.10	Lemma (All vectors are maximized by a non-dominated linear functional)	20
E.11	Corollary (Maximal value is invariant to restriction to non-dominated functionals)	20
E.12	Lemma (How non-domination containment affects optimal value)	20
E.15	Lemma (Invariance of non-domination under positive affine transform)	21
E.16	Lemma (Helper lemma for demonstrating $\geq_{\text{most}: \mathfrak{D}_{\text{any}}}$ )	21
E.17	Lemma (A helper result for expectations of functions)	22
E.20	Proposition (Non-dominated linear functionals and their optimality probability)	22
E.21	Lemma (Expected value of similar linear functional sets)	23
E.22	Lemma (For continuous IID distributions $\mathcal{D}_{X\text{-IID}}, \exists b < c : (b, c)^{ S } \subseteq \text{supp}(\mathcal{D}_{X\text{-IID}})$ )	23
E.24	Lemma (Expectation superiority lemma)	23
E.26	Lemma (Optimality probability inclusion relations)	24
E.27	Lemma (Optimality probability of similar linear functional sets)	25
E.28	Lemma (Optimality probability superiority lemma)	26
E.29	Lemma (Limit probability inequalities which hold for most distributions)	26
E.30	Proposition (How to transfer optimal policy sets across discount rates)	27
E.32	Lemma (Non-domination across $\gamma$ values for expectations of visit distributions)	27
E.33	Lemma ( $\forall \gamma \in (0, 1) : \mathbf{d} \in \mathcal{F}_{\text{nd}}(s, \gamma)$ iff $\mathbf{d} \in \text{ND}(\mathcal{F}(s, \gamma))$ )	28
E.34	Lemma ( $\forall \gamma \in [0, 1) : V_R^*(s, \gamma) = \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \mathbf{f}(\gamma)^\top \mathbf{r}$ )	28
E.35	Lemma (Optimal policy shift bound)	28
E.36	Proposition (Optimality probability's limits exist)	28
E.37	Lemma (Optimality probability identity)	28
E.38	Lemma (POWER identities)	29
E.40	Lemma (Normalized value functions have uniformly bounded derivative)	30
E.41	Lemma (Lower bound on current POWER based on future POWER)	31
E.43	Lemma (Normalized visit distribution functions are continuous)	32
E.44	Lemma (Non-domination of normalized visit distribution functions)	33
E.45	Lemma (POWER limit identity)	33
E.46	Lemma (Lemma for POWER superiority)	34
E.47	Lemma (Non-dominated visit distribution functions never agree with other visit distribution functions at that state)	36

E.48 Corollary (Cardinality of non-dominated visit distributions) . . . . .	36
E.49 Lemma (Optimality probability and state bottlenecks) . . . . .	36
E.50 Lemma (Action optimality probability is a special case of visit distribution optimality probability) . . . . .	37
E.51 Lemma (POWER identity when $\gamma = 1$ ) . . . . .	40
E.52 Proposition (RSD properties) . . . . .	42
E.53 Lemma (When reachable with probability 1, 1-cycles induce non-dominated RSDs) . . . . .	42

## D.1 Contributions of independent interest

We developed new basic MDP theory by exploring the structural properties of visit distribution functions. Echoing Wang et al. [2007, 2008], we believe that this area is interesting and underexplored.

### D.1.1 Optimal value theory

Lemma E.40 shows that  $f(\gamma^*) := \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*)V_R^*(s, \gamma^*)$  is Lipschitz continuous on  $\gamma \in [0, 1]$ , with Lipschitz constant depending only on  $\|R\|_1$ . For all states  $s$  and policies  $\pi \in \Pi$ , corollary E.5 shows that  $V_R^\pi(s, \gamma)$  is rational on  $\gamma$ .

Optimal value has a well-known dual formulation:  $V_R^*(s, \gamma) = \max_{\mathbf{f} \in \mathcal{F}(s)} \mathbf{f}(\gamma)^\top \mathbf{r}$ .

**Lemma E.34** ( $\forall \gamma \in [0, 1] : V_R^*(s, \gamma) = \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \mathbf{f}(\gamma)^\top \mathbf{r}$ ).

In a fixed rewardless MDP, lemma E.34 may enable more efficient computation of optimal value functions for multiple reward functions.

### D.1.2 Optimal policy theory

Proposition E.30 demonstrates how to preserve optimal incentives while changing the discount rate.

**Proposition E.30** (How to transfer optimal policy sets across discount rates). *Suppose reward function  $R$  has optimal policy set  $\Pi^*(R, \gamma)$  at discount rate  $\gamma \in (0, 1)$ . For any  $\gamma^* \in (0, 1)$ , we can construct a reward function  $R'$  such that  $\Pi^*(R', \gamma^*) = \Pi^*(R, \gamma)$ . Furthermore,  $V_{R'}^*(\cdot, \gamma^*) = V_R^*(\cdot, \gamma)$ .*

### D.1.3 Visit distribution theory

While Regan and Boutilier [2010] consider a visit distribution function  $\mathbf{f} \in \mathcal{F}(s)$  to be non-dominated if it is optimal for some reward function in a set  $\mathcal{R} \subseteq \mathbb{R}^{|\mathcal{S}|}$ , our stricter definition 3.6 considers  $\mathbf{f}$  to be non-dominated when  $\exists \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}, \gamma \in (0, 1) : \mathbf{f}(\gamma)^\top \mathbf{r} > \max_{\mathbf{f}' \in \mathcal{F}(s) \setminus \{\mathbf{f}\}} \mathbf{f}'(\gamma)^\top \mathbf{r}$ .

## Appendix E Theoretical results

**Lemma E.1** (A policy is optimal iff it induces an optimal visit distribution at every state). *Let  $\gamma \in (0, 1)$  and let  $R$  be a reward function.  $\pi \in \Pi^*(R, \gamma)$  iff  $\pi$  induces an optimal visit distribution at every state.*

*Proof.* By definition, a policy  $\pi$  is optimal iff  $\pi$  induces the maximal on-policy value at each state, which is true iff  $\pi$  induces an optimal visit distribution at every state (by the dual formulation of optimal value functions).  $\square$

**Definition E.2** (Transition matrix induced by a policy).  $\mathbf{T}^\pi$  is the transition matrix induced by policy  $\pi \in \Pi$ , where  $\mathbf{T}^\pi \mathbf{e}_s := T(s, \pi(s))$ .  $(\mathbf{T}^\pi)^t \mathbf{e}_s$  gives the probability distribution over the states visited at time step  $t$ , after following  $\pi$  for  $t$  steps from  $s$ .

**Proposition E.3** (Properties of visit distribution functions). *Let  $s, s' \in \mathcal{S}, \mathbf{f}^{\pi, s} \in \mathcal{F}(s)$ .*

1.  $\mathbf{f}^{\pi, s}(\gamma)$  is element-wise non-negative and element-wise monotonically increasing on  $\gamma \in [0, 1]$ .

$$2. \forall \gamma \in [0, 1) : \|\mathbf{f}^{\pi, s}(\gamma)\|_1 = \frac{1}{1-\gamma}.$$

*Proof.* Item 1: by examination of definition 3.3,  $\mathbf{f}^{\pi, s} = \sum_{t=0}^{\infty} (\gamma \mathbf{T}^\pi)^t \mathbf{e}_s$ . Since each  $(\mathbf{T}^\pi)^t$  is left stochastic and  $\mathbf{e}_s$  is the standard unit vector, each entry in each summand is non-negative. Therefore,  $\forall \gamma \in [0, 1) : \mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{e}_{s'} \geq 0$ , and this function monotonically increases on  $\gamma$ .

Item 2:

$$\|\mathbf{f}^{\pi, s}(\gamma)\|_1 = \left\| \sum_{t=0}^{\infty} (\gamma \mathbf{T}^\pi)^t \mathbf{e}_s \right\|_1 \quad (6)$$

$$= \sum_{t=0}^{\infty} \gamma^t \left\| (\mathbf{T}^\pi)^t \mathbf{e}_s \right\|_1 \quad (7)$$

$$= \sum_{t=0}^{\infty} \gamma^t \quad (8)$$

$$= \frac{1}{1-\gamma}. \quad (9)$$

Equation (7) follows because all entries in each  $(\mathbf{T}^\pi)^t \mathbf{e}_s$  are non-negative by item 1. Equation (8) follows because each  $(\mathbf{T}^\pi)^t$  is left stochastic and  $\mathbf{e}_s$  is a stochastic vector, and so  $\left\| (\mathbf{T}^\pi)^t \mathbf{e}_s \right\|_1 = 1$ .  $\square$

**Lemma E.4** ( $\mathbf{f} \in \mathcal{F}(s)$  is multivariate rational on  $\gamma$ ).  $\mathbf{f}^\pi \in \mathcal{F}(s)$  is a multivariate rational function on  $\gamma \in [0, 1)$ .

*Proof.* Let  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$  and consider  $\mathbf{f}^\pi \in \mathcal{F}(s)$ . Let  $\mathbf{v}_R^\pi$  be the  $V_R^*(s, \gamma)$  function in column vector form, with one entry per state value.

By the Bellman equations,  $\mathbf{v}_R^\pi = (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{r}$ . Let  $\mathbf{A}_\gamma := (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1}$ , and for state  $s$ , form  $\mathbf{A}_{s, \gamma}$  by replacing  $\mathbf{A}_\gamma$ 's column for state  $s$  with  $\mathbf{r}$ . As noted by Lippman [1968], by Cramer's rule,  $V_R^\pi(s, \gamma) = \frac{\det \mathbf{A}_{s, \gamma}}{\det \mathbf{A}_\gamma}$  is a rational function with numerator and denominator having degree at most  $|\mathcal{S}|$ .

In particular, for each state indicator reward function  $\mathbf{e}_{s_i}$ ,  $V_{s_i}^\pi(s, \gamma) = \mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{e}_{s_i}$  is a rational function of  $\gamma$  whose numerator and denominator each have degree at most  $|\mathcal{S}|$ . This implies that  $\mathbf{f}^\pi(\gamma)$  is multivariate rational on  $\gamma \in [0, 1)$ .  $\square$

**Corollary E.5** (On-policy value is rational on  $\gamma$ ). Let  $\pi \in \Pi$  and  $R$  be any reward function.  $V_R^\pi(s, \gamma)$  is rational on  $\gamma \in [0, 1)$ .

*Proof.*  $V_R^\pi(s, \gamma) = \mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{r}$ , and  $\mathbf{f}$  is a multivariate rational function of  $\gamma$  by lemma E.4. Therefore, for fixed  $\mathbf{r}$ ,  $\mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{r}$  is a rational function of  $\gamma$ .  $\square$

## E.1 Non-dominated visit distribution functions

**Definition E.6** (Continuous reward function distribution). Results with  $\mathcal{D}_{\text{cont}}$  hold for any absolutely continuous reward function distribution.

**Remark.** We assume  $\mathbb{R}^{|\mathcal{S}|}$  is endowed with the standard topology.

**Lemma E.7** (Distinct linear functionals disagree almost everywhere on their domains). Let  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{|\mathcal{S}|}$  be distinct.  $\mathbb{P}_{\mathbf{r} \sim \mathcal{D}_{\text{cont}}}(\mathbf{x}^\top \mathbf{r} = \mathbf{x}'^\top \mathbf{r}) = 0$ .

*Proof.*  $\{\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} \mid (\mathbf{x} - \mathbf{x}')^\top \mathbf{r} = 0\}$  is a hyperplane since  $\mathbf{x} - \mathbf{x}' \neq \mathbf{0}$ . Therefore, it has no interior in the standard topology on  $\mathbb{R}^{|\mathcal{S}|}$ . Since this empty-interior set is also convex, it has zero Lebesgue measure. By the Radon-Nikodym theorem, it has zero measure under any continuous distribution  $\mathcal{D}_{\text{cont}}$ .  $\square$

**Corollary E.8** (Unique maximization of almost all vectors). *Let  $X \subsetneq \mathbb{R}^{|\mathcal{S}|}$  be finite.  $\mathbb{P}_{\mathbf{r} \sim \mathcal{D}_{\text{cont}}} \left( \left| \arg \max_{\mathbf{x}'' \in X} \mathbf{x}''^\top \mathbf{r} \right| > 1 \right) = 0$ .*

*Proof.* Let  $\mathbf{x}, \mathbf{x}' \in X$  be distinct. For any  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$ ,  $\mathbf{x}, \mathbf{x}' \in \arg \max_{\mathbf{x}'' \in X} \mathbf{x}''^\top \mathbf{r}$  iff  $\mathbf{x}^\top \mathbf{r} = \mathbf{x}'^\top \mathbf{r} \geq \max_{\mathbf{x}'' \in X \setminus \{\mathbf{x}, \mathbf{x}'\}} \mathbf{x}''^\top \mathbf{r}$ . By lemma E.7,  $\mathbf{x}^\top \mathbf{r} = \mathbf{x}'^\top \mathbf{r}$  holds with probability 0 under any  $\mathcal{D}_{\text{cont}}$ .  $\square$

### E.1.1 Generalized non-domination results

Our formalism includes both  $\mathcal{F}_{\text{nd}}(s)$  and  $\text{RSD}_{\text{nd}}(s)$ ; we therefore prove results that are applicable to both.

**Definition E.9** (Non-dominated linear functionals). Let  $X \subsetneq \mathbb{R}^{|\mathcal{S}|}$  be finite.  $\text{ND}(X) := \left\{ \mathbf{x} \in X \mid \exists \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} : \mathbf{x}^\top \mathbf{r} > \max_{\mathbf{x}' \in X \setminus \{\mathbf{x}\}} \mathbf{x}'^\top \mathbf{r} \right\}$ .

**Lemma E.10** (All vectors are maximized by a non-dominated linear functional). *Let  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$  and let  $X \subsetneq \mathbb{R}^{|\mathcal{S}|}$  be finite and non-empty.  $\exists \mathbf{x}^* \in \text{ND}(X) : \mathbf{x}^{*\top} \mathbf{r} = \max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r}$ .*

*Proof.* Let  $A(\mathbf{r} \mid X) := \arg \max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . Then

$$\mathbf{x}_1^\top \mathbf{r} = \dots = \mathbf{x}_n^\top \mathbf{r} > \max_{\mathbf{x}' \in X \setminus A(\mathbf{r} \mid X)} \mathbf{x}'^\top \mathbf{r}. \quad (10)$$

In eq. (10), each  $\mathbf{x}^\top \mathbf{r}$  expression is linear on  $\mathbf{r}$ . The max is piecewise linear on  $\mathbf{r}$  since it is the maximum of a finite set of linear functionals. In particular, all expressions in eq. (10) are continuous on  $\mathbf{r}$ , and so we can find some  $\delta > 0$  neighborhood  $B(\mathbf{r}, \delta)$  such that  $\forall \mathbf{r}' \in B(\mathbf{r}, \delta) : \max_{\mathbf{x}_i \in A(\mathbf{r} \mid X)} \mathbf{x}_i^\top \mathbf{r}' > \max_{\mathbf{x}' \in X \setminus A(\mathbf{r} \mid X)} \mathbf{x}'^\top \mathbf{r}'$ .

But almost all  $\mathbf{r}' \in B(\mathbf{r}, \delta)$  are maximized by a unique functional  $\mathbf{x}^*$  by corollary E.8; in particular, at least one such  $\mathbf{r}''$  exists. Formally,  $\exists \mathbf{r}'' \in B(\mathbf{r}, \delta) : \mathbf{x}^{*\top} \mathbf{r}'' > \max_{\mathbf{x}' \in X \setminus \{\mathbf{x}^*\}} \mathbf{x}'^\top \mathbf{r}''$ . Therefore,  $\mathbf{x}^* \in \text{ND}(X)$  by definition E.9.

$\mathbf{x}^{*\top} \mathbf{r}' \geq \max_{\mathbf{x}_i \in A(\mathbf{r} \mid X)} \mathbf{x}_i^\top \mathbf{r}' > \max_{\mathbf{x}' \in X \setminus A(\mathbf{r} \mid X)} \mathbf{x}'^\top \mathbf{r}'$ , with the strict inequality following because  $\mathbf{r}'' \in B(\mathbf{r}, \delta)$ . These inequalities imply that  $\mathbf{x}^* \in A(\mathbf{r} \mid X)$ .  $\square$

**Corollary E.11** (Maximal value is invariant to restriction to non-dominated functionals). *Let  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$  and let  $X \subsetneq \mathbb{R}^{|\mathcal{S}|}$  be finite.  $\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} = \max_{\mathbf{x} \in \text{ND}(X)} \mathbf{x}^\top \mathbf{r}$ .*

*Proof.* If  $X$  is empty, holds trivially. Otherwise, apply lemma E.10.  $\square$

**Lemma E.12** (How non-domination containment affects optimal value). *Let  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$  and let  $X, X' \subsetneq \mathbb{R}^{|\mathcal{S}|}$  be finite.*

1. *If  $\text{ND}(X) \subseteq X'$ , then  $\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \leq \max_{\mathbf{x}' \in X'} \mathbf{x}'^\top \mathbf{r}$ .*
2. *If  $\text{ND}(X) \subseteq X' \subseteq X$ , then  $\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} = \max_{\mathbf{x}' \in X'} \mathbf{x}'^\top \mathbf{r}$ .*

*Proof.* Item 1:

$$\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} = \max_{\mathbf{x} \in \text{ND}(X)} \mathbf{x}^\top \mathbf{r} \quad (11)$$

$$\leq \max_{\mathbf{x}' \in X'} \mathbf{x}'^\top \mathbf{r}. \quad (12)$$

Equation (11) follows by corollary E.11. Equation (12) follows because  $\text{ND}(X) \subseteq X'$ .

Item 2: by item 1,  $\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \leq \max_{\mathbf{x}' \in X'} \mathbf{x}'^\top \mathbf{r}$ . Since  $X' \subseteq X$ , we also have  $\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{x}' \in X'} \mathbf{x}'^\top \mathbf{r}$ , and so equality must hold.  $\square$

**Definition E.13** (Non-dominated vector functions). Let  $I \subseteq \mathbb{R}$  and let  $F \subsetneq \left(\mathbb{R}^{|\mathcal{S}|}\right)^I$  be a finite set of vector-valued functions on  $I$ .  $\text{ND}(F) := \left\{ \mathbf{f} \in F \mid \exists \gamma \in I, \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} : \mathbf{f}(\gamma)^\top \mathbf{r} > \max_{\mathbf{f}' \in F \setminus \{\mathbf{f}\}} \mathbf{f}'(\gamma)^\top \mathbf{r} \right\}$ .

