

Dimension-Free Rates for Natural Policy Gradient in Multi-Agent Reinforcement Learning

Carlo Alfano* and Patrick Rebeschini*

*Department of Statistics, University of Oxford

Abstract

Cooperative multi-agent reinforcement learning is a decentralized paradigm in sequential decision making where agents distributed over a network iteratively collaborate with neighbors to maximize global (network-wide) notions of rewards. Exact computations typically involve a complexity that scales exponentially with the number of agents. To address this curse of dimensionality, we design a scalable algorithm based on the Natural Policy Gradient framework that uses local information and only requires agents to communicate with neighbors within a certain range. Under standard assumptions on the spatial decay of correlations for the transition dynamics of the underlying Markov process and the localized learning policy, we show that our algorithm converges to the globally optimal policy with a dimension-free statistical and computational complexity, incurring a localization error that does not depend on the number of agents and converges to zero exponentially fast as a function of the range of communication.

1 Introduction

Sequential decision-making is a prominent setting in modern statistical theories and applications, where agents sequentially interact with an environment—observing its state, taking actions, and receiving rewards—to maximize notions of reward. Reinforcement learning is the setting where agents do not have complete knowledge of the environment dynamics, and it has received increased attention due to its recent successes on a variety of domains, e.g. games [Silver et al., 2016, 2017] and autonomous driving [Shalev-Shwartz et al., 2016].

Modern applications typically involve high-dimensional state and action spaces, and classical algorithms often lead to a computational complexity that scales exponentially with the number of degrees of freedom in the model. Understanding which structures can be used to design approximate methods that can overcome this curse of dimensionality while retaining near-optimal statistical guarantees is a question of paramount importance.

A class of algorithms that have proven successful to face high-dimensional models is that of Natural Policy Gradient (NPG) methods [Amari, 1998, Kakade, 2002, Peters and Schaal, 2008, Bhatnagar et al., 2009]. It has recently been shown [Agarwal et al., 2020] that NPG converges to an optimal policy with an iteration complexity that scales only logarithmically with the cardinality of the action space and with no explicit dependence on the cardinality of the state space.

Despite the favorable iteration complexity of NPG, NPG still faces the curse of dimensionality in applications where the computational cost per iteration scales exponentially with the number of degrees of freedom. This is the case in the setting of multi-agent reinforcement learning (MARL), for instance, where agents distributed over a network iteratively interact with each other to maximize global notions of reward. In this setting, the computational complexity of NPG—when applied to the entire network of agents—scales exponentially with the dimension of the model, which corresponds to the number of agents (see Section 3.1).

Along with the curse of dimensionality, NPG also faces scalability and implementability issues within the MARL framework. Applying NPG to the entire network of agents requires global communication, i.e. it requires each agent to be able to communicate with every other agent in the network, at every time step. This requirement is unrealistic in many multi-agent applications of interest, where the network topology is typically sparse, often grid-like, and where agents are only allowed to perform computation and communication in a decentralized manner, interacting only with neighboring agents within a certain range. These computational and communicational constraints arise, for instance, in the case of sensor networks, e.g. [Rabbat and Nowak, 2004, Nedic and Ozdaglar, 2009], and in the case of intelligent transportation systems, e.g. [Adler and Blue, 2002].

Over the past decades, various approaches have been proposed to address the curse of dimensionality in high-dimensional reinforcement learning models and, before that, in high-dimensional dynamical programming models, where exact knowledge of the probabilistic structure describing the environment is assumed. A popular approach involves designing algorithms that can exploit notions of *locality*, which encodes the assumption that, in some regimes, information can dissipate when it propagates through the network so that global computation and communication are not required to meet a prescribed level of error accuracy. Exploiting locality prompted the use of ad-hoc approximate factorization and truncation techniques, such as expressing the value function as a linear combination of basis functions that only depend on a small subset of local variables [Koller and Parr, 1999, Guestrin et al., 2001b, Koller and Parr, 2013, Yang and Wang, 2019, Jin et al., 2020]. These ideas have been applied to the MARL setting [Guestrin et al., 2001a, 2002, Sunehag et al., 2017, Rashid et al., 2018, Zhang et al., 2018a,b, Zhang and Zavlanos, 2019] and have proven successful in experiments, but lack theoretical guarantees or non-asymptotic analysis. A recent line of work has formally considered spatial decay of correlation assumptions for nearest-neighbors dynamics and designed decentralized algorithms based on policy gradient and actor-critic methods [Qu and Li, 2019, Qu et al., 2020a, Lin et al., 2020, Qu et al., 2020b], establishing non-asymptotic convergence guarantees towards a *stationary* point, but

not towards an *optimal* policy.¹ An application of the same principles to the setting of *mean-field* MARL [Yang et al., 2018] can be found in Haotian Gu [2021], where the authors show that a neural network based version of the actor-critic algorithm can achieve global convergence. In this setting, however, agents are considered to be indistinguishable and the transition scheme of an agent is only affected by the mean effect from its neighbors.

In this paper, we design a decentralized algorithm for the MARL setting based on the NPG framework that only uses local computations and communication for neighbors of agents within a certain range. We show that our algorithm can provably exploit spatial decay of correlation properties to overcome the curse of dimensionality, establishing non-asymptotic convergence guarantees to a globally *optimal* policy. In particular, we consider a general formulation of the decay of correlation assumption from statistical mechanics and probability theory [Dobrusin, 1970, Föllmer, 1982, Georgii, 2011], whereby agents have an influence on each other that decays exponentially with their distance on the network. This type of assumption has been previously considered in the learning literature, e.g. in Mitliagkas and Mackey [2017], Dagan et al. [2019], Borja Balle and Geumlek [2019], Prasad et al. [2020], Ilias Diakonikolas and Sun [2021], and also in the MARL setting, c.f. discussion in the previous paragraph. Under this assumption, we derive convergence bounds that are the same as those established for (centralized) NPG in [Agarwal et al., 2020], worsened only by a localization error that decreases exponentially with the radius of the communication range. A key feature of our bounds is that they are *dimension-free*, as they do not depend on the number of agents, and depend only logarithmically on the cardinality of the action space of *individual* agents and do not explicitly depend on the state space of individual agents. The localization radius controls the trade-off between statistical accuracy and computational complexity, as the overall computational cost of our algorithm scales only with respect to the number of agents within the local communication radius, and not with the total number of agents in the network.

