
Reinforcement learning for mean-field teams

Jayakumar Subramanian

Department of Electrical and Computer Engineering
McGill University
Montreal, QC H3A0G4
jayakumar.subramanian@mail.mcgill.ca

Raihan Seraj

Department of Electrical and Computer Engineering
McGill University
Montreal, QC H3A0G4
raihan.seraj@mail.mcgill.ca

Aditya Mahajan

Department of Electrical and Computer Engineering
McGill University
Montreal, QC H3A0G4
aditya.mahajan@mcgill.ca

Abstract

We develop reinforcement learning (RL) algorithms for a class of multi-agent systems called mean-field teams (MFT). Teams are multi-agent systems where agents have a common goal and receive a common reward at each time step. The team objective is to maximize the expected cumulative discounted reward over an infinite horizon. MFTs are teams with homogeneous, anonymous agents such that the agents are coupled in their dynamics and rewards through the mean-field (i.e., empirical distribution of the agents' state). In our work, we consider MFTs with a mean-field sharing information structure, i.e., each agent knows its local state and the empirical mean-field at each time step. We obtain a dynamic programming (DP) decomposition for MFTs using a decomposition approach from literature called the common information approach, which splits the decision making process into two parts. The first part is a centralized coordination rule that yields the second part, which are prescriptions to be followed by each agent based on their local information. We develop an RL approach for MFTs under the assumption of parametrized prescriptions. We consider the parameters as actions and use conventional RL algorithms to solve the DP. We illustrate the use of these algorithms through two examples based on stylized models of the demand response problem in smart grids and malware spread in networks.

Keywords: Reinforcement learning, mean-field teams, multi-agent reinforcement learning.

1 Introduction

In this paper, we look at reinforcement learning in cooperative multi-agent systems. Several algorithms for multi-agent reinforcement learning have been proposed in the literature [2–4, 9–12, 18–21]. These algorithms perform well on certain benchmark domains but there is little theoretical analysis on whether these algorithms converge to a team optimal solution. In this paper, we present a different view on multi-agent reinforcement learning. Our central thesis is that multi-agent systems for which the team optimal planning solution can be obtained by dynamic programming [13–15], it should be straightforward to translate these dynamic programs to reinforcement learning algorithms.

2 Model

Consider a multi-agent team with n agents, indexed by the set $N = \{1, \dots, n\}$. The team operates in discrete time for an infinite horizon. Let $X_t^i \in \mathcal{X}$ and $U_t^i \in \mathcal{U}$ denote the state and action of agent $i \in N$ at time t . Note that the state space \mathcal{X} and action space \mathcal{U} are the same for all agents. For ease of exposition, we assume that \mathcal{X} and \mathcal{U} are finite sets. Given a vector $x = (x^1 \dots x^n) \in \mathcal{X}^n$ of length n , let $\xi(x)$ denote the mean-field (or empirical distribution) of x , i.e., $\xi(x) = \frac{1}{n} \sum_{i \in N} \delta_{x^i}$. Let $Z_t = \xi(X_t)$ denote the mean-field of the team at time t and \mathcal{Z} denote the space of space of realizations of Z_t . Note that \mathcal{Z} has at most $(n+1)^{|\mathcal{X}|}$ elements. Let $(\{x_t\}_{t \geq 0}, \{u_t\}_{t \geq 0})$ denote a realization of $(\{X_t\}_{t \geq 0}, \{U_t\}_{t \geq 0})$ and let $z_t = \xi(x_t)$. We assume that the initial states of all agents are independent, i.e., $\mathbb{P}(X_0 = x_0) = \prod_{i \in N} \mathbb{P}(X_0^i = x_0^i) =: \prod_{i \in N} P_0(x_0^i)$, where P_0 denotes the initial state distribution of agents. We assume that the global state of the system evolves in a controlled Markov manner, i.e., $\mathbb{P}(X_{t+1} = x_{t+1} \mid X_{0:t} = x_{1:t}, U_{0:t} = u_{0:t}) = \mathbb{P}(X_{t+1} = x_{t+1} \mid X_t = x_t, U_t = u_t)$. All agents are exchangeable, so the state evolution of a generic agent depends on the states and actions of other agents only through the mean-fields of the states, i.e., for agent i :

$$\mathbb{P}(X_{t+1} = x_{t+1} \mid X_t = x_t, U_t = u_t) = \prod_{i \in N} \mathbb{P}(X_t^i = x_t^i, U_t^i = u_t^i, Z_t = z_t) =: \prod_{i \in N} P(x_{t+1}^i \mid x_t^i, u_t^i, z_t),$$

where P denotes the control transition matrix. Combining all of the above, we have

$$\mathbb{P}(X_{t+1} = x_{t+1} \mid X_{0:t} = x_{0:t}, U_{0:t} = u_{0:t}) = \prod_{i \in N} P(x_{t+1}^i \mid x_t^i, u_t^i, z_t). \quad (1)$$