**Remark.**  $\mathcal{F}_{\text{nd}}(s) = \text{ND}(\mathcal{F}(s))$  by definition 3.6.

**Definition E.14** (Affine transformation of visit distribution sets). For notational convenience, we define set-scalar multiplication and set-vector addition on  $X \subseteq \mathbb{R}^{|\mathcal{S}|}$ : for  $c \in \mathbb{R}$ ,  $cX := \{c\mathbf{x} \mid \mathbf{x} \in X\}$ . For  $\mathbf{a} \in \mathbb{R}^{|\mathcal{S}|}$ ,  $X + \mathbf{a} := \{\mathbf{x} + \mathbf{a} \mid \mathbf{x} \in X\}$ . Similar operations hold when  $X$  is a set of vector functions  $\mathbb{R} \mapsto \mathbb{R}^{|\mathcal{S}|}$ .

**Lemma E.15** (Invariance of non-domination under positive affine transform).

1. Let  $X \subsetneq \mathbb{R}^{|\mathcal{S}|}$  be finite. If  $\mathbf{x} \in \text{ND}(X)$ , then  $\forall c > 0, \mathbf{a} \in \mathbb{R}^{|\mathcal{S}|} : (c\mathbf{x} + \mathbf{a}) \in \text{ND}(cX + \mathbf{a})$ .
2. Let  $I \subseteq \mathbb{R}$  and let  $F \subsetneq \left(\mathbb{R}^{|\mathcal{S}|}\right)^I$  be a finite set of vector-valued functions on  $I$ . If  $\mathbf{f} \in \text{ND}(F)$ , then  $\forall c > 0, \mathbf{a} \in \mathbb{R}^{|\mathcal{S}|} : (c\mathbf{f} + \mathbf{a}) \in \text{ND}(cF + \mathbf{a})$ .

*Proof.* Item 1: Suppose  $\mathbf{x} \in \text{ND}(X)$  is strictly optimal for  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$ . Then let  $c > 0, \mathbf{a} \in \mathbb{R}^{|\mathcal{S}|}$  be arbitrary, and define  $b := \mathbf{a}^\top \mathbf{r}$ .

$$\mathbf{x}^\top \mathbf{r} > \max_{\mathbf{x}' \in X \setminus \{\mathbf{x}\}} \mathbf{x}'^\top \mathbf{r} \quad (13)$$

$$c\mathbf{x}^\top \mathbf{r} + b > \max_{\mathbf{x}' \in X \setminus \{\mathbf{x}\}} c\mathbf{x}'^\top \mathbf{r} + b \quad (14)$$

$$(c\mathbf{x} + \mathbf{a})^\top \mathbf{r} > \max_{\mathbf{x}' \in X \setminus \{\mathbf{x}\}} (c\mathbf{x}' + \mathbf{a})^\top \mathbf{r} \quad (15)$$

$$(c\mathbf{x} + \mathbf{a})^\top \mathbf{r} > \max_{\mathbf{x}'' \in (cX + \mathbf{a}) \setminus \{c\mathbf{x} + \mathbf{a}\}} \mathbf{x}''^\top \mathbf{r}. \quad (16)$$

Equation (14) follows because  $c > 0$ . Equation (15) follows by the definition of  $b$ .

Item 2: If  $\mathbf{f} \in \text{ND}(F)$ , then by definition E.13, there exist  $\gamma \in I, \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$  such that

$$\mathbf{f}(\gamma)^\top \mathbf{r} > \max_{\mathbf{f}' \in F \setminus \{\mathbf{f}\}} \mathbf{f}'(\gamma)^\top \mathbf{r}. \quad (17)$$

Apply item 1 to conclude

$$(c\mathbf{f}(\gamma) + \mathbf{a})^\top \mathbf{r} > \max_{(c\mathbf{f}' + \mathbf{a}) \in (cF + \mathbf{a}) \setminus \{c\mathbf{f} + \mathbf{a}\}} (c\mathbf{f}'(\gamma) + \mathbf{a})^\top \mathbf{r}. \quad (18)$$

Therefore,  $(c\mathbf{f} + \mathbf{a}) \in \text{ND}(cF + \mathbf{a})$ .  $\square$

## E.1.2 Inequalities which hold under most reward function distributions

**Definition 6.5** (Inequalities which hold for most probability distributions). Let  $f_1, f_2 : \Delta(\mathbb{R}^{|\mathcal{S}|}) \rightarrow \mathbb{R}$  be functions from reward function distributions to real numbers and let  $\mathfrak{D} \subseteq \Delta(\mathbb{R}^{|\mathcal{S}|})$ . We write  $f_1(\mathcal{D}) \geq_{\text{most: } \mathfrak{D}} f_2(\mathcal{D})$  when, for all  $\mathcal{D} \in \mathfrak{D}$ , the following cardinality inequality holds:

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

**Lemma E.16** (Helper lemma for demonstrating  $\geq_{\text{most: } \mathfrak{D}_{\text{any}}}$ ). Let  $\mathfrak{D} \subseteq \Delta(\mathbb{R}^{|\mathcal{S}|})$ . If  $\exists \phi \in S_{|\mathcal{S}|}$  such that for all  $\mathcal{D} \in \mathfrak{D}$ ,  $f_1(\mathcal{D}) < f_2(\mathcal{D})$  implies that  $f_1(\phi \cdot \mathcal{D}) > f_2(\phi \cdot \mathcal{D})$ , then  $f_1(\mathcal{D}) \geq_{\text{most: } \mathfrak{D}} f_2(\mathcal{D})$ .

*Proof.* Since  $\phi$  does not belong to the stabilizer of  $S_{|\mathcal{S}|}$ ,  $\phi$  acts injectively on  $S_{|\mathcal{S}|} \cdot \mathcal{D}$ . By assumption on  $\phi$ , the image of  $\{\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D} \mid f_1(\mathcal{D}') < f_2(\mathcal{D}')\}$  under  $\phi$  is a subset of  $\{\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D} \mid f_1(\mathcal{D}') > f_2(\mathcal{D}')\}$ . Since  $\phi$  is injective,  $\left| \{\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D} \mid f_1(\mathcal{D}') < f_2(\mathcal{D}')\} \right| \leq \left| \{\mathcal{D}' \in S_{|\mathcal{S}|} \cdot \mathcal{D} \mid f_1(\mathcal{D}') > f_2(\mathcal{D}')\} \right|$ .  $f_1(\mathcal{D}) \geq_{\text{most: } \mathfrak{D}} f_2(\mathcal{D})$  by definition 6.5.  $\square$



**Lemma E.17** (A helper result for expectations of functions). *Let  $B_1, \dots, B_n \subseteq \mathbb{R}^{|\mathcal{S}|}$  be finite and let  $\mathcal{D} \subseteq \Delta(\mathbb{R}^{|\mathcal{S}|})$ . Suppose  $f$  is a function of the form*

$$f(B_1, \dots, B_n | \mathcal{D}) = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}} \left[ g \left( \max_{\mathbf{b}_1 \in B_1} \mathbf{b}_1^\top \mathbf{r}, \dots, \max_{\mathbf{b}_n \in B_n} \mathbf{b}_n^\top \mathbf{r} \right) \right] \quad (19)$$

for some function  $g$ , and that  $f$  is well-defined for all  $\mathcal{D} \in \mathcal{D}$ . Let  $\phi$  be a state permutation. Then

$$f(B_1, \dots, B_n | \mathcal{D}) = f(\phi \cdot B_1, \dots, \phi \cdot B_n | \phi \cdot \mathcal{D}). \quad (20)$$

*Proof.* Let distribution  $\mathcal{D}$  have probability measure  $F$ , and let  $\phi \cdot \mathcal{D}$  have probability measure  $F_\phi$ .

$$f(B_1, \dots, B_n | \mathcal{D}) \quad (21)$$

$$:= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}} \left[ g \left( \max_{\mathbf{b}_1 \in B_1} \mathbf{b}_1^\top \mathbf{r}, \dots, \max_{\mathbf{b}_n \in B_n} \mathbf{b}_n^\top \mathbf{r} \right) \right] \quad (22)$$

$$:= \int_{\mathbb{R}^{|\mathcal{S}|}} g \left( \max_{\mathbf{b}_1 \in B_1} \mathbf{b}_1^\top \mathbf{r}, \dots, \max_{\mathbf{b}_n \in B_n} \mathbf{b}_n^\top \mathbf{r} \right) dF(\mathbf{r}) \quad (23)$$

$$= \int_{\mathbb{R}^{|\mathcal{S}|}} g \left( \max_{\mathbf{b}_1 \in B_1} \mathbf{b}_1^\top \mathbf{r}, \dots, \max_{\mathbf{b}_n \in B_n} \mathbf{b}_n^\top \mathbf{r} \right) dF_\phi(\mathbf{P}_\phi \mathbf{r}) \quad (24)$$

$$= \int_{\mathbb{R}^{|\mathcal{S}|}} g \left( \max_{\mathbf{b}_1 \in B_1} \mathbf{b}_1^\top (\mathbf{P}_\phi^{-1} \mathbf{r}'), \dots, \max_{\mathbf{b}_n \in B_n} \mathbf{b}_n^\top (\mathbf{P}_\phi^{-1} \mathbf{r}') \right) |\det \mathbf{P}_\phi| dF_\phi(\mathbf{r}') \quad (25)$$

$$= \int_{\mathbb{R}^{|\mathcal{S}|}} g \left( \max_{\mathbf{b}_1 \in B_1} (\mathbf{P}_\phi \mathbf{b}_1)^\top \mathbf{r}', \dots, \max_{\mathbf{b}_n \in B_n} (\mathbf{P}_\phi \mathbf{b}_n)^\top \mathbf{r}' \right) dF_\phi(\mathbf{r}') \quad (26)$$

$$= \int_{\mathbb{R}^{|\mathcal{S}|}} g \left( \max_{\mathbf{b}'_1 \in \phi \cdot B_1} \mathbf{b}'_1{}^\top \mathbf{r}', \dots, \max_{\mathbf{b}'_n \in \phi \cdot B_n} \mathbf{b}'_n{}^\top \mathbf{r}' \right) dF_\phi(\mathbf{r}') \quad (27)$$

$$=: f(\phi \cdot B_1, \dots, \phi \cdot B_n | \phi \cdot \mathcal{D}). \quad (28)$$

Equation (24) follows by the definition of  $F_\phi$  (definition 6.3). Equation (25) follows by substituting  $\mathbf{r}' := \mathbf{P}_\phi \mathbf{r}$ . Equation (26) follows from the fact that all permutation matrices have unitary determinant and are orthogonal (and so  $(\mathbf{P}_\phi^{-1})^\top = \mathbf{P}_\phi$ ).  $\square$

**Definition E.18** (Support of  $\mathcal{D}_{\text{any}}$ ). Let  $\mathcal{D}_{\text{any}}$  be any reward function distribution.  $\text{supp}(\mathcal{D}_{\text{any}})$  is the smallest closed subset of  $\mathbb{R}^{|\mathcal{S}|}$  whose complement has measure zero under  $\mathcal{D}_{\text{any}}$ .

**Definition E.19** (Linear functional optimality probability). For finite  $A, B \subseteq \mathbb{R}^{|\mathcal{S}|}$ , the probability under  $\mathcal{D}_{\text{any}}$  that  $A$  is optimal over  $B$  is  $p_{\mathcal{D}_{\text{any}}}(A \geq B) := \mathbb{P}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}}(\max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \geq \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r})$ .

**Proposition E.20** (Non-dominated linear functionals and their optimality probability). *Let  $A \subseteq \mathbb{R}^{|\mathcal{S}|}$  be finite. If  $\exists b < c : [b, c]^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ , then  $\mathbf{a} \in \text{ND}(A)$  implies that  $\mathbf{a}$  is strictly optimal for a set of reward functions with positive measure under  $\mathcal{D}_{\text{any}}$ .*

*Proof.* Suppose  $\exists b < c : [b, c]^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ . If  $\mathbf{a} \in \text{ND}(A)$ , then let  $\mathbf{r}$  be such that  $\mathbf{a}^\top \mathbf{r} > \max_{\mathbf{a}' \in A \setminus \{\mathbf{a}\}} \mathbf{a}'^\top \mathbf{r}$ . For  $a_1 > 0, a_2 \in \mathbb{R}$ , positively affinely transform  $\mathbf{r}' := a_1 \mathbf{r} + a_2 \mathbf{1}$  (where  $\mathbf{1} \in \mathbb{R}^{|\mathcal{S}|}$  is the all-ones vector) so that  $\mathbf{r}' \in (b, c)^{|\mathcal{S}|}$ .

Note that  $\mathbf{a}$  is still strictly optimal for  $\mathbf{r}'$ :

$$\mathbf{a}^\top \mathbf{r} > \max_{\mathbf{a}' \in A \setminus \{\mathbf{a}\}} \mathbf{a}'^\top \mathbf{r} \iff \mathbf{a}^\top \mathbf{r}' > \max_{\mathbf{a}' \in A \setminus \{\mathbf{a}\}} \mathbf{a}'^\top \mathbf{r}'. \quad (29)$$

Furthermore, by the continuity of both terms on the right-hand side of eq. (29),  $\mathbf{a}$  is strictly optimal for reward functions in some open neighborhood  $N$  of  $\mathbf{r}'$ . Let  $N' := N \cap (b, c)^{|\mathcal{S}|}$ .  $N'$  is still open in  $\mathbb{R}^{|\mathcal{S}|}$  since it is the intersection of two open sets  $N$  and  $(b, c)^{|\mathcal{S}|}$ .

$\mathcal{D}_{\text{any}}$  must assign positive probability measure to all open sets in its support; otherwise, its support would exclude these zero-measure sets by definition E.18. Therefore,  $\mathcal{D}_{\text{any}}$  assigns positive probability to  $N' \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ .  $\square$

**Lemma E.21** (Expected value of similar linear functional sets). *Let  $A, B \subseteq \mathbb{R}^{|\mathcal{S}|}$  be finite, let  $A'$  be such that  $\text{ND}(A) \subseteq A' \subseteq A$ , and let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be an increasing function. If  $B$  contains a copy  $B'$  of  $A'$  via  $\phi$ , then*

$$\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \leq \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]. \quad (30)$$

*If  $\text{ND}(B) \setminus B'$  is empty, then eq. (30) is an equality. If  $\text{ND}(B) \setminus B'$  is non-empty,  $g$  is strictly increasing, and  $\exists b < c : (b, c)^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{\text{bound}})$ , then eq. (30) is strict.*

*Proof.* Because  $g : \mathbb{R} \rightarrow \mathbb{R}$  is increasing, it is measurable (as is  $\max$ ). Therefore, the relevant expectations exist for all  $\mathcal{D}_{\text{bound}}$ .

$$\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A'} \mathbf{a}^\top \mathbf{r} \right) \right] \quad (31)$$

$$= \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in \phi \cdot A'} \mathbf{a}^\top \mathbf{r} \right) \right] \quad (32)$$

$$= \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B'} \mathbf{b}^\top \mathbf{r} \right) \right] \quad (33)$$

$$\leq \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]. \quad (34)$$

Equation (31) holds because  $\forall \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} : \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} = \max_{\mathbf{a} \in A'} \mathbf{a}^\top \mathbf{r}$  by lemma E.12's item 2 with  $X := A, X' := A'$ . Equation (32) holds by lemma E.17. Equation (33) holds by the definition of  $B'$ . Furthermore, our assumption on  $\phi$  guarantees that  $B' \subseteq B$ . Therefore,  $\max_{\mathbf{b} \in B'} \mathbf{b}^\top \mathbf{r} \leq \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r}$ , and so eq. (34) holds by the fact that  $g$  is an increasing function. Then eq. (30) holds.

If  $\text{ND}(B) \setminus B'$  is empty, then  $\text{ND}(B) \subseteq B'$ . By assumption,  $B' \subseteq B$ . Then apply lemma E.12 item 2 with  $X := B, X' := B'$  in order to conclude that eq. (34) is an equality. Then eq. (30) is also an equality.

Suppose that  $g$  is strictly increasing,  $\text{ND}(B) \setminus B'$  is non-empty, and  $\exists b < c : (b, c)^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{\text{bound}})$ . Let  $\mathbf{x} \in \text{ND}(B) \setminus B'$ .

$$\mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B'} \mathbf{b}^\top \mathbf{r} \right) \right] < \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in B' \cup \{\mathbf{x}\}} \mathbf{a}^\top \mathbf{r} \right) \right] \quad (35)$$

$$\leq \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]. \quad (36)$$

$\mathbf{x}$  is strictly optimal for a positive-probability subset of  $\text{supp}(\mathcal{D}_{\text{bound}})$  by proposition E.20. Since  $g$  is strictly increasing, eq. (35) is strict. Therefore, we conclude that eq. (30) is strict.  $\square$

**Lemma E.22** (For continuous IID distributions  $\mathcal{D}_{X\text{-IID}}, \exists b < c : (b, c)^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{X\text{-IID}})$ ).

*Proof.*  $\mathcal{D}_{X\text{-IID}} := X^{|\mathcal{S}|}$ . Since the state reward distribution  $X$  is continuous,  $X$  must have support on some open interval  $(b, c)$ . Since  $\mathcal{D}_{X\text{-IID}}$  is IID across states,  $(b, c)^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{X\text{-IID}})$ .  $\square$

**Definition E.23** (Bounded, continuous IID reward).  $\mathcal{D}_{C/B\text{IID}}$  is the set of  $\mathcal{D}_{X\text{-IID}}$  which equal  $X^{|\mathcal{S}|}$  for some continuous, bounded-support distribution  $X$  over  $\mathbb{R}$ .

**Lemma E.24** (Expectation superiority lemma). *Let  $A, B \subseteq \mathbb{R}^{|\mathcal{S}|}$  be finite and let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be an increasing function. If  $B$  contains a copy  $B'$  of  $\text{ND}(A)$  via  $\phi$ , then*

$$\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \leq_{\text{most: } \mathcal{D}_{\text{bound}}} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]. \quad (37)$$

Furthermore, if  $g$  is strictly increasing and  $\text{ND}(B) \setminus \phi \cdot \text{ND}(A)$  is non-empty, then eq. (37) is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B/\text{IID}}$ . In particular,  $\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \not\leq_{\text{most: } \mathfrak{D}_{\text{bound}}} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]$ .

