

A. Appendix

A.1. Trajectory Distribution Induced by Logistic Stochastic Best Response Equilibrium

Let $\{\pi_{-i}^t(a_{-i}^t|s^t)\}_{t=1}^T$ denote other agents' marginal LSBRE policies, and $\{\hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t)\}_{t=1}^T$ denote agent i 's conditional policy. With chain rule, the induced trajectory distribution is given by:

$$\hat{p}(\tau) = \left[\eta(s^1) \cdot \prod_{t=1}^T P(s^{t+1}|s^t, \mathbf{a}^t) \cdot \pi_{-i}^t(\mathbf{a}_{-i}^t|s^t) \right] \cdot \prod_{t=1}^T \hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t) \quad (14)$$

Suppose the desired distribution is given by:

$$p(\tau) \propto \left[\eta(s^1) \cdot \prod_{t=1}^T P(s^{t+1}|s^t, \mathbf{a}^t) \cdot \pi_{-i}^t(\mathbf{a}_{-i}^t|s^t) \right] \cdot \exp \left(\sum_{t=1}^T r_i(s^t, a_i^t, \mathbf{a}_{-i}^t) \right) \quad (15)$$

Now we will show that the optimal solution to the following optimization problem correspond to the LSBRE conditional policies:

$$\min_{\hat{\pi}_i^{1:T}} D_{\text{KL}}(\hat{p}(\tau)||p(\tau)) \quad (16)$$

The optimization problem in Equation (16) is equivalent to (the partition function of the desired distribution is a constant with respect to optimized policies):

$$\begin{aligned} \max_{\hat{\pi}_i^{1:T}} \mathbb{E}_{\tau \sim \hat{p}(\tau)} & \left[\log \eta(s^1) + \sum_{t=1}^T (\log P(s^{t+1}|s^t, \mathbf{a}^t) + \log \pi_{-i}^t(\mathbf{a}_{-i}^t|s^t) + r_i(s^t, \mathbf{a}^t)) - \right. \\ & \left. \log \eta(s^1) - \sum_{t=1}^T (\log P(s^{t+1}|s^t, \mathbf{a}^t) + \log \pi_{-i}^t(\mathbf{a}_{-i}^t|s^t) + \log \hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t)) \right] \\ &= \mathbb{E}_{\tau \sim \hat{p}(\tau)} \left[\sum_{t=1}^T r_i(s^t, \mathbf{a}^t) - \log \hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t) \right] = \sum_{t=1}^T \mathbb{E}_{(s^t, \mathbf{a}^t) \sim \hat{p}(s^t, \mathbf{a}^t)} [r_i(s^t, \mathbf{a}^t) - \log \hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t)] \end{aligned} \quad (17)$$

To maximize this objective, we can use a dynamic programming procedure. Let us first consider the base case of optimizing $\hat{\pi}_i^T(a_i^T|\mathbf{a}_{-i}^T, s^T)$:

$$\begin{aligned} \mathbb{E}_{(s^T, \mathbf{a}^T) \sim \hat{p}(s^T, \mathbf{a}^T)} [r_i(s^T, \mathbf{a}^T) - \log \hat{\pi}_i^T(a_i^T|\mathbf{a}_{-i}^T)] &= \\ \mathbb{E}_{s^T \sim \hat{p}(s^T), \mathbf{a}_{-i}^T \sim \pi_{-i}^T(\cdot|s^T)} \left[-D_{\text{KL}} \left(\hat{\pi}_i^T(a_i^T|\mathbf{a}_{-i}^T, s^T) \parallel \frac{\exp(r_i(s^T, a_i^T, \mathbf{a}_{-i}^T))}{\exp(V_i(s^T, \mathbf{a}_{-i}^T))} \right) + V_i(s^T, \mathbf{a}_{-i}^T) \right] \end{aligned} \quad (18)$$

where $\exp(V_i(s^T, \mathbf{a}_{-i}^T))$ is the partition function and $V_i(s^T, \mathbf{a}_{-i}^T) = \log \sum_{a'_i} \exp(r_i(s^T, a'_i, \mathbf{a}_{-i}^T))$. The optimal policy is given by:

$$\pi_i^T(a_i^T|\mathbf{a}_{-i}^T, s^T) = \exp(r_i(s^T, a_i^T, \mathbf{a}_{-i}^T) - V_i(s^T, \mathbf{a}_{-i}^T)) \quad (19)$$

