

482 **A Proofs**

483 **Lemma 15** (A one-step good policy is close to optimal). *Let $\Delta(h) := |V_\xi^*(h) - V_\xi^\pi(h)|$ with*
 484 *$h \in (\mathcal{A} \times \mathcal{E})^t$ for $t \geq t_0 \in \mathbb{N}$.*

$$\begin{aligned} & \text{If } \mathbb{E}_\xi^\pi \left| \max_a Q_\xi^\pi(h, a) - V_\xi^\pi(h) \right| < \beta \quad \forall t \geq t_0 \\ & \text{and } \mathbb{E}_\xi^\pi \left[\max_a \sum_e \xi(e|ha) \Delta(hae) \right] \leq (1 + \alpha) \mathbb{E}_\xi^\pi \Delta(hae) \quad \forall t \geq t_0 \\ & \text{then } \mathbb{E}_\xi^\pi \Delta(h) < \frac{\beta}{1 - \gamma(1 + \alpha)} \quad \forall t \geq t_0 \quad \text{provided } 1 + \alpha < 1/\gamma \end{aligned}$$

485 *Proof.* Let $\delta := \sup_{t \geq t_0} \mathbb{E} \Delta(h)$, where $h \in (\mathcal{A} \times \mathcal{E})^t$ and \mathbb{E} is short for \mathbb{E}_ξ^π .

$$\begin{aligned} \mathbb{E} \Delta(h) &= \left| \max_a Q_\xi^*(h, a) - V_\xi^\pi(h) \right| \\ &= \mathbb{E} \left| \max_a Q_\xi^*(h, a) - \max_a Q_\xi^\pi(h, a) + \max_a Q_\xi^\pi(h, a) - V_\xi^\pi(h) \right| \\ &\leq \mathbb{E} \left| \max_a Q_\xi^\pi(h, a) - V_\xi^\pi(h) \right| + \mathbb{E} \left| \max_a Q_\xi^*(h, a) - \max_a Q_\xi^\pi(h, a) \right| \\ &\stackrel{(1)}{<} \beta + \mathbb{E} \left| \max_a \sum_e \xi(e|ha) (r + \gamma V_\xi^*(hae)) - \max_a \sum_e \xi(e|ha) (r + \gamma V_\xi^\pi(hae)) \right| \\ &\leq \beta + \gamma \mathbb{E} \max_a \sum_e \xi(e|ha) |V_\xi^*(hae) - V_\xi^\pi(hae)| \\ &\leq \beta + \gamma(1 + \alpha) \mathbb{E} \Delta(hae) \end{aligned}$$

486 Taking $\sup_{t \geq t_0}$ on both sides implies $\delta < \beta + \gamma(1 + \alpha)\delta$ implies $\delta < \beta/(1 - \gamma(1 + \alpha))$. \square \square

487 **Lemma 17** ($\mathbb{E}_\xi^\pi \rightarrow 0$ implies $\mathbb{E}_\mu^\pi \rightarrow 0$). *If π is such that*

$$\mathbb{E}_\xi^\pi [V_\xi^*(h_{<t}) - V_\xi^\pi(h_{<t})] \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

488 *then for all $\mu \in \mathcal{M}$ we have*

$$\mathbb{E}_\mu^\pi [V_\xi^*(h_{<t}) - V_\xi^\pi(h_{<t})] \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Proof.

$$\mathbb{E}_\mu^\pi [V_\xi^*(h_{<t}) - V_\xi^\pi(h_{<t})] \leq \frac{1}{w(\mu)} \mathbb{E}_\xi^\pi [V_\xi^*(h_{<t}) - V_\xi^\pi(h_{<t})] \rightarrow 0$$

489 by the dominance of $\xi(\cdot) \geq w(\mu)\mu(\cdot)$. \square

490 **Lemma 20** ($V_\xi^{\pi'} \rightarrow V_\xi^\pi$ implies $V_\mu^{\pi'} \rightarrow V_\mu^\pi$ in μ -expectation). *If π is such that for all $\mu \in \mathcal{M}$*

$$\mathbb{E}_\mu^\pi [V_\xi^{\pi'}(h_{<t}) - V_\xi^\pi(h_{<t})] \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

491 *and $D_\infty(\mu^{\pi'}, \xi^{\pi'} | h_{<t}) \rightarrow 0$ μ^π -almost surely then we have*

$$\mathbb{E}_\mu^\pi [V_\mu^{\pi'}(h_{<t}) - V_\mu^\pi(h_{<t})] \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Proof.