*Proof.* Because  $g : \mathbb{R} \rightarrow \mathbb{R}$  is increasing, it is measurable (as is  $\max$ ). Therefore, the relevant expectations exist for all  $\mathcal{D}_{\text{bound}}$ .

Suppose that  $\mathcal{D}_{\text{bound}}$  is such that  $\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right] < \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right]$ .

$$\mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \leq \mathbb{E}_{\mathbf{r} \sim \phi^2 \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right] \quad (38)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right] \quad (39)$$

$$< \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \quad (40)$$

$$\leq \mathbb{E}_{\mathbf{r} \sim \phi \cdot \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]. \quad (41)$$

Equation (38) follows by applying lemma E.21 with permutation  $\phi$  and  $A' := \text{ND}(A)$ . Equation (39) follows because involutions satisfy  $\phi^{-1} = \phi$ , and  $\phi^2$  is therefore the identity. Equation (40) follows because we assumed that  $\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right] < \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right]$ . Equation (41) follows by applying lemma E.21 with permutation  $\phi$  and  $A' := \text{ND}(A)$ . By lemma E.16, eq. (37) holds.

Suppose  $g$  is strictly increasing and  $\text{ND}(B) \setminus B'$  is non-empty. Let  $\phi' \in S_{|S|}$ .

$$\mathbb{E}_{\mathbf{r} \sim \phi' \cdot \mathcal{D}_{X\text{-IID}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{X\text{-IID}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \quad (42)$$

$$< \mathbb{E}_{\mathbf{r} \sim \phi' \cdot \mathcal{D}_{X\text{-IID}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right] \quad (43)$$

$$= \mathbb{E}_{\mathbf{r} \sim \phi' \cdot \mathcal{D}_{X\text{-IID}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]. \quad (44)$$

Equation (42) and eq. (44) hold because  $\mathcal{D}_{X\text{-IID}}$  distributes reward identically across states:  $\forall \phi_x \in S_{|S|} : \phi_x \cdot \mathcal{D}_{X\text{-IID}} = \mathcal{D}_{X\text{-IID}}$ . By lemma E.22,  $\exists b < c : (b, c)^{|S|} \subseteq \text{supp}(\mathcal{D}_{X\text{-IID}})$ . Therefore, apply lemma E.21 with  $A' := \text{ND}(A)$  to conclude that eq. (43) holds.

Therefore,  $\forall \phi' \in S_{|S|} : \mathbb{E}_{\mathbf{r} \sim \phi' \cdot \mathcal{D}_{X\text{-IID}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] < \mathbb{E}_{\mathbf{r} \sim \phi' \cdot \mathcal{D}_{X\text{-IID}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]$ , and so  $\mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{a} \in A} \mathbf{a}^\top \mathbf{r} \right) \right] \not\leq_{\text{most: } \mathfrak{D}_{\text{bound}}} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ g \left( \max_{\mathbf{b} \in B} \mathbf{b}^\top \mathbf{r} \right) \right]$  by definition 6.5.  $\square$

**Definition E.25** (Indicator function). Let  $L$  be a predicate which takes input  $x$ .  $\mathbb{1}_{L(x)}$  is the function which returns 1 when  $L(x)$  is true, and 0 otherwise.

**Lemma E.26** (Optimality probability inclusion relations). Let  $X, Y \subseteq \mathbb{R}^{|S|}$  be finite and suppose  $Y' \subseteq Y$ .

$$p_{\mathcal{D}_{\text{any}}}(X \geq Y) \leq p_{\mathcal{D}_{\text{any}}}(X \geq Y') \leq p_{\mathcal{D}_{\text{any}}}(X \cup (Y \setminus Y') \geq Y). \quad (45)$$

If  $\exists b < c : (b, c)^{|S|} \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ ,  $X \subseteq Y$ , and  $\text{ND}(Y) \cap (Y \setminus Y')$  is non-empty, then the second inequality is strict.

*Proof.*

$$p_{\mathcal{D}_{\text{any}}}(X \geq Y) := \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y} \mathbf{y}^\top \mathbf{r}} \right] \quad (46)$$

$$\leq \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y'} \mathbf{y}^\top \mathbf{r}} \right] \quad (47)$$

$$\leq \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{x} \in X \cup (Y \setminus Y')} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y'} \mathbf{y}^\top \mathbf{r}} \right] \quad (48)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{x} \in X \cup (Y \setminus Y')} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y' \cup (Y \setminus Y')} \mathbf{y}^\top \mathbf{r}} \right] \quad (49)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{x} \in X \cup (Y \setminus Y')} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y} \mathbf{y}^\top \mathbf{r}} \right] \quad (50)$$

$$=: p_{\mathcal{D}_{\text{any}}}(X \cup (Y \setminus Y') \geq Y). \quad (51)$$

Equation (47) follows because  $\forall \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} : \mathbb{1}_{\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y} \mathbf{y}^\top \mathbf{r}} \leq \mathbb{1}_{\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y'} \mathbf{y}^\top \mathbf{r}}$  since  $Y' \subseteq Y$ ; note that eq. (47) equals  $p_{\mathcal{D}_{\text{any}}}(X \geq Y')$ , and so the first inequality of eq. (45) is shown. Equation (48) holds because  $\forall \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} : \mathbb{1}_{\max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y'} \mathbf{y}^\top \mathbf{r}} \leq \mathbb{1}_{\max_{\mathbf{x} \in X \cup (Y \setminus Y')} \mathbf{x}^\top \mathbf{r} \geq \max_{\mathbf{y} \in Y'} \mathbf{y}^\top \mathbf{r}}$ .

Suppose  $\exists b < c : (b, c)^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ ,  $X \subseteq Y$ , and  $\text{ND}(Y) \cap (Y \setminus Y')$  is non-empty. Let  $\mathbf{y}^* \in \text{ND}(Y) \cap (Y \setminus Y')$ . By proposition E.20,  $\mathbf{y}^*$  is strictly optimal on a subset of  $\text{supp}(\mathcal{D}_{\text{any}})$  with positive measure under  $\mathcal{D}_{\text{any}}$ . In particular, for a set of  $\mathbf{r}^*$  with positive measure under  $\mathcal{D}_{\text{any}}$ , we have  $\mathbf{y}^{*\top} \mathbf{r}^* > \max_{\mathbf{y} \in Y'} \mathbf{y}^\top \mathbf{r}^*$ .

Then eq. (48) is strict, and therefore the second inequality of eq. (45) is strict as well.  $\square$

**Lemma E.27** (Optimality probability of similar linear functional sets). *Let  $A, B, C \subseteq \mathbb{R}^{|\mathcal{S}|}$  be finite, and let  $Z \subseteq \mathbb{R}^{|\mathcal{S}|}$  be such that  $\text{ND}(C) \subseteq Z \subseteq C$ . If  $\text{ND}(A)$  is similar to  $B' \subseteq B$  via  $\phi$  such that  $\phi \cdot (Z \setminus (B \setminus B')) = Z \setminus (B \setminus B')$ , then*

$$p_{\mathcal{D}_{\text{any}}}(A \geq C) \leq p_{\phi \cdot \mathcal{D}_{\text{any}}}(B \geq C). \quad (52)$$

*If  $B' = B$ , then eq. (52) is an equality. If  $\exists b < c : (b, c)^{|\mathcal{S}|} \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ ,  $B' \subseteq C$ , and  $\text{ND}(C) \cap (B \setminus B')$  is non-empty, then eq. (52) is strict.*

*Proof.*

$$p_{\mathcal{D}_{\text{any}}}(A \geq C) = p_{\mathcal{D}_{\text{any}}}(A \geq Z) \quad (53)$$

$$= p_{\mathcal{D}_{\text{any}}}(\text{ND}(A) \geq Z) \quad (54)$$

$$\leq p_{\mathcal{D}_{\text{any}}}(\text{ND}(A) \geq Z \setminus (B \setminus B')) \quad (55)$$

$$= p_{\phi \cdot \mathcal{D}_{\text{any}}}(\phi \cdot \text{ND}(A) \geq \phi \cdot Z \setminus (B \setminus B')) \quad (56)$$

$$= p_{\phi \cdot \mathcal{D}_{\text{any}}}(B' \geq Z \setminus (B \setminus B')) \quad (57)$$

$$\leq p_{\phi \cdot \mathcal{D}_{\text{any}}}(B' \cup (B \setminus B') \geq Z) \quad (58)$$

$$= p_{\phi \cdot \mathcal{D}_{\text{any}}}(B \geq C). \quad (59)$$

Equation (53) and eq. (59) follow by lemma E.12's item 2 with  $X := C$ ,  $X' := Z$ . Similarly, eq. (54) follows by lemma E.12's item 2 with  $X := A$ ,  $X' := \text{ND}(A)$ . Equation (55) follows by applying the first inequality of lemma E.26 with  $X := \text{ND}(A)$ ,  $Y := Z$ ,  $Y' := Z \setminus (B \setminus B')$ . Equation (56) follows by applying lemma E.17 to eq. (53) with permutation  $\phi$ .

Equation (57) follows by our assumptions on  $\phi$ . Equation (58) follows because by applying the second inequality of lemma E.26 with  $X := B'$ ,  $Y := \text{ND}(C)$ ,  $Y' := \text{ND}(C) \setminus (B \setminus B')$ .

Suppose  $B' = B$ . Then  $B \setminus B' = \emptyset$ , and so eq. (55) and eq. (58) are trivially equalities. Then eq. (52) is an equality.

Suppose  $\exists b < c : (b, c)^{|S|} \subseteq \text{supp}(\mathcal{D}_{\text{any}})$ ; note that  $(b, c)^{|S|} \subseteq \text{supp}(\phi \cdot \mathcal{D}_{\text{any}})$ , since such support must be invariant to permutation. Further suppose that  $B' \subseteq C$  and that  $\text{ND}(C) \cap (B \setminus B')$  is non-empty. Then letting  $X := B'$ ,  $Y := Z$ ,  $Y' := Z \setminus (B \setminus B')$  and noting that  $\text{ND}(\text{ND}(Z)) = \text{ND}(Z)$ , apply lemma E.26 to eq. (58) to conclude that eq. (52) is strict.  $\square$

**Lemma E.28** (Optimality probability superiority lemma). *Let  $A, B, C \subseteq \mathbb{R}^{|S|}$  be finite, and let  $Z$  satisfy  $\text{ND}(C) \subseteq Z \subseteq C$ . If  $B$  contains a copy  $B'$  of  $\text{ND}(A)$  via  $\phi$  such that  $\phi \cdot (Z \setminus (B \setminus B')) = Z \setminus (B \setminus B')$ , then  $p_{\mathcal{D}_{\text{any}}}(A \geq C) \leq_{\text{most: } \mathfrak{D}_{\text{any}}} p_{\mathcal{D}_{\text{any}}}(B \geq C)$ .*

*If  $B' \subseteq C$  and  $\text{ND}(C) \cap (B \setminus B')$  is non-empty, then the inequality is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B\text{IID}}$  and  $p_{\mathcal{D}_{\text{any}}}(A \geq C) \not\leq_{\text{most: } \mathfrak{D}_{\text{any}}} p_{\mathcal{D}_{\text{any}}}(B \geq C)$ .*

*Proof.* Suppose  $\mathcal{D}_{\text{any}}$  is such that  $p_{\mathcal{D}_{\text{any}}}(B \geq C) < p_{\mathcal{D}_{\text{any}}}(A \geq C)$ .

$$p_{\phi \cdot \mathcal{D}_{\text{any}}}(A \geq C) = p_{\phi^{-1} \cdot \mathcal{D}_{\text{any}}}(A \geq C) \quad (60)$$

$$\leq p_{\mathcal{D}_{\text{any}}}(B \geq C) \quad (61)$$

$$< p_{\mathcal{D}_{\text{any}}}(A \geq C) \quad (62)$$

$$\leq p_{\phi \cdot \mathcal{D}_{\text{any}}}(B \geq C). \quad (63)$$

Equation (60) holds because  $\phi$  is an involution. Equation (61) and eq. (63) hold by applying lemma E.27 with permutation  $\phi$ . Equation (62) holds by assumption. Therefore,  $p_{\mathcal{D}_{\text{any}}}(A \geq C) \leq_{\text{most: } \mathfrak{D}_{\text{any}}} p_{\mathcal{D}_{\text{any}}}(B \geq C)$  by lemma E.16.

Suppose  $B' \subseteq C$  and  $\text{ND}(C) \cap (B \setminus B')$  is non-empty, and let  $\mathcal{D}_{X\text{-IID}}$  be any continuous distribution which distributes reward independently and identically across states. Let  $\phi' \in S_{|S|}$ .

$$p_{\phi' \cdot \mathcal{D}_{X\text{-IID}}}(A \geq C) = p_{\mathcal{D}_{X\text{-IID}}}(A \geq C) \quad (64)$$

$$< p_{\phi' \cdot \mathcal{D}_{X\text{-IID}}}(B \geq C) \quad (65)$$

$$= p_{\phi' \cdot \mathcal{D}_{X\text{-IID}}}(A \geq C). \quad (66)$$

Equation (64) and eq. (66) hold because  $\mathcal{D}_{X\text{-IID}}$  distributes reward identically across states,  $\forall \phi_x \in S_{|S|} : \phi_x \cdot \mathcal{D}_{X\text{-IID}} = \mathcal{D}_{X\text{-IID}}$ . By lemma E.22,  $\exists b < c : (b, c)^{|S|} \subseteq \text{supp}(\mathcal{D}_{X\text{-IID}})$ . Therefore, apply lemma E.27 to conclude that eq. (65) holds.

Therefore,  $\forall \phi' \in S_{|S|} : p_{\phi' \cdot \mathcal{D}_{X\text{-IID}}}(A \geq C) < p_{\phi' \cdot \mathcal{D}_{X\text{-IID}}}(B \geq C)$ . In particular,  $p_{\mathcal{D}_{\text{any}}}(A \geq C) \not\leq_{\text{most: } \mathfrak{D}_{\text{any}}} p_{\mathcal{D}_{\text{any}}}(B \geq C)$  by definition 6.5.  $\square$

**Lemma E.29** (Limit probability inequalities which hold for most distributions). *Let  $I \subseteq \mathbb{R}$ , let  $\mathfrak{D} \subseteq \Delta(\mathbb{R}^{|S|})$  be closed under permutation, and let  $F_A, F_B, F_C$  be finite sets of vector functions  $I \mapsto \mathbb{R}^{|S|}$ . Let  $\gamma$  be a limit point of  $I$  such that  $f_1(\mathfrak{D}) := \lim_{\gamma^* \rightarrow \gamma} p_{\mathfrak{D}}(F_B(\gamma^*) \geq F_C(\gamma^*))$ ,  $f_2(\mathfrak{D}) := \lim_{\gamma^* \rightarrow \gamma} p_{\mathfrak{D}}(F_A(\gamma^*) \geq F_C(\gamma^*))$  are well-defined for all  $\mathfrak{D} \in \mathfrak{D}$ .*

*Let  $F_Z$  satisfy  $\text{ND}(F_C) \subseteq F_Z \subseteq F_C$ . Suppose  $F_B$  contains a copy of  $F_A$  via  $\phi$  such that  $\phi \cdot (F_Z \setminus (F_B \setminus \phi \cdot F_A)) = F_Z \setminus (F_B \setminus \phi \cdot F_A)$ . Then  $f_2(\mathfrak{D}) \leq_{\text{most: } \mathfrak{D}} f_1(\mathfrak{D})$ .*

*Proof.* Suppose  $\mathfrak{D} \in \mathfrak{D}$  is such that  $f_2(\mathfrak{D}) > f_1(\mathfrak{D})$ .

$$f_2(\phi \cdot \mathfrak{D}) = f_2(\phi^{-1} \cdot \mathfrak{D}) \quad (67)$$

$$:= \lim_{\gamma^* \rightarrow \gamma} p_{\phi^{-1} \cdot \mathfrak{D}}(F_A(\gamma^*) \geq F_C(\gamma^*)) \quad (68)$$

$$\leq \lim_{\gamma^* \rightarrow \gamma} p_{\mathfrak{D}}(F_B(\gamma^*) \geq F_C(\gamma^*)) \quad (69)$$

$$< \lim_{\gamma^* \rightarrow \gamma} p_{\mathfrak{D}}(F_A(\gamma^*) \geq F_C(\gamma^*)) \quad (70)$$



$$\leq \lim_{\gamma^* \rightarrow \gamma} p_{\phi \cdot \mathcal{D}} (F_B(\gamma^*) \geq F_C(\gamma^*)) \quad (71)$$

$$=: f_1(\phi \cdot \mathcal{D}). \quad (72)$$

By the assumption that  $\mathfrak{D}$  is closed under permutation and  $f_2$  is well-defined for all  $\mathcal{D} \in \mathfrak{D}$ ,  $f_2(\phi \cdot \mathcal{D})$  is well-defined. Equation (67) follows since  $\phi = \phi^{-1}$  because  $\phi$  is an involution. For all  $\gamma^* \in I$ , let  $A := F_A(\gamma^*)$ ,  $B := F_B(\gamma^*)$ ,  $C := F_C(\gamma^*)$ ,  $Z := F_Z(\gamma^*)$  (by definition E.13,  $\text{ND}(C) \subseteq Z \subseteq C$ ). Since  $\phi \cdot A \subseteq B$  by assumption, and since  $\text{ND}(A) \subseteq A$ ,  $B$  also contains a copy of  $\text{ND}(A)$  via  $\phi$ . Furthermore,  $\phi \cdot (Z \setminus (B \setminus \phi \cdot A)) = Z \setminus (B \setminus \phi \cdot A)$  (by assumption), and so apply lemma E.27 to conclude that  $p_{\phi^{-1} \cdot \mathcal{D}} (F_A(\gamma^*) \geq F_C(\gamma^*)) \leq p_{\mathcal{D}} (F_B(\gamma^*) \geq F_C(\gamma^*))$ . Therefore, the limit inequality eq. (69) holds. Equation (70) follows because we assumed that  $f_1(\mathcal{D}) < f_2(\mathcal{D})$ . Equation (71) holds by reasoning similar to that given for eq. (69).

