

## A. Proof of Theorem 1

Let us first introduce some relevant notation and definitions. Let  $\nu$  be an arbitrary policy over options, and  $c : \mathcal{S} \times \mathcal{O} \rightarrow [0, 1]$  a coefficient function. It will be convenient to define the following:

$$\begin{aligned} p^{c\pi^o}(s'|s) &\stackrel{\text{def}}{=} c(s', o)p^{\pi^o}(s'|s), \\ \mathcal{P}^{c\nu}q(s, o) &\stackrel{\text{def}}{=} \sum_{s'} p^{c\pi^o}(s'|s) \sum_{o'} \nu(o'|s')q(s', o'). \end{aligned}$$

It will also be convenient to introduce the following operators for continuation and termination:

$$\begin{aligned} \mathcal{P}^{(1-\beta)^\iota}q(s, o) &\stackrel{\text{def}}{=} \sum_{s'} p_{ss'}^{\pi^o}(1 - \beta^o(s'))q(s', o), \\ \mathcal{P}^{\beta\mu}q(s, o) &\stackrel{\text{def}}{=} \sum_{s'} p_{ss'}^{\pi^o}\beta^o(s') \sum_{o'} \mu(o'|s')q(s', o'). \end{aligned}$$

That is:  $\iota$  (“iota”) denotes the policy over options that maintains the current (argument) option. Let  $r^\pi$  denote the matrix of  $r^{\pi^o}$  for all options. Using these notations we can rewrite the Bellman operator over options for a discount  $\gamma$ :

$$\mathcal{T}_\gamma^\mu q = (I - \gamma\mathcal{P}^{(1-\beta)^\iota})^{-1}(r^\pi + \gamma\mathcal{P}^{\beta\mu}q). \quad (5)$$

We are now ready to give our proof.

*Proof.* In order to preserve the appealing form of the option equations, it will be convenient to renormalize the *one-step* reward model  $r^{\pi^o}$  to be in terms of  $\gamma_p$ . We do this by setting  $R_{\gamma_r}^o$  and  $R_{\gamma_p}^o$  to be the same:

$$\begin{aligned} R_{\gamma_r}^o &= (I - \gamma_r p^{(1-\beta)\pi^o})^{-1}r^{\pi^o} \\ &= (I - \gamma_p p^{(1-\beta)\pi^o})^{-1}z^{\pi^o} = R_{\gamma_p}^o, \end{aligned}$$

and solving for  $z^{\pi^o}$ :

$$z^{\pi^o} = (I - \gamma_p p^{(1-\beta)\pi^o})(I - \gamma_r p^{(1-\beta)\pi^o})^{-1}r^{\pi^o}.$$

The Bellman operator can be rewritten similarly to (5):  $\mathcal{T}_\Gamma^\mu q \stackrel{\text{def}}{=} R + \mathcal{P}_\Gamma^\mu q = (I - \gamma_p \mathcal{P}^{(1-\beta)^\iota})^{-1}(z^\pi + \gamma_p \mathcal{P}^{\gamma_d \beta \mu} q)$ , where we slightly abuse notation by writing  $\mathcal{P}^{\Gamma\beta\mu}$  for  $\sum_{s'} \gamma_d(s')\beta^o(s')p^{\pi^o}(s'|s) \sum_{o'} \nu(o'|s')q(s', o')$ . Let us now derive the fixed point of  $\mathcal{T}_\Gamma^\mu$ :

$$\begin{aligned} q &= (I - \gamma_p \mathcal{P}^{(1-\beta)^\iota})^{-1}(z^\pi + \gamma_p \mathcal{P}^{\Gamma\beta\mu} q) \\ q - \gamma_p \mathcal{P}^{(1-\beta)^\iota} q &= z^\pi + \gamma_p \mathcal{P}^{\Gamma\beta\mu} q \\ q &= z^\pi + \gamma_p \left( \mathcal{P}^{\Gamma\beta\mu} q + \mathcal{P}^{(1-\beta)^\iota} q \right) \\ &= \left( I - \gamma_p (\mathcal{P}^{\Gamma\beta\mu} + \mathcal{P}^{(1-\beta)^\iota}) \right)^{-1} z^\pi \\ &= \sum_{t=0}^{\infty} \gamma^t \left( \mathcal{P}^{\Gamma\beta\mu} + \mathcal{P}^{(1-\beta)^\iota} \right)^t z^\pi. \end{aligned}$$

