

• Supplementary File •

A Formal Model for Multiagent Q -Learning on Graphs

Jinzhao LIU^{1,2}, Guangchen JIANG¹, Chen CHU^{2,3,4}, Yong LI⁵, Zhen WANG^{2,6*} & Shuyue HU^{7*}

¹Department of School of Software, Yunnan University, Kunming 650091, China;

²Department of School of Artificial Intelligence, OPTics and ElectroNics (iOPEN),
Northwestern Polytechnical University, Xi'an 710072, China;

³Department of Statistics, School of Statistics and Mathematics,
Yunnan University of Finance and Economics, Kunming 650091, China;

⁴Yunnan Key Laboratory of Service Computing,
Yunnan University of Finance and Economics, Kunming 650091, China;

⁵Department of School of Information Science & Engineering, Yunnan University, Kunming 650091, China;

⁶Department of School of Cybersecurity, Northwestern Polytechnical University, Xi'an 710072, China;

⁷Shanghai Artificial Intelligence Laboratory, Shanghai 200232, China

Appendix A Proof of Lemma 1

Proof. First, we analyze the reward of an arbitrary agent i . For the agent's reward $r_t(a_j | \mathbf{a}^{-i})$, we can expand

$$r_t^i(a_j | \mathbf{a}^{-i}) = \frac{1}{k^i} \sum_{\forall z \in \mathcal{N}_i, a_y \in \mathcal{A}} r_t^i(a_j | a_y^z),$$

and obtain:

$$r_t(a_j | \mathbf{a}^{-i}) = \frac{1}{k^i} [k_{a_1,t}^i r_t(a_j | a_1) + \dots + k_{a_m,t}^i r_t(a_j | a_m)]. \quad (\text{A1})$$

According to

$$r(a_j | a_y) = \mathbf{a}_j^\top \mathbf{R} \mathbf{a}_y,$$

we can transform reward in Equation (A1) into a matrix multiplication form:

$$\begin{aligned} r_t(a_j | \mathbf{a}^{-i}) &= \frac{1}{k^i} [k_{a_1,t}^i \mathbf{a}_j^\top \mathbf{R} \mathbf{a}_1 \dots + k_{a_m,t}^i \mathbf{a}_j^\top \mathbf{R} \mathbf{a}_m] \\ &= \frac{1}{k^i} \mathbf{a}_j^\top \mathbf{R} (k_{a_1,t}^i \mathbf{a}_1 + \dots + k_{a_m,t}^i \mathbf{a}_m) \\ &= \frac{1}{k^i} \mathbf{a}_j^\top \mathbf{R} [k_{a_1,t}^i, \dots, k_{a_m,t}^i]^\top \\ &= \frac{1}{k^i} \mathbf{a}_j^\top \mathbf{R} \gamma_t^i \\ &= r_t(a_j | \gamma_t^i). \end{aligned} \quad (\text{A2})$$

■

Appendix B Proof of Lemma 2

Proof. For agent i , the action choices of their neighbors are generally different. With mean-field approximation [1, 2], we approximate the strategy of each neighbor by the same average strategy $\bar{\mathbf{x}}_t$ of the population. However, different from the previous approaches [1–3], we additionally characterize the effects of graph structure by explicitly considering the neighborhood size. Specifically, we observe that the decision-making process of each neighbor can be seen as an independent multinoulli trial:

- the m available actions amounts to the m events,
- each trial leads to success for exactly one event (a particular action choice),
- the success probability of an event is the probability of the corresponding action choice, which can be approximated by the average probability of choosing that action in the population.

* Corresponding author (email: w-zhen@nwpu.edu.cn, hushuyue@pjlab.org.cn)

Thus, $\lambda(\gamma_t^i)$ can be understood as the probability distribution of a multinomial distribution with the number of independent trials being k^i (the number of neighbors), the number of events being m (the number of available actions), and the success probabilities of events being \bar{x}_t (the average strategy).

