

CONSISTENCY OF MULTI-AGENT BATCH REINFORCEMENT LEARNING

Clemens Possnig

Vancouver School of Economics, University of British Columbia

November 30, 2022

Working Paper

[Link to current version](#)

ABSTRACT. This paper provides asymptotic results for a class of actor-critic batch - reinforcement learning algorithms in the multi-agent setting. At each period, each agent faces an estimation problem (the critic, e.g. a value function $Q(s, a)$), and a policy updating problem. The estimation step is done by parametric function estimation based on a batch of past observations. I give sufficient conditions on the environment, growth rate of the batch-size and speed of their stepsizes, so that each agent's parametric function estimator is consistent in the following sense: For large t , the optimal parameter θ_t is close to a true optimal parameter θ_t^* , depending on t only through the current period's policy profile.

This result greatly simplifies the asymptotic analysis of multi-agent learning, e.g. in the application of long-run characterisations using stochastic approximation techniques.

Keywords. Multi-Agent Reinforcement Learning, Batch-Reinforcement Learning, Consistency.

I thank Nima Haghpahanah, Vadim Marmer, and Kevin Song for helpful discussions. I thank the participants at EC 22, GTA22, CORS/INFORMS 22, and CETC 22 for insightful comments.

1. Introduction

This paper develops asymptotic results for the multi-agent reinforcement learning (MARL) setting which will help analyse what behaviors can be learned by algorithms that interact with one another.

Reinforcement Learning (RL) algorithms are updating rules meant for the learning of optimal policies or value functions for a given problem. Such algorithms are commonly used to solve Markov decision problems. In general, RL updating rules move policies towards actions that have performed well in the past (i.e., such actions are *reinforced*), and away from actions that perform poorly, based on some objective function. Commonly, a RL agent estimates a value function, and updates policies based on that value function. If estimates converge to the correct value function, policies commonly also will converge to optimal policies, and learning is successful. For a thorough introduction to RL see Sutton and Barto (2018). Multiple recent surveys on MARL and related theoretical results exist, notably Zhang, Yang, and Başar (2021), and Hernandez-Leal, Kartal, and Taylor (2019).

Recent years have brought significant advancements in the literature on multi-agent reinforcement learning (MARL). Such algorithms have proven successful in various strategic settings, such as the games of Go and Poker, and also autonomous driving. Despite these widespread successes, the performance of multi-agent learning algorithms is commonly verified only empirically, while theoretical results are relatively lacking. The aim of this paper is to provide novel theoretical results that are useful in determining convergence of multi-agent systems.

This paper considers the setting of actor-critic batch RL, which involves the mixture of offline (training performance measures on a batch of observations) and online (updating policies during play) approaches (c.f. Busoniu et al. (2017), Chapter 3).

A main problem when establishing theoretical results for the MARL setting is the inherent non-stationarity of the environment faced by each agent. This comes from the fact that each agent’s observations are drawn from distributions dependent on each other agent’s policies, which themselves are moving over time. At the same time, an agent needs to find an optimal decision at any given period, where for optimality only the current policies of their opponents matter. This introduces the commonly called problem of ‘tracking a moving target’: To find a best response, agents need to estimate a value function based on their opponent’s *current* policy, but can only use data generated from their opponent’s *past* strategies.

The batch-setting we study allows a useful solution to this issue. The name refers to the fact that the historical data used for estimating the performance measure is only a most recent window (the batch) of past observations, not the full available set of observations.

Akin to the idea of two-timescale approaches (c.f. Borkar (2009), Chapter 6), it will be true that the batch each agent uses to train their performance measure grows at a speed that is slower than the convergence rate of each agent’s policy-stepsizes. In that case, one can imagine that the most recent observations made by each agent are generated from distributions that are quite similar. Once that is true, I apply techniques developed and used in econometric theory due to Newey and McFadden (1994) to show that tracking the moving target becomes feasible. The Batch-RL setting is in contrast to more commonly known online-only RL schemes, which at every period t incorporate only the new information that has been accrued to adapt their performance measure estimator (see for example stochastic gradient descent methods for parametric Q -estimation in Sutton and Barto (2018)). The method of batch-learning is computationally more costly at each period, since a separate optimization routine is run at every period. However, we will see that this can put us at an advantage when it comes to estimation given nonstationarity.

The setting I study is one of discrete state spaces but interval action spaces, which implies the requirement of using function approximation in the estimation step. I do not make assumptions on the strategic nature of the interaction each agent faces, i.e. I make no requirement on the game being played to be zero-sum, cooperative or otherwise as is commonly done in the MARL setting. This paper focuses on results on a more fundamental level: it is only concerned with giving guarantees for the function approximator of each agent to be well-behaved in an appropriate sense. Once this can be verified, stochastic approximation techniques can be applied to paint a full picture of the convergence behavior of the policy-profile process implied by the MARL updating scheme. In Possnig (2022), I give such an analysis for Markov games of discrete states and interval actions under the assumption that function approximators are well-behaved in the sense developed in this paper.

To the best of my knowledge, this is the first paper providing MARL convergence analysis for batch RL in such a setting. The paper closest to mine is Perolat, Piot, and Pietquin (2018), which construct a stochastic approximation result for the two-timescale actor-critic scheme in the discrete state-action setting under a similar intuition as mentioned above in my batch-setting. For a more thorough discussion on recent advancements in the MARL literature, consider Zhang, Yang, and Başar (2021).

This paper is organized as follows: Section 2 gives the consistency result in full generality. Section 3 shows how the result applies to common learning rules such as actor-critic Q learning and gradient learning, and finishes with a Corollary that shows how the results developed here apply to Assumption 3 in Possnig (2022). All proofs are in the appendix.

2. Consistency

We begin by giving a general consistency result on the estimation-step for batch-RL algorithms. We assume there are n algorithmic agents, a finite state space S , a compact interval action space $A^i \subset \mathbb{R}$, with $A = \times_i A^i$, a twice differentiable, bounded payoff function $u^i : A \times S \mapsto \mathbb{R}$. We define the state transition probability to be $P_{ss'}[a] \geq 0$ for all $a \in A$ and $s, s' \in S$, where throughout we will maintain an assumption of irreducibility stated below.