So  $\mathcal{B} = \mathcal{P}^{\Gamma\beta\mu} + \mathcal{P}^{(1-\beta)^\iota}$  is the corresponding one-step operator. Let us verify that it induces the new step-discount and termination scheme. Let  $\ell^o(s') = \gamma_d(s')\beta^o(s') + 1 - \beta^o(s')$ . We have:

$$\begin{aligned} \mathcal{B}q(s, o) &= (\mathcal{P}^{\Gamma\beta\mu} + \mathcal{P}^{(1-\beta)^\iota})q(s, o) \\ &= \sum_{s'} p^{\pi^o}(s'|s) \left( \sum_{o'} \mu(o'|s') \gamma_d(s') \right. \\ &\quad \times \beta^o(s')q(s', o') + (1 - \beta^o(s'))q(s', o) \Big) \\ &= \sum_{s'} p^{\pi^o}(s'|s) \ell^o(s') \left( \sum_{o'} \mu(o'|s') \frac{\gamma_d(s')\beta^o(s')}{\ell^o(s')} \right. \\ &\quad \times q(s', o') + \frac{1 - \beta^o(s')}{\ell^o(s')} q(s', o) \Big) \\ &= \sum_{s'} p^{\pi^o}(s'|s) \ell^o(s') \left( \sum_{o'} \mu(o'|s') \right. \\ &\quad \times \zeta^o(s')q(s', o') + (1 - \zeta^o(s'))q(s', o) \Big). \end{aligned}$$

Thus, we have a new termination scheme  $\zeta^o(s) = \frac{\gamma_d(s)\beta^o(s)}{\ell^o(s)}$ , and combining  $\ell^o$  with the step-discount  $\gamma_p$ , we get our new step discount:

$$\begin{aligned} \kappa(s, o, s') &= \gamma_p \ell^o(s') = \gamma_p (\gamma_d(s')\beta^o(s') + 1 - \beta^o(s')) \\ &= \gamma_p (1 - \beta^o(s'))(1 - \gamma_d(s')). \end{aligned}$$

This implies that we have a state-option-dependent contraction factor.

$$\begin{aligned} \eta(s, o) &= \mathbb{E}_{S_1, S_2, \dots \sim p^{\pi^o}} \left[ \prod_{i=1}^{\infty} \gamma(S_i, o, S_{i+1}) \right] \\ &= \mathbb{E}_{S_1, S_2, \dots \sim p^{\pi^o}} \left[ \prod_{i=1}^{\infty} \gamma_p (1 - \beta^o(S_i)(1 - \gamma_d(S_i))) \right] \\ &\leq \gamma_p. \end{aligned}$$

The operator  $\mathcal{T}_\Gamma^\mu$  is a contraction if  $\eta(s, o) < 1$ . This is trivially true if  $\gamma_p < 1$ . Otherwise, if  $\gamma_p = 1$ , in order for  $\eta(s, o) < 1$ , we need  $1 - \beta^o(S_i)(1 - \gamma_d(S_i)) < 1$  for some  $S_i$  along the trajectory, which is the same as  $\beta^o(S_i)(1 - \gamma_d(S_i)) > 0$  and holds if Assumption 5.2 holds: if such an  $S_i$  is reachable by  $\pi^o$ , is terminating in the sense that  $\beta^o(S_i) > 0$ , and  $\gamma_d(S_i) < 1$ .  $\square$

## B. Proof of Lemma 1

Throughout we will refer to the minimum and maximum duration of options by  $d_{\min}^o$  and  $d_{\max}^o$ , where  $d_{\min}^o$  denotes the minimum, and  $d_{\max}^o$  the maximum duration of  $o$  between any  $s$  and  $s'$  in  $\mathcal{J}^o$ . We will also write  $d_{\min} \stackrel{\text{def}}{=} \min_{o \in \mathcal{O}} d_{\min}^o$ ,  $d_{\max} \stackrel{\text{def}}{=} \max_{o \in \mathcal{O}} d_{\max}^o$  for minimum and maximum durations across options.

Before we proceed, let us show two helper bounds.

