

Modelling the dynamics of multi-agent Q-learning: the stochastic effects of local interaction and incomplete information

Chin-wing Leung¹, Shuyue Hu², Ho-fung Leung¹

¹The Chinese University of Hong Kong

² Shanghai Artificial Intelligence Laboratory

cwleung@cse.cuhk.edu.hk, shuyuehu217@gmail.com, lhf@cuhk.edu.hk

A Derivation for equation (10)-(12)

Starting with

$$\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t \equiv \alpha \left(\frac{1}{m} \sum_{j=1}^d \zeta_j^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j - Q_k^t \right) z_k^t \quad (9)$$

The unconditional first and second moments of $\partial_t \mathbf{Q}^t$ are evaluated as

$$\begin{aligned} \mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t] &= \alpha z_k^t \left(\frac{1}{m} \sum_{j=1}^d \mathbb{E}[\zeta_j^t] \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j - Q_k^t \right) \\ &= \alpha z_k^t \left(\frac{1}{m} \sum_{j=1}^d m y_j^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j - Q_k^t \right) \\ &= \alpha z_k^t (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t - Q_k^t) \end{aligned}$$

$$\mu_k^t = \mathbb{E}[\partial_t Q_k^t] = \alpha x_k^t (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t - Q_k^t) \quad (10)$$

$$Var(\partial_t Q_i^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t) = 0$$

$$\begin{aligned} Var(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t] \mid \mathbf{z}^t) &= \alpha^2 (z_k^t)^2 \frac{1}{m^2} Var\left(\sum_{j=1}^d \zeta_j^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j\right) \\ &= \alpha^2 (z_k^t)^2 \frac{1}{m^2} \left[\sum_{j=1}^d Var(\zeta_j^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j) + 2 \sum_{j_1 \neq j_2} Cov(\zeta_{j_1}^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_{j_1}, \zeta_{j_2}^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_{j_2}) \right] \\ &= \alpha^2 (z_k^t)^2 \frac{1}{m^2} \left[\sum_{j=1}^d (\mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j)^2 m y_j^t (1 - y_j^t) - 2 \sum_{j_1 \neq j_2} m y_{j_1}^t y_{j_2}^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_{j_1} \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_{j_2} \right] \\ &= \alpha^2 (z_k^t)^2 \frac{1}{m} \left[\sum_{j=1}^d y_j^t \mathbf{e}_k^\top \mathbf{U} \circ \mathbf{U} \mathbf{e}_j - \sum_{j=1}^d (y_j^t)^2 (\mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j)^2 \right. \\ &\quad \left. - 2 \sum_{j_1 \neq j_2} m y_{j_1}^t y_{j_2}^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_{j_1} \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_{j_2} \right] \\ &= \frac{1}{m} \alpha^2 (z_k^t)^2 [\mathbf{e}_k^\top \mathbf{U} \circ \mathbf{U} \mathbf{y}^t - (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t)^2] \end{aligned}$$

$$E[Var(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t] \mid \mathbf{z}^t)] = \frac{1}{m} \alpha^2 x_k^t [\mathbf{e}_k^\top \mathbf{U} \circ \mathbf{U} \mathbf{y}^t - (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t)^2]$$

$$Var(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t]) = \alpha^2 (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t - Q_k^t)^2 x_k^t (1 - x_k^t)$$

$$\begin{aligned} \sigma_{kk}^t &= Var(\partial_t Q_k^t) \\ &= \mathbb{E}[Var(\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t)] + E[Var(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t] \mid \mathbf{z}^t)] + Var(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t]) \\ &= \alpha^2 (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t - Q_k^t)^2 x_k^t (1 - x_k^t) + \frac{1}{m} \alpha^2 x_k^t [\mathbf{e}_k^\top \mathbf{U} \circ \mathbf{U} \mathbf{y}^t - (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t)^2] \end{aligned} \tag{11}$$

$$Cov(\partial_t Q_k^t, \partial_t Q_l^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t) = 0$$

$$Cov(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t], \mathbb{E}[\partial_t Q_l^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t] \mid \mathbf{z}^t) = -\alpha^2 z_k^t z_l^t \frac{1}{m^2} Cov(\sum_{j=1}^d \zeta_j^t \mathbf{e}_k^\top \mathbf{U} \mathbf{e}_j, \sum_{j=1}^d \zeta_j^t \mathbf{e}_l^\top \mathbf{U} \mathbf{e}_j)$$

$$E[Cov(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t], \mathbb{E}[\partial_t Q_l^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t] \mid \mathbf{z}^t)] = 0$$

$$Cov(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t], \mathbb{E}[\partial_t Q_l^t \mid \mathbf{z}^t]) = -\alpha^2 (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t - Q_k^t)(\mathbf{e}_l^\top \mathbf{U} \mathbf{y}^t - Q_l^t) x_k^t x_l^t$$