The system has mean-field sharing information-structure, i.e., the information available to agent i is given by: $I_t^i = \{X_t^i, Z_t\}$. We assume that all agents use identical (stochastic) control law: $\mu_t: \mathcal{X} \times \mathcal{Z} \rightarrow \Delta(\mathcal{U})$ to choose the control action at time t , i.e., $U_t^i \sim \mu_t(X_t^i, Z_t)$. Let $\mu = (\mu_1, \mu_2, \dots)$ denote the team policy for all times. Note that, in general, restricting attention to identical policies may lead to a loss of optimality. See [1] for an example. Nonetheless, identical policies are attractive for reasons of fairness, simplicity, and robustness.

The team receives a per-step reward given by: $R_t \sim r(X_t, U_t)$. Given strategy $\mu = (\mu_1, \mu_2, \dots)$ the expected total reward incurred by the team is given by:

$$J(\mu) = \mathbb{E}^\mu \left[\sum_{t=0}^{\infty} \gamma^t R_t \right], \quad (2)$$

where $\gamma \in (0, 1)$ is the discount factor. The objective is to choose a policy μ to maximize the performance $J(\mu)$ given by (2).

3 Solution approach

The mean-field team model formulated above is a multi-agent team problem with non classical information structure. A planning solution of this model was presented in [1], which we summarize below for completeness. We then present a framework for using reinforcement learning in such models.

3.1 Planning solution for mean-field teams

Given any policy $\mu = (\mu_1, \mu_2, \dots)$ and any realization, $z = (z_1, z_2, \dots)$ of the mean-field, define *prescriptions* $h_t: \mathcal{X} \rightarrow \Delta(\mathcal{A})$ given by $h_t(x) = \mu_t(x, z_t)$, $\forall x \in \mathcal{X}$. Let \mathcal{H} denote the space of all such prescriptions. When the mean field trajectory is a random process, the prescriptions h_t is a random vector which we denote H_t . The results of [1] relies on the following two key properties. Let $(z_{1:t+1}, h_{1:t})$ denote any realization of $(Z_{1:t+1}, H_{1:t})$. We have:

1. $\{Z_t\}_{t \geq 1}$ is a controlled Markov process with control action h_t , i.e., $\mathbb{P}^\mu(Z_{t+1} = z_{t+1} \mid Z_{1:t} = z_{1:t}, H_{1:t} = h_{1:t}) = \mathbb{P}(Z_{t+1} = z_{t+1} \mid Z_t = z_t, H_t = h_t)$. Note that the right hand side does not depend on the choice of decision rule μ . Furthermore, the right hand side can be simplified as: $\mathbb{P}(Z_{t+1} = z_{t+1} \mid Z_t = z_t, H_t = h_t) = \sum_{x_{t+1}: \xi(x_{t+1}) = z_{t+1}} \prod_{i \in N} P(x_{t+1}^i \mid x_t^i, h_t(x_t^i), z_t)$, where x_t is any state such that $\xi(x_t) = z_t$.

2. The expected per-step reward simplifies as follows. $\mathbb{E}[r(\mathbf{X}_t, \mathbf{U}_t)|Z_{1:t}, H_{1:t}] = \mathbb{E}[r(\mathbf{X}_t, \mathbf{U}_t)|Z_t, H_t] =: \tilde{r}(Z_t, H_t)$.

It is shown in [1] that these two properties imply that the optimal policy μ can be identified as follows.

Theorem 1 Let $V : \mathcal{Z} \rightarrow \mathbb{R}$ be the unique bounded fixed point of the following equation:

$$V(z) := \max_{h \in \mathcal{H}} \mathbb{E}[\tilde{r}(z, h) + \gamma V(Z_{t+1}) | Z_t = z, H_t = h]. \quad (3)$$

Let $\psi(z)$ be an arg max of the right hand side of (3). Then the policy, $\mu(x, z) = \psi(z)(x)$, is an optimal policy for Problem (2). \square

The action space \mathcal{H} of the above dynamic program is all functions from \mathcal{X} to $\Delta(\mathcal{U})$. We assume that \mathcal{H} is approximated by some family of parametrized functions $\mathcal{H}_\Phi = \{h_\phi\}_{\phi \in \Phi}$ (where Φ is a compact and convex set) such as Gibbs/Boltzmann functions or neural networks. With such a parametrization, the dynamic program of (3) may be approximated as:

$$V(z) = \max_{\phi \in \Phi} \mathbb{E}[\tilde{r}(z, h_\phi) + \gamma V(Z_{t+1}) | Z_t = z, H_t = h_\phi] \quad (4)$$

Let $\hat{\psi}(z)$ be an arg max of the right hand side of (4). Then the policy, $\mu(x, z) = h_{\hat{\psi}(z)}(x)$, is the best policy for Problem (2) when $\mu_t(\cdot, z_t)$ is restricted to belong to \mathcal{H}_Φ .