We assume that each agent follows a batch-RL algorithm to update their policies $\rho_t^i : S \mapsto A^i$ over time. Let the resulting compact policy profile space be called Γ .

Assumption 1. *For all $\rho \in \Gamma$, the Markov chain induced by $P_{ss'}[\rho(s)]$ is irreducible and aperiodic.*¹

We maintain this assumption throughout the paper.

The updates are done using a parametric estimator of an underlying performance measure; this can be e.g. a Q -value function as discussed in Section 3. In this section the statement will be given in general terms, agnostic to the exact definition of the performance measure. All we need is that the parametric estimator can be expressed as the minimizer of a loss function.

We will call the parametric estimator $F^i(\rho^i, \theta^i)$. We assume that F^i is continuous in both arguments for all i , and that parameters $\theta \in \Theta$ for a set $\Theta \subset \mathbb{R}^m$ compact.

The consistency result this paper will establish is of the following form: each agent will use data generated from interactions with each other over time to estimate their parameter vector θ^i , while ρ_t^i are being updated concurrently. As a result, it is likely that an optimal θ_t^{i*} moves with time also, generating a moving-target problem. We will then prove:

Under suitable Assumptions, each agent's estimated θ_t^i behaves in the following way:

$$\|\theta_t^i - \theta_t^{i*}\| \rightarrow_P 0,$$

as $t \rightarrow \infty$, and also that in a sense that will become apparent, θ_t^{i*} depends on time t *only* through current period's policy profile ρ_t . This result is desirable as RL agents commonly face an issue of computing policies optimal with respect to the current distributional environment they face, but have only access to data generated from past distributional environments. This issue is absent in the single-agent stationary Markov Decision Problem, but salient in the multi-agent learning of focus here.

First, let $\mathcal{Z} \subset \mathbb{R}^d$ be a space of observations used in the construction of the loss function. Each period, a realization $Z_t \in \mathcal{Z}$ is generated after each algorithm chooses their actions.

¹For Definitions see e.g. Appendix A in Puterman (2014)

For example, $Z_t^i = \langle s_t, a_t^i, u_t^i, s_{t+1} \rangle$ would be the typical tuple of current state s_t , current action, payoff, and next state observed at the end of each period by a model-free² algorithm. We define a bounded function $l(Z, \theta) \in \mathcal{U} \subset \mathbb{R}$, Lipschitz in both arguments as the basic building block of the loss function.

Each algorithm uses only a batch of the most recent observations to construct their empirical loss function. Define a sequence $0 < K_t < t$ with $K_t \in \mathbb{N}$ such that $K_t \rightarrow \infty$ with t , and let

$$W_t = \{k : t - K_t + 1 \leq k \leq t\},$$

be the batch of periods used in the constuction of the loss function. We define $W_t(1) = t - K_t + 1$ as the first period of the batch. Then

$$L_t(\theta) = \frac{1}{K_t} \sum_{k \in W_t} l(z_k, \theta),$$

is the empirical loss. Then we define

$$\theta_t \in \arg \min_{\theta \in \Theta} L_t(\theta),$$

as the empirical parametric minimizer. Let $\boldsymbol{\rho}_t = \{\rho_k\}_{W_t(1) \leq k \leq t}$ denote batch-sequences of variables. Our first assumption is on the smoothness of the loss function and the behavior of its conditional expectation:

Assumption 2. *There exists a function $\phi(\rho, s, \theta) \in \mathbb{R}$ Lipschitz in the first and third arguments with*

(1)

$$\mathbb{E}[l(Z_t, \theta) \mid \boldsymbol{\rho}_t, \mathbf{s}_t] = \phi(\rho_t, s_t, \theta).$$

(2)

$$\lim_{t \rightarrow \infty} \frac{1}{K_t} \sum_{k \in W_t} \mathbb{E} C_1(Z_k) < \infty,$$

$$\lim_{t \rightarrow \infty} \frac{1}{K_t} \sum_{k \in W_t} \mathbb{E} C_2(\rho_k, s_k) < \infty,$$

and

$$\sup_{\theta \in \Theta} \max_{s \in S} C_3(\theta, s) < \infty,$$

²Model-free algorithms estimate their performance measure without a model of their environment. See Sutton and Barto (2018), e.g. Chapter 6.

Where there exist bounded, nonnegative functions $C_1(Z), C_2(\rho, s), C_3(\theta, s)$ by the Lipschitz properties of l, ϕ so that:

$$\begin{aligned} |l(Z, \theta) - l(Z, \theta')| &\leq C_1(Z) \|\theta - \theta'\|, \\ |\phi(\rho, s, \theta) - \phi(\rho, s, \theta')| &\leq C_2(\rho, s) \|\theta - \theta'\|, \\ |\phi(\rho, s, \theta) - \phi(\rho', s, \theta)| &\leq C_3(\theta, s) \|\rho - \rho'\|. \end{aligned}$$

Assumption 2 (1) can be satisfied if Z_t is Markov given current policy profile ρ_t, s_t , as will be discussed in Section 3.

Now we are ready to state the actor-critic updating schemes studied in this paper.

Definition 1. For each agent, ρ_t^i is updated in the following way:

$$\rho_{t+1}^i = \rho_t^i + \alpha_t [F^i(\rho_t^i, \theta_t^i) + M_{t+1}^i],$$

where $F^i(\rho_t^i, \theta_t^i)$ is the bounded parametric function to estimate the population objective, α_t is a decreasing stepsize sequence satisfying the Robbins-Monro condition:

$\alpha_t \rightarrow 0$ with

$$\sum_{t=0}^{\infty} \alpha_t = \infty; \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty,$$

and M_{t+1}^i is an almost surely bounded martingale-difference noise based on an increasing sequence of sigma algebras \mathcal{G}_t .