Our contribution fits into the more general literature that has shown how spatial decay of correlations can be used to establish dimension-free results and are of interest in a variety of settings, such as [Gamarnik, 2013, Gamarnik et al., 2014], mixing times in spin systems [Hayes, 2006, Dyer et al., 2006], particle filtering [Rebeschini and Van Handel, 2015], epidemics [Mei et al., 2017], social networks [Chakrabarti et al., 2008], communication networks [Zocca, 2019], queuing networks [Papadimitriou and Tsitsiklis, 1994], and smart transportation [Zhang and Pavone, 2016].

The paper organization is as follows. In Section 2 we describe the MARL framework and we discuss the model assumptions we work with. In Section 3 we describe NPG as presented in Agarwal et al. [2020] and discuss its limitations when applied to MARL. In Section 4 we design a decentralized version of MARL and state our main results. The Appendix contains all the proofs of our statements and elaborates on the model assumptions.

¹Remark 10 gives a complete comparison of our results against previous findings that exploit the same type of decay of correlation assumption.

2 Setting

Let $\mathcal{G} = (\mathcal{K}, \mathcal{E})$ be an undirected graph describing a network of $|\mathcal{K}| = K$ agents. On this graph, the distance $d(k, k')$ between two agents $k, k' \in \mathcal{K}$ is defined as the length of the shortest path between the two vertices. Let $N_k^r = \{k' \in \mathcal{K} : d(k, k') \leq r\}$ denote the neighborhood of radius r of agent k , with $N_k = N_k^1$ and $N_{-k}^r = \mathcal{K} \setminus N_k^r$. Let \mathcal{S}_k and \mathcal{A}_k be the state and action spaces associated with agent k . We consider a Markov Decision Process (MDP) $(\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$: $\mathcal{S} = \mathcal{S}_1 \times \cdots \times \mathcal{S}_K$ and $\mathcal{A} = \mathcal{A}_1 \times \cdots \times \mathcal{A}_K$ are, respectively, the global state space and the global action space; $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}$, $P_k(s'_k | s, a)$ is the local transition probability, that is the probability that agent k transitions to state s'_k when the global state and action is (s, a) , and $P(s' | s, a) = \prod_{k \in \mathcal{K}} P_k(s'_k | s, a)$ is the global transition probability; $r(s, a) = \frac{1}{K} \sum_{k \in \mathcal{K}} r_k(s_k, a_k)$ is the global (network-wide) reward function that we wish to maximize, where $r_k : \mathcal{S}_k \times \mathcal{A}_k \rightarrow [0, 1]$ is the reward for agent k ; γ is the discount factor and μ is the starting state distribution. At time t denote the current state and action by $s(t)$ and $a(t)$.

To each agent k is assigned a local differentiable policy parameterized by $\theta_k \in \Theta_k$,

$$\pi_{\theta_k}(a_k | s) = \frac{e^{f_{\theta_k}(s, a_k)}}{\sum_{a' \in \mathcal{A}_k} e^{f_{\theta_k}(s, a')}},$$

which depends on the current global state s . Given the current global state s , each agent acts independently of the others. Denote $\theta = (\theta_1, \dots, \theta_K)$ and $\Theta = \Theta_1 \times \cdots \times \Theta_K$, then $\pi_{\theta}(a | s) = \prod_{k \in \mathcal{K}} \pi_{\theta_k}(a_k | s)$.

For a policy π and for each agent k , let $V_k^\pi : \mathcal{S} \rightarrow \mathbb{R}$ be the respective value function, which is defined as the expected discounted cumulative reward with starting state $s(0) = s$, namely,

$$V_k^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_k(s_k(t), a_k(t)) \middle| \pi, s(0) = s \right],$$

where $a(t) \sim \pi(\cdot | s(t))$ and $s(t+1) \sim P(\cdot | s(t), a(t))$. Let V^π be the global value function, defined as $V^\pi(s) = \frac{1}{K} \sum_{k \in \mathcal{K}} V_k^\pi(s)$, and $V^\pi(\mu)$ be the expected global value function when the initial state distribution is μ , i.e. $V^\pi(\mu) = \mathbb{E}_{s \sim \mu} V^\pi(s)$. Our objective is to find an optimal policy $\pi^* \in \arg\max_{\pi} \mathbb{E}_{s \sim \mu} V^\pi(s)$.

For a policy π and for each agent k , let $Q_k^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be the respective Q-function, which is defined as the expected discounted cumulative reward with starting state $s(0) = s$ and starting action $a(0) = a$, namely,

$$Q_k^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_k(s_k(t), a_k(t)) \middle| \pi, s(0) = s, a(0) = a \right],$$

where $a(t) \sim \pi(\cdot | s(t))$ and $s(t+1) \sim P(\cdot | s(t), a(t))$. Let Q^π be the global value function, defined as $Q^\pi(s, a) = \frac{1}{K} \sum_{k \in \mathcal{K}} Q_k^\pi(s, a)$.

Let $A_k^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be the advantage function for policy π and agent k , representing the advantage of taking the action a at step 0 and then following policy π , with respect to following policy π from the start, and defined as

$$A_k^\pi(s, a) = Q_k^\pi(s, a) - V_k^\pi(s).$$

Let $A^\pi(s, a) = \frac{1}{K} \sum_{k \in \mathcal{K}} A_k^\pi(s, a)$ be the global advantage function.

We define the discounted state visitation distribution [Sutton et al., 1999],

$$d_\rho^\pi(s) = (1 - \gamma) \mathbb{E}_{s(0) \sim \rho} \sum_{t=0}^{\infty} \gamma^t P(s(t) = s | \pi, s(0)),$$

and the discounted state-action visitation distribution,

$$d_\nu^\pi(s, a) = (1 - \gamma) \mathbb{E}_{s(0), a(0) \sim \nu} \sum_{t=0}^{\infty} \gamma^t P(s(t) = s, a(t) = a | \pi, s(0), a(0)),$$

where the trajectory $(s(t), a(t))_{t \geq 0}$ is generated by the MDP following policy π . Lastly, a function $f : \Theta \rightarrow \mathbb{R}$ is said to be a δ -smooth function of θ if, $\forall \theta, \theta' \in \Theta$,

$$\|\nabla f(\theta) - \nabla f(\theta')\|_2 \leq \delta \|\theta - \theta'\|_2.$$

2.1 Model Assumptions

We assume that a version of the Dobrushin condition [Georgii, 2011] holds for the transition dynamics of the network of agents. Let $TV(\mu, \nu) = \sup_{A \in \mathcal{F}} |\mu(A) - \nu(A)|$ be the total variation distance between the probability distributions μ and ν defined on the σ -algebra \mathcal{F} .