With the optimal policy in Equation (19), Equation (18) is equivalent to (with the KL divergence being zero):

$$\mathbb{E}_{(s^T, \mathbf{a}^T) \sim \hat{p}(s^T, \mathbf{a}^T)} [r_i(s^T, \mathbf{a}^T) - \log \hat{\pi}_i^T(a_i^T|\mathbf{a}_{-i}^T)] = \mathbb{E}_{s^T \sim \hat{p}(s^T), \mathbf{a}_{-i}^T \sim \pi_{-i}^T(\cdot|s^T)} [V_i(s^T, \mathbf{a}_{-i}^T)] \quad (20)$$

Then recursively, for a given time step t , $\hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t)$ must maximize:

$$\mathbb{E}_{(s^t, \mathbf{a}^t) \sim \hat{p}(s^t, \mathbf{a}^t)} \left[r_i(s^t, \mathbf{a}^t) - \log \hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t) + \mathbb{E}_{s^{t+1} \sim P(\cdot|s^t, \mathbf{a}^t), \mathbf{a}_{-i}^{t+1} \sim \pi_{-i}^{t+1}(\cdot|s^{t+1})} [V_i^{\pi^{t+2:T}}(s^{t+1}, \mathbf{a}_{-i}^{t+1})] \right] = \quad (21)$$

$$\mathbb{E}_{s^t \sim \hat{p}(s^t), \mathbf{a}_{-i}^t \sim \pi_{-i}^t(\cdot|s^t)} \left[-D_{\text{KL}} \left(\hat{\pi}_i^t(a_i^t|\mathbf{a}_{-i}^t, s^t) \parallel \frac{\exp(Q_i^{\pi^{t+1:T}}(s^t, a_i^t, \mathbf{a}_{-i}^t))}{\exp(V_i^{\pi^{t+1:T}}(s^t, \mathbf{a}_{-i}^t))} \right) + V_i^{\pi^{t+1:T}}(s^t, \mathbf{a}_{-i}^t) \right] \quad (22)$$

where we define:

$$Q_i^{\pi^{t+1:T}}(s^t, \mathbf{a}^t) = r_i(s^t, \mathbf{a}^t) + \mathbb{E}_{s^{t+1} \sim p(\cdot | s^t, \mathbf{a}^t)} \left[\mathcal{H}(\pi_i^{t+1}(\cdot | s^{t+1})) + \mathbb{E}_{\mathbf{a}_{-i}^{t+1} \sim \pi_{-i}^{t+1}(\cdot | s^{t+1})} [V_i(s^{t+1}, \mathbf{a}_{-i}^{t+1})] \right] \quad (23)$$

$$V_i^{\pi^{t+1:T}}(s^t, \mathbf{a}_{-i}^t) = \log \sum_{a'_i} \exp(Q_i^{\pi^{t+1:T}}(s^t, a'_i, \mathbf{a}_{-i}^t)) \quad (24)$$

The optimal policy to Equation (22) is given by:

$$\pi_i^t(a_i^t | \mathbf{a}_{-i}^t, s^t) = \exp(Q_i^{\pi^{t+1:T}}(s^t, \mathbf{a}^t) - V_i^{\pi^{t+1:T}}(s^t, \mathbf{a}_{-i}^t)) \quad (25)$$

which is exactly the set of conditional distributions used to produce LSBRE (Definition 2).

A.2. Maximum Pseudolikelihood Estimation for LSBRE

Theorem 2 strictly follows the asymptotic consistency property of maximum pseudolikelihood estimation (Lehmann & Casella, 2006; Dawid & Musio, 2014). For simplicity, we will show the proof for normal form games and similar to Appendix A.1, the extension to Markov games can be proved by induction.