Therefore,  $f_2(\mathcal{D}) > f_1(\mathcal{D})$  implies that  $f_2(\phi \cdot \mathcal{D}) < f_1(\phi \cdot \mathcal{D})$ , and so apply lemma E.16 to conclude that  $f_2(\mathcal{D}) \leq_{\text{most: } \mathfrak{D}} f_1(\mathcal{D})$ .  $\square$

### E.1.3 $\mathcal{F}_{\text{nd}}$ results

**Proposition E.30** (How to transfer optimal policy sets across discount rates). *Suppose reward function  $R$  has optimal policy set  $\Pi^*(R, \gamma)$  at discount rate  $\gamma \in (0, 1)$ . For any  $\gamma^* \in (0, 1)$ , we can construct a reward function  $R'$  such that  $\Pi^*(R', \gamma^*) = \Pi^*(R, \gamma)$ . Furthermore,  $V_{R'}^*(\cdot, \gamma^*) = V_R^*(\cdot, \gamma)$ .*

*Proof.* Let  $R$  be any reward function. Suppose  $\gamma^* \in (0, 1)$  and construct  $R'(s) := V_R^*(s, \gamma) - \gamma^* \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)]$ .

Let  $\pi \in \Pi$  be any policy. By the definition of optimal policies,  $\pi \in \Pi^*(R', \gamma^*)$  iff for all  $s$ :

$$R'(s) + \gamma^* \mathbb{E}_{s' \sim T(s, \pi(s))} [V_{R'}^*(s', \gamma^*)] = R'(s) + \gamma^* \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_{R'}^*(s', \gamma^*)] \quad (73)$$

$$R'(s) + \gamma^* \mathbb{E}_{s' \sim T(s, \pi(s))} [V_R^*(s', \gamma)] = R'(s) + \gamma^* \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)] \quad (74)$$

$$\gamma^* \mathbb{E}_{s' \sim T(s, \pi(s))} [V_R^*(s', \gamma)] = \gamma^* \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)] \quad (75)$$

$$\mathbb{E}_{s' \sim T(s, \pi(s))} [V_R^*(s', \gamma)] = \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)]. \quad (76)$$

By the Bellman equations,  $R'(s) = V_{R'}^*(s, \gamma^*) - \gamma^* \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_{R'}^*(s', \gamma^*)]$ . By the definition of  $R'$ ,  $V_{R'}^*(\cdot, \gamma^*) = V_R^*(\cdot, \gamma)$  must be the unique solution to the Bellman equations for  $R'$  at  $\gamma^*$ . Therefore, eq. (74) holds. Equation (75) follows by plugging in  $R' := V_R^*(s, \gamma) - \gamma^* \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)]$  to eq. (74) and doing algebraic manipulation. Equation (76) follows because  $\gamma^* > 0$ .

Equation (76) shows that  $\pi \in \Pi^*(R', \gamma^*)$  iff  $\forall s : \mathbb{E}_{s' \sim T(s, \pi(s))} [V_R^*(s', \gamma)] = \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)]$ . That is,  $\pi \in \Pi^*(R', \gamma^*)$  iff  $\pi \in \Pi^*(R, \gamma)$ .  $\square$

**Definition E.31** (Evaluating sets of visit distribution functions at  $\gamma$ ). For  $\gamma \in (0, 1)$ , define  $\mathcal{F}(s, \gamma) := \{\mathbf{f}(\gamma) \mid \mathbf{f} \in \mathcal{F}(s)\}$  and  $\mathcal{F}_{\text{nd}}(s, \gamma) := \{\mathbf{f}(\gamma) \mid \mathbf{f} \in \mathcal{F}_{\text{nd}}(s)\}$ . If  $F \subseteq \mathcal{F}(s)$ , then  $F(\gamma) := \{\mathbf{f}(\gamma) \mid \mathbf{f} \in F\}$ .

**Lemma E.32** (Non-domination across  $\gamma$  values for expectations of visit distributions). *Let  $\Delta_d \in \Delta(\mathbb{R}^{|S|})$  be any state distribution and let  $F := \{\mathbb{E}_{s_d \sim \Delta_d} [\mathbf{f}^{\pi, s_d}] \mid \pi \in \Pi\}$ .  $\mathbf{f} \in \text{ND}(F)$  iff  $\forall \gamma^* \in (0, 1) : \mathbf{f}(\gamma^*) \in \text{ND}(F(\gamma^*))$ .*

*Proof.* Let  $\mathbf{f}^\pi \in \text{ND}(F)$  be strictly optimal for reward function  $R$  at discount rate  $\gamma \in (0, 1)$ :

$$\mathbf{f}^\pi(\gamma)^\top \mathbf{r} > \max_{\mathbf{f}^{\pi'} \in F \setminus \{\mathbf{f}^\pi\}} \mathbf{f}^{\pi'}(\gamma)^\top \mathbf{r}. \quad (77)$$

Let  $\gamma^* \in (0, 1)$ . By proposition E.30, we can produce  $R'$  such that  $\Pi^*(R', \gamma^*) = \Pi^*(R, \gamma)$ . Since the optimal policy sets are equal, lemma E.1 implies that

$$\mathbf{f}^\pi(\gamma^*)^\top \mathbf{r}' > \max_{\mathbf{f}^{\pi'} \in F \setminus \{\mathbf{f}^\pi\}} \mathbf{f}^{\pi'}(\gamma^*)^\top \mathbf{r}'. \quad (78)$$

Therefore,  $\mathbf{f}^\pi(\gamma^*) \in \text{ND}(F(\gamma^*))$ .

The reverse direction follows by the definition of  $\text{ND}(F)$ .  $\square$

**Lemma E.33** ( $\forall \gamma \in (0, 1) : \mathbf{d} \in \mathcal{F}_{\text{nd}}(s, \gamma)$  iff  $\mathbf{d} \in \text{ND}(\mathcal{F}(s, \gamma))$ ).

*Proof.* By definition E.31,  $\mathcal{F}_{\text{nd}}(s, \gamma) := \{\mathbf{f}(\gamma) \mid \mathbf{f} \in \text{ND}(\mathcal{F}(s))\}$ . By applying lemma E.32 with  $\Delta_d := \mathbf{e}_s$ ,  $\mathbf{f} \in \text{ND}(\mathcal{F}(s))$  iff  $\forall \gamma \in (0, 1) : \mathbf{f}(\gamma) \in \text{ND}(\mathcal{F}(s, \gamma))$ .  $\square$

**Lemma E.34** ( $\forall \gamma \in [0, 1) : V_R^*(s, \gamma) = \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \mathbf{f}(\gamma)^\top \mathbf{r}$ ).

*Proof.*  $\text{ND}(\mathcal{F}(s, \gamma)) = \mathcal{F}_{\text{nd}}(s, \gamma)$  by lemma E.33, so apply corollary E.11 with  $X := \mathcal{F}(s, \gamma)$ .  $\square$

## E.2 Some actions have greater probability of being optimal

**Lemma E.35** (Optimal policy shift bound). *For fixed  $R$ ,  $\Pi^*(R, \gamma)$  can take on at most  $(2|\mathcal{S}| + 1) \sum_s \binom{|\mathcal{F}(s)|}{2}$  distinct values over  $\gamma \in (0, 1)$ .*

*Proof.* By lemma E.1,  $\Pi^*(R, \gamma)$  changes value iff there is a change in optimality status for some visit distribution function at some state. Lippman [1968] showed that two visit distribution functions can trade off optimality status at most  $2|\mathcal{S}| + 1$  times. At each state  $s$ , there are  $\binom{|\mathcal{F}(s)|}{2}$  such pairs.  $\square$

**Proposition E.36** (Optimality probability's limits exist). *Let  $F \subseteq \mathcal{F}(s)$ .  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(F, 0) = \lim_{\gamma \rightarrow 0} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma)$  and  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(F, 1) = \lim_{\gamma \rightarrow 1} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma)$ .*

*Proof.* First consider the limit as  $\gamma \rightarrow 1$ . Let  $\mathcal{D}_{\text{any}}$  have probability measure  $F_{\text{any}}$ , and define  $\delta(\gamma) := F_{\text{any}}(\{\{R \in \mathbb{R}^{\mathcal{S}} \mid \exists \gamma^* \in [\gamma, 1) : \Pi^*(R, \gamma^*) \neq \Pi^*(R, 1)\}\})$ . Since  $F_{\text{any}}$  is a probability measure,  $\delta(\gamma)$  is bounded  $[0, 1]$ , and  $\delta(\gamma)$  is monotone decreasing. Therefore,  $\lim_{\gamma \rightarrow 1} \delta(\gamma)$  exists.

If  $\lim_{\gamma \rightarrow 1} \delta(\gamma) > 0$ , then there exist reward functions whose optimal policy sets  $\Pi^*(R, \gamma)$  never converge (in the discrete topology on sets) to  $\Pi^*(R, 1)$ , contradicting lemma E.35. So  $\lim_{\gamma \rightarrow 1} \delta(\gamma) = 0$ .

By the definition of optimality probability (definition 4.3) and of  $\delta(\gamma)$ ,  $|\mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma) - \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, 1)| \leq \delta(\gamma)$ . Since  $\lim_{\gamma \rightarrow 1} \delta(\gamma) = 0$ ,  $\lim_{\gamma \rightarrow 1} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma) = \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, 1)$ .

A similar proof shows that  $\lim_{\gamma \rightarrow 0} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma) = \mathbb{P}_{\mathcal{D}_{\text{any}}}(F, 0)$ .  $\square$

**Lemma E.37** (Optimality probability identity). *Let  $\gamma \in (0, 1)$  and let  $F \subseteq \mathcal{F}(s)$ .*

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma) = p_{\mathcal{D}'}(F(\gamma) \geq \mathcal{F}(s, \gamma)) = p_{\mathcal{D}'}(F(\gamma) \geq \mathcal{F}_{\text{nd}}(s, \gamma)). \quad (79)$$

*Proof.* Let  $\gamma \in (0, 1)$ .

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(F, \gamma) := \mathbb{P}_{R \sim \mathcal{D}_{\text{any}}}(F(\gamma) \geq \mathcal{F}(s, \gamma)) \quad (80)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{f} \in F} \mathbf{f}(\gamma)^\top \mathbf{r} = \max_{\mathbf{f}' \in \mathcal{F}(s)} \mathbf{f}'(\gamma)^\top \mathbf{r}} \right] \quad (81)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{any}}} \left[ \mathbb{1}_{\max_{\mathbf{f} \in F} \mathbf{f}(\gamma)^\top \mathbf{r} = \max_{\mathbf{f}' \in \mathcal{F}_{\text{nd}}(s)} \mathbf{f}'(\gamma)^\top \mathbf{r}} \right] \quad (82)$$

$$=: p_{\mathcal{D}'} (F(\gamma) \geq \mathcal{F}_{\text{nd}}(s, \gamma)). \quad (83)$$

Equation (81) follows because lemma E.1 shows that  $\pi$  is optimal iff it induces an optimal visit distribution  $\mathbf{f}$  at every state. Equation (82) follows because  $\forall \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|} : \max_{\mathbf{f}' \in \mathcal{F}(s)} \mathbf{f}'(\gamma)^\top \mathbf{r} = \max_{\mathbf{f}' \in \mathcal{F}_{\text{nd}}(s)} \mathbf{f}'(\gamma)^\top \mathbf{r}$  by lemma E.34.  $\square$

### E.3 Basic properties of POWER

**Lemma E.38** (POWER identities). *Let  $\gamma \in (0, 1)$ .*

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \frac{1 - \gamma}{\gamma} (\mathbf{f}(\gamma) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (84)$$

$$= \frac{1 - \gamma}{\gamma} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} [V_R^*(s, \gamma) - R(s)] \quad (85)$$

$$= \frac{1 - \gamma}{\gamma} \left( V_{\mathcal{D}_{\text{bound}}}^*(s, \gamma) - \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [R(s)] \right) \quad (86)$$

$$= \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \mathbb{E}_{s' \sim T(s, \pi(s))} \left[ (1 - \gamma) V_R^\pi(s', \gamma) \right] \right]. \quad (87)$$

*Proof.*

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) := \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}(s)} \frac{1 - \gamma}{\gamma} (\mathbf{f}(\gamma) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (88)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \frac{1 - \gamma}{\gamma} (\mathbf{f}(\gamma) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (89)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}(s)} \frac{1 - \gamma}{\gamma} (\mathbf{f}(\gamma) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (90)$$

$$= \frac{1 - \gamma}{\gamma} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} [V_R^*(s, \gamma) - R(s)] \quad (91)$$

$$= \frac{1 - \gamma}{\gamma} \left( V_{\mathcal{D}_{\text{bound}}}^*(s, \gamma) - \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [R(s)] \right) \quad (92)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \mathbb{E}_{s' \sim T(s, \pi(s))} \left[ (1 - \gamma) \mathbf{f}^{\pi, s'}(\gamma)^\top \mathbf{r} \right] \right] \quad (93)$$

$$= \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \mathbb{E}_{s' \sim T(s, \pi(s))} \left[ (1 - \gamma) V_R^\pi(s', \gamma) \right] \right]. \quad (94)$$

Equation (89) follows from lemma E.34. Equation (91) follows from the dual formulation of optimal value functions. Equation (92) holds by the definition of  $V_{\mathcal{D}_{\text{bound}}}^*(s, \gamma)$  (definition 5.1). Equation (93) holds because  $\mathbf{f}^{\pi, s}(\gamma) = \mathbf{e}_s + \gamma \mathbb{E}_{s' \sim T(s, \pi(s))} [\mathbf{f}^{\pi, s'}(\gamma)]$  by the definition of a visit distribution function (definition 3.3).  $\square$

**Definition E.39** (Discount-normalized value function). Let  $\pi$  be a policy,  $R$  a reward function, and  $s$  a state. For  $\gamma \in [0, 1]$ ,  $V_{R, \text{norm}}^\pi(s, \gamma) := \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) V_R^\pi(s, \gamma^*)$ .

**Lemma E.40** (Normalized value functions have uniformly bounded derivative). *There exists  $K \geq 0$  such that for all reward functions  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$ ,  $\sup_{s \in \mathcal{S}, \pi \in \Pi, \gamma \in [0,1]} \left| \frac{d}{d\gamma} V_{R, \text{norm}}^\pi(s, \gamma) \right| \leq K \|\mathbf{r}\|_1$ .*

*Proof.* Let  $\pi$  be any policy,  $s$  a state, and  $R$  a reward function. Since  $V_{R, \text{norm}}^\pi(s, \gamma) = \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}^{\pi, s}(\gamma^*)^\top \mathbf{r}$ ,  $\frac{d}{d\gamma} V_{R, \text{norm}}^\pi(s, \gamma)$  is controlled by the behavior of  $\lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}^{\pi, s}(\gamma^*)$ . We show that this function's gradient is bounded in infinity norm.

By lemma E.4,  $\mathbf{f}^{\pi, s}(\gamma)$  is a multivariate rational function on  $\gamma$ . Therefore, for any state  $s'$ ,  $\mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{e}_{s'} = \frac{P(\gamma)}{Q(\gamma)}$  in reduced form. By proposition E.3,  $0 \leq \mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{e}_{s'} \leq \frac{1}{1-\gamma}$ . Thus,  $Q$  may only have a root of multiplicity 1 at  $\gamma = 1$ , and  $Q(\gamma) \neq 0$  for  $\gamma \in [0, 1)$ . Let  $f_{s'}(\gamma) := (1 - \gamma) \mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{e}_{s'}$ .

If  $Q(1) \neq 0$ , then the derivative  $f'_{s'}(\gamma)$  is bounded on  $\gamma \in [0, 1)$  because the polynomial  $(1 - \gamma)P(\gamma)$  cannot diverge on a bounded domain.

If  $Q(1) = 0$ , then factor out the root as  $Q(\gamma) = (1 - \gamma)Q^*(\gamma)$ .

$$f'_{s'}(\gamma) = \frac{d}{d\gamma} \left( \frac{(1 - \gamma)P(\gamma)}{Q(\gamma)} \right) \quad (95)$$

$$= \frac{d}{d\gamma} \left( \frac{P(\gamma)}{Q^*(\gamma)} \right) \quad (96)$$

$$= \frac{P'(\gamma)Q^*(\gamma) - (Q^*)'(\gamma)P(\gamma)}{(Q^*(\gamma))^2}. \quad (97)$$

Since  $Q^*(\gamma)$  is a polynomial with no roots on  $\gamma \in [0, 1]$ ,  $f'_{s'}(\gamma)$  is bounded on  $\gamma \in [0, 1)$ .

Therefore, whether or not  $Q(\gamma)$  has a root at  $\gamma = 1$ ,  $f'_{s'}(\gamma)$  is bounded on  $\gamma \in [0, 1)$ . Furthermore,  $\sup_{\gamma \in [0, 1)} \|\nabla(1 - \gamma) \mathbf{f}^{\pi, s}(\gamma)\|_\infty = \sup_{\gamma \in [0, 1)} \max_{s' \in \mathcal{S}} |f'_{s'}(\gamma)|$  is finite since there are only finitely many states.

There are finitely many  $\pi \in \Pi$ , and finitely many states  $s$ , and so there exists some  $K'$  such that  $\sup_{\substack{s \in \mathcal{S}, \\ \pi \in \Pi, \gamma \in [0, 1)}} \|\nabla(1 - \gamma) \mathbf{f}^{\pi, s}(\gamma)\|_\infty \leq K'$ . Then  $\|\nabla(1 - \gamma) \mathbf{f}^{\pi, s}(\gamma)\|_1 \leq |\mathcal{S}| K' =: K$ .

$$\sup_{\substack{s \in \mathcal{S}, \\ \pi \in \Pi, \gamma \in [0, 1)}} \left| \frac{d}{d\gamma} V_{R, \text{norm}}^\pi(s, \gamma) \right| := \sup_{\substack{s \in \mathcal{S}, \\ \pi \in \Pi, \gamma \in [0, 1)}} \left| \frac{d}{d\gamma} \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) V_R^\pi(s, \gamma^*) \right| \quad (98)$$

$$= \sup_{\substack{s \in \mathcal{S}, \\ \pi \in \Pi, \gamma \in [0, 1)}} \left| \frac{d}{d\gamma} (1 - \gamma) V_R^\pi(s, \gamma) \right| \quad (99)$$

$$= \sup_{\substack{s \in \mathcal{S}, \\ \pi \in \Pi, \gamma \in [0, 1)}} \left| \nabla(1 - \gamma) \mathbf{f}^{\pi, s}(\gamma)^\top \mathbf{r} \right| \quad (100)$$

$$\leq \sup_{\substack{s \in \mathcal{S}, \\ \pi \in \Pi, \gamma \in [0, 1)}} \|\nabla(1 - \gamma) \mathbf{f}^{\pi, s}(\gamma)\|_1 \|\mathbf{r}\|_1 \quad (101)$$

$$\leq K \|\mathbf{r}\|_1. \quad (102)$$

Equation (99) holds because  $V_R^\pi(s, \gamma)$  is continuous on  $\gamma \in [0, 1)$  by corollary E.5. Equation (101) holds by the Cauchy-Schwarz inequality.