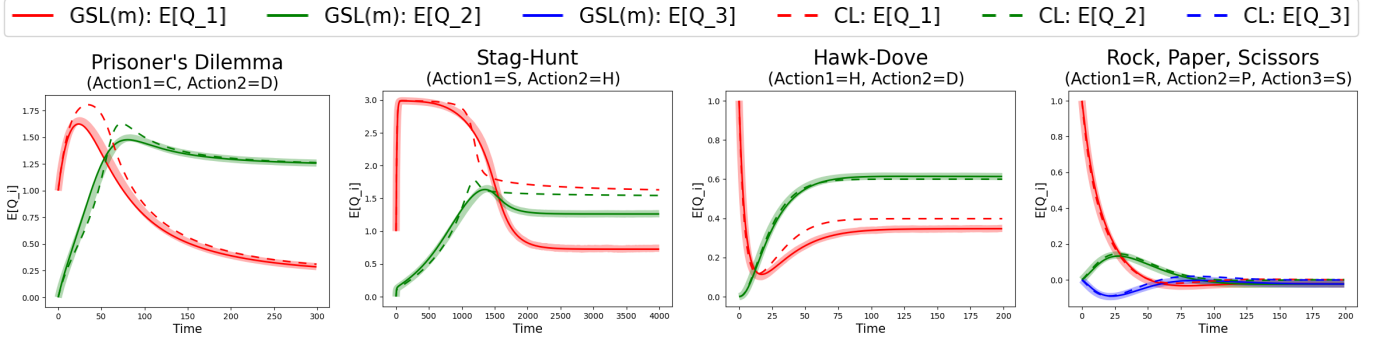
$$\begin{aligned} \sigma_{kl}^t &= Cov(\partial_t Q_k^t, \partial_t Q_l^t) \\ &= \mathbb{E}[Cov(\partial_t Q_k^t, \partial_t Q_l^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t)] + E[Cov(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t], \mathbb{E}[\partial_t Q_l^t \mid \mathbf{z}^t, \boldsymbol{\zeta}^t] \mid \mathbf{z}^t)] \\ &\quad + Cov(\mathbb{E}[\partial_t Q_k^t \mid \mathbf{z}^t], \mathbb{E}[\partial_t Q_l^t \mid \mathbf{z}^t]) \\ &= -\alpha^2 (\mathbf{e}_k^\top \mathbf{U} \mathbf{y}^t - Q_k^t)(\mathbf{e}_l^\top \mathbf{U} \mathbf{y}^t - Q_l^t) x_k^t x_l^t \end{aligned} \tag{12}$$

where $\mathbf{U} \circ \mathbf{U}$ represents the element-wise multiplication of matrix \mathbf{U} and \mathbf{U} .