Assumption 1 (*Spatial Decay of Correlation for the Dynamics*) Let $C \in \mathbb{R}^{K \times K}$ be defined as follows:

$$C_{ij} = \sup_{s_j, s_{-j}, a_j, a_{-j}, s'_j, a'_j} TV(P_i(\cdot | s_j, s_{-j}, a_j, a_{-j}), P_i(\cdot | s'_j, s_{-j}, a'_j, a_{-j})).$$

Assume that there exists $\beta \geq 0$ such that

$$\max_{k \in \mathcal{K}} \sum_{j \in \mathcal{K}} e^{\beta d(k,j)} C_{kj} \leq \rho,$$

with $\rho < 1/\gamma$, where γ is the discount factor of the MDP.

The element (i, j) of the matrix C represents the influence that a perturbation of the state and action of agent j has on the transition probability of agent i . Assumption 1 encodes the

fact that the transition dynamics of each agent is exponentially less sensible to perturbations of the state and action of further away agents. The requirement $\rho < 1/\gamma$ comes as we need the spatial decay of correlation for the dynamics to be strong enough to induce a spatial decay for the Q-function (see Appendix B). A small value of the discount factor γ eases this requirement since it reduces the effect of perturbations through time. When $\beta = 0$ and $\gamma = 1$, Assumption 1 recovers the assumption in Qu et al. [2020a] as a particular case.

Differently from Qu et al. [2020a], we require the Dobrushin condition to hold for the policy as well. This is due to the fact that Assumption 1 is sufficient to prove the decay of correlation for the Q-function, on which the algorithm of Qu et al. [2020a] is based, but it is not sufficient to prove the decay of correlation for the value function, which instead needs an additional assumption on the policy. Since the NPG framework on which we build upon is based on both the Q-function and the value function, we make the following additional assumption.

Assumption 2 (*Spatial Decay of Correlation for the Policy*) Assume that there exist $\xi, \beta \geq 0$ such that, $\forall \theta \in \Theta$,

$$\sup_{s_{N_k^r}, s_{N_{-k}^r}, s'_{N_{-k}^r}} TV(\pi_{\theta_k}(\cdot | s_{N_k^r}, s_{N_{-k}^r}), \pi_{\theta_k}(\cdot | s_{N_k^r}, s'_{N_{-k}^r})) \leq \xi e^{-\beta r}.$$

Assumption 3 (*Local Policy*) Assume that, for any neighborhood radius r , the parameters θ_k can be partitioned in $(\theta_k)_{N_k^r}$ and $(\theta_k)_{N_{-k}^r}$ so that, if $(\theta_k)_{N_{-k}^r} = 0$, then

$$\begin{aligned} \pi_{\theta_k}(a_k | s) &= \pi_{\theta_k}(a_k | s_{N_k^r}), \\ \nabla_{(\theta_k)_{N_k^r}} \log \pi_{\theta_k}(a_k | s) &= \nabla_{(\theta_k)_{N_k^r}} \log \pi_{\theta_k}(a_k | s_{N_k^r}). \end{aligned}$$

Assumptions 2 and 3 impose a design constraint for the policy class $\{\pi_\theta | \theta \in \Theta\}$ rather than being assumptions on the nature of the environment, as the case for Assumption 1. Assumption 2 encodes, for the policy, a type of decay of correlation property that is weaker than Assumption 1. Assumption 2 allows us to consider a policy class that presents properties that are necessary for the optimal policy under Assumption 1, as we show in Appendix C. Assumption 3 is made to address the communication constraints of the network and requires the possibility of computing the policy and its gradient without access to the information coming from distant agents by setting their associated parameters to 0. In practice, we do only need Assumption 3 to hold for the value of r we want Theorem 9 to hold for. In Appendix A, we describe a policy class that satisfies both Assumption 2 and 3 for any value of r .

2.2 Exponential Decay

To take advantage of the local structure of the network, Lin et al. [2020] define a property regarding the dependence of $Q_k^\pi(s, a)$ on the neighbors of k .

Definition 4 (Lin et al. [2020]) The (c, ψ) -exponential decay property for the Q -function holds if, for any agent $k \in \mathcal{K}$ and for any $(s, a), (\tilde{s}, \tilde{a}) \in \mathcal{S} \times \mathcal{A}$ such that $s_{N_k^r} = \tilde{s}_{N_k^r}, a_{N_k^r} = \tilde{a}_{N_k^r}$, we have that

$$|Q_k^\pi(s, a) - Q_k^\pi(\tilde{s}, \tilde{a})| \leq c\psi^{r+1}.$$

In our analysis, we need to define the exponential decay property for the value function as well.

Definition 5 The (c', ϕ) -exponential decay property for the value function holds if, for any agent $k \in \mathcal{K}$ and for any $s, \tilde{s} \in \mathcal{S}$ such that $s_{N_k^r} = \tilde{s}_{N_k^r}$, we have that

$$|V_k^\pi(s) - V_k^\pi(\tilde{s})| \leq c'\phi^{r+1}.$$

These two properties mean that the cumulative discounted rewards of agents have an exponential decaying dependence on the states and actions of distant agents. We show that both these properties hold in our setting.

Proposition 6 If Assumptions 1 and 2 hold, then the exponential decay property holds for both the Q -function and the value function with parameters $(c, \psi) = \left(\frac{\gamma\rho e^\beta}{1-\gamma\rho}, e^{-\beta}\right)$ and $(c', \phi) = \left(\frac{\gamma(\rho+\xi)e^\beta}{1-\gamma(\rho+\xi)}, e^{-\beta}\right)$, respectively.

For clarity of exposition, in the rest of the paper we make the following assumption.

Assumption 7 Assume that the exponential decay property holds for the Q -function with parameters (c, ψ) and that it holds for the value function with parameters (c', ϕ) .

3 Natural Policy Gradient

We consider NPG as presented in Agarwal et al. [2020], which has iteration complexity that scales as $O(\sqrt{\log |\mathcal{A}|/T})$, where T is the number of iterations. We now summarize the algorithm and the results in Agarwal et al. [2020] and show what problems arise in the multi-agent setting that we consider.