Since  $\left| \frac{d}{d\gamma} V_{R, \text{norm}}^\pi(s, \gamma) \right|$  is bounded for all  $\gamma \in [0, 1)$ , eq. (102) also holds for  $\gamma \rightarrow 1$ .  $\square$

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

*Proof.* Let  $b, c$  be such that  $\text{supp}(\mathcal{D}_{\text{bound}}) \subseteq [b, c]^{|\mathcal{S}|}$ . For any  $\mathbf{r} \in \text{supp}(\mathcal{D}_{\text{bound}})$  and  $\pi \in \Pi$ ,  $V_{R, \text{norm}}^\pi(s, \gamma)$  has Lipschitz constant  $K \|\mathbf{r}\|_1 \leq K |\mathcal{S}| \|\mathbf{r}\|_\infty \leq K |\mathcal{S}| \max(|c|, |b|)$  on  $\gamma \in (0, 1)$  by lemma E.40.

For  $\gamma \in (0, 1)$ ,  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) = \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \mathbb{E}_{s' \sim T(s, \pi(s))} \left[ (1 - \gamma) V_R^\pi(s', \gamma) \right] \right]$  by eq. (94). The expectation of the maximum of a set of functions which share a Lipschitz constant, also shares the Lipschitz constant. This shows that  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$  is Lipschitz continuous on  $\gamma \in (0, 1)$ . Thus, its limits are well-defined as  $\gamma \rightarrow 0$  and  $\gamma \rightarrow 1$ . So it is Lipschitz continuous on the closed unit interval.  $\square$

**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 if  $s$  can deterministically reach all states in one step and all states are 1-cycles.

*Proof.* Let  $\gamma \in (0, 1)$ .

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) = \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \mathbb{E}_{s' \sim T(s, \pi(s))} \left[ (1 - \gamma) V_R^*(s', \gamma) \right] \right] \quad (103)$$

$$\leq \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \mathbb{E}_{s' \sim T(s, \pi(s))} \left[ (1 - \gamma) \frac{\max_{s'' \in \mathcal{S}} R(s'')}{1 - \gamma} \right] \right] \quad (104)$$

$$= \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{s'' \in \mathcal{S}} R(s'') \right]. \quad (105)$$

Equation (103) follows from lemma E.38. Equation (104) follows because  $V_R^*(s', \gamma) \leq \frac{\max_{s'' \in \mathcal{S}} R(s'')}{1 - \gamma}$ , as no policy can do better than achieving maximal reward at each time step. Taking limits, the inequality holds for all  $\gamma \in [0, 1]$ .

Suppose that  $s$  can deterministically reach all states in one step and all states are 1-cycles. Then eq. (104) is an equality for all  $\gamma \in (0, 1)$ , since for each  $R$ , the agent can select an action which deterministically transitions to a state with maximal reward. Thus the equality holds for all  $\gamma \in [0, 1]$ .  $\square$

**Lemma E.41 (Lower bound on current POWER based on future POWER).**

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \geq (1 - \gamma) \min_{\substack{a \\ R \sim \mathcal{D}_{\text{bound}}}} \mathbb{E}_{s' \sim T(s, a)} [R(s')] + \gamma \max_{a} \mathbb{E}_{s' \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)]. \quad (106)$$

*Proof.* Let  $\gamma \in (0, 1)$  and let  $a^* \in \arg \max_a \mathbb{E}_{s' \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)]$ .

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \quad (107)$$

$$= (1 - \gamma) \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_a \mathbb{E}_{s' \sim T(s, a)} [V_R^*(s', \gamma)] \right] \quad (108)$$

$$\geq (1 - \gamma) \max_a \mathbb{E}_{s' \sim T(s, a)} \left[ \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [V_R^*(s', \gamma)] \right] \quad (109)$$

$$= (1 - \gamma) \max_a \mathbb{E}_{s' \sim T(s, a)} [V_{\mathcal{D}_{\text{bound}}}^*(s', \gamma)] \quad (110)$$

$$= (1 - \gamma) \max_a \mathbb{E}_{s' \sim T(s, a)} \left[ \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [R(s')] + \frac{\gamma}{1 - \gamma} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \right] \quad (111)$$

$$\geq (1 - \gamma) \mathbb{E}_{s' \sim T(s, a^*)} \left[ \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} [R(s')] + \frac{\gamma}{1 - \gamma} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \right] \quad (112)$$

$$\geq (1 - \gamma) \min_{\substack{a \\ R \sim \mathcal{D}_{\text{bound}}}} \mathbb{E}_{s' \sim T(s, a)} [R(s')] + \gamma \mathbb{E}_{s' \sim T(s, a^*)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)]. \quad (113)$$

Equation (108) holds by lemma E.38. Equation (109) follows because  $\mathbb{E}_{x \sim X} [\max_a f(a, x)] \geq \max_a \mathbb{E}_{x \sim X} [f(a, x)]$  by Jensen's inequality, and eq. (111) follows by lemma E.38.

The inequality also holds when we take the limits  $\gamma \rightarrow 0$  or  $\gamma \rightarrow 1$ .  $\square$



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

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

*Proof.* Suppose  $\gamma \in [0, 1]$ . First consider the case where  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \geq \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)$ .

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \geq (1 - \gamma) \min_{\substack{a \\ s_x \sim T(s', a), \\ R \sim \mathcal{D}_{\text{bound}}}} \mathbb{E} [R(s_x)] + \gamma \max_{\substack{a \\ s_x \sim T(s', a)}} \mathbb{E} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_x, \gamma)] \quad (114)$$

$$\geq (1 - \gamma)b + \gamma \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma). \quad (115)$$

Equation (114) follows by lemma E.41. Equation (115) follows because reward is lower-bounded by  $b$  and because  $s'$  can reach  $s$  in one step with probability 1.

$$|\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) - \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)| = \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) - \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \quad (116)$$

$$\leq \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) - ((1 - \gamma)b + \gamma \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)) \quad (117)$$

$$= (1 - \gamma) (\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) - b) \quad (118)$$

$$\leq (1 - \gamma) \left( \mathbb{E}_{R \sim \mathcal{D}_{\text{bound}}} \left[ \max_{s'' \in \mathcal{S}} R(s'') \right] - b \right) \quad (119)$$

$$\leq (1 - \gamma)(c - b). \quad (120)$$

Equation (116) follows because  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \geq \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)$ . Equation (117) follows by eq. (115). Equation (119) follows by proposition 5.4. Equation (120) follows because reward under  $\mathcal{D}_{\text{bound}}$  is upper-bounded by  $c$ .

The case where  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \leq \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)$  is similar, leveraging the fact that  $s$  can also reach  $s'$  in one step with probability 1.  $\square$

## E.4 Seeking POWER is often more probable under optimality

### E.4.1 Keeping options open tends to be POWER-seeking and tends to be optimal

**Definition E.42** (Normalized visit distribution function). Let  $\mathbf{f} : [0, 1] \rightarrow \mathbb{R}^{|\mathcal{S}|}$  be a vector function. For  $\gamma \in [0, 1]$ ,  $\text{NORM}(\mathbf{f}, \gamma) := \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}(\gamma^*)$  (this limit need not exist for arbitrary  $\mathbf{f}$ ). If  $F$  is a set of such  $\mathbf{f}$ , then  $\text{NORM}(F, \gamma) := \{\text{NORM}(\mathbf{f}, \gamma) \mid \mathbf{f} \in F\}$ .

**Remark.**  $\text{RSD}(s) = \text{NORM}(\mathcal{F}(s), 1)$ .

**Lemma E.43** (Normalized visit distribution functions are continuous). *Let  $\Delta_s \in \Delta(\mathcal{S})$  be a state probability distribution, let  $\pi \in \Pi$ , and let  $\mathbf{f}^* := \mathbb{E}_{s \sim \Delta_s} [\mathbf{f}^{\pi, s}]$ .  $\text{NORM}(\mathbf{f}^*, \gamma)$  is continuous on  $\gamma \in [0, 1]$ .*

*Proof.*

$$\text{NORM}(\mathbf{f}^*, \gamma) := \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbb{E}_{s \sim \Delta_s} [\mathbf{f}^{\pi, s}(\gamma^*)] \quad (121)$$

$$= \mathbb{E}_{s \sim \Delta_s} \left[ \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}^{\pi, s}(\gamma^*) \right] \quad (122)$$

$$=: \mathbb{E}_{s \sim \Delta_s} [\text{NORM}(\mathbf{f}^{\pi, s}, \gamma)]. \quad (123)$$

Equation (122) follows because the expectation is over a finite set. Each  $\mathbf{f}^{\pi, s} \in \mathcal{F}(s)$  is continuous on  $\gamma \in [0, 1]$  by lemma E.4, and  $\lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}^{\pi, s}(\gamma^*)$  exists because RSDs are well-defined [Puterman, 2014]. Therefore, each  $\text{NORM}(\mathbf{f}^{\pi, s}, \gamma)$  is continuous on  $\gamma \in [0, 1]$ . Lastly, eq. (123)'s expectation over finitely many continuous functions is itself continuous.  $\square$

**Lemma E.44** (Non-domination of normalized visit distribution functions). *Let  $\Delta_s \in \Delta(\mathcal{S})$  be a state probability distribution and let  $F := \{\mathbb{E}_{s \sim \Delta_s}[\mathbf{f}^{\pi, s}] \mid \pi \in \Pi\}$ . For all  $\gamma \in [0, 1]$ ,  $\text{ND}(\text{NORM}(F, \gamma)) \subseteq \text{NORM}(\text{ND}(F), \gamma)$ , with equality when  $\gamma \in (0, 1)$ .*

*Proof.* Suppose  $\gamma \in (0, 1)$ .

$$\text{ND}(\text{NORM}(F, \gamma)) = \text{ND}((1 - \gamma)F(\gamma)) \quad (124)$$

$$= (1 - \gamma)\text{ND}(F(\gamma)) \quad (125)$$

$$= (1 - \gamma)(\text{ND}(F)(\gamma)) \quad (126)$$

$$= \text{NORM}(\text{ND}(F), \gamma). \quad (127)$$

Equation (124) and eq. (127) follow by the continuity of  $\text{NORM}(\mathbf{f}, \gamma)$  (lemma E.43). Equation (125) follows by lemma E.15 item 1. Equation (126) follows by lemma E.32.

Let  $\gamma = 1$ . Let  $\mathbf{d} \in \text{ND}(\text{NORM}(F, 1))$  be strictly optimal for  $\mathbf{r}^* \in \mathbb{R}^{|\mathcal{S}|}$ . Then let  $F_{\mathbf{d}} \subseteq F$  be the subset of  $\mathbf{f} \in F$  such that  $\text{NORM}(\mathbf{f}, 1) = \mathbf{d}$ .

$$\max_{\mathbf{f} \in F_{\mathbf{d}}} \text{NORM}(\mathbf{f}, 1)^\top \mathbf{r}^* > \max_{\mathbf{f}' \in F \setminus F_{\mathbf{d}}} \text{NORM}(\mathbf{f}', 1)^\top \mathbf{r}^*. \quad (128)$$

Since  $\text{NORM}(\mathbf{f}, 1)$  is continuous at  $\gamma = 1$  (lemma E.43),  $\mathbf{x}^\top \mathbf{r}^*$  is continuous on  $\mathbf{x} \in \mathbb{R}^{|\mathcal{S}|}$ , and  $F$  is finite, eq. (128) holds for some  $\gamma^* \in (0, 1)$  sufficiently close to  $\gamma = 1$ . By lemma E.10, at least one  $\mathbf{f} \in F_{\mathbf{d}}$  is an element of  $\text{ND}(F(\gamma^*))$ . Then by lemma E.32,  $\mathbf{f} \in \text{ND}(F)$ . We conclude that  $\text{ND}(\text{NORM}(F, 1)) \subseteq \text{NORM}(\text{ND}(F), 1)$ .

The case for  $\gamma = 0$  proceeds similarly.  $\square$

**Lemma E.45** (POWER limit identity). *Let  $\gamma \in [0, 1]$ .*

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \lim_{\gamma^* \rightarrow \gamma} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}(\gamma^*) - \mathbf{e}_s)^\top \mathbf{r} \right]. \quad (129)$$

*Proof.* Let  $\gamma \in [0, 1]$ .

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) = \lim_{\gamma^* \rightarrow \gamma} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma^*) \quad (130)$$

$$= \lim_{\gamma^* \rightarrow \gamma} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}(\gamma^*) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (131)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \lim_{\gamma^* \rightarrow \gamma} \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}(\gamma^*) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (132)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \lim_{\gamma^* \rightarrow \gamma} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}(\gamma^*) - \mathbf{e}_s)^\top \mathbf{r} \right]. \quad (133)$$

Equation (130) follows because  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$  is continuous on  $\gamma \in [0, 1]$  by lemma 5.3. Equation (131) follows by lemma E.38.

For  $\gamma^* \in (0, 1)$ , let  $f_{\gamma^*}(\mathbf{r}) := \max_{\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}(\gamma^*) - \mathbf{e}_s)^\top \mathbf{r}$ . For any sequence  $\gamma_n \rightarrow \gamma$ ,  $(f_{\gamma_n})_{n=1}^\infty$  is a sequence of functions which are piecewise linear on  $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$ , which means they are continuous and therefore measurable. Since lemma E.4 shows that each  $\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)$  is multivariate rational on  $\gamma^*$  (and therefore continuous on  $\gamma^*$ ),  $\{f_{\gamma_n}\}_{n=1}^\infty$  converges pointwise to limit function  $f_\gamma$ . Furthermore,  $|V_R^*(s, \gamma_n) - R(s)| \leq \frac{\gamma}{1 - \gamma_n} \|R\|_\infty$ , and so  $|f_{\gamma_n}(\mathbf{r})| = \left| \frac{1 - \gamma_n}{\gamma_n} (V_R^*(s, \gamma_n) - R(s)) \right| \leq g(\mathbf{r}) \leq \|g\|_\infty =: g(\mathbf{r})$ , which is measurable. Therefore, apply Lebesgue's dominated convergence theorem to conclude that eq. (132) holds. Equation (133) holds because  $\max$  is a continuous function.  $\square$

**Lemma E.46** (Lemma for POWER superiority). *Let  $\Delta_1, \Delta_2 \in \Delta(\mathcal{S})$  be state probability distributions. For  $i = 1, 2$ , let  $F_{\Delta_i} := \{\gamma^{-1} \mathbb{E}_{s_i \sim \Delta_i} [\mathbf{f}^{\pi, s_i} - \mathbf{e}_{s_i}] \mid \pi \in \Pi\}$ . Suppose  $F_{\Delta_2}$  contains a copy of  $\text{ND}(F_{\Delta_1})$  via  $\phi$ . Then  $\forall \gamma \in [0, 1] : \mathbb{E}_{s_1 \sim \Delta_1} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_1, \gamma)] \leq_{\text{most: } \mathcal{D}_{\text{bound}}} \mathbb{E}_{s_2 \sim \Delta_2} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_2, \gamma)]$ .*

*If  $\text{ND}(F_{\Delta_2}) \setminus \phi \cdot \text{ND}(F_{\Delta_1})$  is non-empty, then for all  $\gamma \in (0, 1)$ , the inequality is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathcal{D}_{C/B/\text{IID}}$  and  $\mathbb{E}_{s_1 \sim \Delta_1} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_1, \gamma)] \not\leq_{\text{most: } \mathcal{D}_{\text{bound}}} \mathbb{E}_{s_2 \sim \Delta_2} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_2, \gamma)]$ .*

*These results also hold when replacing  $F_{\Delta_i}$  with  $F_{\Delta_i}^* := \{\mathbb{E}_{s_i \sim \Delta_i} [\mathbf{f}^{\pi, s_i}] \mid \pi \in \Pi\}$  for  $i = 1, 2$ .*

*Proof.*

$$\phi \cdot \text{ND}(\text{NORM}(F_{\Delta_1}, \gamma)) \subseteq \phi \cdot \text{NORM}(\text{ND}(F_{\Delta_1}), \gamma) \quad (134)$$

$$:= \left\{ \mathbf{P}_\phi \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}(\gamma^*) \mid \mathbf{f} \in \text{ND}(F_{\Delta_1}) \right\} \quad (135)$$

$$= \left\{ \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{P}_\phi \mathbf{f}(\gamma^*) \mid \mathbf{f} \in \text{ND}(F_{\Delta_1}) \right\} \quad (136)$$

$$= \left\{ \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}(\gamma^*) \mid \mathbf{f} \in F'_{\text{sub}} \right\} \quad (137)$$

$$\subseteq \left\{ \lim_{\gamma^* \rightarrow \gamma} (1 - \gamma^*) \mathbf{f}(\gamma^*) \mid \mathbf{f} \in F_{\Delta_2} \right\} \quad (138)$$

$$=: \text{NORM}(F_{\Delta_2}, \gamma). \quad (139)$$

Equation (134) follows by lemma E.44. Equation (136) follows because  $\mathbf{P}_\phi$  is a continuous linear operator. Equation (138) follows by assumption.

$$\mathbb{E}_{s_1 \sim \Delta_1} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_1, \gamma)] := \mathbb{E}_{\substack{s_1 \sim \Delta_1, \\ \mathbf{r} \sim \mathcal{D}_{\text{bound}}}} \left[ \max_{\pi \in \Pi} \lim_{\gamma^* \rightarrow \gamma} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}^{\pi, s_1}(\gamma^*) - \mathbf{e}_{s_1})^\top \mathbf{r} \right] \quad (140)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \lim_{\gamma^* \rightarrow \gamma} \frac{1 - \gamma^*}{\gamma^*} \mathbb{E}_{s_1 \sim \Delta_1} [\mathbf{f}^{\pi, s_1}(\gamma^*) - \mathbf{e}_{s_1}]^\top \mathbf{r} \right] \quad (141)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{NORM}(F_{\Delta_1}, \gamma)} \mathbf{d}^\top \mathbf{r} \right] \quad (142)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{ND}(\text{NORM}(F_{\Delta_1}, \gamma))} \mathbf{d}^\top \mathbf{r} \right] \quad (143)$$

$$\leq_{\text{most: } \mathcal{D}_{\text{bound}}} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{NORM}(F_{\Delta_2}, \gamma)} \mathbf{d}^\top \mathbf{r} \right] \quad (144)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\pi \in \Pi} \lim_{\gamma^* \rightarrow \gamma} \frac{1 - \gamma^*}{\gamma^*} \mathbb{E}_{s_2 \sim \Delta_2} [\mathbf{f}^{\pi, s_2}(\gamma^*) - \mathbf{e}_{s_2}]^\top \mathbf{r} \right] \quad (145)$$

$$= \mathbb{E}_{\substack{s_2 \sim \Delta_2, \\ \mathbf{r} \sim \mathcal{D}_{\text{bound}}}} \left[ \max_{\pi \in \Pi} \lim_{\gamma^* \rightarrow \gamma} \frac{1 - \gamma^*}{\gamma^*} (\mathbf{f}^{\pi, s_2}(\gamma^*) - \mathbf{e}_{s_2})^\top \mathbf{r} \right] \quad (146)$$

$$=: \mathbb{E}_{s_2 \sim \Delta_2} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_2, \gamma)]. \quad (147)$$

Equation (140) and eq. (147) follow by lemma E.45. Equation (141) and eq. (146) follow because each  $R$  has a stationary deterministic optimal policy  $\pi \in \Pi^*(R, \gamma) \subseteq \Pi$  which simultaneously achieves optimal value at all states. Equation (143) follows by corollary E.11.