**Lemma 3.** For each option  $o \in \mathcal{O}$ :

$$|P_{\Gamma}^o(\cdot|s)|_1 \leq \|\gamma_d\| \gamma_p^{d_{\min}^o}.$$

*Proof.* For each  $s$  and  $s'$ , the transition model  $P_{\Gamma}^o(s'|s) \leq \gamma_d(s') \gamma_p^{d_{\min}^o}$ , the minimum duration. Taking a max over the states  $s'$  yields the result.  $\square$

**Lemma 4.** The value function is bounded:

$$\|q_{\Gamma}^{\mu}\| \leq \frac{r_{\max}}{1 - \gamma_r} \frac{1}{1 - \|\gamma_d\| \gamma_p^{d_{\min}^o}}.$$

*Proof.* From the definition of  $q_{\Gamma}^{\mu}$ , Lemma 3 and the definition of the reward model  $R$ :

$$\begin{aligned} \|q_{\Gamma}^{\mu}\|_{\infty} &= \|(I - \mathcal{P}_{\Gamma}^{\mu})^{-1}R\|_{\infty} \\ &\leq \frac{1}{1 - \max_{s \in \mathcal{S}, o \in \mathcal{O}} |P_{\Gamma}^o(\cdot|s)|_1} \|R\|_{\infty} \\ &\leq \frac{1}{1 - \|\gamma_d\| \gamma_p^{d_{\min}^o}} \frac{r_{\max}}{1 - \gamma_r}, \end{aligned}$$

since  $\|\mathcal{P}_{\Gamma}^{\mu}q\|_{\infty} \leq \max_{s \in \mathcal{S}, o \in \mathcal{O}} |P_{\Gamma}^o(\cdot|s)|_1 \|q\|$ .  $\square$

Recall the value functions w.r.t. the exact and approximate transition models and some policy  $\mu$ :  $q_{\Gamma}^{\mu} \stackrel{\text{def}}{=} (I - \mathcal{P}_{\Gamma}^{\mu})^{-1}R$  and  $\widehat{q}_{\Gamma}^{\mu} \stackrel{\text{def}}{=} (I - \widehat{\mathcal{P}}_{\Gamma}^{\mu})^{-1}R$ . We will show Lemma 1 in two steps.

**1) Bounding  $\mathcal{E}_{\text{estim}}$  in terms of the one-step error.** Similarly to Lemma 4 from (Jiang et al., 2015), we can relate the error in the value functions due to the approximate model to the maximum one-step error:

**Lemma 5.** For any policy over options  $\mu$ :

$$\begin{aligned} \|q_{\Gamma}^{\mu} - \widehat{q}_{\Gamma}^{\mu}\| &\leq \frac{1}{1 - \|\gamma_d\| \gamma_p^{d_{\min}^o}} \\ &\quad \times \max_{s,o} \left| R^o(s) + \widehat{\mathcal{P}}_{\Gamma}^{\mu} q_{\Gamma}^{\mu}(s,o) - q_{\Gamma}^{\mu}(s,o) \right|. \end{aligned}$$

*Proof.* Consider the evolution

$$q_m(s,o) = R^o(s) + \widehat{\mathcal{P}}_{\Gamma}^{\mu} q_{m-1}(s,o). \quad (6)$$

We can bound the difference between successive estimates:

$$\begin{aligned} \|q_m - q_{m-1}\|_{\infty} &= \|\widehat{\mathcal{P}}_{\Gamma}^{\mu}(q_{m-1} - q_{m-2})\|_{\infty} \\ &\leq \max_{s \in \mathcal{S}, o \in \mathcal{O}} \widehat{\mathcal{P}}_{\Gamma}^{\mu}(\cdot|s)_1 \|q_{m-1} - q_{m-2}\|_{\infty} \\ &= \|\gamma_d\| \gamma_p^{d_{\min}^o} \|q_{m-1} - q_{m-2}\|_{\infty}, \end{aligned}$$

due to Lemma 3 (which of course still applies to the approximate model.) Thus

$$\begin{aligned} \|q_m - q_0\| &\leq \sum_{k=0}^{m-1} \|q_{k+1} - q_k\|_{\infty} \\ &\leq \|q_1 - q_0\|_{\infty} \sum_{k=1}^{m-1} (\|\gamma_d\| \gamma_p^{d_{\min}^o})^{k-1}. \end{aligned}$$

Since as  $m \rightarrow \infty$ ,  $q_m = q_{\Gamma}^{\mu}$ , we have that  $\|q_{\Gamma}^{\mu} - q_0\|_{\infty} \leq \frac{1}{1 - \|\gamma_d\| \gamma_p^{d_{\min}^o}} \|q_1 - q_0\|_{\infty}$ . Finally, since  $q_0$  can be initialized to  $q_{\Gamma}^{\mu}$ , and from Eq. (6) for  $m = 1$ , we have our result.  $\square$

**2) Bounding the one-step error with the Hoeffding's bound.** Now let us bound the one-step error in terms of the number of samples. The following Lemma is similar to Lemma 2 from (Jiang et al., 2015).