3.2 Reinforcement learning for mean-field teams (MFT-RL)

In this section, we present a reinforcement learning algorithm for the special case where the reward is a cumulative reward, i.e., $R_t = \frac{1}{n} \sum_{i \in N} R_t^i$, where $R_t^i \sim \hat{r}(X_t^i, U_t^i, Z_t)$. We assume that we have access to a simulator for $P(\cdot | x_t^i, u_t^i, z_t)$ and $\hat{r}(x_t^i, u_t^i, z_t)$. This simulator is for a generic agent and takes the current local state, current local action and current mean-field as input and generates a sample of the local next state and the total reward as output. Using n copies of this simulator, we create a simulator for the mean-field dynamics. We start with n agents with initial local state sampled according to P_0 . We assume that all these agents use a common stochastic policy $\hat{\psi} : \mathcal{Z} \rightarrow \Phi$ to generate prescription parameters $\phi_t \sim \hat{\psi}(z_t)$. Given this sampled value of ϕ_t , each agent independently samples a control action $u_t^i \sim h_{\phi_t}(x_t^i)$. The actions u_t^i of agent i and the current mean-field z_t are given as input to the i^{th} simulator and the sampled output (X_{t+1}^i, R_t^i) are averaged to obtain (Z_{t+1}, R_t) . Thus, we have a simulator with internal state z_t . This simulator takes ϕ_t as an input and gives (Z_{t+1}, R_t) as sampled next mean-field state and reward. Thus, this is a simulator for $P(z_{t+1} | z_t, h_{\phi_t})$ and $\tilde{r}(z_t, h_{\phi_t})$. We can use this simulator with any standard RL algorithm to find the optimal policy for the dynamic program (4). In our experiments below, we use TRPO [16], PPO [17] and NAFDQN [5].

4 Numerical experiments

4.1 Benchmark domains

We consider the following domains to illustrate different decentralized reinforcement learning algorithms.

4.1.1 Demand response in smart grids

This is a stylized model for demand response in smart grids [1]. The system consists of n agents, where $\mathcal{X} = \{0, 1\}$, $\mathcal{U} = \{\emptyset, 0, 1\}$,

$$P(\cdot | \cdot, \emptyset, z) = M \quad (5)$$

$$P(\cdot | \cdot, 0, z) = (1 - \varepsilon_1) \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} + \varepsilon_1 M \quad (6)$$

$$P(\cdot | \cdot, 1, z) = (1 - \varepsilon_2) \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} + \varepsilon_2 M, \quad (7)$$

where M denotes the ‘‘natural’’ dynamics of the systems and ε_1 and ε_2 are small positive constants.

The per-step reward is given by: $R_t = - \left(\frac{1}{n} \sum_{i \in N} \left(c_0 \mathbb{1}_{\{U_t^i=0\}} + c_1 \mathbb{1}_{\{U_t^i=1\}} \right) + KL(Z_t || \zeta) \right)$, where c_0 and c_1 are costs for taking actions 0 and 1 respectively, ζ is a given target distribution and $KL(Z_t || \zeta)$ denotes the Kullback-Leibler divergence between Z_t and ζ . In our experiments, we consider we consider a system with $n = 100$ agents, initial state distribution $P_0 = [1/3, 2/3]$, $M = \begin{bmatrix} 0.25 & 0.75 \\ 0.375 & 0.625 \end{bmatrix}$, $c_0 = 0.1$, $c_1 = 0.2$, $\zeta = [0.7, 0.3]$, $\varepsilon_1 = \varepsilon_2 = 0.2$ and discount factor $\gamma = 0.9$.