The next Assumption will ensure that the data used by the loss functions appropriately adjusts for the fact that a moving target has to be followed:

Assumption 3. Assume that

$$K_t \alpha_{t-K_t} \rightarrow 0,$$

as $t \rightarrow \infty$.

Define $\mu_{\rho_t}(s) \in (0, 1)$ for every s as the unique invariant state distribution given ρ_t , which exists by our irreducibility assumption 1, and let

$$\lambda_k(s, \boldsymbol{\rho}_t, s_{W_t(1)}) = \mathbf{P}(s_k = s \mid \boldsymbol{\rho}_t, s_{W_t(1)}),$$

be the likelihood of reaching state s in period $k \in W_t$, if over periods $W_t(1), \dots, k$, $\boldsymbol{\rho}_t$ is the policy profile sequence, and $s_{W_t(1)}$ is the initial state in the first batch-period. Also let $\lambda_k(s, \rho, s_t)$ be the counterpart where $\rho_l = \rho$ in all periods $W_t(1), \dots, k$.

Assumption 4.

- (1) Assume for all t , $\lambda_k(s, \rho, s_t)$ and μ_ρ are Lipschitz in ρ with Lipschitz constants bounded uniformly over S .

(2) There exists $c_P > 0$ and $1 \leq k < \infty$ such that for all $s', s \in S$

$$\inf_{\rho \in \Gamma} \mathbf{P}[s_k = s' \mid s_0 = s, \rho] \geq c_P.$$

Assumption 4 (2) is slightly stronger than our irreducibility assumption on the Markov chain over s . It ensures that in the asymptotic analysis we can safely assume $\lambda_k > 0$ for t large enough.

Next, define $\Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)}) = \frac{1}{K_t} \sum_{k \in W_t} \lambda_k(s, \boldsymbol{\rho}_k, s_{W_t(1)})$, and let $l(Z_k, \theta)_s = l(Z_k, \theta) \mathbf{1}\{s_k = s\}$ for all $s \in S$.

The population counterpart to $L_t(\theta)$ is then defined as

$$L_t^*(\theta, \boldsymbol{\rho}_t) = \sum_{s \in S} \Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)}) \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \boldsymbol{\rho}_k]}{\Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)})}.$$

Now we can define

$$\theta^*(\boldsymbol{\rho}_t) \in \arg \min_{\theta \in \Theta} L_t^*(\theta, \boldsymbol{\rho}_t), \quad (1)$$

as the best parameter for the parametric function approximation problem. Note that θ^* is a random variable due to the random trajectories $\boldsymbol{\rho}_t$. Define

$$L_t^*(\theta, \rho_t) = \sum_{s \in S} \Lambda_t(s, \rho_t, s_{W_t(1)}) \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \rho_k]}{\Lambda_t(s, \rho_t, s_{W_t(1)})},$$

as the population loss in the case where in all periods $k \in W_t$, ρ_t is the policy profile played. The t -limit of this loss function will play an important role in our results:

Lemma 1. Suppose Assumptions 1, 2, and 4 hold. Then for any sequence ρ_t ,

$$\lim_{t \rightarrow \infty} \sup_{\theta \in \Theta} \|L_t^*(\theta, \rho_t) - \sum_{s \in S} \mu_{\rho_t}(s) \phi(\rho_t, s, \theta)\| = 0.$$

Proof. All proofs can be found in the Appendix. □

From now on, we define

$$L_\infty^*(\theta, \rho_t) = \sum_{s \in S} \mu_{\rho_t}(s) \phi(\rho_t, s, \theta),$$

and

$$\theta^*(\rho_t) = \arg \min_{\theta \in \Theta} L_\infty^*(\theta, \rho_t).$$

Also define

$$\bar{\theta}(\boldsymbol{\rho}_t) = \arg \min_{\theta \in \Theta} L_t^*(\theta, \boldsymbol{\rho}_t).$$

$$\bar{\theta}(\rho_t) = \arg \min_{\theta \in \Theta} L_t^*(\theta, \rho_t).$$

The next Assumption ensures that for any trajectories, there is a unique minimizer θ^* . Define $B(x, \varepsilon)$ as the ε -ball centered at x .

Assumption 5 (Identification). *For any sequence ρ_t , any $\varepsilon > 0$ and $\theta_1 \notin B(\theta^*(\rho_t), \varepsilon)$, $\theta_2 \notin B(\bar{\theta}(\rho_t), \varepsilon)$, $\theta_3 \notin B(\bar{\theta}(\boldsymbol{\rho}_t), \varepsilon)$ there exists $\delta > 0$ such that for all $t \geq 1$:*

$$L_\infty^*(\theta_1, \rho_t) \geq L_\infty^*(\theta^*(\rho_t), \rho_t) + \delta,$$

$$L_t^*(\theta_2, \rho_t) \geq L_t^*(\bar{\theta}(\rho_t), \rho_t) + \delta,$$

$$L_t^*(\theta_3, \boldsymbol{\rho}_t) \geq L_t^*(\bar{\theta}(\boldsymbol{\rho}_t), \boldsymbol{\rho}_t) + \delta.$$

We can prove the following result:

Theorem 1. *Impose Assumptions 1 - 5. Then for any sequence ρ_t satisfying Definition 1, and for any $\varepsilon > 0$,*

$$\mathbf{P}(\|\theta_t - \theta^*(\rho_t)\| > \varepsilon) \rightarrow 0,$$

as $t \rightarrow \infty$.

This result is useful in the following sense: in general the function approximation parameter vector θ_t will depend on the whole policy profile trajectory $\boldsymbol{\rho}_t$. Given that opponent's policies are moving over time, this can result in a quite hard to interpret estimator and can lead to bad performance of the iteration ρ_t^i . However, the Assumptions taken in the Theorem ensure that in fact, θ_t will, for large enough t , depend on the trajectory of policy profiles only through the most *current* period t . Thus, the resulting loss function behaves as if each agent knew their opponent's current policy, and sampled observations from that policy to estimate their loss function.

