

A PROOF OF PROPOSITION 4.1

Proof. For all $(s, a^i) \in \mathcal{S} \times \mathcal{A}^i$, taking the derivative of the Lagrangian of the optimization problem with simplex constraints over π^i ,

$$\begin{aligned} & \alpha_{t,(i)} \cdot \langle \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s), \pi^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle \\ & - \text{KL}(\pi^i(\cdot | s) \| \pi_{\theta_t}^i(\cdot | s)) + C \cdot \left(\sum_{a^i \in \mathcal{A}^i} \pi^i(a^i | s) - 1 \right) + C' \cdot \sum_{a^i \in \mathcal{A}^i} \pi^i(a^i | s) \end{aligned}$$

with respect to $\pi^i(a^i | s)$, we obtain

$$\alpha_{t,(i)} \cdot \left(\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i) - \lambda_i \cdot \log \pi_{\theta_t}^i(a^i | s) \right) - \log \pi^i(a^i | s) + \log \pi_{\theta_t}^i(a^i | s) + C + C' - 1, \quad (\text{A.1})$$

where C is a constant. Setting (A.1) to zero we further have

$$\begin{aligned} \pi_{t+1}^i(a^i | s) &= \left(\pi_{\theta_t}^i(a^i | s) \right)^{1 - \lambda_i \alpha_{t,(i)}} \cdot \exp \{ \alpha_{t,(i)} \cdot \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i) + C + C' - 1 \} \\ &= \exp \{ \alpha'_{t,(i)} \cdot \mathcal{E}_{\theta_t,(i)}(s, a^i) + \alpha_{t,(i)} \cdot \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i) + C + C' - 1 \}, \end{aligned}$$

where the second equality follows from $\lambda_i = (1 - \alpha'_{t,(i)}) / \alpha_{t,(i)}$. Thus, we know that (4.2) is $\pi_{t+1}^i(\cdot | s) \propto \exp \{ \alpha'_{t,(i)} \cdot \mathcal{E}_{\theta_t,(i)}(s, \cdot) + \alpha_{t,(i)} \cdot \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) \}$, which coincides with (3.3). Therefore, we finish the proof. \square

B DISCUSSION ON ASSUMPTION 4.4

In this section, we make a detailed discussion on Assumption 4.4. First, we give the following condition on the mean squared error of the estimations $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}$ and $\mathcal{E}_{\theta_{t+1},(i)}$.

Condition B.1 (Approximation Error). For all $0 \leq t \leq T - 1$, the estimator $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i)$ of the marginalized $Q_{(i)}$ -function $\widetilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i)$ satisfies

$$\mathbb{E}_{\sigma_t} \left[\left(\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i) - \widetilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i) \right)^2 \right] \leq \varepsilon_t, \quad (\text{B.1})$$

and the energy function update $\mathcal{E}_{\theta_{t+1},(i)}(s, a^i)$ satisfy

$$\mathbb{E}_{\sigma_t} \left[\left(\mathcal{E}_{\theta_{t+1},(i)}(s, a^i) - \widehat{\mathcal{E}}_{\theta_{t+1},(i)}(s, a^i) \right)^2 \right] \leq \varepsilon'_t. \quad (\text{B.2})$$

Condition B.1 can be satisfied for arbitrary small errors ε_t and ε'_t if the following conditions are satisfied: (I) the representation powers of parametrization $\mathcal{E}_{\theta_{t+1},(i)}(s, a^i)$ and the parametrization of $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}$ are strong enough. (II) the algorithms for learning $\widehat{\mathcal{E}}_{\theta_{t+1},(i)}(s, a^i)$ and $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}$ attain stationary points after sufficiently many iterations. In this section, as an example, we take the most commonly used neural network parameterization for $\mathcal{E}_{\theta_{t+1},(i)}(s, a^i)$ and $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}$.

Representation Power. There is a line of literature discussing the representation power of neural networks (see, e.g., Daniely et al. (2016); Khurikov et al. (2017)) showing that the overparameterized neural networks possess strong representation power. Specifically, the representation power of neural networks can be approximated as a subset of the reproducing kernel Hilbert space (RKHS) with neural tangent kernel (Jacot et al., 2018; Chizat & Bach, 2018; Allen-Zhu et al., 2018; Lee et al., 2019; Arora et al., 2019), which is a sufficiently rich function class. Moreover, we present the following lemma justifying the choice of the truncation parameter $\mathcal{E}_{(i)}^{\max}$ in Theorem 4.5.

Lemma B.2 (Bounded Energy Function). Let the regularization parameter λ_i and the stepsizes $\alpha_{t,(i)}, \alpha'_{t,(i)}$ be chosen as in Theorem 4.5. In Algorithm 1, setting $\mathcal{E}_{(i)}^{\max} = Q_{(i)}^{\max} / (\lambda_i - M_i)$ makes the function class $\mathcal{F}_{\mathcal{E}_{(i)}^{\max}}$ always cover the range of the estimated energy function update $\widehat{\mathcal{E}}_{t,(i)}$ obtained in Line 6 of Algorithm 1.

Proof. See Appendix H for a detailed proof. \square

Lemma B.2 states that, truncating the energy functions $\mathcal{E}_{\theta_t, (i)}$ within the function class $\mathcal{F}_{Q_{(i)}^{\max}/(\lambda_i - M_i)}$ does not compromise the MSE ε'_t defined in (B.2).

Learning Algorithms. There are some recent advances (Cai et al., 2019; Liu et al., 2019) showing that, when equipped with neural network parameterization, temporal-difference (TD) learning converges to the stationary point. Also, as discussed in Liu et al. (2019), under neural network parameterization, stochastic gradient descent also converges to the optima at a sublinear rate.

Next, by Lemmas 4.7-4.8 of Liu et al. (2019), we lay out error bounds in the form of those in Assumption 4.4.

Lemma B.3. ε_t and ε'_t in Assumption 4.4 take the forms of

$$\varepsilon_t = \max_{i \in \{1, 2\}} \{|\mathcal{A}^i|\} \cdot (\varepsilon'_t)^2, \quad \varepsilon'_t = \varepsilon'_t \cdot \max_{i \in \{1, 2\}} \{\phi_{\pi_{\theta_t}^*}^i\} + \varepsilon_t \cdot \psi_t^*,$$

where

$$\psi_t^* = \mathbb{E}_{\sigma_t} \left[\left| \frac{d\sigma^*}{d\sigma_t} - \frac{d\nu^*}{d\nu_t} \right| \right], \quad \phi_{\pi_{\theta_t}^*}^i = \mathbb{E}_{\sigma_t} \left[\left| \frac{d\pi_{*}^i}{d\pi_0^i} - \frac{d\pi_{\theta_t}^i}{d\pi_0^i} \right| \right],$$

Here the density ratios are Radon-Nikodym derivatives.

Thus, as long as the learning algorithms run sufficiently many iterations such that the errors ε_t and ε'_t are sufficiently small, for σ defined in (4.13), $\sigma = \tilde{O}(1)$ can be achieved. Finally, we remark that, with some recent advances in the variance reduced techniques for policy optimization (Papini et al., 2018; Xu et al., 2020; Shen et al., 2019; Xu et al., 2019; Huang et al., 2020), we expect the MSEs in (B.1) and (B.2) being further reduced, which could possibly allow us to use relatively smaller neural networks for $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}$ and $\mathcal{E}_{\theta_{t+1}, (i)}$. This can help to boost practicality of smooth FSP due to the reduced computational cost associated with the reduced network size. We leave this direction to our future research.

C PROOF OF THEOREM 4.5

In this section, we lay out the proof of Theorem 4.5. We have the following proposition on the regularization bias of $\mathcal{J}_{(i)}(\pi^i, \pi^{-i})$ defined in (4.3).

Proposition C.1 (Regularization Bias). The regularized performance function $\mathcal{J}_{(i)}(\pi^i, \pi^{-i})$ satisfies

$$\mathcal{J}(\pi^1, \pi^2) < \mathcal{J}_{(1)}(\pi^1, \pi^2) \leq \mathcal{J}(\pi^1, \pi^2) + \frac{\lambda_1 \cdot \log |\mathcal{A}^1|}{1 - \gamma}, \quad (\text{C.1})$$

$$-\mathcal{J}(\pi^1, \pi^2) < \mathcal{J}_{(2)}(\pi^2, \pi^1) \leq -\mathcal{J}(\pi^1, \pi^2) + \frac{\lambda_2 \cdot \log |\mathcal{A}^2|}{1 - \gamma} \quad (\text{C.2})$$

for all policy pairs $[\pi^1, \pi^2]$.

Proof. See Appendix D for a detailed proof. \square

Based on Assumptions 4.2 and 4.3, we have the following lemma on the Lipschitz continuity of the marginalized $Q_{(i)}$ -function. Recall that ζ is the concentrability coefficient in Assumption 4.2 and ι_i is the Lipschitz coefficient in Assumption 4.3.