Let π_θ be a differentiable policy and define the Fisher information matrix induced by π_θ as

$$F_\mu(\theta) = \mathbb{E}_{s \sim d_\mu^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[\nabla_\theta \log \pi_\theta(a|s) (\nabla_\theta \log \pi_\theta(a|s))^\top \right].$$

The NPG update, with step-size η , is defined as

$$\theta^{(t+1)} = \theta^{(t)} + \eta F_\mu(\theta^{(t)})^{-1} \nabla_\theta V^{\theta^{(t)}}(\mu), \quad (1)$$

where $\theta^{(t)}$ is the set of parameters at iteration t , $\nabla_{\theta} V^{\theta}(\mu)$ is the gradient of the value function with respect to the policy parameters, and $F_{\mu}(\theta^{(t)})^{-1}$ is the Moore-Penrose pseudo-inverse of the Fisher information matrix. As discussed in Agarwal et al. [2020], the update in (1) is equivalent to solving the problem

$$w_{\star} \in \operatorname{argmin}_w \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot|s)} [(A^{\pi_{\theta}}(s, a) - w \cdot \nabla_{\theta} \log \pi_{\theta}(\cdot|s))]^2 \quad (2)$$

and then performing the following update:

$$\theta^{(t+1)} = \theta^{(t)} + \frac{\eta}{1 - \gamma} w_{\star}. \quad (3)$$

Define

$$L(w, \theta, \nu) := \mathbb{E}_{s, a \sim \nu} [(A^{\pi_{\theta}}(s, a) - w \cdot \nabla_{\theta} \log \pi_{\theta}(a|s))]^2.$$

Assume that $\log \pi_{\theta}(a|s)$ is a δ -smooth function of θ and that $\pi^{(0)}$ is the uniform distribution. Let $d^{(t)} = d_{\nu}^{\pi^{(t)}}(s, a)$ and $d^{\star}(s, a) = d_{\mu}^{\pi^{\star}}(s) \pi^{\star}(a|s)$. Let ν be a distribution of s, a such that

$$\sup_{w \in \mathbb{R}^d} \frac{w^{\top} \Sigma_{d^{\star}}^{(t)} w}{w^{\top} \Sigma_{\nu}^{(t)} w} \leq \kappa,$$

where

$$\Sigma_{\nu}^{\theta} = \mathbb{E}_{s, a \sim \nu} [\nabla_{\theta} \log \pi_{\theta}(a|s) (\nabla_{\theta} \log \pi_{\theta}(a|s))^{\top}]$$

and $\Sigma^{(t)} = \Sigma^{\theta^{(t)}}$. Lastly, assume that

$$\mathbb{E} [L(w_{\star}^{(t)}, \theta^{(t)}, d^{\star})] \leq \varepsilon_{\text{bias}},$$

$$\mathbb{E} [L(w^{(t)}, \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}, \theta^{(t)}, d^{(t)}) | \theta^{(t)}] \leq \varepsilon_{\text{stat}},$$

where

$$w_{\star}^{(t)} \in \operatorname{argmin}_{\|w\|_2 \leq W} L(w, \theta^{(t)}, d^{(t)})$$

and the expectations are taken w.r.t. the sequence $(w^{(t)})_{t=0, \dots, T-1}$. Then, running algorithm (3) for T time steps with $\eta = \sqrt{2 \frac{\log |\mathcal{A}|}{\delta W^2 T}}$ for a given parameter W , we have the following guarantee: [Agarwal et al. [2020], Theorem 6.2]

$$\mathbb{E} \left[\min_{t \leq T} \{V^{\pi^{\star}}(\mu) - V^{(t)}(\mu)\} \right] \leq \frac{W}{1 - \gamma} \sqrt{\frac{2\delta \log |\mathcal{A}|}{T}} + \sqrt{\frac{\kappa \varepsilon_{\text{stat}}}{(1 - \gamma)^3}} + \frac{\sqrt{\varepsilon_{\text{bias}}}}{1 - \gamma}. \quad (4)$$

As we highlighted in the introduction, the guarantee (4) is particularly suitable for high-dimensional settings, as there is no explicit dependence on $|\mathcal{S}|$ and the dependence on $|\mathcal{A}|$ is only logarithmic. As to implicit dependencies, Agarwal et al. [2020] state that it is reasonable to expect that κ is not a quantity related to $|\mathcal{S}|$. On the other hand, $\varepsilon_{\text{stat}}$ and $\varepsilon_{\text{bias}}$ are constants related to a minimization problem that depends on both \mathcal{S} and \mathcal{A} .

3.1 Curse of Dimensionality and Scalability/Implementability Issues in MARL

When applied to the MARL setting, with $\mathcal{S} = \mathcal{S}_1 \times \dots \times \mathcal{S}_K$ and $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_K$, NPG would incur a curse of dimensionality or scalability and implementability issues, depending on the approach used for the minimization problem in (6). The guarantees for the algorithm in Agarwal et al. [2020] would yield:

$$\mathbb{E} \left[\min_{t < T} \{V^{\pi^*}(\rho) - V^{(t)}(\rho)\} \right] \leq \frac{WK}{1-\gamma} \sqrt{\frac{2\delta \log \max_{k \in \mathcal{K}} |\mathcal{A}_k|}{T}} + \sqrt{\frac{\kappa \varepsilon_{\text{stat}}}{(1-\gamma)^3}} + \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma}, \quad (5)$$

where

$$w_{\star}^{(t)} \in \underset{\|w\|_2 \leq \sqrt{KW}}{\operatorname{argmin}} L(w, \theta^{(t)}, d^{(t)}), \quad (6)$$

under the assumptions

$$\mathbb{E} [L(w^{(t)}, \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}, \theta^{(t)}, d^{(t)}) | \theta^{(t)}] \leq \varepsilon_{\text{stat}},$$

$$\mathbb{E} [L(w_{\star}^{(t)}, \theta^{(t)}, d^{\star})] \leq \varepsilon_{\text{bias}}.$$

In the original analysis in Agarwal et al. [2020], W is a parameter set by the user to control the norm of $w_{\star}^{(t)}$. In our setting, we normalize this parameter to \sqrt{KW} , which is analogous to requiring, for each agent k , the maximum norm of its optimal update $w_{k,\star}^{(t)}$ to be W . Not doing so would mean keeping a constant parameter W despite increases in the dimensions, i.e. agents, of problem (6), incurring increases in the bias term. In the multi-agent setting, the iteration complexity given by the bound in (5) is worse by a factor K , compared to the single-agent setting. The curse of dimensionality appears when solving the minimization problem in (6), e.g. with gradient descent because the computation of exact gradients involves a sum/integral over $\mathcal{S} \times \mathcal{A}$, which has a dimension that grows exponentially with the number of agents. If we solve the minimization problem with stochastic projected gradient descent, then the problem of computing gradients disappears as we *assume* access to samples to estimate gradients; however, the statistical guarantee becomes $O(K^2/\sqrt{N})$, from being $O(1/\sqrt{N})$ in the single-agent setting, due to the increase in dimensionality of the update w , in particular due to the scaling bound $\|w\|_2 \leq \sqrt{KW}$ that is used by a classical convergence result of stochastic projected gradient descent [Bubeck, 2014]. Implementing a sampler could, in turn, involve a curse of dimensionality.