Furthermore, the limiting population loss L_∞^* can represent desirable population loss functions commonly used in the literature, as will be shown in the next section. This will allow to make much more accurate predictions about future behavior of opponents, and therefore better performance of the algorithm as will be seen in Section 3.

3. Applications

Given the setup defined in the previous section, a valid performance measure would be based on the commonly used action-value function $Q^* : S \times A \mapsto \mathbb{R}$. Given a reward

function $u : A \times S \mapsto \mathbb{R}$, it is defined implicitly as

$$Q^*(s, a) = u(a, s) + \delta \mathbb{E} \left[\max_{a' \in A} Q^*(s', a') \mid a, s \right]. \quad (2)$$

An extensive literature of reinforcement learning theory has focused on estimating this function. In the single agent setting, where states evolve according to a controlled markov chain, many convergence results exist for estimators of Q^* , starting with the seminal results in Watkins (1989). See Sutton and Barto (2018) for a thorough exposition of learning algorithms related to Q^* .

Q -learning algorithms are algorithms used to estimate Q^* , and compute an optimal policy based on it. A class of algorithms fitting our Batch-RL framework would be versions of *Fitted Q-Iteration* (FQI), (Ernst, Geurts, and Wehenkel (2005), and Busoniu et al. (2017) Chapter 3 for a general discussion) as will be introduced below. Such iterations are meant to minimize parameters based on what is often called the squared Bellman-loss:

$$l(Z_t, \theta) = \left[u_t + \delta \max_{a'} Q(s_{t+1}, a', \theta) - Q(s_t, a_t, \theta) \right]^2, \quad (3)$$

where we let $Z_t = \langle s_t, a_t, u_t, s_{t+1} \rangle$. Suppose that each agent samples actions using a randomized policy $\bar{\rho}_t^i$ based on their iteration policy ρ_t^i . We will assume for simplicity here that $\mathbb{E} \bar{\rho}_t^i = \rho_t^i$, with full support on A^i for all states. In that case, we have

$$\mathbb{E}[l(Z_t, \theta) \mid \boldsymbol{\rho}_t, \mathbf{s}_t] = \mathbb{E}[l(Z_t, \theta) \mid \rho_t, s_t],$$

by the Markov property, and thus Assumption 2 (1) is satisfied. If we then impose Assumptions 1- 5, we get that empirical minimizer θ_t approaches

$$\theta^*(\rho_t) \in \arg \min_{\theta \in \Theta} L_\infty^*(\theta, \rho_t),$$

where $L_\infty^*(\theta, \rho_t)$ is the mean-squared Bellman loss,

$$\sum_{s \in S} \mu_{\rho_t}(s) \mathbb{E} \left[\left(u_t + \delta \max_{a'} Q(s_{t+1}, a', \theta) - Q(s_t, a_t, \theta) \right)^2 \mid \rho_t, s_t = s \right], \quad (4)$$

which represents a desirable population loss commonly studied in the literature (see for example Sutton and Barto (2018), e.g. Chapters 9, 11).

Two important examples of algorithms based on such a loss function then are Actor-Critic Q- learning, for which

$$F^i(\rho_t^i, \theta_t^i) = \left\{ \arg \max_{a \in A_i} Q^i(s, a, \theta_t^i) \right\}_{s \in S}, \quad (5)$$

and actor-critic gradient learning, where gradient here refers to a gradient in policies, not parameters as more commonly done in the literature:

$$F^i(\rho_t^i, \theta_t^i) = \left\{ \frac{\partial Q^i(s, a, \theta_t^i)}{\partial a} \right\}_{s \in S}, \quad (6)$$

where $\{\}_{s \in S}$ is to be understood as stacking a vector over $s \in S$.

Now we can verify that the result given in Theorem 1 can be applied to show that a given algorithm falls into the class developed in Possnig (2022). For a fixed opponent profile ρ_t^{-i} , we can define

$$Q^*(s, a, \rho_t^{-i}) = u(a, \rho_t^{-i}(s), s) + \delta \mathbb{E} \left[\max_{a' \in A} Q^*(s', a', \rho_t^{-i}) \mid a, s, \rho_t^{-i}(s) \right],$$

the action-value function in a repeated game where opponents play ρ_t^{-i} forever. The next result will show that for sufficient conditions given in this paper, Assumption 3 in Possnig (2022) holds. For convenience, we re-state this assumption here:

Assumption 6 (Assumption 3 in Possnig (2022)). *For each agent i there exists a bounded function $g^i(s, a, \rho^{-i})$, \mathcal{C}^2 in a, ρ^{-i} , such that*

$$\lim_{t \rightarrow \infty} \mathbf{P} \left[\sup_{(s, a) \in S \times A} \|Q_t^i(s, a) - Q^{i*}(s, a, \rho_t^{-i}) - g^i(s, a, \rho_t^{-i})\| \right] = 0.$$

Theorem 1 then allows us to conclude the following:

Corollary 1. *Suppose Assumptions 1 - 5 are satisfied, and algorithms update according to Definition 1. Write $Q_t(s, q) = Q(s, q, \theta_t)$, and*

$$g(s, a, \rho_t) = Q(s, a, \theta^*(\rho_t)) - Q^*(s, a, \rho_t^{-i}).$$

Then Assumption 3 in Possnig (2022) holds if $g(s, a, \rho)$ is twice differentiable in a, ρ .

The differentiability of g in ρ needs to be verified, however for well-behaved payoff and transition functions, what this condition comes down to is the twice-differentiability of $\theta^*(\rho)$. This requires on top of Assumption 5, that L_∞^* be smooth enough at its minimizer θ^* . For the mean-squared Bellman loss (4) as discussed in this section, this comes down to a condition on the differentiability of payoffs, transition probabilities, and the stationary distribution μ_ρ .

4. Conclusion

This paper gives sufficient conditions on the payoff structures, state evolution, and hyper-parameters of batch-RL algorithms so that their batch-estimation procedure has a tractable

analytical interpretation. The setting studied here is one of discrete states and interval action spaces. However, it is likely that an extension can be constructed for more general state spaces, which is subject of further study here.