Lemma C.2 (Lipschitz Marginalized $Q_{(i)}$ -Function). Suppose that Assumption 4.2 holds. We choose the regularization parameter λ_i , the truncation parameter $\mathcal{E}_{(i)}^{\max}$, and the stepsizes $\alpha_{t, (i)}, \alpha'_{t, (i)}$ as in Theorem 4.5. We have for all $[\pi^i; \pi^{-i}] = [\pi_{\theta_t}^i; \pi_{\theta_t}^{-i}]$ generated by the policy update in Line 4 of

Algorithm 1 that,

$$\begin{aligned} & \left| \mathbb{E}_{\nu^*} \left[\langle \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s, \cdot) - \tilde{Q}_{(i)}^{\pi^i, \pi_*^{-i}}(s, \cdot), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \rangle \right] \right| \\ & \leq \sqrt{2l_i \cdot l_{-i}} \cdot \left\{ \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s))^{1/2} \cdot \text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s))^{1/2} \right] \right. \\ & \quad \left. + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2} \right\}. \end{aligned}$$

Proof. See Appendix E for a detailed proof. \square

Recall that $\mathcal{J}_{(i)}(\pi^i, \pi^{-i})$ is defined in (4.3). We present the following extended performance difference lemma, which extends the performance difference lemma of Kakade & Langford (2002) to the two-agent setting with entropy regularization.

Lemma C.3 (Extended Performance Difference). We have for all $[\pi^i; \pi^{-i}]$ that,

$$\begin{aligned} & [\mathcal{J}_{(i)}(\pi_*^i, \pi_*^{-i}) - \mathcal{J}_{(i)}(\pi^i, \pi^{-i})] + \frac{\lambda_i}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right] \\ & = \frac{1}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\langle \tilde{Q}_{(i)}^{\pi^i, \pi_*^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi^i(\cdot | s), \pi_*^i(\cdot | s) - \pi^i(\cdot | s) \rangle \right]. \end{aligned}$$

Proof. See Appendix F for a detailed proof. \square

The following lemma establishes the one-step descent of the KL-divergence between a Nash equilibrium $[\pi_*^1; \pi_*^2]$ and the policy sequence $\{[\pi_{\theta_t}^1; \pi_{\theta_t}^2]\}_{0 \leq t \leq T-1}$ generated by Line 4 of Algorithm 1 in the infinite-dimensional policy space, which extends the analysis of mirror descent (Nemirovski & Yudin, 1983; Nesterov, 2013). Recall that the energy function $\mathcal{E}_{\theta_t, (i)}$ is obtained in Line 7 of Algorithm 1 and the ideal energy function update $\bar{\mathcal{E}}_{t+1, (i)}$ is defined in (4.10).

Lemma C.4 (One-Step Descent). Suppose that the stepsizes satisfy $\alpha'_{t, (i)} = 1 - \lambda_i \alpha_{t, (i)}$. For the policy sequence $\{\pi_{\theta_t}^i\}_{0 \leq t \leq T-1}$ generated by the policy update in Line 4 of Algorithm 1, we have for all $s \in \mathcal{S}$ that,

$$\begin{aligned} & \text{KL}(\pi_*^i(\cdot | s) \| \pi_{\theta_{t+1}}^i(\cdot | s)) - \text{KL}(\pi_*^i(\cdot | s) \| \pi_{\theta_t}^i(\cdot | s)) \tag{C.3} \\ & \geq \langle \mathcal{E}_{\theta_{t+1}, (i)}(s, \cdot) - \bar{\mathcal{E}}_{t+1, (i)}(s, \cdot), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle - 1/2 \cdot \|\mathcal{E}_{\theta_{t+1}, (i)}(s, \cdot) - \mathcal{E}_{\theta_t, (i)}(s, \cdot)\|_{\infty}^2 \\ & \quad + \alpha_{t, (i)} \cdot \langle -\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) + \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle. \end{aligned}$$

Proof. See Appendix G for a detailed proof. \square

Proof of Theorem 4.5. For notational simplicity, we write $\text{KL}(\pi_*^i(\cdot | s) \| \pi_{\theta_{t+1}}^i(\cdot | s))$ as $\text{KL}_{t, (i)}(s)$ throughout this proof. By the choices of the stepsizes in (4.15) of Theorem 4.5, we have $\alpha'_{t, (i)} = 1 - \lambda_i \alpha_{t, (i)}$. By Lemma C.4, we have under (4.12) of Assumption 4.4 that,

$$\begin{aligned} & \mathbb{E}_{\nu^*} [\text{KL}_{t+1, (i)}(s)] - \mathbb{E}_{\nu^*} [\text{KL}_{t, (i)}(s)] \tag{C.4} \\ & \leq \underbrace{\epsilon'_t + \alpha_{t, (i)} \cdot \mathbb{E}_{\nu^*} \left[\langle -\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_*^{-i}}(s, \cdot) + \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle \right]}_{\text{(I)}} \\ & \quad + \underbrace{\alpha_{t, (i)} \cdot \mathbb{E}_{\nu^*} \left[\langle -\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) + \tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_*^{-i}}(s, \cdot), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle \right]}_{\text{(II)}} \\ & \quad + \underbrace{\mathbb{E}_{\nu^*} \left[1/2 \cdot \|\mathcal{E}_{\theta_{t+1}, (i)}(s, \cdot) - \mathcal{E}_{\theta_t, (i)}(s, \cdot)\|_{\infty}^2 \right]}_{\text{(III)}}. \end{aligned}$$

For (I), by Lemma C.3, we have

$$(I) = (1 - \gamma) \cdot [\mathcal{J}_{(i)}(\pi_{\theta_t}^i, \pi_{*}^{-i}) - \mathcal{J}_{(i)}(\pi_{*}^i, \pi_{*}^{-i})] - \lambda_i \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)]. \quad (\text{C.5})$$

For (II), by Lemma C.2, we have

$$\begin{aligned} (\text{II}) &\leq \sqrt{2\iota_i \cdot \iota_{-i}} \cdot \left\{ \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}^{1/2}(s) \cdot \text{KL}_{t,(i)}^{1/2}(s)] \right. \\ &\quad \left. + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)]^{1/2} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(-i)}(s)]^{1/2} \right\} \\ &\leq \left(1 + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \right) \cdot \left(\iota_i \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)] + \iota_{-i} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(-i)}(s)] \right). \end{aligned} \quad (\text{C.6})$$

For (III), by (4.11) of Assumption 4.4 and the definition of $\widehat{\mathcal{E}}_{t+1,(i)}$ in (3.8), we have

$$\begin{aligned} (\text{III}) &\leq \mathbb{E}_{\nu^*} \left[\left\| \mathcal{E}_{\theta_{t+1,(i)}}(s, \cdot) - \widehat{\mathcal{E}}_{t+1,(i)}(s, \cdot) \right\|_{\infty}^2 + \left\| \widehat{\mathcal{E}}_{t+1,(i)}(s, \cdot) - \mathcal{E}_{\theta_{t,(i)}}(s, \cdot) \right\|_{\infty}^2 \right] \\ &\leq \epsilon_t + \alpha_{t,(i)}^2 \cdot \mathbb{E}_{\nu^*} \left[\left\| -\lambda_i \cdot \mathcal{E}_{\theta_{t,(i)}}(s, \cdot) + \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) \right\|_{\infty}^2 \right] \\ &\leq \epsilon_t + \alpha_{t,(i)}^2 \cdot \mathbb{E}_{\nu^*} \left[2\lambda_i^2 \cdot \left\| \mathcal{E}_{\theta_{t,(i)}}(s, \cdot) \right\|_{\infty}^2 + 2 \cdot \left\| \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) \right\|_{\infty}^2 \right] \\ &\leq \epsilon_t + [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2 \cdot \alpha_{t,(i)}^2. \end{aligned} \quad (\text{C.7})$$

Here the last inequality follows from the truncations $\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}} \in \mathcal{F}_{Q_{(i)}^{\max}}$ and $\mathcal{E}_{\theta_{t,(i)}} \in \mathcal{F}_{\mathcal{E}_{(i)}^{\max}}$, where $\mathcal{E}_{(i)}^{\max} = Q_{(i)}^{\max} / (\lambda_i - M_i)$.