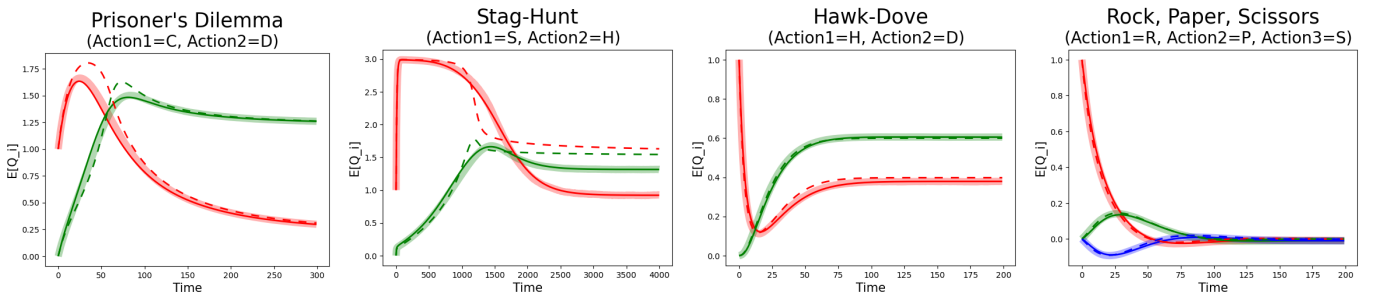
B Additional experiments results

We provide additional experiments results in the following. Figure 1 and 2 contrast the Q-learning dynamics predicted by our model (GSL) with the prediction made by the previous model (CL), the experimental settings are described in the section 4.1. We take the agent-based simulation results as the benchmark. Figure 1 presents the comparisons under the situations with a *small* value of m .¹ Figure 2 presents the comparisons under the situations where $m \rightarrow \infty$. It is clear that our model always provides more accurate descriptions on the dynamics of the expected Q-values $E[Q_k]$, $\forall k$ in a population. Figure 3 and 4 shows the dynamics of Q-values and population policy for different m , it is clear that the value of m plays an important role on the *outcome* of multiagent Q-learning in the GSL protocol.

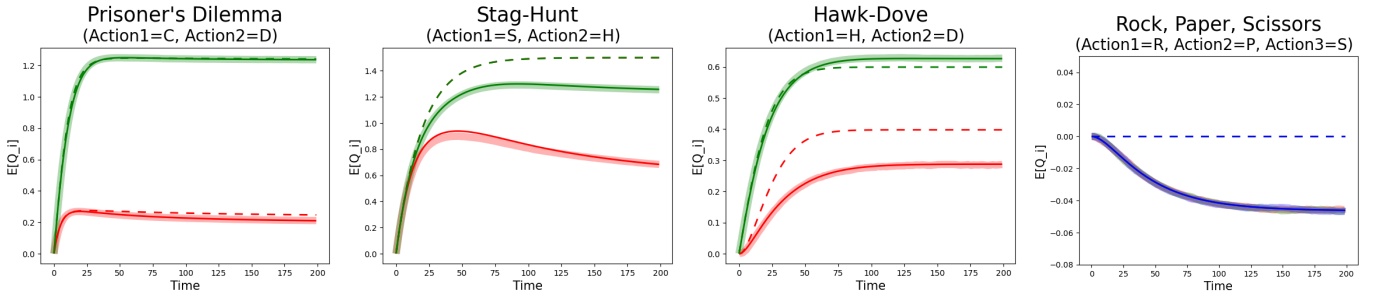
¹When m is large, the dynamics are close to the case of $m \rightarrow \infty$.



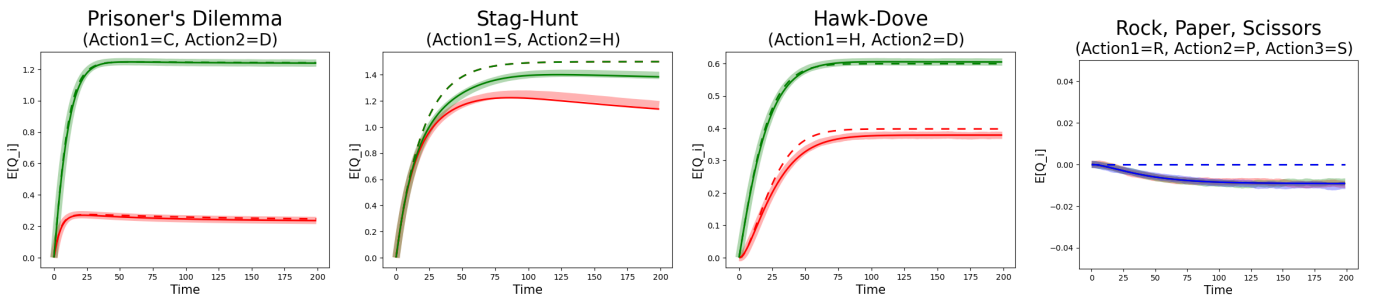
(a) $m = 2$, the initial Q-value of action 1 is 1 and the initial Q-values of other actions are 0 for all the agents



(b) $m = 5$, the initial Q-value of action 1 is 1 and the initial Q-values of other actions are 0 for all the agents