The assumption throughout this paper is that each agent uses parametric function estimation in the classical sense, where the number of parameters is finite and *smaller* than the number of observations. This precludes the analysis of Deep RL methods, which by definition overparametrize. However, recent advancements in the convergence analysis of Deep RL for function approximation (Ramaswamy and Hullermeier (2021)) make me optimistic that an extension to this paper can be made that appropriately adjusts to Deep RL methods.

References

- Borkar, Vivek S (2009). *Stochastic approximation: a dynamical systems viewpoint*. Vol. 48. Springer.
- Busoniu, Lucian et al. (2017). *Reinforcement learning and dynamic programming using function approximators*. CRC press.
- Ernst, Damien, Pierre Geurts, and Louis Wehenkel (2005). “Tree-based batch mode reinforcement learning”. In: *Journal of Machine Learning Research* 6, pp. 503–556.
- Freedman, Ari (2017). “Convergence theorem for finite markov chains”. In: *Proc. REU*.
- Hernandez-Leal, Pablo, Bilal Kartal, and Matthew E Taylor (2019). “A survey and critique of multiagent deep reinforcement learning”. In: *Autonomous Agents and Multi-Agent Systems* 33.6, pp. 750–797.
- Newey, Whitney K and Daniel McFadden (1994). “Large sample estimation and hypothesis testing”. In: *Handbook of econometrics* 4, pp. 2111–2245.
- Perolat, Julien, Bilal Piot, and Olivier Pietquin (2018). “Actor-critic fictitious play in simultaneous move multistage games”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 919–928.
- Possnig, Clemens (2022). “Reinforcement Learning and Collusion”. URL: https://cjmpossnig.github.io/papers/jmp_CPossnig.pdf.
- Puterman, Martin L (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Ramaswamy, Arunselvan and Eyke Hullermeier (2021). “Deep Q-Learning: Theoretical Insights from an Asymptotic Analysis”. In: *IEEE Transactions on Artificial Intelligence*.
- Sutton, Richard S and Andrew G Barto (2018). *Reinforcement learning: An introduction*. MIT press.

Watkins, Christopher John Cornish Hellaby (1989). “Learning from delayed rewards”. In: Zhang, Kaiqing, Zhuoran Yang, and Tamer Başar (2021). “Multi-agent reinforcement learning: A selective overview of theories and algorithms”. In: *Handbook of Reinforcement Learning and Control*, pp. 321–384.

Appendix A. Appendix

All proofs are given here.

A.1. Proof of Lemma 1

We can write

$$\begin{aligned}
& \|L_t^*(\theta, \rho_t) - \sum_{s \in S} \mu_{\rho_t}(s) \phi(\rho_t, s, \theta)\| \\
& \leq \sum_{s \in S} \left\| \Lambda_t(s, \rho_t, s_{W_t(1)}) - \mu_{\rho_t}(s) \right\| \left\| \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \rho_t]}{\Lambda_t(s, \rho_t, s_{W_t(1)})} \right\| \\
& \quad + \sum_{s \in S} \mu_{\rho_t}(s) \left\| \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \rho_t]}{\Lambda_t(s, \rho_t, s_{W_t(1)})} - \phi(\rho_t, s, \theta) \right\| \\
& \leq \sum_{s \in S} R_{1,s,\rho_t,t} \left\| \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \rho_t]}{\Lambda_t(s, \rho_t, s_{W_t(1)})} \right\| + \max_{s \in S} R_{2,s,\rho_t,t},
\end{aligned}$$

where

$$\begin{aligned}
R_{1,s,\rho_t,t} &= \left\| \Lambda_t(s, \rho_t, s_{W_t(1)}) - \mu_{\rho_t}(s) \right\|, \\
R_{2,s,\rho_t,t} &= \left\| \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \rho_t]}{\Lambda_t(s, \rho_t, s_{W_t(1)})} - \phi(\rho_t, s, \theta) \right\|.
\end{aligned}$$

First note that for all fixed $\rho \in \Gamma$,

$$|\Lambda_t(s, \rho_t, s_{W_t(1)}) - \mu_\rho(s)| \rightarrow 0$$

as $t \rightarrow \infty$, and independently from initial state $s_{W_t(1)}$, which follows from irreducibility (c.f. Freedman (2017), Theorem 4.9). Since ρ_t can move over time, we prove uniform convergence in ρ . Firstly, $\rho_t \in \Gamma$ compact. So for any fixed $\delta > 0$, we can find an open cover of Γ of cardinality J_δ , using δ -balls centered at θ_j with $1 \leq j \leq J_\delta$. Now write for all $s \in S$

$$H_t(\rho, s) = \frac{1}{K_t} \sum_{k \in W_t} \lambda_k(s, \rho, s_{W_t(1)}) - \mu_\rho(s).$$

Then

$$\begin{aligned}
\sup_{\rho \in \Gamma} |H_t(\rho, s)| &\leq \max_{1 \leq j \leq J_\delta} \sup_{\rho \in B(\rho_j, \delta)} \{ |H_t(\rho, s) - H_t(\rho_j, s)| + |H_t(\rho_j, s)| \} \\
&\leq \sup_{\rho \in \Gamma} \sup_{\rho_1 \in B(\rho, \delta)} |H_t(\rho, s) - H_t(\rho_1, s)| + \max_{1 \leq j \leq J_\delta} |H_t(\rho_j, s)| \\
&= A_t + B_t,
\end{aligned}$$

where

$$A_t = \sup_{\rho \in \Gamma} \sup_{\rho_1 \in B(\rho, \delta)} |H_t(\rho, s) - H_t(\rho_1, s)|,$$

and

$$B_t = \max_{1 \leq j \leq J_\delta} |H_t(\rho_j, s)|.$$

Pointwise convergence of $H_t(\rho, s)$ implies that $B_t \rightarrow 0$ as $t \rightarrow \infty$. Then,