Then plugging (C.5), (C.6), and (C.7) into (C.4), we obtain

$$\begin{aligned} \mathbb{E}_{\nu^*} [\text{KL}_{t+1,(i)}(s)] &- \left\{ 1 - \left[\lambda_i - \left(1 + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \right) \cdot \iota_i \right] \cdot \alpha_{t,(i)} \right\} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)] \\ &\leq (\epsilon_t + \epsilon'_t) + \left(1 + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \right) \cdot \iota_{-i} \cdot \alpha_{t,(i)} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(-i)}(s)]^{1/2} \\ &\quad + (1 - \gamma) \alpha_{t,(i)} \cdot [\mathcal{J}_{(i)}(\pi_{\theta_t}^i, \pi_{*}^{-i}) - \mathcal{J}_{(i)}(\pi_{*}^i, \pi_{*}^{-i})] + [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2 \cdot \alpha_{t,(i)}^2. \end{aligned} \quad (\text{C.8})$$

Summing (C.8) for $i \in \{1, 2\}$ and setting $\alpha_{t,(1)} = \alpha_{t,(2)} = \eta_t$, we obtain

$$\begin{aligned} (1 - \gamma) \eta_t \cdot \sum_{i \in \{1, 2\}} [\mathcal{J}_{(i)}(\pi_{*}^i, \pi_{*}^{-i}) - \mathcal{J}_{(i)}(\pi_{\theta_t}^i, \pi_{*}^{-i})] \\ \leq \sum_{i \in \{1, 2\}} \left\{ 1 - \left\{ \lambda_i - \left[2 + \frac{\sum_{i \in \{1, 2\}} (V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta)}{1 - \gamma} \right] \cdot \iota_i \right\} \cdot \eta_t \right\} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)] \\ - \sum_{i \in \{1, 2\}} \mathbb{E}_{\nu^*} [\text{KL}_{t+1,(i)}(s)] + \left[2(\epsilon_t + \epsilon'_t) + \sum_{i \in \{1, 2\}} [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2 \cdot \eta_t^2 \right] \\ = \sum_{i \in \{1, 2\}} [1 - (\lambda_i - M_i) \cdot \eta_t] \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)] - \sum_{i \in \{1, 2\}} \mathbb{E}_{\nu^*} [\text{KL}_{t+1,(i)}(s)] \\ + \left\{ 2(\epsilon_t + \epsilon'_t) + \sum_{i \in \{1, 2\}} [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2 \cdot \eta_t^2 \right\}. \end{aligned} \quad (\text{C.9})$$

Multiplying the both sides of (C.9) by $t + 1$ and setting $\eta_t = 1/[(t + 1) \cdot \min_{i \in \{1,2\}} \{\lambda_i - M_i\}]$, we obtain for all $0 \leq t \leq T - 1$,

$$\begin{aligned} & \frac{1 - \gamma}{\min_{i \in \{1,2\}} \{\lambda_i - M_i\}} \cdot \sum_{i \in \{1,2\}} [\mathcal{J}_{(i)}(\pi_*^i, \pi_*^{-i}) - \mathcal{J}_{(i)}(\pi_{\theta_t}^i, \pi_*^{-i})] \\ & \leq t \cdot \sum_{i \in \{1,2\}} \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)] - (t + 1) \cdot \sum_{i \in \{1,2\}} \mathbb{E}_{\nu^*} [\text{KL}_{t+1,(i)}(s)] \\ & \quad + \left\{ 2(t + 1) \cdot (\epsilon_t + \epsilon'_t) + \frac{\sum_{i \in \{1,2\}} [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2}{(t + 1) \cdot \min_{i \in \{1,2\}} \{\lambda_i - M_i\}^2} \right\}. \end{aligned} \quad (\text{C.10})$$

Telescoping (C.10) over $0 \leq t \leq T - 1$, by $\sum_{t=1}^T 1/t \leq \log T$ and the nonnegativity of the KL-divergence, we obtain

$$\begin{aligned} & \frac{1 - \gamma}{T \cdot \min_{i \in \{1,2\}} \{\lambda_i - M_i\}} \cdot \sum_{t=0}^{T-1} \sum_{i \in \{1,2\}} [\mathcal{J}_{(i)}(\pi_*^i, \pi_*^{-i}) - \mathcal{J}_{(i)}(\pi_{\theta_t}^i, \pi_*^{-i})] \\ & \leq \frac{1}{T} \cdot \left\{ 2 \cdot \sum_{t=0}^{T-1} (t + 1) \cdot (\epsilon_t + \epsilon'_t) + \frac{\sum_{i \in \{1,2\}} [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2}{\min_{i \in \{1,2\}} \{\lambda_i - M_i\}^2} \cdot \sum_{t=1}^T \frac{1}{t} \right\} \\ & \quad - \sum_{i \in \{1,2\}} \mathbb{E}_{\nu^*} [\text{KL}_{T,(i)}(s)] \\ & \leq \frac{2\sigma}{T} + \frac{\sum_{i \in \{1,2\}} [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2}{\min_{i \in \{1,2\}} \{\lambda_i - M_i\}^2} \cdot \frac{\log T}{T}, \end{aligned} \quad (\text{C.11})$$

where σ is defined in (4.13). Finally, by Proposition C.1, we have

$$\sum_{i \in \{1,2\}} [\mathcal{J}_{(i)}(\pi_*^i, \pi_*^{-i}) - \mathcal{J}_{(i)}(\pi_{\theta_t}^i, \pi_*^{-i})] \geq \mathcal{J}(\pi_*^1, \pi_{\theta_t}^2) - \mathcal{J}(\pi_{\theta_t}^1, \pi_*^2) - \sum_{i \in \{1,2\}} \lambda_i \cdot \log |\mathcal{A}^i|,$$

combining which with (C.11), we obtain

$$\begin{aligned} & \frac{1}{T} \cdot \sum_{t=0}^{T-1} [\mathcal{J}(\pi_*^1, \pi_{\theta_t}^2) - \mathcal{J}(\pi_{\theta_t}^1, \pi_*^2)] \\ & \leq \frac{2\sigma \cdot \min_{i \in \{1,2\}} \{\lambda_i - M_i\}}{(1 - \gamma) \cdot T} + \frac{\sum_{i \in \{1,2\}} [2 + 2\lambda_i^2 / (\lambda_i - M_i)^2] \cdot (Q_{(i)}^{\max})^2}{(1 - \gamma) \cdot \min_{i \in \{1,2\}} \{\lambda_i - M_i\}^2} \cdot \frac{\log T}{T} + \sum_{i \in \{1,2\}} \lambda_i \cdot \log |\mathcal{A}^i|. \end{aligned}$$

Thus, we conclude the proof of Theorem 4.5. \square

D PROOF OF PROPOSITION C.1

Proof. By the definition of $\mathcal{J}_{(1)}(\pi^1, \pi^2)$ and $\mathcal{J}(\pi^1, \pi^2)$, we have

$$\begin{aligned} \mathcal{J}_{(1)}(\pi^1, \pi^2) - \mathcal{J}(\pi^1, \pi^2) &= \mathbb{E}_{\nu^*} [V_{(1)}^{\pi^1, \pi^2}(s) - V_1^{\pi^1, \pi^2}(s)] \\ &= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{\nu^*}^{\pi^1, \pi^2}} [r_{(1)}^{\pi^1, \pi^2}(s') - r_1^{\pi^1, \pi^2}(s')] \\ &= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{\nu^*}^{\pi^1, \pi^2}} [\lambda_1 \cdot H(\pi^1(\cdot | s))]. \end{aligned}$$

Since $0 < H(\pi^1(\cdot | s)) \leq \log |\mathcal{A}^1|$, we further obtain

$$\mathcal{J}(\pi^1, \pi^2) < \mathcal{J}_{(1)}(\pi^1, \pi^2) \leq \mathcal{J}(\pi^1, \pi^2) + \frac{\lambda_1 \cdot \log |\mathcal{A}^1|}{1 - \gamma},$$

which finishes the proof of (C.1). Using the same argument, we can also prove (C.2). Thus, we conclude the proof of Proposition C.1. \square

E PROOF OF LEMMA C.2

In the subsequent analysis, using the notion of the state-transition operator $\mathcal{P}^{\pi^i, \pi^{-i}}$ in (4.7), we write $\rho_{s, a^i, \pi^{-i}}^{\pi^i, \pi^{-i}}$, which is defined in (4.5), as

$$\rho_{s, a^i, \pi^{-i}}^{\pi^i, \pi^{-i}} = \left[(1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot (\mathcal{P}^{\pi^i, \pi^{-i}})^t \right] \circ \mathcal{P}^{a^i, \pi^{-i}} \circ \delta_s, \quad (\text{E.1})$$

where δ_s is the Dirac delta function.