**Lemma 6.** Let  $\widehat{P}_{\Gamma}^o$  denote the modified transition model of an option  $o$  estimated i.i.d. from  $n$  samples, and  $\widehat{\mathcal{P}}_{\Gamma}^{\mu}$  the corresponding operator w.r.t. some policy over options  $\mu$ . Let  $\mathcal{F}^o$  denote the set of possible terminating states of an option  $o$ . We have, with probability  $1 - \delta$ :

$$\begin{aligned} \|R + \widehat{\mathcal{P}}_{\Gamma}^{\mu} q_{\Gamma}^{\mu} - q_{\Gamma}^{\mu}\| &\leq ((\gamma_p^{d_{\min}^o} - \gamma_p^{d_{\max}^o})w + \gamma_p^{d_{\min}^o} \Delta v) \\ &\quad \times \sqrt{\frac{1}{2n} \log \frac{2|\mathcal{S}||\mathcal{O}|}{\delta}}, \end{aligned}$$

where  $\Delta v = \max_{o \in \mathcal{O}} \left( \max_{s \in \mathcal{F}^o} v_{\Gamma}^{\mu}(s) - \min_{s \in \mathcal{F}^o} v_{\Gamma}^{\mu}(s) \right)$  is the maximum variation of value in terminating states, and  $w = \max_{o \in \mathcal{O}} \min_{s \in \mathcal{F}^o} \gamma_d(s) v_{\Gamma}^{\mu}(s)$ .

*Proof.* Let us fix  $s$  and  $o$ , and consider a sample  $Y$  of  $\sum_{s'} \widehat{P}_{\Gamma}^o(s'|s) v(s')$ , where we write  $v(s) = v_{\Gamma}^{\mu}(s) = \sum_o \mu(o|s) q_{\Gamma}^{\mu}(s,o)$  throughout. Because each  $\widehat{P}_{\Gamma}^o(s'|s)$  is an average of zeros and samples of  $\gamma_d(s') \gamma_p^D$ , where  $D$  is the random variable corresponding to option duration, we have:

$$\underbrace{\gamma_p^{d_{\max}^o} \min_{s' \in \mathcal{F}^o} \gamma_d(s') v(s')}_{v_{\min}^o} \leq Y \leq \underbrace{\gamma_p^{d_{\min}^o} \max_{s' \in \mathcal{F}^o} \gamma_d(s') v(s')}_{v_{\max}^o}.$$

Now let  $X_o = R^o(s) + \sum_{s'} \widehat{P}_{\Gamma}^o(s'|s) v(s')$ . We have that:

$$R^o(s) + \gamma_p^{d_{\max}^o} v_{\min}^o \leq X_o \leq R^o(s) + \gamma_p^{d_{\min}^o} v_{\max}^o.$$

Thus the range  $a = X_{\min} - X_{\max}$  of  $X$  is:

$$\begin{aligned} a &= \gamma_p^{d_{\min}^o} v_{\max}^o - \gamma_p^{d_{\max}^o} v_{\min}^o \\ &= (\gamma_p^{d_{\min}^o} - \gamma_p^{d_{\max}^o}) v_{\min}^o + \gamma_p^{d_{\min}^o} \underbrace{(v_{\max}^o - v_{\min}^o)}_{\Delta v^o}. \end{aligned}$$