$$\begin{aligned}
A_t &\leq \sup_{\rho \in \Gamma} \sup_{\rho_1 \in B(\rho, \delta)} \frac{1}{K_t} \sum_{k \in W_t} |\lambda_k(s, \rho, s_{W_t(1)}) - \lambda_k(s, \rho_1, s_{W_t(1)})| + \sup_{\rho \in \Gamma} \sup_{\rho_1 \in B(\rho, \delta)} |\mu_\rho(s) - \mu_{\rho_1}(s)| \\
&\leq \sup_{\rho \in \Gamma} \sup_{\rho_1 \in B(\rho, \delta)} D_1 \frac{1}{K_t} \sum_{k \in W_t} \|\rho - \rho_1\| + \sup_{\rho \in \Gamma} \sup_{\rho_1 \in B(\rho, \delta)} D_2 \|\rho - \rho_1\| \\
&\leq \delta(D_1 + D_2),
\end{aligned}$$

where $0 < D_1, D_2 < \infty$ are the Lipschitz constants existing by Assumption 4. Thus for all $t \geq 1$, $A_t \rightarrow 0$ as $\delta \rightarrow 0$ (and recall that δ is picked arbitrarily), and the result follows:

$$\sup_{\rho \in \Gamma} |R_{1,s,\rho,t}| \rightarrow 0,$$

as $t \rightarrow \infty$.

Next, note that

$$\phi(\rho_t, s, \theta) = \frac{\mathbb{E}[l(Z_k, \theta)_s | \rho_t]}{P(s_k = s, \rho_t)},$$

with $P(s_k = s, \rho_t) = \sum_{s' \in S} \mu_{\rho_t}(s') \lambda_k(s, \rho_t, s')$ is the stationary expected value of λ_k over initial states s' . By Assumption 1, for all fixed ρ , $\lim_{k \rightarrow \infty} \lambda_k(s, \rho, s') = \lim_{k \rightarrow \infty} P(s_k = s, \rho) = \mu_\rho(s)$. Then we get

$$R_{2,s,\rho_t,t} \leq \|\mathbb{E}[l(Z_t, \theta)_s | \rho_t]\| \left\| \frac{1}{\Lambda_t(s, \rho_t, s_{W_t(1)})} - \frac{1}{P(s_k = s, \rho_t)} \right\| \quad (7)$$

$$\leq D_3 D_4 \|\Lambda_t(s, \rho_t, s_{W_t(1)}) - P(s_k = s, \rho_t)\|, \quad (8)$$

where $0 < D_3 < \infty$ is an upper bound on $\|\phi(\rho_t, s, \theta)\|$ following from the boundedness of the loss function, and $0 < D_4 < \infty$ is an upper bound on $\frac{1}{\Lambda_t(s, \rho_t, s_{W_t(1)})} \frac{1}{P(s_k=s, \rho_t)}$ which follows from irreducibility and Assumption 4, which implies that both fractions cannot diverge. Finally, the last term in (7) converges to zero uniformly over ρ by an argument analogous to the convergence of $H_t(\rho, s)$.

Finally, since

$$\left\| \frac{\frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}[l(Z_k, \theta)_s \mid \rho_t]}{\Lambda_t(s, \rho_t, s_{W_t(1)})} \right\| \leq D_3 D_4,$$

where the last bound is independent of θ , convergence of $\|L_t^*(\theta, \rho_t) - \sum_{s \in S} \mu_{\rho_t}(s) \phi(\rho_t, s, \theta)\|$ is uniform over $\theta \in \Theta$. The conclusion follows. ■

A.2. Proof of Theorem 1

The following Lemma will help prove the result. From now on, we drop the i -superscript whenever possible.

Lemma 2. *Impose Assumptions 1, 2, and 4.*

For all $\varepsilon > 0$,

$$\mathbf{P}\left(\sup_{\theta \in \Theta} \|L_t(\theta) - L_t^*(\theta, \boldsymbol{\rho}_t)\| > \varepsilon\right) \rightarrow 0,$$

as $t \rightarrow \infty$.

Proof. We first show pointwise convergence of $\|L_t(\theta) - L_t^*(\theta, \boldsymbol{\rho}_t)\|$.

We can write

$$L_t(\theta) = \sum_{s \in S} \frac{n_t(s)}{K_t} \frac{\sum_{k \in W_t} l(Z_k, \theta) \mathbf{1}\{s_k = s\}}{n_t(s)},$$

where

$$n_t(s) = \sum_{k \in W_t} \mathbf{1}\{s_k = s\}.$$

First we show, for all $s \in S$

$$\left| \frac{n_t(s)}{K_t} - \Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)}) \right| \rightarrow_P 0,$$

as $t \rightarrow \infty$. For this, define

$$V_{1,t} = \mathbb{E}\left[\left(\frac{n_t(s)}{K_t} - \Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)})\right)^2\right],$$

and let $d_t = \mathbf{1}\{s_k = s\} - \lambda_k(s, \boldsymbol{\rho}_k, s_{W_t(1)})$.

$$V_{1,t} = \frac{1}{K_t^2} \sum_{k \in W_t} \mathbb{E}[d_k^2] + \frac{1}{K_t^2} \sum_{k, k' \in W_t | k \neq k'} \mathbb{E}[d_k d_{k'}],$$

where the second term

$$\frac{1}{K_t^2} \sum_{k, k' \in W_t | k \neq k'} \mathbb{E}[d_k d_{k'}] = \frac{1}{K_t^2} \mathbb{E} \sum_{k, k' \in W_t | k \neq k'} \mathbb{E}[d_k d_{k'} \mid \mathbf{Z}_{k \vee k'}, \boldsymbol{\rho}_{k \wedge k'}, s_{k \wedge k'}] = 0,$$

since by definition of d_t and Assumption 1, states form a controlled markov chain and thus

$$\begin{aligned} \mathbb{E}[\mathbb{E}[d_k d_{k'} \mid \mathbf{Z}_{k \vee k'}, \boldsymbol{\rho}_{k \wedge k'}, s_{k \wedge k'}]] &= \mathbb{E}[d_{k \vee k'} \mathbb{E}[d_{k \wedge k'} \mid \mathbf{Z}_{k \vee k'}, \boldsymbol{\rho}_{k \wedge k'}, s_{k \wedge k'}]] \\ &= \mathbb{E}[d_{k \vee k'} \mathbb{E}[d_{k \wedge k'} \mid \rho_{k \wedge k'}, s_{k \wedge k'}]] = 0. \end{aligned}$$