Lemma E.1. Under Assumption 4.3, we have

$$\left\| \sum_{t=0}^{\infty} \gamma^t \cdot [(\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t] \right\|_{\text{op}} \leq \frac{\gamma l_i}{(1 - \gamma)^2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}.$$

Proof. Since $\|\mathcal{P}^{\pi^1, \pi^2}\|_{\text{op}} \leq 1$ (Lasota & Mackey, 2013), we have

$$\begin{aligned} & \left\| (\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t \right\|_{\text{op}} \\ &= \left\| (\mathcal{P}^{\pi^i, \pi^{-i}})^{t-1} \circ (\mathcal{P}^{\pi^i, \pi^{-i}} - \mathcal{P}^{\pi^i, \pi_*^{-i}}) + [(\mathcal{P}^{\pi^i, \pi^{-i}})^{t-1} - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^{t-1}] \circ \mathcal{P}^{\pi^i, \pi_*^{-i}} \right\|_{\text{op}} \\ &\leq \left\| (\mathcal{P}^{\pi^i, \pi^{-i}})^{t-1} \circ (\mathcal{P}^{\pi^i, \pi^{-i}} - \mathcal{P}^{\pi^i, \pi_*^{-i}}) \right\|_{\text{op}} + \left\| [(\mathcal{P}^{\pi^i, \pi^{-i}})^{t-1} - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^{t-1}] \circ \mathcal{P}^{\pi^i, \pi_*^{-i}} \right\|_{\text{op}} \\ &\leq \|\mathcal{P}^{\pi^i, \pi^{-i}} - \mathcal{P}^{\pi^i, \pi_*^{-i}}\|_{\text{op}} + \left\| (\mathcal{P}^{\pi^i, \pi^{-i}})^{t-1} - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^{t-1} \right\|_{\text{op}}. \end{aligned} \quad (\text{E.2})$$

Recursively applying (E.2) gives

$$\begin{aligned} \left\| (\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t \right\|_{\text{op}} &\leq t \cdot \|\mathcal{P}^{\pi^i, \pi^{-i}} - \mathcal{P}^{\pi^i, \pi_*^{-i}}\|_{\text{op}} \\ &\leq t \cdot l_i \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s))^{1/2} \right], \end{aligned}$$

where the second inequality follows from Assumption 4.3. Thus, we have

$$\begin{aligned} \left\| \sum_{t=0}^{\infty} \gamma^t \cdot [(\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t] \right\|_{\text{op}} &\leq \sum_{t=0}^{\infty} \gamma^t \cdot \left\| (\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t \right\|_{\text{op}} \\ &\leq \left(\sum_{t=0}^{\infty} t \gamma^t \right) \cdot l_i \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2} \\ &= \frac{\gamma l_i}{(1 - \gamma)^2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}, \end{aligned}$$

which concludes the proof of Lemma E.1. \square

Now we are ready to prove Lemma C.2.

Proof of Lemma C.2. By the definition of the $Q_{(i)}$ -function in (2.5), we have

$$Q_{(i)}^{\pi^i, \pi^{-i}}(s, a^i, a^{-i}) = r_i(s, a^i, a^{-i}) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a^i, a^{-i})} [V_{(i)}^{\pi^i, \pi^{-i}}(s')],$$

which gives

$$\begin{aligned}
& \left| \mathbb{E}_{\nu^*} \left[\left\langle \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s, \cdot) - \tilde{Q}_{(i)}^{\pi^i, \pi_*^{-i}}(s, \cdot), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\rangle \right] \right| & (E.3) \\
&= \left| \mathbb{E}_{\nu^*} \left[\left\langle r_i^{\cdot, \pi^{-i}}(s) - r_i^{\cdot, \pi_*^{-i}}(s), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\rangle \right. \right. \\
&\quad \left. \left. + \left\langle [\mathcal{P}^{\cdot, \pi^{-i}} \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) - [\mathcal{P}^{\cdot, \pi_*^{-i}} \circ V_{(i)}^{\pi^i, \pi_*^{-i}}](s), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\rangle \right] \right| \\
&= \underbrace{\left| \mathbb{E}_{\nu^*} \left[\left\langle r_i^{\cdot, \pi^{-i}}(s) - r_i^{\cdot, \pi_*^{-i}}(s), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\rangle \right] \right|}_{(I)} \\
&\quad + \underbrace{\left| \mathbb{E}_{\nu^*} \left[\left\langle [\mathcal{P}^{\cdot, \pi^{-i}} \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) - [\mathcal{P}^{\cdot, \pi_*^{-i}} \circ V_{(i)}^{\pi^i, \pi_*^{-i}}](s), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\rangle \right] \right|}_{(II)}.
\end{aligned}$$

Upper Bounding (I): By the Cauchy-Schwartz inequality, we have

$$\begin{aligned}
(I) &\leq \mathbb{E}_{\nu^*} \left[\left\| r_i^{\cdot, \pi^{-i}}(s) - r_i^{\cdot, \pi_*^{-i}}(s) \right\|_{\infty} \cdot \left\| \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\|_1 \right] & (E.4) \\
&\leq \sqrt{2} \cdot l_i \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^{-i}(\cdot | s))^{1/2} \cdot \text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s))^{1/2} \right],
\end{aligned}$$

where the second inequality follows from Assumption 4.3 and the Pinsker's inequality. Meanwhile, we have

$$\begin{aligned}
(I) &= \left| \mathbb{E}_{\nu^*} \left[\left\langle r_i^{\pi^i, \cdot}(s) - r_i^{\pi_*^i, \cdot}(s), \pi^{-i}(\cdot | s) - \pi_*^{-i}(\cdot | s) \right\rangle \right] \right| & (E.5) \\
&\leq \sqrt{2} \cdot l_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s))^{1/2} \cdot \text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s))^{1/2} \right],
\end{aligned}$$

where the inequality follows from the same argument as (E.4). Combining (E.4) and (E.5), we obtain

$$(I) \leq \sqrt{2} l_i \cdot l_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s))^{1/2} \cdot \text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s))^{1/2} \right]. \quad (E.6)$$

Upper Bounding (II): Recall that we have

$$\begin{aligned}
V_{(i)}^{\pi^i, \pi^{-i}}(s) &= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{s, \pi^i, \pi^{-i}}} \left[r_{(i)}^{\pi^i, \pi^{-i}}(s') \right] \\
&= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{s, \pi^i, \pi^{-i}}} \left[r_i^{\pi^i, \pi^{-i}}(s') + \lambda_i \cdot H(\pi^i(\cdot | s')) \right],
\end{aligned}$$

which gives

$$\begin{aligned}
(1 - \gamma) \cdot & \left| [\mathcal{P}^{a^i, \pi^{-i}} \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) - [\mathcal{P}^{a^i, \pi_*^{-i}} \circ V_{(i)}^{\pi^i, \pi_*^{-i}}](s) \right| & (E.7) \\
&= \left| \mathbb{E}_{s' \sim \rho_{s, a^i, \pi^{-i}}} \left[r_i^{\pi^i, \pi^{-i}}(s') + \lambda_i \cdot H(\pi^i(\cdot | s')) \right] - \mathbb{E}_{s' \sim \rho_{s, a^i, \pi_*^{-i}}} \left[r_i^{\pi^i, \pi_*^{-i}}(s') + \lambda_i \cdot H(\pi^i(\cdot | s')) \right] \right| \\
&\leq \left| \mathbb{E}_{s' \sim \rho_{s, a^i, \pi^{-i}}} \left[r_i^{\pi^i, \pi^{-i}}(s') + \lambda_i \cdot H(\pi^i(\cdot | s')) \right] - \mathbb{E}_{s' \sim \rho_{s, a^i, \pi_*^{-i}}} \left[r_i^{\pi^i, \pi^{-i}}(s') + \lambda_i \cdot H(\pi^i(\cdot | s')) \right] \right| \\
&\quad + l_i \cdot \mathbb{E}_{s' \sim \rho_{s, a^i, \pi_*^{-i}}} \left[\left\| \pi^{-i}(\cdot | s') - \pi_*^{-i}(\cdot | s') \right\|_1 \right] \\
&\leq (1 + \lambda_i \cdot \log |\mathcal{A}^i|) \cdot \underbrace{\left\| \rho_{s, a^i, \pi^{-i}}^{\pi^i, \pi^{-i}} - \rho_{s, a^i, \pi_*^{-i}}^{\pi^i, \pi_*^{-i}} \right\|_{L_1}(S)}_{(III)} + l_i \cdot \underbrace{\mathbb{E}_{s' \sim \rho_{s, a^i, \pi_*^{-i}}} \left[\text{KL}(\pi_*^{-i}(\cdot | s') \| \pi^{-i}(\cdot | s'))^{1/2} \right]}_{(IV)},
\end{aligned}$$

where the first inequality follows from (4.9) of Assumption 4.3, and $\rho_{s, a^i, \pi^{-i}}^{\pi^i, \pi^{-i}}$ is defined in (4.5).