The dependencies of the minimization problem in (6) cause the algorithm to incur additional scalability and implementability issues. As the projection step, the advantage function and the policy gradient depend on the states and actions of the entire network, which do not factorize, each agent would have to communicate to every other agent in the network at each iteration to solve the problem. As mentioned in the introduction, requiring such a level of communication is rarely viable in real-world applications in the decentralized MARL setting.

Remark 8 *These aforementioned problems do not arise in case of independent agents, where, for each agent k , the local transition probabilities satisfy $P_k(s'_k|s, a) = P_k(s'_k|s_k, a_k)$ and policies satisfy $\pi_{\theta_k}(a_k|s) = \pi_{\theta_k}(a_k|s_k)$. In this setting, as we show in Appendix D, it is possible to show that applying NPG to the whole network of K agents corresponds to running K independent runs of NPG applied to individual agents and to recover the same results of Agarwal et al. [2020] for the individual agents.*

4 Decentralized NPG

We design a decentralized version of NPG that is capable of exploiting the spatial decay of correlation properties that we assume and of avoiding the curse of dimensionality while still approximately converging to a globally optimal policy. We do so by limiting the communication range to agents that are at most at distance r and defining, for each agent k , the localized advantage function, localized value function, localized Q-function, as follows:

$$\begin{aligned}\tilde{A}_k^\pi(s_{N_k^r}, a_{N_k^r}) &= \tilde{Q}_k^\pi(s_{N_k^r}, a_{N_k^r}) - \tilde{V}_k^\pi(s_{N_k^r}), \\ \tilde{V}_k^\pi(s_{N_k^r}) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_k(s_k(t), a_k(t)) \mid \pi, s_{N_k^r}(0) = s_{N_k^r} \right], \\ \tilde{Q}_k^\pi(s_{N_k^r}, a_{N_k^r}) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_k(s_k(t), a_k(t)) \mid \pi, s_{N_k^r}(0) = s_{N_k^r}, a_{N_k^r}(0) = a_{N_k^r} \right].\end{aligned}$$

Following Assumption 3, we set $(\theta_k)_{N_{-k}^r} = 0, \forall k \in \mathcal{K}$ and we never update these parameters, so that the policy of an agent and its gradient do not depend on the states of agents whose distance is greater than r . For each agent k , define the loss function

$$\tilde{L}_k^r(w, \theta, \nu) = \mathbb{E}_{s, a \sim \nu} \left[\left(\tilde{A}_k^{\pi_\theta}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi_{\theta_k}(a_k | s_{N_k^r}) \cdot w \right)^2 \right]. \quad (7)$$

The minimization problem that each agent k aims to solve at each step becomes

$$w^* \in \underset{\|w\|_2 \leq W}{\operatorname{argmin}} \tilde{L}_k^r(w, \theta^{(t)}, d^{(t)}). \quad (8)$$

In Appendix F we show how to solve this minimization problem in a decentralized manner and that, even if $d^{(t)}$ is a global distribution, it is possible to build a decentralized sampler of it assuming only access to a global clock.

Exploiting decay of correlation properties and the policy design constraints in Assumption 3, Algorithm 1 removes the dependence on K from the iteration complexity bound and addresses the curse of dimensionality and scalability and implementability issues outlined in Section 3.1.

Algorithm 1: Decentralized NPG

Input: Learning rate η ; numbers of iterations T ; an initialized policy $\pi^{(0)}$.

Set $(\theta_k)_{N_k^r} = 0, \forall k \in \mathcal{K}$;

for $t = 0, \dots, T - 1$; $k \in \mathcal{K}$ **do**

 Compute approximately $w_k^{(t)} \in \operatorname{argmin}_{\|w\|_2 \leq W} \tilde{L}_k^r(w, \theta^{(t)}, d^{(t)})$;

 Compute the update $(\theta_k)_{N_k^r}^{(t+1)} = (\theta_k)_{N_k^r}^{(t)} + \frac{\eta}{1-\gamma} w_k^{(t)}$.

end

Theorem 9 *Assume that Assumption 3 and Assumption 7 hold. Assume that $\log \pi_\theta(a|s)$ is a δ -smooth function of θ and that $\pi^{(0)}$ is the uniform distribution. Let $d^{(t)} = d_\nu^{\pi^{(t)}}(s, a)$ and $d^*(s, a) = d_\nu^{\pi^*}(s) \pi^*(a|s)$. Let ν be a distribution of s, a for which there exists $\kappa' \geq 0$ such that*

$$\max_{k \in \mathcal{K}} \sup_{w \in \mathbb{R}^d} \frac{w^\top \Sigma_{d^*, k}^{(t)} w}{w^\top \Sigma_{\nu, k}^{(t)} w} \leq \kappa',$$

where, $\forall \theta, \nu, k$

$$\Sigma_{\nu, k}^\theta = \mathbb{E}_{s, a \sim \nu} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi_\theta(a_k | s_{N_k^r}) (\nabla_{(\theta_k)_{N_k^r}} \log \pi_\theta(a_k | s_{N_k^r}))^\top \right]$$

and $\Sigma_{\nu, k}^{(t)} \equiv \Sigma_{\nu, k}^{\theta^{(t)}}$. Let

$$\max_{k \in \mathcal{K}} \mathbb{E} \left[\tilde{L}_k(w_{k, \star}^{(t)}, \theta^{(t)}, d^*) \right] \leq \varepsilon_{bias},$$

$$\max_{k \in \mathcal{K}} \mathbb{E} \left[\tilde{L}_k(w_k^{(t)}, \theta^{(t)}, d^{(t)}) - \tilde{L}_k(w_{k, \star}^{(t)}, \theta^{(t)}, d^{(t)}) \mid \theta^{(t)} \right] \leq \varepsilon_{stat},$$

where

$$w_{k, \star}^{(t)} \in \operatorname{argmin}_{\|w\|_2 \leq W} \tilde{L}_k(w, \theta^{(t)}, d^{(t)}).$$