It follows that $V_{1,t} \rightarrow 0$ as $K_t \rightarrow \infty$. We can then apply Chebyshev's inequality and the first result follows:

$$\left| \frac{n_t(s)}{K_t} - \Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)}) \right| = o_P(1).$$

Similarly, let $h_t(\theta, s) = l(Z_t, \theta) \mathbf{1}\{s_t = s\} - \mathbb{E}[l(Z_t, \theta) \mid \boldsymbol{\rho}_t]$. Define

$$V_{2,t,s} = \mathbb{E}\left[\left(\frac{1}{K_t} \sum_{k \in W_t} h_k(\theta, s)\right)^2\right],$$

then by an argument analogous to above, using Assumption 2 and boundedness of l we can conclude that $V_{2,t,s} \rightarrow 0$ as $t \rightarrow \infty$ for all s . By Assumption 4 we have that $\frac{n_t(s)}{K_t} > 0$ with probability approaching 1 with t . Thus we can apply the continuous mapping theorem to arrive at the result: for all $\theta \in \Theta$,

$$\|L_t(\theta) - L_t^*(\theta, \boldsymbol{\rho}_t)\| = o_P(1).$$

The rest of the proof is based on Newey and McFadden (1994)'s proof of their Theorem 2.1, but we have to adapt to the fact that we face a random population objective L^* due to the randomness of $\boldsymbol{\rho}_t$.

The proof follows a similar logic as the proof of Lemma 1. First define

$$H_t(\theta, \boldsymbol{\rho}_t) = \frac{1}{K_t} \sum_{k \in W_t} h_k(\theta),$$

where we drop the dependence on state s since the statement holds for any s and there are finitely many.

Take any $\varepsilon > 0$ and any $\delta > 0$. Let $B(x, \delta)$ denote the δ ball centered at x . Then by compactness of Θ , we can construct a finite open cover of Θ with cardinality $J_\delta < \infty$ using open balls $B(\theta_j, \delta)$. Now note that

$$\begin{aligned}
& \mathbf{P}\left(\sup_{\theta \in \Theta} \|H_t(\theta, \boldsymbol{\rho}_t)\| > 2\varepsilon\right) \\
& \leq \mathbf{P}\left(\max_{1 \leq j \leq J_\delta} \sup_{\theta \in B(\theta_j, \delta)} \{\|H_t(\theta, \boldsymbol{\rho}_t) - H_t(\theta_j, \boldsymbol{\rho}_t)\| + \|H_t(\theta_j, \boldsymbol{\rho}_t)\|\} > 2\varepsilon\right) \\
& \leq \mathbf{P}\left(\sup_{\theta \in \Theta} \sup_{\theta_1 \in B(\theta, \delta)} \|H_t(\theta, \boldsymbol{\rho}_t) - H_t(\theta_1, \boldsymbol{\rho}_t)\| + \max_{1 \leq j \leq J_\delta} \|H_t(\theta_j, \boldsymbol{\rho}_t)\| > 2\varepsilon\right) \\
& \leq A_t + B_t,
\end{aligned}$$

where

$$A_t = \mathbf{P}\left(\sup_{\theta \in \Theta} \sup_{\theta_1 \in B(\theta, \delta)} \|H_t(\theta, \boldsymbol{\rho}_t) - H_t(\theta_1, \boldsymbol{\rho}_t)\| > \varepsilon\right),$$

and

$$B_t = \mathbf{P}\left(\max_{1 \leq j \leq J_\delta} \|H_t(\theta_j, \boldsymbol{\rho}_t)\| > \varepsilon\right).$$

The second term must converge to zero by pointwise convergence as proved before, since

$$B_t \leq \sum_{1 \leq j \leq J_\delta} \mathbf{P}\left(\|H_t(\theta_j, \boldsymbol{\rho}_t)\| > \varepsilon\right) \rightarrow 0$$

as $t \rightarrow \infty$. Now define

$$Y_\delta = \sup_{\theta \in \Theta} \sup_{\theta_1 \in B(\theta, \delta)} \frac{1}{K_t} \sum_{k \in W_t} \|l(Z_k, \theta) - l(Z_k, \theta_1)\|,$$

and

$$\tilde{Y}_\delta = \sup_{\theta \in \Theta} \sup_{\theta_1 \in B(\theta, \delta)} \frac{1}{K_t} \sum_{k \in W_t} \|\mathbb{E}[(l(Z_k, \theta) - l(Z_k, \theta_1)) \mid \boldsymbol{\rho}_k]\|.$$

Then note that

$$A_t \leq \mathbf{P}(Y_\delta + \tilde{Y}_\delta > \varepsilon) \leq \frac{1}{\varepsilon} \mathbb{E}[Y_\delta + \tilde{Y}_\delta], \quad (9)$$

by Markov's inequality. Finally, note that

$$\begin{aligned}
\mathbb{E}Y_\delta & \leq \frac{1}{K_t} \sum_{k \in W_t} \mathbb{E} \sup_{\theta \in \Theta} \sup_{\theta_1 \in B(\theta, \delta)} \|l(Z_k, \theta) - l(Z_k, \theta_1)\| \\
& \leq \frac{1}{K_t} \sum_{k \in W_t} \mathbb{E} \sup_{\theta \in \Theta} \sup_{\theta_1 \in B(\theta, \delta)} C_1(Z_k) \|\theta - \theta_1\| \leq \frac{1}{K_t} \sum_{k \in W_t} \mathbb{E} C_1(Z_k) \delta,
\end{aligned}$$

where the second to last inequality follows from Assumption 2 and the Lipschitz property of $l(Z, \theta)$. Thus, we get

$$\lim_{t \rightarrow \infty} \mathbb{E}Y_\delta \leq \lim_{t \rightarrow \infty} \frac{1}{K_t} \sum_{k \in W_t} \mathbb{E}C_1(Z_k)\delta,$$

where the right hand side vanishes as $\delta \rightarrow 0$ by Assumption 2. We can make an analogous argument to show that $\lim_{t \rightarrow \infty} \mathbb{E}\tilde{Y}_\delta \rightarrow 0$ as $\delta \rightarrow 0$. It follows that $A_t \rightarrow 0$ as $t \rightarrow \infty$ and $\delta \rightarrow 0$ by the bound given in (9). The result follows, since $H_t(\theta)$ is the only factor in $L_t(\theta) - L_t^*(\theta, \boldsymbol{\rho}_t)$ that depends on θ :