By (E.1), we have

$$\begin{aligned}
\text{(III)} &= (1 - \gamma) \cdot \left\| \left[\sum_{t=0}^{\infty} \gamma^t \cdot (\mathcal{P}^{\pi^i, \pi^{-i}})^t \right] \circ \mathcal{P}^{a^i, \pi^{-i}} \circ \delta_s - \left[\sum_{t=0}^{\infty} \gamma^t \cdot (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t \right] \circ \mathcal{P}^{a^i, \pi_*^{-i}} \circ \delta_s \right\|_{L_1(\mathcal{S})} \\
&\leq (1 - \gamma) \cdot \left\| \left[\sum_{t=0}^{\infty} \gamma^t \cdot [(\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t] \right] \circ \mathcal{P}^{a^i, \pi^{-i}} \circ \delta_s \right\|_{L_1(\mathcal{S})} \\
&\quad + \left\| \left[(1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot (\mathcal{P}^{\pi^i, \pi^{-i}})^t \right] \circ (\mathcal{P}^{a^i, \pi^{-i}} - \mathcal{P}^{a^i, \pi_*^{-i}}) \circ \delta_s \right\|_{L_1(\mathcal{S})}. \tag{E.8}
\end{aligned}$$

Since $\|\mathcal{P}^{\pi^i, \pi^{-i}}\|_{\text{op}} \leq 1$, we have

$$\left\| (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot (\mathcal{P}^{\pi^i, \pi^{-i}})^t \right\|_{\text{op}} \leq (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot \|\mathcal{P}^{\pi^i, \pi^{-i}}\|_{\text{op}}^t \leq 1,$$

plugging which and Lemma E.1 into (E.8) gives

$$\begin{aligned}
\text{(III)} &\leq (1 - \gamma) \cdot \left\| \sum_{t=0}^{\infty} \gamma^t \cdot [(\mathcal{P}^{\pi^i, \pi^{-i}})^t - (\mathcal{P}^{\pi^i, \pi_*^{-i}})^t] \right\|_{\text{op}} \cdot \|\mathcal{P}^{a^i, \pi^{-i}}\|_{\text{op}} \cdot \|\delta_s\|_{L_1(\mathcal{S})} \\
&\quad + \left\| (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot (\mathcal{P}^{\pi^i, \pi^{-i}})^t \right\| \cdot \|\mathcal{P}^{a^i, \pi^{-i}} - \mathcal{P}^{a^i, \pi_*^{-i}}\|_{\text{op}} \cdot \|\delta_s\|_{L_1(\mathcal{S})} \\
&\leq \frac{l_i}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}. \tag{E.9}
\end{aligned}$$

By the Cauchy-Schwartz inequality, we have

$$\begin{aligned}
\text{(IV)} &= \mathbb{E}_{s' \sim \nu^*} \left[\frac{d\rho_{s, a^i, \pi_*^{-i}}^{\pi^i, \pi_*^{-i}}}{d\nu^*} \cdot \text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s))^{1/2} \right] \\
&\leq \mathbb{E}_{\nu^*} \left[\left| \frac{d\rho_{s, a^i, \pi_*^{-i}}^{\pi^i, \pi_*^{-i}}}{d\nu^*} \right|^2 \right]^{1/2} \cdot \mathbb{E}_{s' \sim \nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s') \| \pi^{-i}(\cdot | s')) \right]^{1/2} \\
&\leq \zeta \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}, \tag{E.10}
\end{aligned}$$

where the last inequality follows from Assumption 4.2 and the Pinsker's inequality. Plugging (E.9) and (E.10) into (E.7), we obtain for all $s \in \mathcal{S}$ and $a^i \in \mathcal{A}^i$ that,

$$\begin{aligned}
&|[\mathcal{P}^{a^i, \pi^{-i}} \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) - [\mathcal{P}^{a^i, \pi_*^{-i}} \circ V_{(i)}^{\pi^i, \pi_*^{-i}}](s)| \\
&\leq \frac{V_{(i)}^{\max} + \zeta}{1 - \gamma} \cdot l_i \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2},
\end{aligned}$$

which further gives

$$\begin{aligned}
\text{(II)} &\leq \mathbb{E}_{\nu^*} \left[\left\| [\mathcal{P}^{\pi^i, \pi^{-i}} \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) - [\mathcal{P}^{\pi^i, \pi_*^{-i}} \circ V_{(i)}^{\pi^i, \pi_*^{-i}}](s) \right\|_{\infty} \cdot \|\pi^i(\cdot | s) - \pi_*^i(\cdot | s)\|_1 \right] \tag{E.11} \\
&\leq \frac{V_{(i)}^{\max} + \zeta}{1 - \gamma} \cdot l_i \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\|\pi^i(\cdot | s) - \pi_*^i(\cdot | s)\|_1 \right] \\
&\leq \frac{\sqrt{2} \cdot (V_{(i)}^{\max} + \zeta) \cdot l_i}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2},
\end{aligned}$$

where the first inequality follows from the Cauchy-Schwartz inequality and the last inequality follows from the Pinsker's inequality. Alternatively, we can also write

$$\begin{aligned}
\text{(II)} &= \left| \mathbb{E}_{\nu^*} \left[\left\langle [(\mathcal{P}^{\pi^i, \cdot} - \mathcal{P}^{\pi_*^i, \cdot}) \circ V_{(i)}^{\pi^i, \pi^{-i}}](s), \pi^{-i}(\cdot | s) - \pi_*^{-i}(\cdot | s) \right\rangle \right] \right. \\
&\quad \left. + \mathbb{E}_{\nu^*} \left[\left\langle [\mathcal{P}^{\cdot, \pi_*^{-i}} \circ (V_{(i)}^{\pi^i, \pi^{-i}} - V_{(i)}^{\pi_*^i, \pi_*^{-i}})](s), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \right\rangle \right] \right| \\
&\leq \underbrace{\left| \mathbb{E}_{\nu^*} \left[\left\langle [(\mathcal{P}^{\pi^i, \cdot} - \mathcal{P}^{\pi_*^i, \cdot}) \circ V_{(i)}^{\pi^i, \pi^{-i}}](s), \pi^{-i}(\cdot | s) - \pi_*^{-i}(\cdot | s) \right\rangle \right] \right|}_{\text{(V)}} \\
&\quad + \underbrace{\left| \mathbb{E}_{\nu^*} \left[\left\langle [(\mathcal{P}^{\pi^i, \pi_*^{-i}} - \mathcal{P}^{\pi_*^i, \pi_*^{-i}}) \circ (V_{(i)}^{\pi^i, \pi^{-i}} - V_{(i)}^{\pi_*^i, \pi_*^{-i}})](s) \right\rangle \right] \right|}_{\text{(VI)}}. \tag{E.12}
\end{aligned}$$

By (4.8) of Assumption 4.3, we have for all $s \in \mathcal{S}$ and $a^{-i} \in \mathcal{A}^{-i}$ that,

$$\begin{aligned}
\left| [(\mathcal{P}^{\pi^i, a^{-i}} - \mathcal{P}^{\pi_*^i, a^{-i}}) \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) \right| &\leq V_{(i)}^{\max} \cdot \|\mathcal{P}^{\pi^i, a^{-i}} - \mathcal{P}^{\pi_*^i, a^{-i}}\|_{\text{op}} \\
&\leq V_{(i)}^{\max} \cdot \iota_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2},
\end{aligned}$$

which gives

$$\begin{aligned}
\text{(V)} &\leq \mathbb{E}_{\nu^*} \left[\left\| [(\mathcal{P}^{\pi^i, \cdot} - \mathcal{P}^{\pi_*^i, \cdot}) \circ V_{(i)}^{\pi^i, \pi^{-i}}](s) \right\|_{\infty} \cdot \|\pi^{-i}(\cdot | s) - \pi_*^{-i}(\cdot | s)\|_1 \right] \tag{E.13} \\
&\leq \sqrt{2} \cdot V_{(i)}^{\max} \cdot \iota_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2} \\
&\leq \frac{\sqrt{2} \cdot V_{(i)}^{\max}}{1 - \gamma} \cdot \iota_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}.
\end{aligned}$$

Here the second inequality follows from the Pinsker's inequality.