Then, Algorithm 1, with $\eta = \sqrt{2 \frac{\log |\mathcal{A}|}{(\delta K W^2 T)}}$, has the following guarantee:

$$\begin{aligned} \mathbb{E} \left[\min_{t < T} \{V^{\pi^*}(\mu) - V^{(t)}(\mu)\} \right] &\leq \frac{W}{1-\gamma} \sqrt{\frac{2\delta \log \max_{k \in \mathcal{K}} |\mathcal{A}_k|}{T}} + \sqrt{\frac{\kappa' \varepsilon_{stat}}{(1-\gamma)^3}} + \frac{\sqrt{\varepsilon_{bias}}}{1-\gamma} \\ &+ \underbrace{\frac{c\psi^{r+1} + c'\phi^{r+1}}{1-\gamma}}_{\text{localization error}}. \end{aligned} \tag{9}$$

Agent-wise, the assumptions in Theorem 9 correspond to the assumptions made in Agarwal et al. [2020], in a setting where, for each agent k , the state space and the action space are $\mathcal{S}_{N_k^r}$ and

\mathcal{A}_k , respectively, where the policy is defined as $\pi_\theta(a|s) = \pi_{\theta_k}(a_k|s_{N_k^r})$ and where the update $w_k^{(t)}$ is bounded by W . Theorem 9 shows that Decentralized NPG recovers the iteration complexity of the algorithm in Agarwal et al. [2020], worsened only by the fourth term on the RHS of (9). This localization error is exponentially small in r . Theorem 9 provides a dimension-free guarantee on the average expected cumulative rewards of the whole network, as the upper bound in (9) does not depend on the number of agents K in the network, and it depends only on the logarithm of the cardinality of the action space of an individual agent, with no explicit dependence on the state space of agents.

Theorem 9 shows that, under the assumption on spatial decay of correlation, Decentralized NPG solves the curse of dimensionality and the scalability and implementability issues outlined in Section 3.1. The minimization problem in (8) can be approximately solved in a decentralized manner through stochastic projected gradient descent, as we show in Appendix F, which leads to computational savings as we manage to eliminate the dependency on K from the statistical guarantee of the algorithm, obtaining the same computational complexity of the single-agent setting, i.e. $O(1/\sqrt{N})$, where N is the number of gradient steps, which is not surprising because problem (8) regards only the advantage function and the policy of an individual agent. Then, using stochastic projected gradient descent and the sampler in Algorithm 2, we recover, for each agent, the same expected sample complexity of the single agent setting ($2NT/(1-\gamma)$, where $2/(1-\gamma)$ is the expected length of a sampler episode). Decentralized NPG can be run locally by each agent and only requires information from neighbors within distance r .

The role of the term $\varepsilon_{\text{bias}}$ in (9) has a difference from the role that $\varepsilon_{\text{bias}}$ has in (4) [Agarwal et al., 2020]. They both represent the worst-case error that is made by the agents when they approximate their current advantage function with a linear combination of the elements of the gradient of their current policy and encode the *transfer* error that we make shifting the distribution to d^* . In addition to that, $\varepsilon_{\text{bias}}$ in (9) also encodes the localization error that we make in Algorithm 1 by using the localized loss defined in (7). In Appendix G we give a bound for this localization error of the bias term, showing that the localized bias is at most the non-localized bias, i.e. the bias associated with an infinite range parameter r , plus a quantity that decreases to 0 exponentially fast in r .

Remark 10 *The works in [Qu and Li, 2019, Qu et al., 2020a, Lin et al., 2020, Qu et al., 2020b] are closely related to our contribution, as they also use decay of correlation assumptions to provably avoid the curse of dimensionality in MARL. Our contribution differs from these works in the following main ways:*

1. *(Decay of correlation) We consider a more general version of the decay of correlation property (Assumption 1) and, differently from them, we also require a decay of correlation property to hold for the policy (Assumption 2). Assumption 1 recovers the version they consider in Qu et al. [2020a] in the case $\beta = 0$ and $\gamma = 1$. The generality of our*

condition allows us to consider transition dynamics that are not truncated, as they do, and to control the truncation of the policy.

2. (Methodology) Our method is based on NPG framework, while their methods is based on policy gradient and actor-critic methods.
3. (Optimality) We present statistical error bounds w.r.t. to the optimal policy, while the bounds they give are w.r.t. a stationary policy.
4. (Computational complexity) Our method has a computational complexity that does not depend, for any agent k , on the number of agents K or the number of neighbors $|N_k^r|$. The method in [Qu et al., 2020b] is shown to have a computational complexity that scales as $O(\log |\mathcal{S}||\mathcal{A}|)$, hence depending linearly on K , using additional assumptions on the minimum local exploration. The method in [Lin et al., 2020] is shown to have computational complexity that scales as $O(\log \max_{k \in \mathcal{K}} |\mathcal{S}_{N_k^r}| |\mathcal{A}_{N_k^r}|)$, hence depending linearly on $|N_k^r|$, using additional assumptions on the stationarity and on the mixing rates of the MDP.
5. (Statistical/Iteration Complexity) Under the only assumptions on decay of correlation and local policy, our method has an iteration complexity that scales as $O(\sqrt{\log \max_{k \in \mathcal{K}} |\mathcal{A}_k|})$. The methods in [Qu et al., 2020b, Lin et al., 2020] have an iteration complexity that does not depend on the state or action spaces.

5 Conclusion

We have investigated applications of the NPG framework to MARL, showing how a standard assumption on the spatial decay of correlation for the dynamics and for the policy on a network of agents, expressed through a form of Dobrushin condition, induces a form of exponential decay in the cumulative rewards that can be exploited by a localized version of NPG to avoid the curse of dimensionality. The version of NPG that we design scales to large networks and yields convergence guarantees to the optimal policy that are analogous to the ones in Agarwal et al. [2020], worsened only by a localization error that decreases exponentially with the communication radius. Our analysis does not consider regularization, which has been shown to accelerate convergence for NPG methods [Geist et al., 2019] and yield linear convergence rates [Cen et al., 2020].

References

Jeffrey L Adler and Victor J Blue. A cooperative multi-agent transportation management and route guidance system. *Transportation Research Part C: Emerging Technologies*, pages 433–454, 2002.

- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. Optimality and approximation with policy gradient methods in markov decision processes. *Conference on Learning Theory*, pages 64–66, 2020.
- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10: 251–276, 1998.
- Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor–critic algorithms. *Automatica*, pages 2471–2482, 2009.
- Marco Gaboardi Borja Balle, Gilles Barthe and Joseph Geumlek. Privacy amplification by mixing and diffusion mechanisms. In *Advances in Neural Information Processing Systems*, 2019.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *arXiv preprint arXiv:1405.4980*, 2014.
- Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *arXiv preprint arXiv:1912.02906v2*, 2020.
- Deepayan Chakrabarti, Yang Wang, Chenxi Wang, Jurij Leskovec, and Christos Faloutsos. Epidemic thresholds in real networks. *ACM Transactions on Information and System Security*, pages 1–26, 2008.
- Yuval Dagan, Constantinos Daskalakis, Nishanth Dikkala, and Siddhartha Jayanti. Learning from weakly dependent data under dobrushin’s condition. In *Conference on Learning Theory*, pages 914–928. PMLR, 2019. ISBN 2640-3498.
- RL Dobrusin. Definition of a system of random variables by means of conditional distributions. *Theory of Probability and its Applications*, pages 458–486, 1970.
- Martin Dyer, Leslie Ann Goldberg, and Mark Jerrum. *Dobrushin conditions and systematic scan*, pages 327–338. 2006.
- Hans Föllmer. A covariance estimate for gibbs measures. *Journal of Functional Analysis*, pages 387–395, 1982.
- David Gamarnik. *Correlation decay method for decision, optimization, and inference in large-scale networks*, pages 108–121. 2013.
- David Gamarnik, David A Goldberg, and Theophane Weber. Correlation decay in random decision networks. *Mathematics of Operations Research*, pages 229–261, 2014.
- Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A theory of regularized markov decision processes. In *International Conference on Machine Learning*, pages 2160–2169, 2019.

- Hans-Otto Georgii. *Gibbs measures and phase transitions*. 2011.
- Alison L Gibbs and Francis Edward Su. On choosing and bounding probability metrics. *International Statistical Review*, pages 419–435, 2002.
- Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. In *Advances in Neural Information Processing Systems*, pages 1523–1530, 2001a.
- Carlos Guestrin, Daphne Koller, and Ronald Parr. Max-norm projections for factored mdps. In *International Joint Conference on Artificial Intelligence*, pages 673–680, 2001b.
- Carlos Guestrin, Michail Lagoudakis, and Ronald Parr. Coordinated reinforcement learning. In *International Conference on Machine Learning*, pages 227–234, 2002.
- Xiaoli Wei Renyuan Xu Haotian Gu, Xin Guo. Mean-field multi-agent reinforcement learning: A decentralized network approach. *arXiv preprint arXiv: 2108.02731*, 2021.
- Thomas P Hayes. A simple condition implying rapid mixing of single-site dynamics on spin systems. In *Annual IEEE Symposium on Foundations of Computer Science*, pages 39–46, 2006.
- Alistair Stewart Ilias Diakonikolas, Daniel M. Kane and Yuxin Sun. Outlier-robust learning of ising models under dobrushins condition. In *Conference on Learning Theory*, 2021.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In Abernethy Jacob and Agarwal Shivani, editors, *Conference on Learning Theory*, pages 2137–2143, 2020.
- S Kakade. A natural policy gradient. *Advances in Neural Information Processing Systems*, 2002.
- Daphne Koller and Ron Parr. Policy iteration for factored mdps. *arXiv preprint arXiv:1301.3869*, 2013.
- Daphne Koller and Ronald Parr. Computing factored value functions for policies in structured mdps. In *International Joint Conference on Artificial Intelligence*, pages 1332–1339, 1999.
- Yiheng Lin, Guannan Qu, Longbo Huang, and Adam Wierman. Distributed reinforcement learning in multi-agent networked systems. *arXiv preprint arXiv:2006.06555*, 2020.
- Wenjun Mei, Shadi Mohagheghi, Sandro Zampieri, and Francesco Bullo. On the dynamics of deterministic epidemic propagation over networks. *Annual Reviews in Control*, pages 116–128, 2017.
- Ioannis Mitliagkas and Lester Mackey. Improving gibbs sampler scan quality with dogs. In *International Conference on Machine Learning*, pages 2469–2477. PMLR, 2017. ISBN 2640-3498.

- Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, pages 48–61, 2009.
- Christos H Papadimitriou and John N Tsitsiklis. The complexity of optimal queueing network control. In *IEEE Conference on Structure in Complexity Theory*, pages 318–322, 1994.
- Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, pages 1180–1190, 2008.
- Adarsh Prasad, Vishwak Srinivasan, Sivaraman Balakrishnan, and Pradeep Ravikumar. On learning ising models under huber’s contamination model. *Advances in Neural Information Processing Systems*, 33, 2020.
- Guannan Qu and Na Li. Exploiting fast decaying and locality in multi-agent mdp with tree dependence structure. In *IEEE Conference on Decision and Control*, 2019.
- Guannan Qu, Yiheng Lin, Adam Wierman, and Na Li. Scalable multi-agent reinforcement learning for networked systems with average reward. *Advances in Neural Information Processing Systems*, pages 2074–2086, 2020a.
- Guannan Qu, Adam Wierman, and Na Li. Scalable reinforcement learning of localized policies for multi-agent networked systems. In *Conference on Learning for Dynamics and Control*, 2020b.
- Michael Rabbat and Robert Nowak. Distributed optimization in sensor networks. In *International Symposium on Information Processing in Sensor Networks*, pages 20–27, 2004.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 4295–4304, 2018.
- Patrick Rebeschini and Ramon Van Handel. Can local particle filters beat the curse of dimensionality? *Annals of Applied Probability*, pages 2809–2866, 2015.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, pages 484–489, 2016.

- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George Van Den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. pages 354–359, 2017.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, pages 1057–1063, 1999.
- Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. In Chaudhuri Kamalika and Salakhutdinov Ruslan, editors, *International Conference on Machine Learning*, pages 6995–7004, 2019.
- Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field multi-agent reinforcement learning. In *International Conference on Machine Learning*. PMLR, 2018.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Basar. Networked multi-agent reinforcement learning in continuous spaces. In *IEEE Conference on Decision and Control*, pages 2771–2776, 2018a.
- Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*, pages 5872–5881, 2018b.
- Rick Zhang and Marco Pavone. Control of robotic mobility-on-demand systems: a queueing-theoretical perspective. *The International Journal of Robotics Research*, pages 186–203, 2016.
- Yan Zhang and Michael M Zavlanos. Distributed off-policy actor-critic reinforcement learning with policy consensus. In *IEEE Conference on Decision and Control*, pages 4674–4679, 2019.
- Alessandro Zocca. Temporal starvation in multi-channel csma networks: an analytical framework. *Queueing Systems*, pages 241–263, 2019.