Then, since  $\widehat{P}_\Gamma^o(s'|s)$  is sampled i.i.d., and  $q_\Gamma^\mu(s, o)$  is the average of  $R^o(s) + \sum_{s'} \widehat{P}_\Gamma^o(s'|s)v_\Gamma^\mu(s')$ , we have by the Hoeffding's bound:

$$\begin{aligned} \Pr(|X_{s,o} - \mathbb{E}[X_{s,o}]| \geq t) &\leq 2 \exp\left(-\frac{2nt^2}{a^2}\right) \\ &= 2 \exp\left(-\frac{2nt^2}{((\gamma_p^{d_{\min}} - \gamma_p^{d_{\max}})v_{\min}^o + \gamma_p^{d_{\min}} \Delta v^o)^2}\right) \\ &= p_o. \end{aligned}$$

Using the union bound we have:

$$\Pr(|X - \mathbb{E}[X]| \geq t) \leq |S| \sum_o p_o \leq |S||\mathcal{O}|p_{\max},$$

where

$$\begin{aligned} p_{\max} &= \max_o p_o \\ &= 2 \exp\left(-\frac{2nt^2}{((\gamma_p^{d_{\min}} - \gamma_p^{d_{\max}})w + \gamma_p^{d_{\min}} \Delta v)^2}\right), \end{aligned}$$

in which  $\Delta v$  is the maximum variation of terminal values over all options, and  $w = \max_{o \in \mathcal{O}} \min_{s \in \mathcal{S}^o} \gamma_d(s)v(s)$ . Solving for  $t$  we have our result:

$$\frac{\delta}{|S||\mathcal{O}|} = 2 \exp\left(-\frac{2nt^2}{((\gamma_p^{d_{\min}} - \gamma_p^{d_{\max}})w + \gamma_p^{d_{\min}} \Delta v)^2}\right)$$

$$\log \frac{2|S||\mathcal{O}|}{\delta} = \frac{2nt^2}{((\gamma_p^{d_{\min}} - \gamma_p^{d_{\max}})w + \gamma_p^{d_{\min}} \Delta v)^2}$$

$$t = ((\gamma_p^{d_{\min}} - \gamma_p^{d_{\max}})w + \gamma_p^{d_{\min}} \Delta v) \sqrt{\frac{1}{2n} \log \frac{2|S||\mathcal{O}|}{\delta}}.$$

□

Lemma 1 then follows from combining Lemmas 5 and 6

### C. Proof of Lemma 2

*Proof.* We have

$$\begin{aligned} q_\Gamma^\mu - q_{\gamma_r}^\mu &= \mathcal{P}_\Gamma^\mu q_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu \\ &= \mathcal{P}_\Gamma^\mu q_\Gamma^\mu - \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu + \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu \\ &= \mathcal{P}_\Gamma^\mu (q_\Gamma^\mu - q_{\gamma_r}^\mu) + (\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu) q_{\gamma_r}^\mu \\ &= (I - \mathcal{P}_\Gamma^\mu)^{-1} \underbrace{(\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu)}_A q_{\gamma_r}^\mu. \end{aligned} \quad (7)$$

Let us now bound this expression. We will start with the inner term first. Noticing that  $\mathcal{P}^{\Gamma\beta\mu} = \mathcal{P}^{\beta\mu}\gamma_d$ , and from

the definitions of the operators, we can expand:

$$\begin{aligned} \|\mathcal{A}q_{\gamma_r}^\mu\| &= \|(\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu)q_{\gamma_r}^\mu\| \\ &= \left\| \gamma_p \mathcal{P}^{\beta\mu} \gamma_d q_{\gamma_r}^\mu + \gamma_p \mathcal{P}^{(1-\beta)\iota} \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \left( \gamma_r \mathcal{P}^{\beta\mu} q_{\gamma_r}^\mu \right. \right. \\ &\quad \left. \left. + \gamma_r \mathcal{P}^{(1-\beta)\iota} \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu \right) \right\| \\ &\leq \left\| \mathcal{P}^{\beta\mu} (\gamma_p \gamma_d - \gamma_r) q_{\gamma_r}^\mu + \mathcal{P}^{(1-\beta)\iota} \right. \\ &\quad \left. (\gamma_p \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \gamma_r \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu) \right\| \\ &= \left\| \mathcal{P}^{\beta\mu} (\gamma_p \gamma_d - \gamma_r) q_{\gamma_r}^\mu + \mathcal{P}^{(1-\beta)\iota} \right. \\ &\quad \left. (\gamma_p \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \gamma_r \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu + \gamma_r \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \gamma_r \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu) \right\| \\ &= \left\| \mathcal{P}^{\beta\mu} (\gamma_p \gamma_d - \gamma_r) q_{\gamma_r}^\mu + \mathcal{P}^{(1-\beta)\iota} \right. \\ &\quad \left. (\gamma_p - \gamma_r) \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu + \gamma_r (\mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu) \right\| \\ &\leq \|\gamma_d\| \gamma_p - \gamma_r \|q_{\gamma_r}^\mu\| + (\gamma_p - \gamma_r) \|\mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu\| \\ &\quad + \gamma_r \|\mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu - \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu\| \\ &\leq \|\gamma_d\| \gamma_p - \gamma_r \|q_{\gamma_r}^\mu\| + (\gamma_p - \gamma_r) \gamma_p^{d_{\min}} \|q_{\gamma_r}^\mu\| \\ &\quad + \gamma_r \|(\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu) q_{\gamma_r}^\mu\|, \end{aligned}$$