Apply lemma E.24 with  $A := \text{NORM}(F_{\Delta_1}, \gamma)$ ,  $B := \text{NORM}(F_{\Delta_2}, \gamma)$ ,  $g$  the identity function, and involution  $\phi$  (satisfying  $\phi \cdot \text{ND}(A) \subseteq B$  by eq. (139)) in order to conclude that eq. (144) holds.

Suppose that  $\text{ND}(F_{\Delta_2}) \setminus \phi \cdot \text{ND}(F_{\Delta_1})$  is non-empty; let  $F'_{\text{sub}} := \phi \cdot \text{ND}(F_{\Delta_1})$ . Lemma E.32 shows that for all  $\gamma \in (0, 1)$ ,  $\text{ND}(F_{\Delta_2}(\gamma)) \setminus F'_{\text{sub}}(\gamma)$  is non-empty. Lemma E.15 item 1 then implies that  $\text{ND}(B) \setminus \phi \cdot A = \frac{1-\gamma}{\gamma} \left( \text{ND}(F_{\Delta_2}(\gamma)) - \mathbf{e}_s \right) \setminus \left( \frac{1-\gamma}{\gamma} F'_{\text{sub}}(\gamma) \right)$  is non-empty. Then lemma E.24 implies that for all  $\gamma \in (0, 1)$ , eq. (144) is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B/\text{IID}}$  and  $\mathbb{E}_{s_1 \sim \Delta_1} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_1, \gamma)] \not\leq_{\text{most: } \mathfrak{D}_{\text{bound}}} \mathbb{E}_{s_2 \sim \Delta_2} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_2, \gamma)]$ .

We show that this result's preconditions holding for  $F_{\Delta_i}^*$  implies the  $F_{\Delta_i}$  preconditions. Suppose  $F_{\Delta_i}^* := \{\mathbb{E}_{s_i \sim \Delta_i} [\mathbf{f}^{\pi, s_i}] \mid \pi \in \Pi\}$  for  $i = 1, 2$  are such that  $F_{\text{sub}}^* := \phi \cdot \text{ND}(F_{\Delta_1}^*) \subseteq F_{\Delta_2}^*$ . In the following, the  $\Delta_i$  are represented as vectors in  $\mathbb{R}^{|S|}$ , and  $\gamma$  is a variable.

$$\phi \cdot \{\gamma \mathbf{f} \mid \mathbf{f} \in \text{ND}(F_{\Delta_1})\} = \phi \cdot \left( \text{ND}(F_{\Delta_1}^* - \Delta_1) \right) \quad (148)$$

$$= \phi \cdot \left( \text{ND}(F_{\Delta_1}^*) - \Delta_1 \right) \quad (149)$$

$$= \left\{ \mathbf{P}_\phi \mathbf{f} - \mathbf{P}_\phi \Delta_1 \mid \mathbf{f} \in \text{ND}(F_{\Delta_1}^*) \right\} \quad (150)$$

$$\subseteq \left\{ \mathbf{f} - \Delta_2 \mid \mathbf{f} \in F_{\Delta_2}^* \right\} \quad (151)$$

$$= \left\{ \gamma \mathbf{f} \mid \mathbf{f} \in F_{\Delta_2} \right\}. \quad (152)$$

Equation (149) follows from lemma E.15 item 2. Since we assumed that  $\phi \cdot \text{ND}(F_{\Delta_1}^*) \subseteq F_{\Delta_2}^*$ ,  $\phi \cdot \{\Delta_1\} = \phi \cdot \left( \text{ND}(F_{\Delta_1}^*)(0) \right) \subseteq F_{\Delta_2}^*(0) = \{\Delta_2\}$ . This implies that  $\mathbf{P}_\phi \Delta_1 = \Delta_2$  and so eq. (151) follows.

Equation (152) shows that  $\phi \cdot \{\gamma \mathbf{f} \mid \mathbf{f} \in \text{ND}(F_{\Delta_1})\} \subseteq \{\gamma \mathbf{f} \mid \mathbf{f} \in F_{\Delta_2}\}$ . But we then have  $\phi \cdot \{\gamma \mathbf{f} \mid \mathbf{f} \in \text{ND}(F_{\Delta_1})\} := \{\gamma \mathbf{P}_\phi \mathbf{f} \mid \mathbf{f} \in \text{ND}(F_{\Delta_1})\} = \{\gamma \mathbf{f} \mid \mathbf{f} \in \phi \cdot \text{ND}(F_{\Delta_1})\} \subseteq \{\gamma \mathbf{f} \mid \mathbf{f} \in F_{\Delta_2}\}$ . Thus,  $\phi \cdot \text{ND}(F_{\Delta_1}) \subseteq F_{\Delta_2}$ .

Suppose  $\text{ND}(F_{\Delta_2}^*) \setminus \phi \cdot \text{ND}(F_{\Delta_1}^*)$  is non-empty, which implies that

$$\phi \cdot \{\gamma \mathbf{f} \mid \mathbf{f} \in \text{ND}(F_{\Delta_1})\} = \left\{ \mathbf{P}_\phi \mathbf{f} - \mathbf{P}_\phi \Delta_1 \mid \mathbf{f} \in \text{ND}(F_{\Delta_1}^*) \right\} \quad (153)$$

$$= \left\{ \mathbf{f} - \mathbf{P}_\phi \Delta_1 \mid \mathbf{f} \in \phi \cdot \text{ND}(F_{\Delta_1}^*) \right\} \quad (154)$$

$$\subsetneq \left\{ \mathbf{f} - \Delta_2 \mid \mathbf{f} \in \text{ND}(F_{\Delta_2}^*) \right\} \quad (155)$$

$$= \left\{ \gamma \mathbf{f} \mid \mathbf{f} \in \text{ND}(F_{\Delta_2}) \right\}. \quad (156)$$

Then  $\text{ND}(F_{\Delta_2}) \setminus \phi \cdot \text{ND}(F_{\Delta_1})$  must be non-empty. Therefore, if the preconditions of this result are met for  $F_{\Delta_i}^*$ , they are met for  $F_{\Delta_i}$ .  $\square$

**Proposition 6.6** (States with “more options” have more POWER). *If  $\mathcal{F}(s)$  contains a copy of  $\mathcal{F}_{\text{nd}}(s')$  via  $\phi$ , then  $\forall \gamma \in [0, 1] : \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma) \geq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma)$ . If  $\mathcal{F}_{\text{nd}}(s) \setminus \phi \cdot \mathcal{F}_{\text{nd}}(s')$  is non-empty, then for all  $\gamma \in (0, 1)$ , the converse  $\leq_{\text{most}}$  statement does not hold.*

*Proof.* Let  $F_{\text{sub}} := \phi \cdot \mathcal{F}_{\text{nd}}(s') \subseteq \mathcal{F}(s)$ . Let  $\Delta_1 := \mathbf{e}_{s'}$ ,  $\Delta_2 := \mathbf{e}_s$ , and define  $F_{\Delta_i}^* := \{\mathbb{E}_{s_i \sim \Delta_i} [\mathbf{f}^{\pi, s_i}] \mid \pi \in \Pi\}$  for  $i = 1, 2$ . Then  $\mathcal{F}_{\text{nd}}(s') = \text{ND}(F_{\Delta_1}^*)$  is similar to  $F_{\text{sub}} = F_{\text{sub}}^* \subseteq F_{\Delta_2}^* = \mathcal{F}(s)$  via involution  $\phi$ . Apply lemma E.46 to conclude that  $\forall \gamma \in [0, 1] : \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \leq_{\text{most: } \mathfrak{D}_{\text{bound}}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$ .

Furthermore,  $\mathcal{F}_{\text{nd}}(s) = \text{ND}(F_{\Delta_2}^*)$ , and  $F_{\text{sub}} = F_{\text{sub}}^*$ , and so if  $\mathcal{F}_{\text{nd}}(s) \setminus \phi \cdot \mathcal{F}_{\text{nd}}(s') := \mathcal{F}_{\text{nd}}(s) \setminus F_{\text{sub}} = \text{ND}(F_{\Delta_2}^*) \setminus F_{\text{sub}}^*$  is non-empty, then lemma E.46 shows that for all  $\gamma \in (0, 1)$ , the inequality is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B/\text{IID}}$  and  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', \gamma) \not\leq_{\text{most: } \mathfrak{D}_{\text{bound}}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, \gamma)$ .  $\square$

**Lemma E.47** (Non-dominated visit distribution functions never agree with other visit distribution functions at that state). *Let  $\mathbf{f} \in \mathcal{F}_{\text{nd}}(s), \mathbf{f}' \in \mathcal{F}(s) \setminus \{\mathbf{f}\}$ .  $\forall \gamma \in (0, 1) : \mathbf{f}(\gamma) \neq \mathbf{f}'(\gamma)$ .*

*Proof.* Let  $\gamma \in (0, 1)$ . Since  $\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)$ , there exists a  $\gamma^* \in (0, 1)$  at which  $\mathbf{f}$  is strictly optimal for some reward function. Then by proposition E.30, we can produce another reward function for which  $\mathbf{f}$  is strictly optimal at discount rate  $\gamma$ ; in particular, proposition E.30 guarantees that the policies which induce  $\mathbf{f}'$  are not optimal at  $\gamma$ . So  $\mathbf{f}(\gamma) \neq \mathbf{f}'(\gamma)$ .  $\square$

**Corollary E.48** (Cardinality of non-dominated visit distributions). *Let  $F \subseteq \mathcal{F}(s)$ .  $\forall \gamma \in (0, 1) : |F \cap \mathcal{F}_{\text{nd}}(s)| = |F(\gamma) \cap \mathcal{F}_{\text{nd}}(s, \gamma)|$ .*

*Proof.* Let  $\gamma \in (0, 1)$ . By applying lemma E.32 with  $\Delta_d := \mathbf{e}_s$ ,  $\mathbf{f} \in \mathcal{F}_{\text{nd}}(s) = \text{ND}(\mathcal{F}(s))$  iff  $\mathbf{f}(\gamma) \in \text{ND}(\mathcal{F}(s, \gamma))$ . By lemma E.33,  $\text{ND}(\mathcal{F}(s, \gamma)) = \mathcal{F}_{\text{nd}}(s, \gamma)$ . So all  $\mathbf{f} \in F \cap \mathcal{F}_{\text{nd}}(s)$  induce  $\mathbf{f}(\gamma) \in F(\gamma) \cap \mathcal{F}_{\text{nd}}(s, \gamma)$ , and  $|F \cap \mathcal{F}_{\text{nd}}(s)| \geq |F(\gamma) \cap \mathcal{F}_{\text{nd}}(s, \gamma)|$ .

Lemma E.47 implies that for all  $\mathbf{f}, \mathbf{f}' \in \mathcal{F}_{\text{nd}}(s)$ ,  $\mathbf{f} = \mathbf{f}'$  iff  $\mathbf{f}(\gamma) = \mathbf{f}'(\gamma)$ . Therefore,  $|F \cap \mathcal{F}_{\text{nd}}(s)| \leq |F(\gamma) \cap \mathcal{F}_{\text{nd}}(s, \gamma)|$ . So  $|F \cap \mathcal{F}_{\text{nd}}(s)| = |F(\gamma) \cap \mathcal{F}_{\text{nd}}(s, \gamma)|$ .  $\square$

**Lemma E.49** (Optimality probability and state bottlenecks). *Suppose that  $s$  can reach  $\text{REACH}(s', a') \cup \text{REACH}(s', a)$ , but only by taking actions equivalent to  $a'$  or  $a$  at state  $s'$ .  $F_{\text{nd}, a'} := \mathcal{F}_{\text{nd}}(s \mid \pi(s') = a')$ ,  $F_a := \mathcal{F}(s \mid \pi(s') = a)$ . Suppose  $F_a$  contains a copy of  $F_{\text{nd}, a'}$  via  $\phi$  which fixes all states not belonging to  $\text{REACH}(s', a') \cup \text{REACH}(s', a)$ . Then  $\forall \gamma \in [0, 1] : \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd}, a'}, \gamma) \leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma)$ .*

*If  $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{\text{nd}, a'})$  is non-empty, then for all  $\gamma \in (0, 1)$ , the inequality is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathcal{D}_{\text{C/B/IID}}$ , and  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd}, a'}, \gamma) \not\leq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma)$ .*

*Proof.* Let  $F_{\text{sub}} := \phi \cdot F_{\text{nd}, a'}$ . Let  $F^* := \bigcup_{\substack{a'' \in \mathcal{A}: \\ (a'' \neq_{s'} a) \wedge (a'' \neq_{s'} a')}} \mathcal{F}(s \mid \pi(s') = a'') \cup F_{\text{nd}, a'} \cup F_{\text{sub}}$ .

$$\phi \cdot F^* := \phi \cdot \left( \bigcup_{\substack{a'' \in \mathcal{A}: \\ (a'' \neq_{s'} a) \wedge (a'' \neq_{s'} a')}} \mathcal{F}(s \mid \pi(s') = a'') \cup F_{\text{nd}, a'} \cup F_{\text{sub}} \right) \quad (157)$$

$$= \bigcup_{\substack{a'' \in \mathcal{A}: \\ (a'' \neq_{s'} a) \wedge (a'' \neq_{s'} a')}} \phi \cdot \mathcal{F}(s \mid \pi(s') = a'') \cup (\phi \cdot F_{\text{nd}, a'}) \cup (\phi \cdot F_{\text{sub}}) \quad (158)$$

$$= \bigcup_{\substack{a'' \in \mathcal{A}: \\ (a'' \neq_{s'} a) \wedge (a'' \neq_{s'} a')}} \phi \cdot \mathcal{F}(s \mid \pi(s') = a'') \cup F_{\text{sub}} \cup F_{\text{nd}, a'} \quad (159)$$

$$= \bigcup_{\substack{a'' \in \mathcal{A}: \\ (a'' \neq_{s'} a) \wedge (a'' \neq_{s'} a')}} \mathcal{F}(s \mid \pi(s') = a'') \cup F_{\text{sub}} \cup F_{\text{nd}, a'} \quad (160)$$

$$=: F^*. \quad (161)$$

Equation (159) follows because the involution  $\phi$  ensures that  $\phi \cdot F_{\text{sub}} = F_{\text{nd}, a'}$ . By assumption,  $\phi$  fixes all  $s' \notin \text{REACH}(s', a') \cup \text{REACH}(s', a)$ . Suppose  $\mathbf{f} \in \mathcal{F}(s) \setminus (F_{\text{nd}, a'} \cup F_a)$ . By the bottleneck assumption,  $\mathbf{f}$  does not visit states in  $\text{REACH}(s', a') \cup \text{REACH}(s', a)$ . Therefore,  $\mathbf{P}_\phi \mathbf{f} = \mathbf{f}$ , and so eq. (160) follows.

Let  $F_Z := (\mathcal{F}(s) \setminus (\mathcal{F}(s \mid \pi(s) = a') \cup F_a)) \cup F_{\text{nd}, a'} \cup F_a$ . By definition,  $F_Z \subseteq \mathcal{F}(s)$ . Furthermore,  $\mathcal{F}_{\text{nd}}(s) = \bigcup_{a'' \in \mathcal{A}} \mathcal{F}_{\text{nd}}(s \mid \pi(s') = a'') \subseteq (\mathcal{F}(s) \setminus (\mathcal{F}(s \mid \pi(s) = a') \cup F_a)) \cup \mathcal{F}_{\text{nd}}(s \mid \pi(s) = a') \cup F_a =: F_Z$ , and so  $\mathcal{F}_{\text{nd}}(s) \subseteq F_Z$ . Note that  $F^* = F_Z \setminus (F_a \setminus F_{\text{sub}})$ .



**Case:**  $\gamma \in (0, 1)$ .

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}, \gamma) = p_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}(\gamma) \geq \mathcal{F}(s, \gamma)) \quad (162)$$

$$\leq_{\text{most: } \mathfrak{D}_{\text{any}}} p_{\mathcal{D}_{\text{any}}}(F_a(\gamma) \geq \mathcal{F}(s, \gamma)) \quad (163)$$

$$= \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}, \gamma). \quad (164)$$

Equation (162) and eq. (164) follow from lemma E.37. Equation (163) follows by applying lemma E.28 with  $A := F_{\text{nd},a'}(\gamma)$ ,  $B' := F_{\text{sub}}(\gamma)$ ,  $B := F_a(\gamma)$ ,  $C := \mathcal{F}(s, \gamma)$ ,  $Z := F_Z(\gamma)$  which satisfies  $\text{ND}(C) = \mathcal{F}_{\text{nd}}(s, \gamma) \subseteq F_Z(\gamma) \subseteq \mathcal{F}(s, \gamma) = C$ , and involution  $\phi$  which satisfies  $\phi \cdot F^*(\gamma) = \phi \cdot (Z \setminus (B \setminus B')) = Z \setminus (B \setminus B') = F^*(\gamma)$ .