$$\begin{aligned} & \mathbb{E}_\mu^\pi \left[|V_\mu^{\pi'}(h_{<t}) - V_\mu^\pi(h_{<t})| \right] \\ &= \mathbb{E}_\mu^\pi \left[|V_\mu^{\pi'}(h_{<t}) - V_\xi^{\pi'}(h_{<t}) + V_\xi^{\pi'}(h_{<t}) - V_\xi^\pi(h_{<t}) + V_\xi^\pi(h_{<t}) - V_\mu^\pi(h_{<t})| \right] \\ &\leq \mathbb{E}_\mu^\pi \left[|V_\mu^{\pi'}(h_{<t}) - V_\xi^{\pi'}(h_{<t})| \right] + \mathbb{E}_\mu^\pi \left[|V_\xi^{\pi'}(h_{<t}) - V_\xi^\pi(h_{<t})| \right] + \mathbb{E}_\mu^\pi \left[|V_\xi^\pi(h_{<t}) - V_\mu^\pi(h_{<t})| \right] \end{aligned}$$

492 The second and third term go to 0 as $t \rightarrow \infty$ by the assumptions and Lemma 3 with Lemma 13. The
 493 first term goes to 0 as $D_\infty(\mu^{\pi'}, \xi^{\pi'} | h_{<t}) \rightarrow 0$ μ^π -almost surely implies $\mathbb{E}_\mu^\pi \left[D_\infty(\mu^{\pi'}, \xi^{\pi'} | h_{<t}) \right] \rightarrow 0$
 494 and we have $\mathbb{E}_\mu^\pi \left[|V_\mu^{\pi'}(h_{<t}) - V_\xi^{\pi'}(h_{<t})| \right] \leq \mathbb{E}_\mu^\pi \left[D_\infty(\mu^{\pi'}, \xi^{\pi'} | h_{<t}) \right]$.

495

□

496 **Theorem 22** (Self-AIXI is Self-optimizing). *Let μ be some environment. If there is a policy*
 497 *π and a sequence of policies $\bar{\pi}_1, \bar{\pi}_2, \dots$ all contained within \mathcal{P} such that for all $t, h_{<t}$ we have*
 498 *$V_\xi^\zeta(h_{<t}) \geq V_\xi^{\bar{\pi}_t}(h_{<t}) - \epsilon_t$ with $\epsilon_t \rightarrow 0$, and for all $\nu \in \mathcal{M}$*

$$V_\nu^*(h_{<t}) - V_\nu^{\bar{\pi}_t}(h_{<t}) \rightarrow 0 \text{ as } t \rightarrow \infty \text{ } \mu^\pi\text{-almost surely} \quad (4)$$

499 then

$$V_\nu^*(h_{<t}) - V_\nu^{\pi_S}(h_{<t}) \rightarrow 0 \text{ as } t \rightarrow \infty \text{ } \mu^\pi\text{-almost surely}$$

500 If $\pi = \pi_S$ and Equation 4 holds for all $\mu \in \mathcal{M}$, then π_S is strongly asymptotically optimal in the
 501 class \mathcal{M} .

Proof.

$$0 \leq w(\mu | h_{<t}) (V_\mu^*(h_{<t}) - V_\mu^{\pi_S}(h_{<t})) \quad (5)$$

$$\leq \sum_{\nu \in \mathcal{M}} w(\nu | h_{<t}) (V_\nu^*(h_{<t}) - V_\nu^{\pi_S}(h_{<t})) \quad (6)$$

$$= \sum_{\nu \in \mathcal{M}} w(\nu | h_{<t}) V_\nu^*(h_{<t}) - V_\xi^{\pi_S}(h_{<t}) \quad (7)$$

$$\leq \sum_{\nu \in \mathcal{M}} w(\nu | h_{<t}) V_\nu^*(h_{<t}) - V_\xi^\zeta(h_{<t}) \quad (8)$$

$$\leq \sum_{\nu \in \mathcal{M}} w(\nu | h_{<t}) V_\nu^*(h_{<t}) - V_\xi^{\bar{\pi}_t}(h_{<t}) + \epsilon_t \quad (9)$$

$$= \sum_{\nu \in \mathcal{M}} w(\nu | h_{<t}) (V_\nu^*(h_{<t}) - V_\nu^{\bar{\pi}_t}(h_{<t})) + \epsilon_t \quad (10)$$

$$\rightarrow 0 \quad (11)$$

502 Equation 6 comes from adding positive terms. Equations 7 and 10 comes from the linearity of the
 503 value function. Equation 8 comes from π_S being one step optimal then following ζ and . Equation 9
 504 comes from the assumptions. Lastly, 11 comes from Equation 4 and [14, Lem.5.28ii].

505 $w(\mu | h_{<t}) \rightarrow 0$ as $h_{<t}$ is generated from μ^π (for more details see Self-Optimizing proof in [14]).
 506 Therefore $V_\mu^*(h_{<t}) - V_\mu^{\pi_S}(h_{<t}) \rightarrow 0$ μ^π -almost surely.

507

□