Meanwhile, we have

$$\begin{aligned}
\text{(VI)} &\leq \mathbb{E}_{\nu^*} \left[\left\| [(\mathcal{P}^{\pi^i, \pi_*^{-i}} - \mathcal{P}^{\pi_*^i, \pi_*^{-i}}) \circ (V_{(i)}^{\pi^i, \pi^{-i}} - V_{(i)}^{\pi_*^i, \pi_*^{-i}})](s) \right\| \right] \tag{E.14} \\
&\leq \|\mathcal{P}^{\pi^i, \pi_*^{-i}} - \mathcal{P}^{\pi_*^i, \pi_*^{-i}}\|_{\text{op}} \cdot \mathbb{E}_{\nu^*} \left[\left| V_{(i)}^{\pi^i, \pi^{-i}}(s) - V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s) \right| \right] \\
&\leq \iota_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\left| V_{(i)}^{\pi^i, \pi^{-i}}(s) - V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s) \right| \right].
\end{aligned}$$

By the same argument in the proof of Lemma C.3, we have

$$V_{(i)}^{\pi^i, \pi^{-i}}(s) - V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s) = \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{\nu^*, \pi_*^i, \pi_*^{-i}}^{\pi^i, \pi^{-i}}} \left[\left\langle \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s', \cdot), \pi_*^{-i}(\cdot | s') - \pi^{-i}(\cdot | s') \right\rangle \right],$$

which gives

$$\begin{aligned}
\mathbb{E}_{\nu^*} \left[\left| V_{(i)}^{\pi^i, \pi^{-i}}(s) - V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s) \right| \right] &\leq \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{\nu^*, \pi_*^i, \pi_*^{-i}}^{\pi^i, \pi^{-i}}} \left[\left\| \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s', \cdot) \right\|_{\infty} \cdot \|\pi_*^{-i}(\cdot | s') - \pi^{-i}(\cdot | s')\|_1 \right] \\
&\leq \frac{Q_{(i)}^{\max}}{1 - \gamma} \cdot \mathbb{E}_{s' \sim \rho_{\nu^*, \pi_*^i, \pi_*^{-i}}^{\pi^i, \pi^{-i}}} \left[\|\pi_*^{-i}(\cdot | s') - \pi^{-i}(\cdot | s')\|_1 \right] \\
&\leq \frac{\sqrt{2} \cdot Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}, \tag{E.15}
\end{aligned}$$

where the last inequality follows from the same arguments in (E.10). Taking (E.15) into (E.14), we obtain

$$\text{(VI)} \leq \frac{\sqrt{2} \cdot Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \iota_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2},$$

pugging which and (E.13) into (E.12) gives

$$(II) \leq \frac{\sqrt{2} \cdot (V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta)}{1 - \gamma} \cdot \iota_{-i} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}. \quad (\text{E.16})$$

Combining (E.11) and (E.16), we obtain

$$(II) \leq \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \sqrt{2\iota_i \cdot \iota_{-i}} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2}. \quad (\text{E.17})$$

Finally, taking (E.6) and (E.17) in to (E.3), we obtain

$$\begin{aligned} & \left| \mathbb{E}_{\nu^*} \left[\langle \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s, \cdot) - \tilde{Q}_{(i)}^{\pi_*^i, \pi_*^{-i}}(s, \cdot), \pi^i(\cdot | s) - \pi_*^i(\cdot | s) \rangle \right] \right| \\ & \leq \sqrt{2\iota_i \cdot \iota_{-i}} \cdot \left\{ \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right\} \\ & \quad + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)) \right]^{1/2} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^{-i}(\cdot | s) \| \pi^{-i}(\cdot | s)) \right]^{1/2} \Bigg\}, \end{aligned}$$

which concludes the proof of Lemma C.2. \square

F PROOF OF LEMMA C.3

Proof. The proof extends that of Lemma 6.1 in Kakade & Langford (2002) to zero-sum Markov games. By the definition of $V_{(i)}^{\pi^i, \pi^{-i}}(s)$ in (2.1), we have

$$\begin{aligned} & V_{(i)}^{\pi^i, \pi^{-i}}(s) \quad (\text{F.1}) \\ & = \mathbb{E}_{\pi_*^i, \pi_*^{-i}} \left[\sum_{t=0}^{\infty} \gamma^t \cdot (r_{(i)}^{\pi_*^i, \pi_*^{-i}}(s_t, a^i, a^{-i}) + V_{(i)}^{\pi^i, \pi^{-i}}(s_t) - V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s_t)) \mid s_0 = s \right] \\ & = \mathbb{E}_{\pi_*^i, \pi_*^{-i}} \left[\sum_{t=0}^{\infty} \gamma^t \cdot (r_{(i)}^{\pi_*^i, \pi_*^{-i}}(s_t, a^i, a^{-i}) + \gamma \cdot V_{(i)}^{\pi^i, \pi^{-i}}(s_{t+1}) - V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s_t)) \mid s_0 = s \right] + V_{(i)}^{\pi^i, \pi^{-i}}(s). \end{aligned}$$

By the definition of $Q_{(i)}$ -function, we have

$$Q_{(i)}^{\pi^i, \pi^{-i}}(s, a^i, a^{-i}) = r(s, a^i, a^{-i}) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a^i, a^{-i})} [V_{(i)}^{\pi^i, \pi^{-i}}(s')],$$

which gives

$$\begin{aligned} & \mathbb{E}_{\pi_*^i, \pi_*^{-i}} \left[r_{(i)}^{\pi_*^i, \pi_*^{-i}}(s) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a^i, a^{-i})} [V_{(i)}^{\pi^i, \pi^{-i}}(s')] - V_{(i)}^{\pi^i, \pi^{-i}}(s) \right] \quad (\text{F.2}) \\ & = \mathbb{E}_{\pi_*^i, \pi_*^{-i}} \left[r(s, a^1, a^{-i}) - \lambda_i \cdot \log \pi_*^i(a^i | s) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a^i, a^{-i})} [V_{(i)}^{\pi^i, \pi^{-i}}(s')] - V_{(i)}^{\pi^i, \pi^{-i}}(s) \right] \\ & = \mathbb{E}_{\pi_*^i, \pi_*^{-i}} [Q_{(i)}^{\pi^i, \pi^{-i}}(s, a^i, a^{-i}) - V_{(i)}^{\pi^i, \pi^{-i}}(s)] + \lambda_i \cdot H(\pi_*^i(\cdot | s)) \\ & = \langle \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s, \cdot), \pi_*^i(\cdot | s) - \pi^i(\cdot | s) \rangle + \lambda_i \cdot H(\pi_*^i(\cdot | s)) - \lambda_i \cdot H(\pi^i(\cdot | s)). \end{aligned}$$

Here the last equality follows from

$$\begin{aligned} V_{(i)}^{\pi^i, \pi^{-i}}(s) & = \mathbb{E}_{\pi_*^i, \pi_*^{-i}} [Q_{(i)}^{\pi^i, \pi^{-i}}(s, a^i, a^{-i}) - \lambda_i \cdot \log \pi^i(a^i | s)] \\ & = \langle \tilde{Q}_{(i)}^{\pi^i, \pi^{-i}}(s, \cdot), \pi^i(\cdot | s) \rangle + \lambda_i \cdot H(\pi^i(\cdot | s)). \end{aligned}$$

For the entropy terms in (F.2), we have

$$\begin{aligned} H(\pi_*^i(\cdot | s)) & = \langle -\log \pi_*^i(\cdot | s), \pi_*^i(\cdot | s) \rangle \\ & = \langle -\log \pi^i(\cdot | s), \pi_*^i(\cdot | s) \rangle - \langle \log(\pi_*^i(\cdot | s)/\pi^i(\cdot | s)), \pi_*^i(\cdot | s) \rangle \\ & = \langle -\log \pi^i(\cdot | s), \pi_*^i(\cdot | s) \rangle - \text{KL}(\pi_*^i(\cdot | s) \| \pi^i(\cdot | s)), \end{aligned}$$

which gives

$$\begin{aligned} & H(\pi_*^i(\cdot | s)) - H(\pi^i(\cdot | s)) \\ &= \langle -\log \pi^i(\cdot | s), \pi_*^i(\cdot | s) - \pi^i(\cdot | s) \rangle - \text{KL}(\pi_*^i(\cdot | s) \parallel \pi^i(\cdot | s)). \end{aligned} \quad (\text{F.3})$$

Taking (F.2) and (F.3) into (F.1), we have

$$\begin{aligned} & V_{(i)}^{\pi_*^i, \pi_*^{-i}}(s) - V_{(i)}^{\pi^i, \pi^{-i}}(s) \\ &= \mathbb{E}_{\pi_*^i, \pi_*^{-i}} \left[\sum_{t=0}^{\infty} \gamma^t \cdot \left(\langle \tilde{Q}_{(i)}^{\pi^i, \pi_*^{-i}}(s_t, \cdot) - \lambda_i \cdot \log \pi^i(\cdot | s_t), \pi_*^i(\cdot | s_t) - \pi^i(\cdot | s_t) \rangle \right. \right. \\ &\quad \left. \left. - \lambda_i \cdot \text{KL}(\pi_*^i(\cdot | s_t) \parallel \pi^i(\cdot | s_t)) \right) \Big| s_0 = s \right] \\ &= \frac{1}{1-\gamma} \cdot \mathbb{E}_{s' \sim \rho_{s'}^{\pi_*^i, \pi_*^{-i}}} \left[\langle \tilde{Q}_{(i)}^{\pi^i, \pi_*^{-i}}(s', \cdot) - \lambda_i \cdot \log \pi^i(\cdot | s'), \pi_*^i(\cdot | s') - \pi^i(\cdot | s') \rangle \right. \\ &\quad \left. - \lambda_i \cdot \text{KL}(\pi_*^i(\cdot | s') \parallel \pi^i(\cdot | s')) \right], \end{aligned} \quad (\text{F.4})$$

where $\rho_{s'}^{\pi_*^i, \pi_*^{-i}}$ is defined in (4.4) as the visitation measure of the policy pair $[\pi_*^i; \pi_*^{-i}]$ starting from state s . Taking $\mathbb{E}_{\nu^*}[\cdot]$ on both sides of (F.4) and recalling the definition of $\mathcal{J}_{(i)}(\pi^i, \pi^{-i})$ in (4.3), we have

$$\begin{aligned} & \mathcal{J}_{(i)}(\pi_*^i, \pi_*^{-i}) - \mathcal{J}_{(i)}(\pi^i, \pi^{-i}) + \frac{\lambda_i}{1-\gamma} \cdot \mathbb{E}_{\nu^*} \left[\text{KL}(\pi_*^i(\cdot | s) \parallel \pi^i(\cdot | s)) \right] \\ &= \frac{1}{1-\gamma} \cdot \mathbb{E}_{\nu^*} \left[\langle \tilde{Q}_{(i)}^{\pi^i, \pi_*^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi^i(\cdot | s), \pi_*^i(\cdot | s) - \pi^i(\cdot | s) \rangle \right]. \end{aligned}$$