$$\lambda(\gamma_t^i) = \frac{k^i!}{\prod_{l=1}^m k_{a_l}^i!} \prod_{j=1}^m \bar{x}_t(a_j)^{k_{a_j}^i}$$

is the probability mass function of the above multinomial distribution. ■

Appendix C Proof of Lemma 3

Proof. The expected change in the Q -value on the graph caused by the action a_j of agent i with degree k^i is the product of the probability that these agents have neighbor configuration γ_t^i , the probability of adopting action, and the change in the Q -value vector \mathbf{q}_t in the j -th dimension, $\Delta q_{t,j}^i(\gamma_t^i, \mathbf{q}_t^i)$:

$$\begin{aligned} f_j(\mathbf{q}_t^i) &= x_t^i(a_j) \sum_{\gamma_t^i \in \Gamma^i} \lambda(\gamma_t^i) \Delta q_{t,j}^i(\gamma_t^i, \mathbf{q}_t^i) \\ &= x_t^i(a_j) \sum_{\gamma_t^i \in \Gamma^i} \lambda(\gamma_t^i) [q_{t+1}^i(a_j) - q_t^i(a_j)] \\ &= x_t^i(a_j) \sum_{\gamma_t^i \in \Gamma^i} \lambda(\gamma_t^i) \alpha [r_t^i(a_j | \gamma_t^i) - q_t^i(a_j)] \\ &= \alpha x_t^i(a_j) \left[\sum_{\gamma_t^i \in \Gamma^i} [\lambda(\gamma_t^i) r_t^i(a_j | \gamma_t^i)] - q_t^i(a_j) \right]. \end{aligned}$$

■

Appendix D Proof of Lemma 4

Proof. The Lemma 3 corresponds to the case of an agent with a known degree k^i . However, in a population where agents can have different degrees, we must consider the proportion of these agents, i.e. $\rho(k)$. In this case, we can express the average reward of the entire population in round t as follows:

$$\sum_{k=1}^{\bar{k}} \sum_{\gamma \in \Gamma(k)} [\rho(k) \lambda(\gamma) r_t(a_j | \gamma_t)]$$

Therefore,

$$f_j^i(\mathbf{q}_t^i) = \alpha x_t^i(a_j) \left[\sum_{\gamma_t^i \in \Gamma^i} [\lambda(\gamma_t^i) r_t(a_j | \gamma_t^i)] - q_t^i(a_j) \right].$$

can be expanded to:

$$v_j(\mathbf{q}_t) = \alpha x_t(a_j) \left[\sum_{k=1}^{\bar{k}} \sum_{\gamma \in \Gamma(k)} [\rho(k) \lambda(\gamma) r_t(a_j | \gamma_t)] - q_t(a_j) \right]$$

■

Appendix E Proof of Corollary 1

Proof. Intuitively, the probability of a neighbor configuration $\lambda(\gamma_t)$ occurring follows a multinomial distribution Multinom(k, \bar{x}_t) with

$$\bar{\mathbf{x}}_t := [\bar{x}_t(a_1), \dots, \bar{x}_t(a_m)]^\top, \bar{\mathbf{x}}_t \in \mathbb{R}^m, \bar{x}_t \geq 0, \sum_{j=1}^m \bar{x}_t(a_j) = 1,$$

where $\bar{x}_t(a_j)$ represents the average probability of all agents in the population adopting action a_j at time step t .

Since the limit of a multinomial distribution with m random variables can be approximated by a $(m-1)$ -dimensional normal distribution, we can restrict $\gamma \in \mathbb{R}^m$ to its first $m-1$ rows, denote it as $\hat{\gamma}_t \in \mathbb{R}^{m-1}$. Then, the agents' average reward by adopting action a_j , denoted as $\bar{r}(a_j)$, can be expressed as follows in this population at time step t .

$$\sum_{\gamma_t \in \Gamma(k)} \int \cdots \int \varphi(\hat{\gamma}_t) r_t(a_j | \gamma_t) dk_{a_1} \cdots dk_{a_{m-1}}, \quad (\text{E1})$$

where $\varphi(\hat{\gamma}_t)$ is density function

$$\varphi(\hat{\gamma}_t) = \frac{\exp\left(-1/2(\hat{\gamma}_t - \hat{\mathbf{x}}_t)^\top \hat{\Sigma}_t^{-1}(\hat{\gamma}_t - \hat{\mathbf{x}}_t)\right)}{\sqrt{(2\pi)^{m-1} |\hat{\Sigma}_t|}}$$