where the last inequality is due to Lemma 3. Let us simplify the first coefficient:

$$\begin{aligned} \|\gamma_d\| \gamma_p - \gamma_r &= \|\gamma_d\| \gamma_p - \gamma_p + \gamma_p - \gamma_r \\ &\leq \gamma_p (1 - \|\gamma_d\|) + \gamma_p - \gamma_r. \end{aligned}$$

Solving for  $\|(\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu)q_{\gamma_r}^\mu\|$ , and by Lemma 4, we get:

$$\begin{aligned} \|(\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu)q_{\gamma_r}^\mu\| &\leq \frac{1}{1 - \gamma_r} (\gamma_p (1 - \|\gamma_d\|) + \gamma_p - \gamma_r \\ &\quad + (\gamma_p - \gamma_r) \gamma_p^{d_{\min}}) \|q_{\gamma_r}^\mu\| \\ &\leq \frac{r_{\max}}{(1 - \gamma_r)^2} ((\gamma_p - \gamma_r) (\gamma_p^{d_{\min}} + 1) + \\ &\quad \gamma_p (1 - \|\gamma_d\|)). \end{aligned}$$

Finally, from Eq. (7) and using the bound on  $\mathcal{P}_\Gamma^\mu$  from Lemma 3:

$$\mathcal{E}_{targ} \leq \frac{r_{\max}}{(1 - \gamma_r)^2} \frac{(\gamma_p - \gamma_r) (\gamma_p^{d_{\min}} + 1) + \gamma_p (1 - \|\gamma_d\|)}{1 - \|\gamma_d\| \gamma_p^{d_{\min}}}.$$

□

### D. Proof of Proposition 1

*Proof.* From Eq. (7):

$$\mathcal{E}_{estim} = (I - \mathcal{P}_\Gamma^\mu)^{-1} (\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu) q_{\gamma_r}^\mu.$$

If the inner term is zero, the bias will be zero as well:

$$\begin{aligned} (\mathcal{P}_\Gamma^\mu - \mathcal{P}_{\gamma_r}^\mu) q_{\gamma_r}^\mu(s, o) &= 0, \\ \mathcal{P}_\Gamma^\mu q_{\gamma_r}^\mu(s, o) &= \mathcal{P}_{\gamma_r}^\mu q_{\gamma_r}^\mu(s, o), \end{aligned}$$

and hence:

$$\begin{aligned} \sum_{s'} P_{\gamma_p}^o(s'|s) \gamma_d(s, s') \sum_{o'} \mu(o'|s') q_{\gamma_r}^\mu(s', o') \\ = \sum_{s'} P_{\gamma_r}^o(s'|s) \sum_{o'} \mu(o'|s') q_{\gamma_r}^\mu(s', o'). \end{aligned}$$

This equality can be achieved, if for each option  $o$ :

$$\begin{aligned} P_{\gamma_p}^o(s'|s) \gamma_d(s, s') \sum_{o'} \mu(o'|s') q_{\gamma_r}^\mu(s', o') \\ = P_{\gamma_r}^o(s'|s) \sum_{o'} \mu(o'|s') q_{\gamma_r}^\mu(s', o'), \end{aligned}$$

$$P_{\gamma_p}^o(s'|s) \gamma_d(s, s') = P_{\gamma_r}^o(s'|s),$$

$$\gamma_d(s, s') = \frac{P_{\gamma_r}^o(s'|s)}{P_{\gamma_p}^o(s'|s)}.$$

□