Here we use the fact that $\mathbb{E}_{s' \sim \nu^*}[\cdot] = \mathbb{E}_{s' \sim \rho_{s'}^{\pi_*^i, \pi_*^{-i}}, s \sim \nu^*}[\cdot]$. Thus, we finish the proof of Lemma C.3. \square

G PROOF OF LEMMA C.4

Proof. First, we have for any $s \in \mathcal{S}$ that,

$$\begin{aligned} & \text{KL}(\pi_*^i(\cdot | s) \parallel \pi_{\theta_t}^i(\cdot | s)) - \text{KL}(\pi_*^i(\cdot | s) \parallel \pi_{\theta_{t+1}}^i(\cdot | s)) \\ &= \text{KL}(\pi_{\theta_{t+1}}^i(\cdot | s) \parallel \pi_{\theta_t}^i(\cdot | s)) + \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right], \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \right\rangle \\ &\quad + \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right], \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle \\ &= \text{KL}(\pi_{\theta_{t+1}}^i(\cdot | s) \parallel \pi_{\theta_t}^i(\cdot | s)) + \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right], \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \right\rangle \\ &\quad + \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right] - \alpha_{t,(i)} \cdot (\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s)), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle \\ &\quad + \alpha_{t,(i)} \cdot \langle \tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle. \end{aligned} \quad (\text{G.1})$$

Recall that

$$\begin{aligned} \bar{\pi}_{\theta_{t+1}}^i(\cdot | s) &\propto \exp\{(1 - \lambda_i \alpha_{t,(i)}) \cdot \mathcal{E}_{\theta_t,(i)}(s, \cdot) + \alpha_{t,(i)} \cdot \tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot)\} \\ &\propto \exp\{\mathcal{E}_{\theta_t,(i)}(s, \cdot) + \alpha_{t,(i)} \cdot (\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s))\}. \end{aligned}$$

Let

$$Z_{t+1,(i)}(s) = \sum_{a^i \in \mathcal{A}^i} \exp\{\mathcal{E}_{\theta_t,(i)}(s, a^i) + \alpha_{t,(i)} \cdot (\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i) - \lambda_i \cdot \log \pi_{\theta_t}^i(a^i | s))\},$$

and

$$Z_{\theta_{t,(i)}}(s) = \sum_{a^i \in \mathcal{A}^i} \exp\{\mathcal{E}_{\theta_{t+1,(i)}}(s, a^i)\}.$$

where are only dependent on the state s . It can be verified that $\langle \log Z_{\theta_{t,(i)}}(s), \pi^i(\cdot | s) - \pi^{i'}(\cdot | s) \rangle = \langle \log Z_{t,(i)}(s), \pi^i(\cdot | s) - \pi^{i'}(\cdot | s) \rangle = 0$ for all t , π^i , and $\pi^{i'}$, which implies that, on the right-hand-side of (G.1),

$$\begin{aligned} & \left\langle \log \pi_{\theta_t}^i(\cdot | s) + \alpha_{t,(i)} \cdot (\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s)), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle \\ &= \left\langle \mathcal{E}_{\theta_{t,(i)}}(s, \cdot) - \log Z_{\theta_{t,(i)}}(s) + \alpha_{t,(i)} \cdot (\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s)), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle \\ &= \left\langle \mathcal{E}_{\theta_{t,(i)}}(s, \cdot) + \alpha_{t,(i)} \cdot (\tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s)) - \log Z_{t+1,(i)}(s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle \\ &= \langle \log \bar{\pi}_{t+1}^i(\cdot | s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \rangle, \end{aligned} \quad (\text{G.2})$$

and

$$\begin{aligned} & \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right], \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \right\rangle \\ &= \langle \mathcal{E}_{\theta_{t+1,(i)}}(s, \cdot) - \mathcal{E}_{\theta_{t,(i)}}(s, \cdot), \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \rangle \\ &\quad - \langle \log Z_{\theta_{t+1,(i)}}(s), \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \rangle + \langle \log Z_{\theta_{t,(i)}}(s), \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \rangle \\ &= \langle \mathcal{E}_{\theta_{t+1,(i)}}(s, \cdot) - \mathcal{E}_{\theta_{t,(i)}}(s, \cdot), \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \rangle. \end{aligned} \quad (\text{G.3})$$

Plugging (G.2) and (G.3) into (G.1), we obtain

$$\text{KL}(\pi_*^i(\cdot | s) \| \pi_{\theta_t}^i(\cdot | s)) - \text{KL}(\pi_*^i(\cdot | s) \| \pi_{\theta_{t+1}}^i(\cdot | s)) \quad (\text{G.4})$$

$$\begin{aligned} &= \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right], \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle + \langle \mathcal{E}_{\theta_{t+1,(i)}}(s, \cdot) - \mathcal{E}_{\theta_{t,(i)}}(s, \cdot), \pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s) \rangle \\ &\quad + \alpha_{t,(i)} \cdot \left\langle \tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle + \text{KL}(\pi_{\theta_{t+1}}^i(\cdot | s) \| \pi_{\theta_t}^i(\cdot | s)) \\ &\geq \left\langle \log \left[\frac{\pi_{\theta_{t+1}}^i(\cdot | s)}{\pi_{\theta_t}^i(\cdot | s)} \right], \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle - \|\mathcal{E}_{\theta_{t+1,(i)}}(s, \cdot) - \mathcal{E}_{\theta_{t,(i)}}(s, \cdot)\|_\infty \cdot \|\pi_{\theta_t}^i(\cdot | s) - \pi_{\theta_{t+1}}^i(\cdot | s)\|_1 \\ &\quad + \alpha_{t,(i)} \cdot \left\langle \tilde{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, \cdot) - \lambda_i \cdot \log \pi_{\theta_t}^i(\cdot | s), \pi_*^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s) \right\rangle + 1/2 \cdot \|\pi_{\theta_{t+1}}^i(\cdot | s) - \pi_{\theta_t}^i(\cdot | s)\|_1^2, \end{aligned}$$

where in the last inequality we use the Cauchy-Schwartz inequality and the Pinsker's inequality. Rearranging the terms in (G.4), we finish the proof of Lemma C.4. \square

H PROOF OF LEMMA B.2

Proof. We prove the lemma by induction. First, since $|\widehat{Q}_{(i)}^{\pi_{\theta_0}^i, \pi_{\theta_0}^{-i}}(s, a^i)| \leq Q_{(i)}^{\max}$, we have

$$\begin{aligned} |\widehat{\mathcal{E}}_{1,(i)}(s, a^i)| &= |(1 - \lambda_i \alpha_{0,(i)}) \cdot \mathcal{E}_{\theta_0,(i)}(s, a^i) + \alpha_{0,(i)} \cdot \widehat{Q}_{(i)}^{\pi_{\theta_0}^i, \pi_{\theta_0}^{-i}}(s, a^i)| \\ &= \alpha_{0,(i)} \cdot |\widehat{Q}_{(i)}^{\pi_{\theta_0}^i, \pi_{\theta_0}^{-i}}(s, a^i)| \leq Q_{(i)}^{\max} / (\lambda_i - M_i), \end{aligned}$$

where the last inequality follows from $\alpha_{0,(i)} = 1/(\lambda_i - M_i)$. This means that setting $|\mathcal{E}_{\theta_1,(i)}(s, a^i)| \leq Q_{(i)}^{\max} / (\lambda_i - M_i)$ covers the range of $\widehat{\mathcal{E}}_{1,(i)}(s, a^i)$.