and $\hat{\gamma}_t \sim \mathcal{N}_{m-1}(\hat{\mathbf{x}}_t, \hat{\Sigma}_t)$ with $(m-1)$ -dimensional mean vector

$$\hat{\mathbf{x}}_t := [\bar{x}_t(a_1), \dots, \bar{x}_t(a_{m-1})]^\top, \hat{\mathbf{x}}_t \in \mathbb{R}^{m-1},$$

and $(m-1) \times (m-1)$ covariance matrix

$$\begin{aligned}\hat{\Sigma}_{t,i,j} &= \mathbb{E}[(\hat{\gamma}_{t,i} - \bar{x}_{t,i}(a_i))(\hat{\gamma}_{t,j} - \bar{x}_{t,j}(a_j))] \\ &= \text{Cov}[\hat{\gamma}_{t,i}, \hat{\gamma}_{t,j}]\end{aligned}$$

such that $1 \leq i \leq (m-1)$ and $1 \leq j \leq (m-1)$.

$v_j(\mathbf{q}_t)$ can be viewed as a function of the expected reward, i.e.

$$v_j(\mathbf{q}_t) = g(\mathbb{E}[r_t(a_j)]),$$

where

$$\mathbb{E}[r_t(a_j)] = \sum_{k=1}^{\bar{k}} \sum_{\gamma \in \Gamma(k)} [\rho(k) \lambda(\gamma) r_t(a_j | \gamma)]. \quad (\text{E2})$$

For a fixed k ,

$$\mathbb{E}[r_t(a_j)] = \sum_{\gamma \in \Gamma(k)} [\lambda(\gamma) r_t(a_j | \gamma)].$$

For the case of the mean-field method [1,2], we can assume that the actions of agents' neighbors are the average strategy of the entire population. In other words, we consider the average payoff across the population for an agent's payoff, i.e.

$$g(\bar{r}_t(a_j)),$$

where

$$\bar{r}_t(a_j) = \mathbf{a}_j^T \mathbf{R} \bar{\mathbf{x}}_t, \quad (\text{E3})$$

which is the agents' average reward by adopting action a_j .

Please note that Equation (E2) and (E3) are different and should not be considered identical.

Then, we take the Taylor series expansion of $v_j(\mathbf{q}_t) = g(\mathbb{E}[r_t(a_j)])$ at $\bar{r}_t(a_j)$, and obtain

$$\begin{aligned}g(\mathbb{E}[r_t(a_j)]) &= g(\bar{r}_t(a_j)) + \frac{g'(\bar{r}_t(a_j))}{1!} (\mathbb{E}[r_t(a_j)] - \bar{r}_t(a_j)) \\ &\quad + \frac{g''(\bar{r}_t(a_j))}{2!} (\mathbb{E}[r_t(a_j)] - \bar{r}_t(a_j))^2 + \dots \\ &= g(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t) + \frac{g'(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t)}{1!} \left(\sum_{\gamma \in \Gamma(k)} \left[\lambda(\gamma) \frac{1}{k} \mathbf{a}_j \mathbf{R} \gamma \right] - \mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t \right) \\ &\quad + \frac{g''(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t)}{2!} \left(\sum_{\gamma \in \Gamma(k)} \left[\lambda(\gamma) \frac{1}{k} \mathbf{a}_j \mathbf{R} \gamma \right] - \mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t \right)^2 \\ &\quad + \dots \\ &= g(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t) + \sum_{\gamma_t \in \Gamma(k)} \lambda(\gamma_t) \sum_{h=1}^{\infty} \frac{\Delta_{\bar{\gamma}_t}^h}{h!} g^{(h)}(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t)\end{aligned} \quad (\text{E4})$$

where $\Delta_{\bar{\gamma}_t} = \bar{\gamma}_t/k - \bar{\mathbf{x}}_t$, which is the difference between the actual neighborhood configuration and the average policy. Note that $\Delta_{\bar{\gamma}_t}^h$ is the h -th central moment of the normal distribution; therefore, $\Delta_{\bar{\gamma}_t}^h$ is 0 if and only if h is odd. For this reason, we focus only on even h , meaning that Equation (E4) can be reduced to

$$\underbrace{g(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t)}_{\text{term1}} + \underbrace{\sum_{\gamma_t \in \Gamma(k)} \lambda(\gamma_t) \sum_{l=1}^{\infty} \frac{\Delta_{\bar{\gamma}_t}^{2l}}{2l!} g^{(2l)}(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t)}_{\text{term2}} \quad (\text{E5})$$

It is obvious that the bias of the mean-field method [1,2] is mainly caused by the central moments $\Delta_{\bar{\gamma}_t}^{2l}$ of the distribution in term2 of Equation (E5).