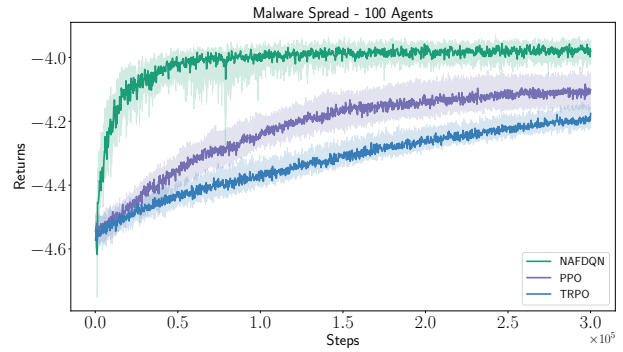
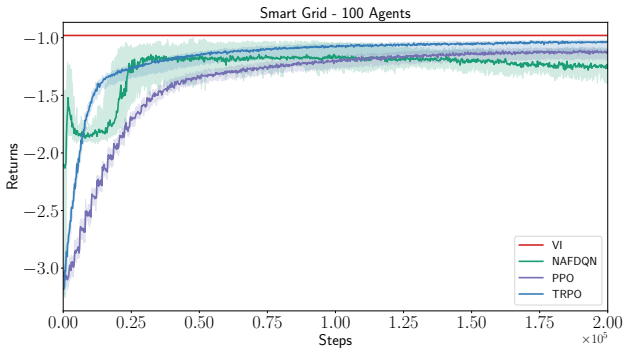


Figure 1: Demand response domain (25 independent runs). Figure 2: Malware spread domain (15 independent runs).

4.1.2 Malware spread in networks

This is a stylized model for malware spread in networks [6–8]. The system consists of n agents where $\mathcal{X} = [0, 1], \mathcal{U} = \{0, 1\}$. The dynamics are given by:

$$X_{t+1}^i = \begin{cases} X_t^i + (1 - X_t^i)\omega_t, & \text{for } U_t = 0, \\ 0 & \text{for } U_t = 1, \end{cases}$$

where $\omega_t \sim \text{Uniform}[0, 1]$. The per-step reward is given by: $R_t = -\left(\frac{1}{n} \sum_{i \in N} (k + \langle Z_t \rangle) X_t^i + \lambda U_t^i\right)$, where $\langle Z_t \rangle$ denotes the average of Z_t , and λ is the cost of taking action 1. In our experiments, we consider $k = 0.2$, initial state distribution $P_0 = \text{Uniform}(\mathcal{X})$, $\lambda = 0.5$ and discount factor $\gamma = 0.9$. For the simulation, we discretize the state space into 11 bins— $0, 0.1, \dots, 1$.

4.2 Simulation results

We consider three variants of MFT-RL algorithms, which use different RL algorithms for the mean-field system—TRPO, PPO and NAFDQN. Figure 1 shows the result for the demand response domain and Figure 2 shows the result for the malware spread domain. For each of the MFT-RL algorithms, the dark line shows the median performance and the shaded region shows the region between the first and third quartiles across multiple independent runs. For the demand response domain we also show the optimal performance obtained using the value iteration algorithm presented in [1]. All these variants of MFT-RL algorithms converge almost to the optimal value.

4.3 Mean-field approximations

Mean-field approximations are a common approach to simplify the planning solution of mean-field coupled systems. The main idea is to approximate a large population system with an infinite population system, find the optimal policy for the infinite population system and use that policy in the finite population system. Under appropriate regularity conditions, it can be shown that such an approximate policy is ε -optimal where ε is $O(1/n)$ or $O(1/\sqrt{n})$. Such approximations rely on the system model and are not appropriate in the learning setup. However, the mean-field approximation results suggest some form of continuity in the optimal policy as the number of agents becomes large. This motivates us to investigate the reverse question. Can we find an approximate policy for a n -agent mean-field team by running MFT-RL on m agents, where $m < n$?

We investigate this idea in the demand response domain. We use MFT-RL for $m = 100$ agents using TRPO and PPO, and use the resultant policy in the systems with $n > 100$ agents. We compare this performance with optimal planning solution obtained using value iteration. The results are shown in Figure 3. This shows that the policy obtained for the 100 agent RL environment performs reasonably well in environments with larger number of agents as well.

5 Conclusion

There are many results in the Dec-POMDP/decentralized control literature where a team optimal solution can be obtained using dynamic programming. Our central thesis is that for such models one can easily translate the dynamic program to a reinforcement learning algorithm. We illustrate this point by using mean-field teams as an example. This allows us to use standard off-the-shelf RL algorithms to obtain solutions for some MARL setups. The numerical results show that standard single agent RL algorithms work for RL for MFTs.