Consider a normal form game with N players and reward functions $\{r_i(\mathbf{a}; \omega_i)\}_{i=1}^N$. Suppose the expert demonstrations $\mathcal{D} = \{(a_1, \dots, a_N)^m\}_{m=1}^M$ are generated by $\pi(\mathbf{a}; \omega^*)$, where ω^* denotes the true value of the parameters. The pseudolikelihood objective we want to maximize is given by:

$$\ell_{\text{PL}}(\omega) = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^N \log \pi_i(a_i^m | \mathbf{a}_{-i}^m; \omega_i) = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^N \log \frac{\exp(r_i(a_i^m, \mathbf{a}_{-i}^m; \omega_i))}{\sum_{a'_i} \exp(r_i(a'_i, \mathbf{a}_{-i}^m; \omega_i))} \quad (26)$$

$$= \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^N r_i(a_i^m, \mathbf{a}_{-i}^m; \omega_i) - \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^N \log Z(\mathbf{a}_{-i}^m; \omega_i) \quad (27)$$

$$= \sum_{i=1}^N \sum_{\mathbf{a}} p_{\mathcal{D}}(\mathbf{a}) r_i(a_i, \mathbf{a}_{-i}; \omega_i) - \sum_{i=1}^N \sum_{\mathbf{a}_{-i}} p_{\mathcal{D}}(\mathbf{a}_{-i}) \log Z(\mathbf{a}_{-i}; \omega_i) \quad (28)$$

where $p_{\mathcal{D}}$ is the empirical data distribution and $Z(\mathbf{a}_{-i}; \omega_i)$ is the partition function.

Take derivatives of $\ell_{\text{PL}}(\omega)$:

$$\frac{\partial}{\partial \omega} \ell_{\text{PL}}(\omega) = \sum_{i=1}^N \sum_{\mathbf{a}} p_{\mathcal{D}}(\mathbf{a}) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) - \sum_{i=1}^N \sum_{\mathbf{a}_{-i}} p_{\mathcal{D}}(\mathbf{a}_{-i}) \frac{1}{Z(\mathbf{a}_{-i}; \omega_i)} \frac{\partial}{\partial \omega} Z(\mathbf{a}_{-i}; \omega_i) \quad (29)$$

$$= \sum_{i=1}^N \sum_{\mathbf{a}} p_{\mathcal{D}}(\mathbf{a}) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) - \sum_{i=1}^N \sum_{\mathbf{a}_{-i}} p_{\mathcal{D}}(\mathbf{a}_{-i}) \sum_{a_i} \frac{\exp(r_i(a_i, \mathbf{a}_{-i}; \omega_i))}{Z(\mathbf{a}_{-i}; \omega_i)} \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) \quad (30)$$

$$= \sum_{i=1}^N \sum_{\mathbf{a}} p_{\mathcal{D}}(\mathbf{a}) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) - \sum_{i=1}^N \sum_{\mathbf{a}_{-i}} p_{\mathcal{D}}(\mathbf{a}_{-i}) \sum_{a_i} \pi_i(a_i | \mathbf{a}_{-i}; \omega_i) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) \quad (31)$$

When the sample size $m \rightarrow \infty$, Equation (31) is equivalent to:

$$\frac{\partial}{\partial \omega} \ell_{\text{PL}}(\omega) = \sum_{i=1}^N \sum_{\mathbf{a}} p(\mathbf{a}; \omega^*) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) - \sum_{i=1}^N \sum_{\mathbf{a}_{-i}} p(\mathbf{a}_{-i}; \omega^*) \sum_{a_i} \pi_i(a_i | \mathbf{a}_{-i}; \omega_i) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) \quad (32)$$

$$= \sum_{i=1}^N \sum_{\mathbf{a}_{-i}} p(\mathbf{a}_{-i}; \omega^*) \sum_{a_i} (p(a_i | \mathbf{a}_{-i}; \omega^*) - \pi_i(a_i | \mathbf{a}_{-i}; \omega_i)) \frac{\partial}{\partial \omega} r_i(a_i, \mathbf{a}_{-i}; \omega_i) \quad (33)$$

When $\omega = \omega^*$, the gradients in Equation (33) will be zero.