In the scenario we are considering, where the actions are independent of each other, we can obtain

$$\Delta_{\bar{\gamma}_t}^{2l} = \frac{\sum_{j=1}^m [\bar{x}_t(a_j)(1 - \bar{x}_t(a_j))]^l}{k^l} (2l-1)!!$$

[4] and

$$\underbrace{g(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t)}_{\text{term1}} + \underbrace{\sum_{\gamma_t \in \Gamma(k)} \lambda(\gamma_t) \sum_{l=1}^{\infty} \left[\frac{(2l-1)!!}{(2l)!} g^{(2l)}(\mathbf{a}_j \mathbf{R} \bar{\mathbf{x}}_t) \frac{\sum_{j=1}^m [\bar{x}_t(a_j)(1 - \bar{x}_t(a_j))]^l}{k^l} \right]}_{\text{term2}}. \quad (\text{E6})$$

It is evident that the difference $\mathbb{E}[r_t(a_j)]$ and $\bar{r}_t(a_j)$ is only related to $\sum_{j=1}^m [\bar{x}_t(a_j)(1 - \bar{x}_t(a_j))]^l / k^l$ in term2 of Equation (E6). For Equation (E6), term2 $\rightarrow 0$ with $k \rightarrow 0$. Hence, we can obtain

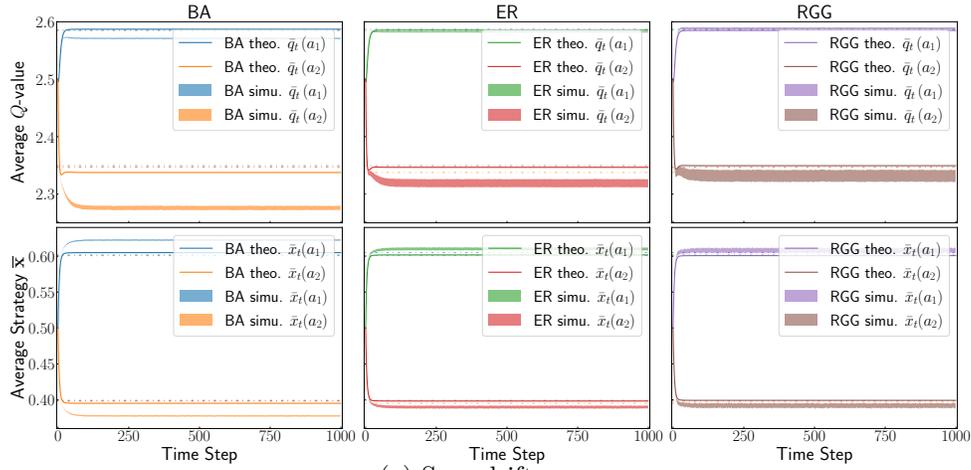
$$\lim_{k \rightarrow \infty} \frac{\partial p(\mathbf{q}_t, t, k)}{\partial t} = \frac{\partial \hat{p}(\mathbf{q}_t, t)}{\partial t},$$