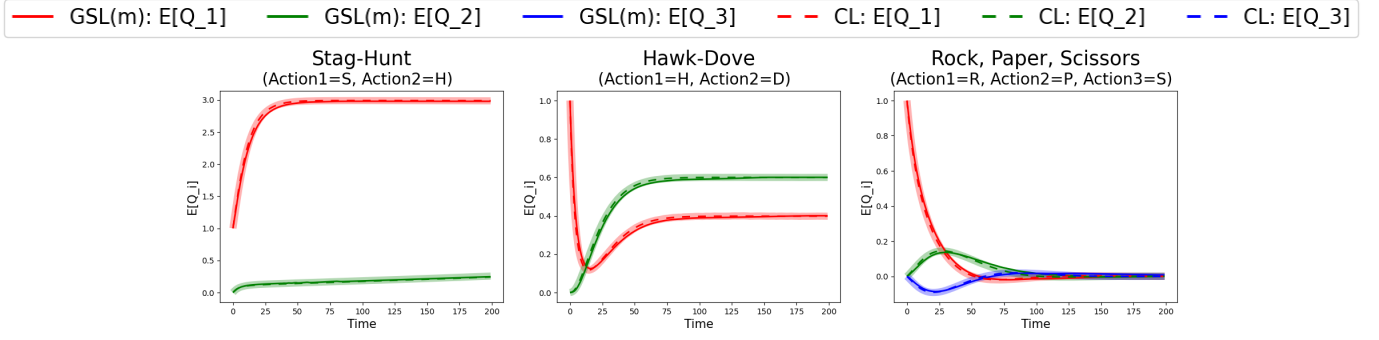


(c) $m = 1$, the initial Q-value of each action is 0 for all the agents

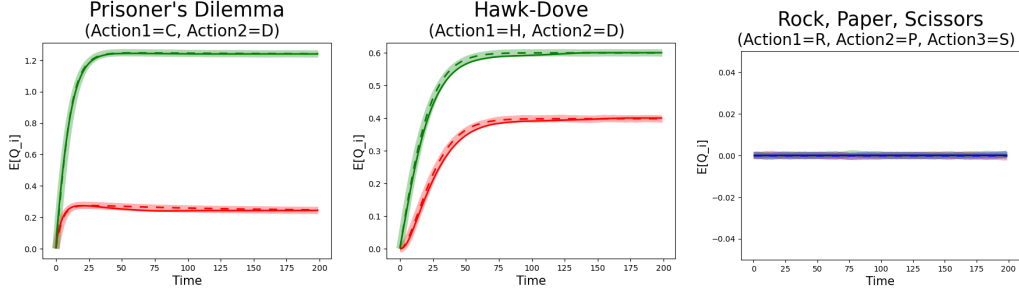


(d) $m = 5$, the initial Q-value of each action is 0 for all the agents

Figure 1: With a small value of m , comparison among the dynamics of average Q-values predicted by our model (solid line) and the previous model (dashed line), and the actual dynamics averaged over 100 runs of agent-based simulations (shaded line). In all these settings, our model better captures the qualitative and quantitative dynamics of the populations.

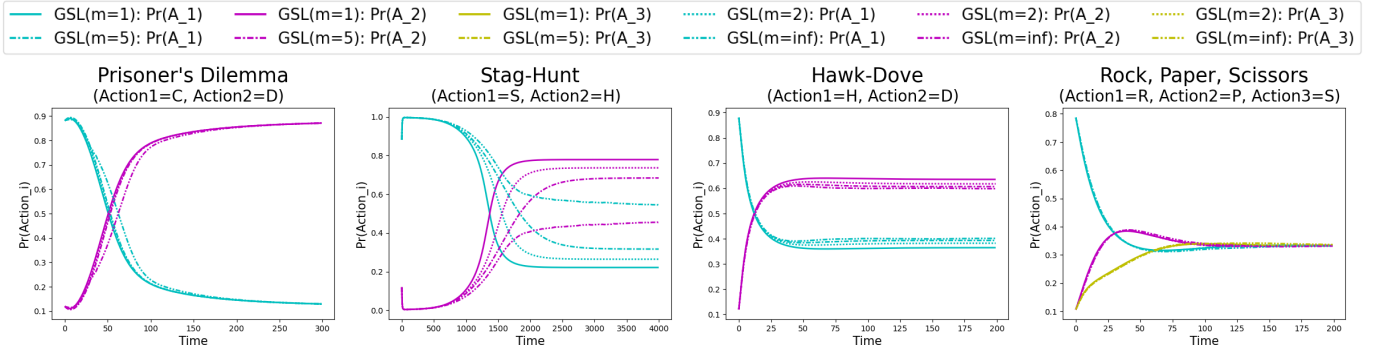


(a) $m \rightarrow \infty$, the initial Q-value of action 1 is 1 and the initial Q-values of other actions are 0 for all the agents