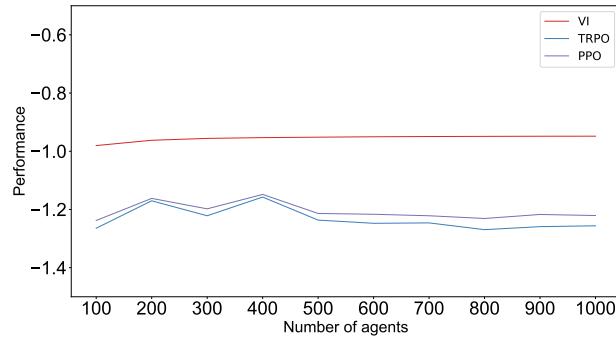


Figure 3: Performance of policy obtained in 100 agent system in systems with larger number of agents.

References

- [1] ARABNEYDI, J., AND MAHAJAN, A. Team optimal control of coupled subsystems with mean-field sharing. In *IEEE Conference on Decision and Control* (2014), IEEE, pp. 1669–1674.
- [2] BUŞONIU, L., BABUŞKA, R., AND DE SCHUTTER, B. Multi-agent reinforcement learning: An overview. In *Innovations in multi-agent systems and applications-1*. Springer, 2010, pp. 183–221.
- [3] FOERSTER, J., NARDELLI, N., FARQUHAR, G., TORR, P., KOHLI, P., WHITESON, S., ET AL. Stabilising experience replay for deep multi-agent reinforcement learning. *arXiv preprint arXiv:1702.08887* (2017).
- [4] FOERSTER, J. N., SONG, F., HUGHES, E., BURCH, N., DUNNING, I., WHITESON, S., BOTVINICK, M., AND BOWLING, M. Bayesian action decoder for deep multi-agent reinforcement learning. *CoRR* (2018).
- [5] GU, S., LILLICRAP, T., SUTSKEVER, I., AND LEVINE, S. Continuous deep Q-learning with model-based acceleration. In *ICML* (2016).
- [6] HUANG, M., AND MA, Y. Mean field stochastic games: Monotone costs and threshold policies. In *2016 IEEE 55th Conference on Decision and Control (CDC)* (Dec 2016), pp. 7105–7110.
- [7] HUANG, M., AND MA, Y. Mean field stochastic games with binary action spaces and monotone costs. *ArXiv e-prints* (Jan. 2017).
- [8] HUANG, M., AND MA, Y. Mean field stochastic games with binary actions: Stationary threshold policies. In *2017 IEEE 56th Annual Conference on Decision and Control (CDC)* (Dec 2017), pp. 27–32.
- [9] LEIBO, J. Z., ZAMBALDI, V., LANCTOT, M., MARECKI, J., AND GRAEPEL, T. Multi-agent reinforcement learning in sequential social dilemmas. In *Conference on Autonomous Agents and MultiAgent Systems* (2017).
- [10] LITTMAN, M. L. Markov games as a framework for multi-agent reinforcement learning. In *ICML*. (1994).
- [11] LITTMAN, M. L. Friend-or-foe q-learning in general-sum games. In *ICML*. (2001).
- [12] LITTMAN, M. L. Value-function reinforcement learning in markov games. *Cogn. Sys. Research* 2, 1 (2001), 55–66.
- [13] MAHAJAN, A., AND MANNAN, M. Decentralized stochastic control. *Ann Oper Res.*, 241 (June 2016), 109–126.
- [14] NAYYAR, A., MAHAJAN, A., AND TENEKETZIS, D. Decentralized stochastic control with partial history sharing: A common information approach. *IEEE Transactions on Automatic Control* 58, 7 (2013), 1644–1658.
- [15] OLIEHOEK, F. A., AND AMATO, C. *A concise introduction to decentralized POMDPs*, vol. 1. Springer, 2015.
- [16] SCHULMAN, J., LEVINE, S., ABBEEL, P., JORDAN, M., AND MORITZ, P. Trust region policy optimization. In *ICML* (June 2015).
- [17] SCHULMAN, J., WOLSKI, F., DHARIWAL, P., RADFORD, A., AND KLIMOV, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [18] YANG, J., YE, X., TRIVEDI, R., XU, H., AND ZHA, H. Deep mean field games for learning optimal behavior policy of large populations. In *ICLR* (2018).
- [19] YANG, Y., LUO, R., LI, M., ZHOU, M., ZHANG, W., AND WANG, J. Mean field multi-agent reinforcement learning. In *ICML* (2018).
- [20] YIN, H., MEHTA, P. G., MEYN, S. P., AND SHANBHAG, U. V. Learning in mean-field games. *IEEE TAC* 59, 3 (Mar 2014).
- [21] ZHANG, K., YANG, Z., LIU, H., ZHANG, T., AND BASAR, T. Fully decentralized multi-agent reinforcement learning with networked agents. In *ICML* (2018).