and

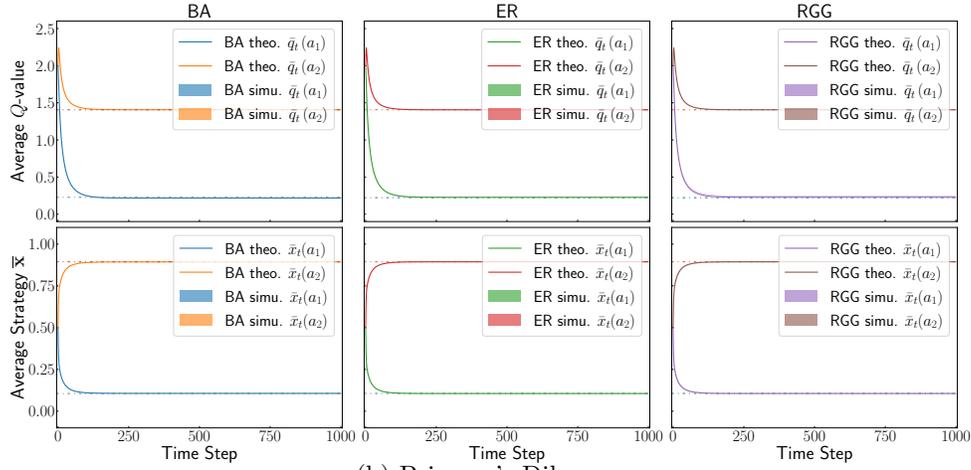
$$p(\mathbf{q}_t, t, k) \xrightarrow[k \rightarrow \infty]{\mathcal{D}} \hat{p}(\mathbf{q}_t, t)$$

■

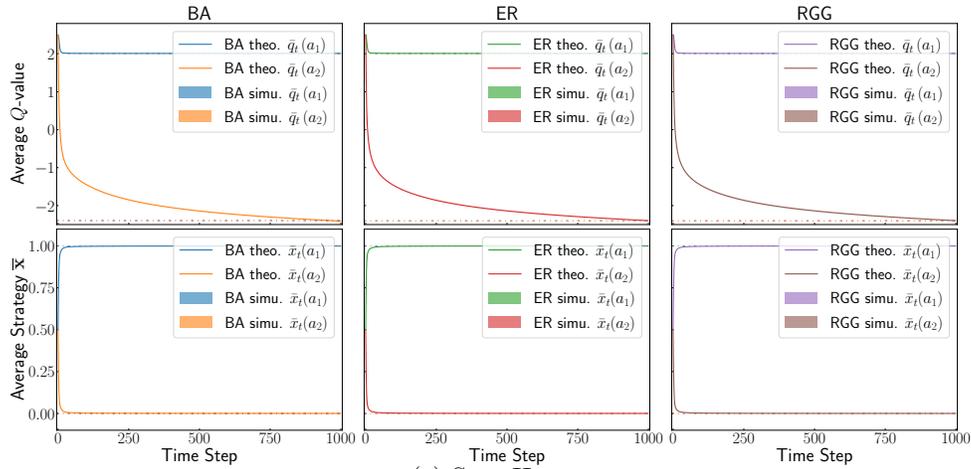
Appendix F Extended Validation of Figure 2. (a): Results from 100 to 1000 Steps



(a) Snowdrift game



(b) Prisoner's Dilemma



(c) Stag Hunt

Figure F1 Extended experimental results from the Snowdrift (SD), Prisoner's Dilemma (PD), and Stag Hunt (SH) games for populations in three graph structures — Barabási–Albert (BA), Erdős–Rényi (ER), and Random Geometric Graph (RGG) - extend from 100-steps to 1,000-steps. The results clearly demonstrate that the discrepancies between theoretical predictions and experimental outcomes stabilize and do not increase beyond a certain point. Specifically, in SD, which shows the largest difference among the games tested, the discrepancies remain relatively small even in the BA network, with errors ≈ 0.016 for $\bar{q}(a_1)$ and ≈ 0.061 for $\bar{q}(a_2)$. Similarly, in the ER network and RGG, these errors are about ≈ 0.0033 and ≈ 0.02 , respectively.

References

- 1 Hu S, Leung C, Leung H. Modelling the dynamics of multiagent q-learning in repeated symmetric games: a mean field theoretic approach. In: *Advances in Neural Information Processing Systems*, 2019, 32
- 2 Hu S, Leung C W, Leung H, et al. The Dynamics of Q-learning in Population Games: a Physics-Inspired Continuity Equation Model. In: *21st International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2022, Auckland, New Zealand, 2022*: 615-623
- 3 Leung C W, Hu S, Leung H F. Formal Modeling of Reinforcement Learning with Many Agents through Repeated Local Interactions. In: *2021 IEEE 33rd International Conference on Tools with Artificial Intelligence (ICTAI), IEEE, 2021*: 714-718
- 4 Shiryaev A N. *Probability-1*. Springer, 2016