Now suppose that $|\widehat{\mathcal{E}}_{\theta_t,(i)}(s, a^i)| \leq Q_{(i)}^{\max} / (\lambda_i - M_i)$. By $\lambda_i \geq 2M_i$, we have for all $t \geq 1$ that,

$$\lambda_i \alpha_{t,(i)} = \frac{\lambda_i}{(t+1) \cdot (\lambda_i - M_i)} \leq \frac{\lambda_i}{2(\lambda_i - M_i)} \leq 1.$$

Thus, we have that, for all $t \geq 1$ and $(s, a^i) \in \mathcal{S} \times \mathcal{A}^i$,

$$\begin{aligned} |\widehat{\mathcal{E}}_{t+1,(i)}(s, a^i)| &= |(1 - \lambda_i \alpha_{t,(i)}) \cdot \mathcal{E}_{\theta_{t,(i)}}(s, a^i) + \alpha_{t,(i)} \cdot \widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i)| \\ &\leq \left(1 - \frac{\lambda_i}{(t+1) \cdot (\lambda_i - M_i)}\right) \cdot |\mathcal{E}_{\theta_{t,(i)}}(s, a^i)| + \frac{1}{(t+1) \cdot (\lambda_i - M_i)} \cdot |\widehat{Q}_{(i)}^{\pi_{\theta_t}^i, \pi_{\theta_t}^{-i}}(s, a^i)| \\ &\leq \frac{t \cdot (\lambda_i - M_i) - M_i}{(t+1) \cdot (\lambda_i - M_i)} \cdot \frac{Q_{(i)}^{\max}}{\lambda_i - M_i} + \frac{Q_{(i)}^{\max}}{(t+1) \cdot (\lambda_i - M_i)} \\ &\leq \frac{t}{t+1} \cdot \frac{Q_{(i)}^{\max}}{\lambda_i - M_i} + \frac{Q_{(i)}^{\max}}{(t+1) \cdot (\lambda_i - M_i)} = \frac{Q_{(i)}^{\max}}{\lambda_i - M_i}, \end{aligned}$$

which means that setting $|\mathcal{E}_{\theta_{t+1,(i)}}(s, a^i)| \leq Q_{(i)}^{\max}/(\lambda_i - M_i)$ covers the range of $\widehat{\mathcal{E}}_{t+1,(i)}(s, a^i)$. By induction, we conclude the proof of Lemma B.2. \square

I IMBALANCED INFLUENCE

When the two players have imbalanced influence to the game (without loss of generality, we assume Player 2 has a dominant influence to the game), we let $\iota_1/\iota_2 = z$. In this case, we can replace (C.6) by

$$\begin{aligned} \text{(II)} &\leq \sqrt{2z} \cdot \iota_2 \cdot \left\{ \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}^{1/2}(s) \cdot \text{KL}_{t,(i)}^{1/2}(s)] \right. \\ &\quad \left. + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)]^{1/2} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(-i)}(s)]^{1/2} \right\} \\ &\leq \sqrt{2z} \cdot \left(1 + \frac{V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta}{1 - \gamma}\right) \cdot \left(\iota_2 \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(i)}(s)] + \iota_2 \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(-i)}(s)]\right). \end{aligned}$$

As a consequence, with M_i in (4.14) replaced by

$$M_i = \sqrt{2z} \cdot \left[2 + \sum_{i \in \{1,2\}} (V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta)/(1 - \gamma)\right] \cdot \iota_2, \quad i \in \{1,2\},$$

the convergence guarantee established in Theorem 4.5 remains valid. Taking $\lambda_1 = \lambda_2 = x \cdot \iota_2 \sqrt{2z}$ into $\lambda_i \geq 2M_i$, we obtain

$$x \geq 2 + \sum_{i \in \{1,2\}} (V_{(i)}^{\max} + Q_{(i)}^{\max} \cdot \zeta)/(1 - \gamma).$$

For the above inequality to hold, we have a sufficient requirement for the ratio z as

$$z \leq (1 - \gamma)^4 / \left[16(1 + \gamma)\iota_2 \cdot \log(|\mathcal{A}^1| \cdot |\mathcal{A}^2|)\right]^2.$$

J STRONGER ASSUMPTION: HANNAN CONSISTENCY

In the proof of Theorem 4.5 in Appendix C, we replace $\pi_{\theta_t}^2$ and π_*^2 with $\widetilde{\pi}_t^2$. Also, we replace π_*^1 with

$$\pi_{\dagger}^1 = \operatorname{argmax}_{\pi^1} \left\{ \frac{1}{T} \cdot \sum_{t=0}^{T-1} \mathcal{J}(\pi^1, \widetilde{\pi}_t^2) \right\}.$$

With stronger assumptions stated in 4.2, corresponding to (C.8), we have

$$\begin{aligned} &(1 - \gamma)\alpha_{t,(1)} \cdot \left[\mathcal{J}_{(1)}(\pi_{\dagger}^1, \widetilde{\pi}_t^2) - \mathcal{J}_{(1)}(\pi_{\theta_t}^1, \widetilde{\pi}_t^2)\right] \\ &\leq \mathbb{E}_{\nu^*} [\text{KL}_{t+1,(1)}^{\dagger}(s)] - \left\{1 - \left\{\lambda_1 - \left[1 + \frac{V_{(1)}^{\max} + Q_{(1)}^{\max} \cdot \zeta}{1 - \gamma}\right] \cdot \iota_1\right\} \cdot \alpha_{t,(1)}\right\} \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(1)}^{\dagger}(s)] \\ &\quad (\epsilon_t + \epsilon'_t) + [2 + 2\lambda_1^2/(\lambda_1 - M_1)^2] \cdot (Q_{(1)}^{\max})^2 \cdot \alpha_{t,(1)}^2, \end{aligned} \tag{J.1}$$

where $\text{KL}_{t+1,(1)}^{\dagger}(s) = \text{KL}(\pi_{\dagger}^1(\cdot | s) \| \pi_{\theta_t}^1(\cdot | s))$. Here we drop the KL-divergence term for Player 2 since π_*^2 and $\pi_{\theta_t}^2$ are replaced by the same policy $\widetilde{\pi}_t^2$. Also, the regularization parameter λ_2 is dropped since now Player 2 no longer uses such a regularization parameter to update its policies. Multiplying

both sides of (J.2) and setting the stepsize $\alpha_{t,(1)}$ as the stepsize choice in Theorem 4.5, we obtain

$$\begin{aligned} & \frac{1-\gamma}{(\lambda_1 - M_1)} \cdot [\mathcal{J}_{(1)}(\pi_{\dagger}^1, \tilde{\pi}_t^2) - \mathcal{J}_{(1)}(\pi_{\theta_t}^1, \tilde{\pi}_t^2)] \\ & \leq t \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t,(1)}^{\dagger}(s)] - (t+1) \cdot \mathbb{E}_{\nu^*} [\text{KL}_{t+1,(1)}^{\dagger}(s)] \\ & \quad + (t+1) \cdot (\epsilon_t + \epsilon'_t) + \frac{[2 + 2\lambda_1^2/(\lambda_1 - M_1)^2] \cdot (Q_{(1)}^{\max})^2}{(t+1) \cdot (\lambda_1 - M_1)}, \end{aligned} \quad (\text{J.2})$$

applying same argument in the proof of Theorem 4.5, we obtain

$$\begin{aligned} & \sup_{\pi^1} \left\{ \frac{1}{T} \cdot \sum_{t=0}^{T-1} [\mathcal{J}(\pi^1, \tilde{\pi}_t^2) - \mathcal{J}(\pi_{\theta_t}^1, \tilde{\pi}_t^2)] \right\} \\ & = \frac{1}{T} \cdot \sum_{t=0}^{T-1} [\mathcal{J}(\pi_{\dagger}^1, \tilde{\pi}_t^2) - \mathcal{J}(\pi_{\theta_t}^1, \tilde{\pi}_t^2)] \\ & \leq \frac{\sigma \cdot (\lambda_1 - M_1)}{(1-\gamma) \cdot T} + \frac{[2 + 2\lambda_1^2/(\lambda_1 - M_1)^2] \cdot (Q_{(1)}^{\max})^2}{(1-\gamma) \cdot (\lambda_1 - M_1)} \cdot \frac{\log T}{T} + \lambda_1 \cdot \log |\mathcal{A}^1|. \end{aligned}$$

Thus, we conclude the proof of (4.17). Using the same argument as above, we can also prove the same result for Player 2 when Player 1 does not update its policies according to Algorithm 1. Setting $\tilde{\pi}_t^1 = \pi_{\dagger}^1 = \pi_*^1$ and $\tilde{\pi}_t^2 = \pi_{\dagger}^2 = \pi_*^2$, and combining such result for both players, we can replace the left-hand side of the convergence guarantee established in Theorem 4.5 by (4.18). Thus, the proposed Algorithm 1 is Hannan consistent.