Suppose  $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus F_{\text{sub}})$  is non-empty.  $0 < \left| \mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus F_{\text{sub}}) \right| = \left| \mathcal{F}_{\text{nd}}(s, \gamma) \cap (F_a(\gamma) \setminus F_{\text{sub}}(\gamma)) \right| =: \left| \text{ND}(C) \cap (B \setminus B') \right|$  (with the first equality holding by corollary E.48), and so  $\text{ND}(C) \cap (B \setminus B')$  is non-empty. We also have  $B := F_a(\gamma) \subseteq \mathcal{F}(s, \gamma) =: C$ . Then reapplying lemma E.28, eq. (163) is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B/\text{IID}}$ , and  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}, \gamma) \not\leq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma)$ .

**Case:**  $\gamma = 1, \gamma = 0$ .

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}, 1) = \lim_{\gamma^* \rightarrow 1} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}, \gamma^*) \quad (165)$$

$$= \lim_{\gamma^* \rightarrow 1} p_{\mathcal{D}_{\text{any}}}(F_{\text{nd},a'}(\gamma^*) \geq \mathcal{F}(s, \gamma^*)) \quad (166)$$

$$\leq_{\text{most: } \mathfrak{D}_{\text{any}}} \lim_{\gamma^* \rightarrow 1} p_{\mathcal{D}_{\text{any}}}(F_a(\gamma^*) \geq \mathcal{F}(s, \gamma^*)) \quad (167)$$

$$= \lim_{\gamma^* \rightarrow 1} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma^*) \quad (168)$$

$$= \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, 1). \quad (169)$$

Equation (165) and eq. (169) hold by proposition E.36. Equation (166) and eq. (168) follow by lemma E.37. Applying lemma E.29 with  $\gamma := 1$ ,  $I := (0, 1)$ ,  $F_A := F_{\text{nd},a'}$ ,  $F_B := F_a$ ,  $F_C := \mathcal{F}(s)$ ,  $F_Z$  as defined above, and involution  $\phi$  (for which  $\phi \cdot (F_Z \setminus (F_B \setminus \phi \cdot F_A)) = F_Z \setminus (F_B \setminus \phi \cdot F_A)$ ), we conclude that eq. (167) follows.

The  $\gamma = 0$  case proceeds similarly to  $\gamma = 1$ .  $\square$

**Lemma E.50** (Action optimality probability is a special case of visit distribution optimality probability).  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma) = \mathbb{P}_{\mathcal{D}_{\text{any}}}(\mathcal{F}(s \mid \pi(s) = a), \gamma)$ .

*Proof.* Let  $F_a := \mathcal{F}(s \mid \pi(s) = a)$ . For  $\gamma \in (0, 1)$ ,

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma) := \mathbb{P}_{R \sim \mathcal{D}_{\text{any}}}( \exists \pi^* \in \Pi^*(R, \gamma) : \pi^*(s) = a ) \quad (170)$$

$$= \mathbb{P}_{r \sim \mathcal{D}_{\text{any}}}\left( \exists \mathbf{f}^{\pi^*, s} \in F_a : \mathbf{f}^{\pi^*, s}(\gamma)^\top \mathbf{r} = \max_{\mathbf{f} \in \mathcal{F}(s)} \mathbf{f}(\gamma)^\top \mathbf{r} \right) \quad (171)$$

$$= \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma). \quad (172)$$

By lemma E.1, if  $\exists \pi^* \in \Pi^*(R, \gamma) : \pi^*(s) = a$ , then it induces some optimal  $\mathbf{f}^{\pi^*, s} \in F_a$ . Conversely, if  $\mathbf{f}^{\pi^*, s} \in F_a$  is optimal at  $\gamma \in (0, 1)$ , then  $\pi^*$  chooses optimal actions on the support of  $\mathbf{f}^{\pi^*, s}(\gamma)$ . Let  $\pi'$  agree with  $\pi^*$  on that support and let  $\pi'$  take optimal actions at all other states. Then  $\pi' \in \Pi^*(R, \gamma)$  and  $\pi'(s) = a$ . So eq. (171) follows.

Suppose  $\gamma = 0$  or  $\gamma = 1$ . Consider any sequence  $(\gamma_n)_{n=1}^\infty$  converging to  $\gamma$ , and let  $\mathcal{D}_{\text{any}}$  induce probability measure  $F$ .

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma) := \lim_{\gamma^* \rightarrow \gamma} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma^*) \quad (173)$$

$$= \lim_{\gamma^* \rightarrow \gamma} \mathbb{P}_{R \sim \mathcal{D}_{\text{any}}} (\exists \pi^* \in \Pi^*(R, \gamma^*) : \pi^*(s) = a) \quad (174)$$

$$= \lim_{n \rightarrow \infty} \mathbb{P}_{R \sim \mathcal{D}_{\text{any}}} (\exists \pi^* \in \Pi^*(R, \gamma_n) : \pi^*(s) = a) \quad (175)$$

$$= \lim_{n \rightarrow \infty} \int_{\mathbb{R}^S} \mathbb{1}_{\exists \pi^* \in \Pi^*(R, \gamma_n) : \pi^*(s) = a} dF(R) \quad (176)$$

$$= \int_{\mathbb{R}^S} \lim_{n \rightarrow \infty} \mathbb{1}_{\exists \pi^* \in \Pi^*(R, \gamma_n) : \pi^*(s) = a} dF(R) \quad (177)$$

$$= \int_{\mathbb{R}^S} \mathbb{1}_{\exists \pi^* \in \Pi^*(R, \gamma) : \pi^*(s) = a} dF(R) \quad (178)$$

$$=: \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma). \quad (179)$$

Equation (174) follows by eq. (172). for  $\gamma^* \in [0, 1]$ , let  $f_{\gamma^*}(R) := \mathbb{1}_{\exists \pi^* \in \Pi^*(R, \gamma^*) : \pi^*(s) = a}$ . For each  $R \in \mathbb{R}^S$ , lemma E.35 exists  $\gamma_x \approx \gamma$  such that for all intermediate  $\gamma'_x$  between  $\gamma_x$  and  $\gamma$ ,  $\Pi^*(R, \gamma'_x) = \Pi^*(R, \gamma)$ . Since  $\gamma_n \rightarrow \gamma$ , this means that  $(f_{\gamma_n})_{n=1}^{\infty}$  converges pointwise to  $f_{\gamma}$ . Furthermore,  $\forall n \in \mathbb{N}, R \in \mathbb{R}^S : |f_{\gamma_n}(R)| \leq 1$  by definition. Therefore, eq. (177) follows by Lebesgue's dominated convergence theorem.  $\square$

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

Suppose  $F_a := \mathcal{F}(s \mid \pi(s) = a)$  contains a copy of  $F_{a'} := \mathcal{F}(s \mid \pi(s) = a')$  via  $\phi$ .

1. If  $s \notin \text{REACH}(s, a')$ , then  $\forall \gamma \in [0, 1] : \mathbb{E}_{s_a \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_a, \gamma)] \geq_{\text{most}} \mathfrak{D}_{\text{bound}} \mathbb{E}_{s_{a'} \sim T(s, a')} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_{a'}, \gamma)]$ .
2. If  $s$  can only reach the states of  $\text{REACH}(s, a') \cup \text{REACH}(s, a)$  by taking actions equivalent to  $a'$  or  $a$  at state  $s$ , then  $\forall \gamma \in [0, 1] : \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma) \geq_{\text{most}} \mathfrak{D}_{\text{any}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a', \gamma)$ .

If  $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'})$  is non-empty, then  $\forall \gamma \in (0, 1)$ , the converse  $\leq_{\text{most}}$  statements do not hold.

*Proof.* Note that by definition 3.3,  $F_{a'}(0) = \{\mathbf{e}_s\} = F_a(0)$ . Since  $\phi \cdot F_{a'} \subseteq F_a$ , in particular we have  $\phi \cdot F_{a'}(0) = \{\mathbf{P}_{\phi} \mathbf{e}_s\} \subseteq \{\mathbf{e}_s\} = F_a(0)$ , and so  $\phi(s) = s$ .

**Item 1.** For state probability distribution  $\Delta_s \in \Delta(\mathcal{S})$ , let  $F_{\Delta_s}^* := \left\{ \mathbb{E}_{s' \sim \Delta_s} [\mathbf{f}^{\pi, s'}] \mid \pi \in \Pi \right\}$ .

Unless otherwise stated, we treat  $\gamma$  as a variable in this item; we apply element-wise vector addition, constant multiplication, and variable multiplication via the conventions outlined in definition E.14.

$$F_{a'} = \left\{ \mathbf{e}_s + \gamma \mathbb{E}_{s_{a'} \sim T(s, a')} [\mathbf{f}^{\pi, s_{a'}}] \mid \pi \in \Pi : \pi(s) = a' \right\} \quad (180)$$

$$= \left\{ \mathbf{e}_s + \gamma \mathbb{E}_{s_{a'} \sim T(s, a')} [\mathbf{f}^{\pi, s_{a'}}] \mid \pi \in \Pi \right\} \quad (181)$$

$$= \mathbf{e}_s + \gamma F_{T(s, a')}^*. \quad (182)$$

Equation (180) follows by definition 3.3, since each  $\mathbf{f} \in \mathcal{F}(s)$  has an initial term of  $\mathbf{e}_s$ . Equation (181) follows because  $s \notin \text{REACH}(s, a')$ , and so for all  $s_{a'} \in \text{supp}(T(s, a'))$ ,  $\mathbf{f}^{\pi, s_{a'}}$  is unaffected by the choice of action  $\pi(s)$ . Note that similar reasoning implies that  $F_a \subseteq \mathbf{e}_s + \gamma F_{T(s, a)}^*$  (because eq. (181) is a containment relation in general).

Since  $F_{a'} = \mathbf{e}_s + \gamma F_{T(s, a')}^*$ , if  $F_a$  contains a copy of  $F_{a'}$  via  $\phi$ , then  $F_{T(s, a)}^*$  contains a copy of  $F_{T(s, a')}^*$  via  $\phi$ . Then  $\phi \cdot \text{ND}(F_{T(s, a')}^*) \subseteq \phi \cdot F_{T(s, a')}^* \subseteq F_{T(s, a)}^*$ , and so  $F_{T(s, a)}^*$  contains a copy of  $\text{ND}(F_{T(s, a')}^*)$ . Then apply lemma E.46 with  $\Delta_1 := T(s, a')$  and  $\Delta_2 := T(s, a)$  to conclude that  $\forall \gamma \in [0, 1] : \mathbb{E}_{s_{a'} \sim T(s, a')} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_{a'}, \gamma)] \leq_{\text{most}} \mathfrak{D}_{\text{bound}} \mathbb{E}_{s_a \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_a, \gamma)]$ .

Suppose  $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'})$  is non-empty. To apply the second condition of lemma E.46, we want to demonstrate that  $\text{ND} \left( F_{T(s,a)}^* \right) \setminus \phi \cdot \text{ND} \left( F_{T(s,a')}^* \right)$  is also non-empty.

First consider  $\mathbf{f} \in \mathcal{F}_{\text{nd}}(s) \cap F_a$ . Because  $F_a \subseteq \mathbf{e}_s + \gamma F_{T(s,a)}^*$ , we have that  $\gamma^{-1}(\mathbf{f} - \mathbf{e}_s) \in F_{T(s,a)}^*$ . Because  $\mathbf{f} \in \mathcal{F}_{\text{nd}}(s)$ , by definition 3.6,  $\exists \mathbf{r} \in \mathbb{R}^{|\mathcal{S}|}$ ,  $\gamma_x \in (0, 1)$  such that

$$\mathbf{f}(\gamma_x)^\top \mathbf{r} > \max_{\mathbf{f}' \in \mathcal{F}(s) \setminus \{\mathbf{f}\}} \mathbf{f}'(\gamma_x)^\top \mathbf{r}. \quad (183)$$

Then since  $\gamma_x \in (0, 1)$ ,

$$\gamma_x^{-1}(\mathbf{f}(\gamma_x) - \mathbf{e}_s)^\top \mathbf{r} > \max_{\mathbf{f}' \in \mathcal{F}(s) \setminus \{\mathbf{f}\}} \gamma_x^{-1}(\mathbf{f}'(\gamma_x) - \mathbf{e}_s)^\top \mathbf{r} \quad (184)$$

$$= \max_{\mathbf{f}' \in \gamma_x^{-1}((\mathcal{F}(s) \setminus \{\mathbf{f}\}) - \mathbf{e}_s)} \mathbf{f}'(\gamma_x)^\top \mathbf{r} \quad (185)$$

$$\geq \max_{\mathbf{f}' \in \gamma_x^{-1}((F_a \setminus \{\mathbf{f}\}) - \mathbf{e}_s)} \mathbf{f}'(\gamma_x)^\top \mathbf{r} \quad (186)$$

$$= \max_{\mathbf{f}' \in F_{T(s,a)}^* \setminus \{\gamma_x^{-1}(\mathbf{f} - \mathbf{e}_s)\}} \mathbf{f}'(\gamma_x)^\top \mathbf{r}. \quad (187)$$

Equation (186) holds because  $F_a \subseteq \mathcal{F}(s)$ . By assumption, action  $a$  is optimal for  $\mathbf{r}$  at state  $s$  and at discount rate  $\gamma_x$ . Equation (181) shows that  $F_{T(s,a)}^*$  potentially allows the agent a non-stationary policy choice at  $s$ , but non-stationary policies cannot increase optimal value [Puterman, 2014]. Therefore, eq. (187) holds.

We assumed that  $\gamma^{-1}(\mathbf{f} - \mathbf{e}_s) \in \gamma^{-1}(\mathcal{F}_{\text{nd}}(s) - \mathbf{e}_s)$ . Furthermore, since we just showed that  $\gamma^{-1}(\mathbf{f} - \mathbf{e}_s) \in F_{T(s,a)}^*$  is strictly optimal over the other elements of  $F_{T(s,a)}^*$  for reward function  $\mathbf{r}$  at discount rate  $\gamma_x \in (0, 1)$ , we conclude that it is an element of  $\text{ND} \left( F_{T(s,a)}^* \right)$  by definition E.13.

Then we conclude that  $\gamma^{-1}(\mathcal{F}_{\text{nd}}(s) - \mathbf{e}_s) \cap F_{T(s,a)}^* \subseteq \text{ND} \left( F_{T(s,a)}^* \right)$ .

We now show that  $\text{ND} \left( F_{T(s,a)}^* \right) \setminus \phi \cdot \text{ND} \left( F_{T(s,a')}^* \right)$  is non-empty.

$$0 < \left| \mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'}) \right| \quad (188)$$

$$= \left| \gamma^{-1} \left( \mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'}) - \mathbf{e}_s \right) \right| \quad (189)$$

$$\leq \left| \gamma^{-1} \left( \mathcal{F}_{\text{nd}}(s) - \mathbf{e}_s \right) \cap \left( F_{T(s,a)}^* \setminus \phi \cdot F_{T(s,a')}^* \right) \right| \quad (190)$$

$$= \left| \left( \gamma^{-1} \left( \mathcal{F}_{\text{nd}}(s) - \mathbf{e}_s \right) \cap F_{T(s,a)}^* \right) \setminus \phi \cdot F_{T(s,a')}^* \right| \quad (191)$$

$$\leq \left| \text{ND} \left( F_{T(s,a)}^* \right) \setminus \phi \cdot F_{T(s,a')}^* \right| \quad (192)$$

$$\leq \left| \text{ND} \left( F_{T(s,a)}^* \right) \setminus \phi \cdot \text{ND} \left( F_{T(s,a')}^* \right) \right|. \quad (193)$$

Equation (188) follows by the assumption that  $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'})$  is non-empty. Let  $\mathbf{f}, \mathbf{f}' \in \mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'})$  be distinct. Then we must have that for some  $\gamma_x \in (0, 1)$ ,  $\mathbf{f}(\gamma_x) \neq \mathbf{f}'(\gamma_x)$ . This holds iff  $\gamma_x^{-1}(\mathbf{f}(\gamma_x) - \mathbf{e}_s) \neq \gamma_x^{-1}(\mathbf{f}'(\gamma_x) - \mathbf{e}_s)$ , and so eq. (189) holds.

Equation (190) holds because  $F_a \subseteq \mathbf{e}_s + \gamma F_{T(s,a)}^*$  and  $F_{a'} = \mathbf{e}_s + \gamma F_{T(s,a')}^*$  by eq. (182). Equation (192) holds because we showed above that  $\gamma^{-1}(\mathcal{F}_{\text{nd}}(s) - \mathbf{e}_s) \cap F_{T(s,a)}^* \subseteq \text{ND} \left( F_{T(s,a)}^* \right)$ .

Equation (193) holds because  $\text{ND} \left( F_{T(s,a')}^* \right) \subseteq F_{T(s,a')}^*$  by definition E.13.

Therefore,  $\text{ND} \left( F_{T(s,a)}^* \right) \setminus \phi \cdot \text{ND} \left( F_{T(s,a')}^* \right)$  is non-empty, and so apply the second condition of lemma E.46 to conclude that for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B\text{/IID}}$ ,  $\forall \gamma \in (0, 1)$  :

$\mathbb{E}_{s_{a'} \sim T(s, a')} [\text{POWER}_{\mathcal{D}_{X\text{-IID}}}(s_{a'}, \gamma)] < \mathbb{E}_{s_a \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{X\text{-IID}}}(s_a, \gamma)]$ , and that  $\forall \gamma \in (0, 1) : \mathbb{E}_{s_{a'} \sim T(s, a')} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_{a'}, \gamma)] \not\leq_{\text{most: } \mathfrak{D}_{\text{bound}}} \mathbb{E}_{s_a \sim T(s, a)} [\text{POWER}_{\mathcal{D}_{\text{bound}}}(s_a, \gamma)]$ .