We can write

$$\begin{aligned} & \|L_t(\theta) - L_t^*(\theta, \boldsymbol{\rho}_t)\| \\ \leq & \sum_{s \in S} \left\| \frac{n_t(s)}{K_t} - \Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)}) \right\| \left\| \frac{K_t}{n_t(s)} \frac{1}{K_t} \sum_{k \in W_t} l(Z_k, \theta)_s \right\| \\ & + \max_{s \in S} \left\| \frac{K_t}{n_t(s)} - \frac{1}{\Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)})} \right\| \left\| \frac{1}{K_t} \sum_{k \in W_t} l(Z_k, \theta)_s \right\| \\ & + \max_{s \in S} \left\| \frac{1}{\Lambda_t(s, \boldsymbol{\rho}_t, s_{W_t(1)})} \frac{1}{K_t} \sum_{k \in W_t} h_k(\theta, s) \right\|. \end{aligned}$$

There first two terms converge uniformly in θ to zero by our first arguments in this proof, due to the boundedness assumption on l and Assumption 4. Only the last term depends on $h_t(\theta, s)$, the uniform convergence of which has been shown above. □

Now we show that for any trajectory $\boldsymbol{\rho}_t$,

$$\sup_{\theta \in \Theta} \|L_t^*(\theta, \boldsymbol{\rho}_t) - L_t^*(\theta, \rho_t)\| \rightarrow 0,$$

as $t \rightarrow \infty$. For any $\theta \in \Theta$ we can write

$$\begin{aligned} \|L_t^*(\theta, \boldsymbol{\rho}_t) - L_t^*(\theta, \rho_t)\| & \leq \frac{1}{K_t} \sum_{k \in W_t} \|\mathbb{E}[l(Z_k, \theta) \mid \rho_k] - \mathbb{E}[l(Z_k, \theta) \mid \rho_t]\| \\ & \leq C_4 \frac{1}{K_t} \sum_{k \in W_t} \|\rho_k - \rho_t\| \leq C_4 C_5 \frac{1}{K_t} \sum_{k \in W_t} \sum_{l=k}^t \alpha_l, \end{aligned}$$

with $0 < C_4 < \infty$ being the bound on C_3 given by Assumption 2 and $0 < C_5 < \infty$ being the bound resulting from $F(\rho_t, \theta_t) + M_{t+1}$ being almost surely bounded given \mathcal{G}_t for all t . Since α_t is decreasing, we have

$$\frac{1}{K_t} \sum_{k \in D_t} \sum_{l=k}^t \alpha_l \leq K_t \alpha_{t-K_t},$$

and the last term vanishes by assumption of the Theorem. Since the last term is independent of θ , convergence is uniform.

As a last step we prove the convergence of θ_t . By Assumption 5,

$$\begin{aligned} \mathbf{P}(\theta_t \notin B(\theta^*(\rho_t), \varepsilon)) &\leq \mathbf{P}(L_\infty^*(\theta_t, \boldsymbol{\rho}_t) - L_\infty^*(\theta^*(\rho_t), \rho_t) \geq \delta) \\ &= \mathbf{P}\left(L_\infty^*(\theta_t, \rho_t) - L_t(\theta_t) + L_t(\theta_t) - L_t^*(\bar{\theta}_t(\boldsymbol{\rho}_t), \boldsymbol{\rho}_t) \right. \\ &\quad \left. + L_t^*(\bar{\theta}_t(\boldsymbol{\rho}_t), \boldsymbol{\rho}_t) - L_t^*(\bar{\theta}_t(\rho_t), \rho_t) + L_t^*(\bar{\theta}_t(\rho_t), \rho_t) - L_\infty^*(\theta^*(\rho_t), \rho_t) \geq \delta\right) \\ &\leq \mathbf{P}\left(L_\infty^*(\theta_t, \rho_t) - L_t(\theta_t) + L_t(\bar{\theta}_t(\boldsymbol{\rho}_t)) - L_t^*(\bar{\theta}_t(\boldsymbol{\rho}_t), \boldsymbol{\rho}_t) \right. \\ &\quad \left. + L_t^*(\bar{\theta}_t(\rho_t), \boldsymbol{\rho}_t) - L_t^*(\bar{\theta}_t(\rho_t), \rho_t) + L_t^*(\theta^*(\rho_t), \rho_t) - L_\infty^*(\theta^*(\rho_t), \rho_t) \geq \delta\right) \\ &\leq \mathbf{P}\left(\sup_{\theta \in \Theta} \|L_t(\theta) - L_t^*(\theta, \boldsymbol{\rho}_t)\| + \sup_{\theta \in \Theta} \|L_t^*(\theta, \rho_t) - L_t^*(\theta, \boldsymbol{\rho}_t)\| \right. \\ &\quad \left. + \sup_{\theta \in \Theta} \|L_t^*(\theta, \rho_t) - L_\infty^*(\theta, \rho_t)\| \geq \delta\right), \end{aligned}$$

where the second-to-last inequality follows from Assumption 5. The last bound vanishes the uniform convergence results given here and in Lemmas 1, 2, and the result follows. It follows that we can write $F(\rho_t, \theta_t) = F(\rho_t, \theta^*(\rho_t)) + o_P(1)$ as a function approximator that depends on policy profiles only through the *current period's* profile ρ_t , and not some weighted average of past profiles. ■