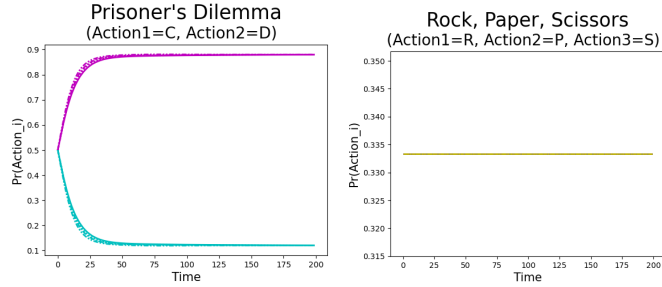


(b) $m \rightarrow \infty$, the initial Q-value of each action is 0 for all the agents

Figure 2: With a large value of m , comparison among the dynamics of average Q-values predicted by our model (solid line) and the previous model(dashed line), and the actual dynamics averaged over 100 runs of simulations (shaded line).

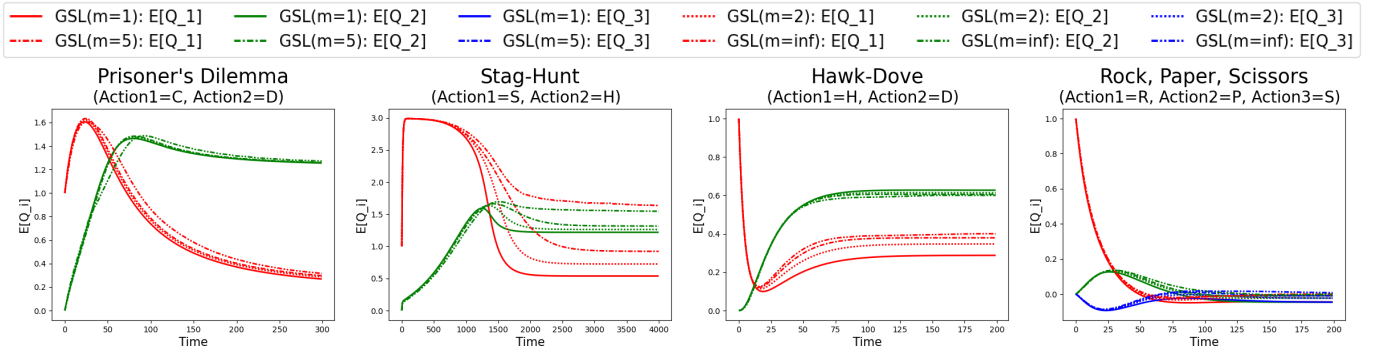


(a) the initial Q-value of action 1 is 1 and the initial Q-values of other actions are 0 for all the agents

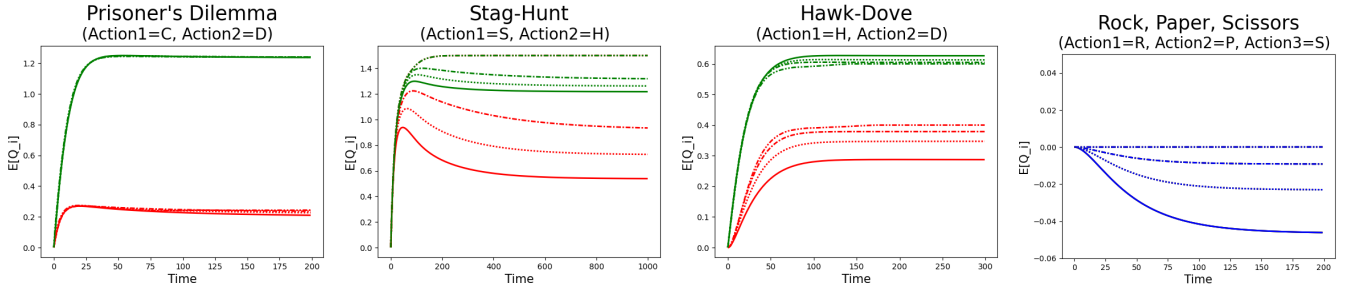


(b) the initial Q-value of each action is 0 for all the agents

Figure 3: The effects of local interactions and incomplete information on multiagent Q-learning. Our model shows that as the value of m varies, the population can stabilize at significantly different proportions of agents using each action.



(a) the initial Q-value of action 1 is 1 and the initial Q-values of other actions are 0 for all the agents



(b) the initial Q-value of each action is 0 for all the agents

Figure 4: The effects of local interactions and incomplete information on multiagent Q-learning. Our model shows that as the value of m varies, the population will establish different Q-values dynamics.