**Item 2.** Let  $\phi'(s_x) := \phi(s_x)$  when  $s_x \in \text{REACH}(s, a') \cup \text{REACH}(s, a)$ , and equal  $s_x$  otherwise. Since  $\phi$  is an involution, so is  $\phi'$ .

$$\phi' \cdot F_{a'} := \left\{ \mathbf{P}_{\phi'} \left( \mathbf{e}_s + \gamma \mathbb{E}_{s_{a'} \sim T(s, a')} [\mathbf{f}^{\pi, s_{a'}}] \right) \mid \pi \in \Pi, \pi(s) = a' \right\} \quad (194)$$

$$= \left\{ \mathbf{e}_s + \gamma \mathbb{E}_{s_{a'} \sim T(s, a')} [\mathbf{P}_{\phi'} \mathbf{f}^{\pi, s_{a'}}] \mid \pi \in \Pi, \pi(s) = a' \right\} \quad (195)$$

$$= \left\{ \mathbf{P}_{\phi} \mathbf{e}_s + \gamma \mathbb{E}_{s_{a'} \sim T(s, a')} [\mathbf{P}_{\phi} \mathbf{f}^{\pi, s_{a'}}] \mid \pi \in \Pi, \pi(s) = a' \right\} \quad (196)$$

$$=: \phi \cdot F_{a'} \quad (197)$$

$$\subseteq F_a. \quad (198)$$

Equation (195) follows because if  $s \in \text{REACH}(s, a') \cup \text{REACH}(s, a)$ , then we already showed that  $\phi$  fixes  $s$ . Otherwise,  $\phi'(s) = s$  by definition. Equation (196) follows by the definition of  $\phi'$  on  $\text{REACH}(s, a') \cup \text{REACH}(s, a)$  and because  $\mathbf{e}_s = \mathbf{P}_{\phi} \mathbf{e}_s$ . Next, we assumed that  $\phi \cdot F_{a'} \subseteq F_a$ , and so eq. (198) holds.

Therefore,  $F_a$  contains a copy of  $F_{a'}$  via  $\phi'$  fixing all  $s_x \notin \text{REACH}(s, a') \cup \text{REACH}(s, a)$ . Therefore,  $F_a$  contains a copy of  $F_{\text{nd}, a'} := \mathcal{F}_{\text{nd}}(s) \cap F_{a'}$  via the same  $\phi'$ . Then apply lemma E.49 with  $s' := s$  to conclude that  $\forall \gamma \in [0, 1] : \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{a'}, \gamma) \leq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma)$ . By lemma E.50,  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a', \gamma) = \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_{a'}, \gamma)$  and  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma) = \mathbb{P}_{\mathcal{D}_{\text{any}}}(F_a, \gamma)$ . Therefore,  $\forall \gamma \in [0, 1] : \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a', \gamma) \leq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma)$ .

If  $\mathcal{F}_{\text{nd}}(s) \cap (F_a \setminus \phi \cdot F_{a'})$  is non-empty, then apply the second condition of lemma E.49 to conclude that for all  $\gamma \in (0, 1)$ , the inequality is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B\text{IID}}$ , and  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a', \gamma) \not\leq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(s, a, \gamma)$ .  $\square$

#### E.4.2 When $\gamma = 1$ , optimal policies tend to navigate towards “larger” sets of cycles

**Lemma E.51** (POWER identity when  $\gamma = 1$ ).

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1) = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{RSD}(s)} \mathbf{d}^\top \mathbf{r} \right] = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{RSD}_{\text{nd}}(s)} \mathbf{d}^\top \mathbf{r} \right]. \quad (199)$$

*Proof.*

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1) = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{f}^{\pi, s} \in \mathcal{F}(s)} \lim_{\gamma \rightarrow 1} \frac{1 - \gamma}{\gamma} (\mathbf{f}^{\pi, s}(\gamma) - \mathbf{e}_s)^\top \mathbf{r} \right] \quad (200)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{RSD}(s)} \mathbf{d}^\top \mathbf{r} \right] \quad (201)$$

$$= \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{RSD}_{\text{nd}}(s)} \mathbf{d}^\top \mathbf{r} \right]. \quad (202)$$

Equation (200) follows by lemma E.45. Equation (201) follows by the definition of  $\text{RSD}(s)$  (definition 6.10). Equation (202) follows because for all  $\mathbf{r} \in \mathbb{R}^{|S|}$ , corollary E.11 shows that  $\max_{\mathbf{d} \in \text{RSD}(s)} \mathbf{d}^\top \mathbf{r} = \max_{\mathbf{d} \in \text{ND}(\text{RSD}(s))} \mathbf{d}^\top \mathbf{r} =: \max_{\mathbf{d} \in \text{RSD}_{\text{nd}}(s)} \mathbf{d}^\top \mathbf{r}$ .  $\square$

**Proposition 6.12** (When  $\gamma = 1$ , RSDs control POWER). *If  $\text{RSD}(s)$  contains a copy of  $\text{RSD}_{\text{nd}}(s')$  via  $\phi$ , then  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1) \geq_{\text{most}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s', 1)$ . If  $\text{RSD}_{\text{nd}}(s) \setminus \phi \cdot \text{RSD}_{\text{nd}}(s')$  is non-empty, then the converse  $\leq_{\text{most}}$  statement does not hold.*

*Proof.* Suppose  $\text{RSD}_{\text{nd}}(s')$  is similar to  $D \subseteq \text{RSD}(s)$  via involution  $\phi$ .

$$\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', 1) = \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{RSD}_{\text{nd}}(s')} \mathbf{d}^\top \mathbf{r} \right] \quad (203)$$

$$\leq_{\text{most: } \mathcal{D}_{\text{bound}}} \mathbb{E}_{\mathbf{r} \sim \mathcal{D}_{\text{bound}}} \left[ \max_{\mathbf{d} \in \text{RSD}_{\text{nd}}(s)} \mathbf{d}^\top \mathbf{r} \right] \quad (204)$$

$$= \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1) \quad (205)$$

Equation (203) and eq. (205) follow from lemma E.51. By applying lemma E.24 with  $A := \text{RSD}(s')$ ,  $B' := D$ ,  $B := \text{RSD}(s)$  and  $g$  the identity function, eq. (204) follows.

Suppose  $\text{RSD}_{\text{nd}}(s) \setminus D$  is non-empty. By the same result, eq. (204) is a strict inequality for all  $\mathcal{D}_{X\text{-IID}} \in \mathcal{D}_{C/B\text{-IID}}$ , and we conclude that  $\text{POWER}_{\mathcal{D}_{\text{bound}}}(s', 1) \not\leq_{\text{most: } \mathcal{D}_{\text{bound}}} \text{POWER}_{\mathcal{D}_{\text{bound}}}(s, 1)$ .  $\square$

**Theorem 6.13** (Average-optimal policies tend to end up in “larger” sets of RSDs). *Let  $D, D' \subseteq \text{RSD}(s)$ . Suppose that  $D$  contains a copy of  $D'$  via  $\phi$ , and that the sets  $D \cup D'$  and  $\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)$  have pairwise orthogonal vector elements (i.e. pairwise disjoint vector support). Then  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(D, \text{average}) \geq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(D', \text{average})$ . If  $\text{RSD}_{\text{nd}}(s) \cap (D \setminus \phi \cdot D')$  is non-empty, the converse  $\leq_{\text{most}}$  statement does not hold.*

*Proof.* Let  $D_{\text{sub}} := \phi \cdot D'$ , where  $D_{\text{sub}} \subseteq D$  by assumption. Let  $X := \{s_i \in \mathcal{S} \mid \max_{\mathbf{d} \in D' \cup D} \mathbf{d}^\top \mathbf{e}_{s_i} > 0\}$ . Define

$$\phi'(s_i) := \begin{cases} \phi(s_i) & \text{if } s_i \in X \\ s_i & \text{else.} \end{cases} \quad (206)$$

Since  $\phi$  is an involution,  $\phi'$  is also an involution. Furthermore, by the definition of  $X$ ,  $\phi' \cdot D' = D_{\text{sub}}$  and  $\phi' \cdot D_{\text{sub}} = D'$  (because we assumed that both equalities hold for  $\phi$ ).

Let  $D^* := D' \cup D_{\text{sub}} \cup (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D))$ .

$$\phi' \cdot D^* := \phi' \cdot (D' \cup D_{\text{sub}} \cup (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D))) \quad (207)$$

$$= (\phi' \cdot D') \cup (\phi' \cdot D_{\text{sub}}) \cup \phi' \cdot (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)) \quad (208)$$

$$= D_{\text{sub}} \cup D' \cup (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)) \quad (209)$$

$$=: D^*. \quad (210)$$

In eq. (209), we know that  $\phi' \cdot D' = D_{\text{sub}}$  and  $\phi' \cdot D_{\text{sub}} = D'$ . We just need to show that  $\phi' \cdot (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)) = \text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)$ .

Suppose  $\exists s_i \in X, \mathbf{d}' \in \text{RSD}_{\text{nd}}(s) \setminus (D' \cup D) : \mathbf{d}'^\top \mathbf{e}_{s_i} > 0$ . By the definition of  $X, \exists \mathbf{d} \in D' \cup D : \mathbf{d}^\top \mathbf{e}_{s_i} > 0$ . Then

$$\mathbf{d}^\top \mathbf{d}' = \sum_{j=1}^{|\mathcal{S}|} \mathbf{d}^\top (\mathbf{d}' \odot \mathbf{e}_{s_j}) \quad (211)$$

$$\geq \mathbf{d}^\top (\mathbf{d}' \odot \mathbf{e}_{s_i}) \quad (212)$$

$$= \mathbf{d}^\top ((\mathbf{d}'^\top \mathbf{e}_{s_i}) \mathbf{e}_{s_i}) \quad (213)$$

$$= (\mathbf{d}'^\top \mathbf{e}_{s_i}) \cdot (\mathbf{d}^\top \mathbf{e}_{s_i}) \quad (214)$$

$$> 0. \quad (215)$$

Equation (211) follows from the definitions of the dot and Hadamard products. Equation (212) follows because  $\mathbf{d}$  and  $\mathbf{d}'$  have non-negative entries. Equation (215) follows because  $\mathbf{d}^\top \mathbf{e}_{s_i}$  and  $\mathbf{d}'^\top \mathbf{e}_{s_i}$  are both positive. But eq. (215) shows that  $\mathbf{d}^\top \mathbf{d}' > 0$ , contradicting our assumption that  $\mathbf{d}$  and  $\mathbf{d}'$  are orthogonal.

Therefore, such an  $s_i$  cannot exist, and  $X' := \left\{ s'_i \in \mathcal{S} \mid \max_{\mathbf{d}' \in \text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)} \mathbf{d}'^\top \mathbf{e}_{s_i} > 0 \right\} \subseteq (\mathcal{S} \setminus X)$ . By eq. (206),  $\forall s'_i \in X' : \phi'(s'_i) = s'_i$ . Thus,  $\phi' \cdot (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)) = \text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)$ , and eq. (209) follows. We conclude that  $\phi' \cdot D^* = D^*$ .

Consider  $Z := (\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D)) \cup D \cup D'$ . First,  $Z \subseteq \text{RSD}(s)$  by definition. Second,  $\text{RSD}_{\text{nd}}(s) = \text{RSD}_{\text{nd}}(s) \setminus (D' \cup D) \cup (\text{RSD}_{\text{nd}}(s) \cap D') \cup (\text{RSD}_{\text{nd}}(s) \cap D) \subseteq Z$ . Note that  $D^* = Z \setminus (D \setminus D_{\text{sub}})$ .

$$\mathbb{P}_{\mathcal{D}_{\text{any}}}(D', \text{average}) = p_{\mathcal{D}_{\text{any}}}(D' \geq \text{RSD}(s)) \quad (216)$$

$$\leq_{\text{most: } \mathfrak{D}_{\text{any}}} p_{\mathcal{D}_{\text{any}}}(D \geq \text{RSD}(s)) \quad (217)$$

$$= \mathbb{P}_{\mathcal{D}_{\text{any}}}(D, \text{average}). \quad (218)$$

Since  $\phi \cdot D' \subseteq D$  and  $\text{ND}(D') \subseteq D'$ ,  $\phi \cdot \text{ND}(D') \subseteq D$ . Then eq. (217) holds by applying lemma E.28 with  $A := D'$ ,  $B' := D_{\text{sub}}$ ,  $B := D$ ,  $C := \text{RSD}(s)$ , and the previously defined  $Z$  which we showed satisfies  $\text{ND}(C) \subseteq Z \subseteq C$ . Furthermore, involution  $\phi'$  satisfies  $\phi' \cdot B^* = \phi' \cdot (Z \setminus (B \setminus B')) = Z \setminus (B \setminus B') = B^*$  by eq. (210).

When  $\text{RSD}_{\text{nd}}(s) \cap (D \setminus D_{\text{sub}})$  is non-empty, since  $B' \subseteq C$  by assumption, lemma E.28 also shows that eq. (217) is strict for all  $\mathcal{D}_{X\text{-IID}} \in \mathfrak{D}_{C/B/\text{IID}}$ , and that  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(D', \text{average}) \not\leq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(D, \text{average})$ .  $\square$

**Proposition E.52** (RSD properties). *Let  $\mathbf{d} \in \text{RSD}(s)$ .  $\mathbf{d}$  is element-wise non-negative and  $\|\mathbf{d}\|_1 = 1$ .*

*Proof.*  $\mathbf{d}$  has non-negative elements because it equals the limit of  $\lim_{\gamma \rightarrow 1} (1 - \gamma)\mathbf{f}(\gamma)$ , whose elements are non-negative by proposition E.3 item 1.

$$\|\mathbf{d}\|_1 = \left\| \lim_{\gamma \rightarrow 1} (1 - \gamma)\mathbf{f}(\gamma) \right\|_1 \quad (219)$$

$$= \lim_{\gamma \rightarrow 1} (1 - \gamma) \|\mathbf{f}(\gamma)\|_1 \quad (220)$$

$$= 1. \quad (221)$$

Equation (219) follows because the definition of RSDs (definition 6.10) ensures that  $\exists \mathbf{f} \in \mathcal{F}(s) : \lim_{\gamma \rightarrow 1} (1 - \gamma)\mathbf{f}(\gamma) = \mathbf{d}$ . Equation (220) follows because  $\|\cdot\|_1$  is a continuous function. Equation (221) follows because  $\|\mathbf{f}(\gamma)\|_1 = \frac{1}{1 - \gamma}$  by proposition E.3 item 2.  $\square$

**Lemma E.53** (When reachable with probability 1, 1-cycles induce non-dominated RSDs). *If  $\mathbf{e}_{s'} \in \text{RSD}(s)$ , then  $\mathbf{e}_{s'} \in \text{RSD}_{\text{nd}}(s)$ .*

*Proof.* If  $\mathbf{d} \in \text{RSD}(s)$  is distinct from  $\mathbf{e}_{s'}$ , then  $\|\mathbf{d}\|_1 = 1$  and  $\mathbf{d}$  has non-negative entries by proposition E.52. Since  $\mathbf{d}$  is distinct from  $\mathbf{e}_{s'}$ , then its entry for index  $s'$  must be strictly less than 1:  $\mathbf{d}^\top \mathbf{e}_{s'} < 1 = \mathbf{e}_{s'}^\top \mathbf{e}_{s'}$ . Therefore,  $\mathbf{e}_{s'} \in \text{RSD}(s)$  is strictly optimal for the reward function  $\mathbf{r} := \mathbf{e}_{s'}$ , and so  $\mathbf{e}_{s'} \in \text{RSD}_{\text{nd}}(s)$ .  $\square$

**Corollary 6.14** (Average-optimal policies tend not to end up in any given 1-cycle). *Suppose  $\mathbf{e}_{s_x}, \mathbf{e}_{s'} \in \text{RSD}(s)$  are distinct. Then  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(\text{RSD}(s) \setminus \{\mathbf{e}_{s_x}\}, \text{average}) \geq_{\text{most}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(\{\mathbf{e}_{s_x}\}, \text{average})$ . If there is a third  $\mathbf{e}_{s''} \in \text{RSD}(s)$ , the converse  $\leq_{\text{most}}$  statement does not hold.*

*Proof.* Suppose  $\mathbf{e}_{s_x}, \mathbf{e}_{s'} \in \text{RSD}(s)$  are distinct. Let  $\phi := (s_x \ s')$ ,  $D' := \{\mathbf{e}_{s_x}\}$ ,  $D := \text{RSD}(s) \setminus \{\mathbf{e}_{s_x}\}$ .  $\phi \cdot D' = \{\mathbf{e}_{s'}\} \subseteq \text{RSD}(s) \setminus \{\mathbf{e}_{s_x}\} =: D$  since  $s_x \neq s'$ .  $D' \cup D = \text{RSD}(s)$  and  $\text{RSD}_{\text{nd}}(s) \setminus (D' \cup D) = \text{RSD}_{\text{nd}}(s) \setminus \text{RSD}(s) = \emptyset$  trivially have pairwise orthogonal vector elements. Then apply theorem 6.13 to conclude that  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(\{\mathbf{e}_{s_x}\}, \text{average}) \leq_{\text{most: } \mathfrak{D}_{\text{any}}} \mathbb{P}_{\mathcal{D}_{\text{any}}}(\text{RSD}(s) \setminus \{\mathbf{e}_{s_x}\}, \text{average})$ .

Suppose there exists another  $\mathbf{e}_{s''} \in \text{RSD}(s)$ . By lemma E.53,  $\mathbf{e}_{s''} \in \text{RSD}_{\text{nd}}(s)$ . Furthermore, since  $s'' \notin \{s', s_x\}$ ,  $\mathbf{e}_{s''} \in (\text{RSD}(s) \setminus \{\mathbf{e}_{s_x}\}) \setminus \{\mathbf{e}_{s'}\} = D \setminus \phi \cdot D'$ . Therefore,



$\mathbf{e}_{s''} \in \text{RSD}_{\text{nd}}(s) \cap (D \setminus \phi \cdot D')$ . Then apply the second condition of theorem 6.13 to conclude that  $\mathbb{P}_{\mathcal{D}_{\text{any}}}(\{\mathbf{e}_{s_x}\}, \text{average}) \not\stackrel{\text{most: } \mathcal{D}_{\text{bound}}}{\leq} \mathbb{P}_{\mathcal{D}_{\text{any}}}(\text{RSD}(s) \setminus \{\mathbf{e}_{s_x}\}, \text{average})$ .  $\square$