A Policy Class Example

Let $\tilde{r} = \max_{k, k' \in \mathcal{K}} d(k, k')$ be the maximum distance between two agents. Define a set of parameterized differentiable functions $\{f_{(\theta_k)_r} : \mathcal{S}_{N_k^r} \times \mathcal{A}_k \rightarrow \mathcal{C} \mid 0 \leq r \leq \tilde{r}\}$, where $\mathcal{C} \subset [-C, C]$ and $C > 0$, a set of parameters $\{\alpha_r \geq 0 \mid 0 \leq r \leq \tilde{r}\}$ and let, for each agent k ,

$$f_{\theta_k}(s, a_k) = \sum_{r=0}^{\tilde{r}} \alpha_r f_{(\theta_k)_r}(s_{N_k^r}, a_k),$$

$$\pi_{\theta_k}(a_k | s) = \frac{\exp(f_{\theta_k}(s, a_k))}{\sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(s, a'))}.$$

By tuning the parameters α_r , we can make any policy belonging to this policy class respect Assumptions 2 and 3, as we show in the following. Let $r \in \{0, \dots, \tilde{r}\}$, let $s, \tilde{s} \in \mathcal{S}$ be such that $s_{N_k^r} = \tilde{s}_{N_k^r}$, then

$$\begin{aligned} TV(\pi_{\theta_k}(\cdot | s), \pi_{\theta_k}(\cdot | \tilde{s})) &= \frac{1}{2} \sum_{a \in \mathcal{A}_k} |\pi_{\theta_k}(a | s) - \pi_{\theta_k}(a | \tilde{s})| \\ &= \frac{1}{2} \sum_{a \in \mathcal{A}_k} \left| \frac{\exp(f_{\theta_k}(s, a))}{\sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(s, a'))} - \frac{\exp(f_{\theta_k}(\tilde{s}, a))}{\sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(\tilde{s}, a'))} \right| \\ &= \frac{\sum_{a \in \mathcal{A}_k} \left| \sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(s, a)) \exp(f_{\theta_k}(\tilde{s}, a')) - \exp(f_{\theta_k}(\tilde{s}, a)) \exp(f_{\theta_k}(s, a')) \right|}{2 \sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(\tilde{s}, a')) \sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(s, a'))} \\ &\leq \frac{\sum_{a \in \mathcal{A}_k} \sum_{a' \in \mathcal{A}_k} |\exp(f_{\theta_k}(s, a)) \exp(f_{\theta_k}(\tilde{s}, a')) - \exp(f_{\theta_k}(\tilde{s}, a)) \exp(f_{\theta_k}(s, a'))|}{2 \sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(\tilde{s}, a')) \sum_{a' \in \mathcal{A}_k} \exp(f_{\theta_k}(s, a'))} \\ &\leq \frac{\sum_{a \in \mathcal{A}_k} |\exp(f_{\theta_k}(\tilde{s}, a)) - \exp(f_{\theta_k}(s, a))|}{\sum_{a \in \mathcal{A}_k} \exp(f_{\theta_k}(\tilde{s}, a))} \\ &\leq \frac{\sum_{a \in \mathcal{A}_k} |f_{\theta_k}(\tilde{s}, a) - f_{\theta_k}(s, a)| \exp(\sup_{s' \in \{s, \tilde{s}\}} f_{\theta_k}(s', a))}{\sum_{a \in \mathcal{A}_k} \exp(f_{\theta_k}(\tilde{s}, a))} \\ &\leq e^{2C(\tilde{r}-r)} \frac{\sum_{a \in \mathcal{A}_k} |f_{\theta_k}(\tilde{s}, a) - f_{\theta_k}(s, a)| \exp(f_{\theta_k}(\tilde{s}, a))}{\sum_{a \in \mathcal{A}_k} \exp(f_{\theta_k}(\tilde{s}, a))} \\ &= e^{2C(\tilde{r}-r)} \mathbb{E}_{\pi_{\theta_k}} \left| \sum_{r'=r+1}^{\tilde{r}} \alpha_{r'} \left(f_{(\theta_k)_{r'}}(\tilde{s}_{N_k^{r'}}, a) - f_{(\theta_k)_{r'}}(s_{N_k^{r'}}, a) \right) \right| \end{aligned}$$

$$\begin{aligned}
&\leq e^{2C(\tilde{r}-r)} \sum_{r'=r+1}^{\tilde{r}} \alpha_{r'} \mathbb{E}_{\pi_{\theta_k}} \left| f_{(\theta_k)_{r'}}(\tilde{s}_{N_k^{r'}}, a) - f_{(\theta_k)_{r'}}(s_{N_k^{r'}}, a) \right| \\
&\leq 2C e^{2C(\tilde{r}-r)} \sum_{r'=r+1}^{\tilde{r}} \alpha_{r'}.
\end{aligned}$$

Setting the parameters $\{\alpha_{r'}\}_{r' \in \{r+1, \dots, \tilde{r}\}}$ small enough ensures that the policy respects Assumption 2. Similarly, Assumption 3 is satisfied for a value r of the range parameter if $\alpha_{r'} = 0 \forall r' \in \{r+1, \dots, \tilde{r}\}$.

B Proof of Proposition 6

B.1 Preliminary Lemmas

To prove Proposition 6, we need a series of intermediate results, which we state and prove for completeness. Results similar to Lemmas 11 and 12 can be found in Chapter 8 of Georgii [2011], Lemma 13 is an extension of results from Qu et al. [2020a].

Lemma 11 *Let $f : \mathcal{Z} \rightarrow [m, M]$, where $\mathcal{Z} = \prod_{k \in \mathcal{K}} \mathcal{Z}_k$ and $m, M \in \mathbb{R}$. For every $k \in \mathcal{K}$, let μ_k and ν_k be two distributions on \mathcal{Z}_k . Let μ and ν be the respective product distributions. Let $\delta_k(f(z)) = \sup_{z_k, z_{-k}, z'_k} |f(z_k, z_{-k}) - f(z'_k, z_{-k})|$. Then:*

$$|\mathbb{E}_{z \sim \mu} f(z) - \mathbb{E}_{z \sim \nu} f(z)| \leq \sum_{k \in \mathcal{K}} TV(\mu_k, \nu_k) \delta_k(f).$$

Proof: We prove Lemma 11 by induction. Note that

$$TV(\mu, \nu) = \frac{1}{2} \max_{|h| \leq 1} |\mathbb{E}_\mu(h) - \mathbb{E}_\nu(h)|$$

is an equivalent formulation of the total variation distance [Gibbs and Su, 2002]. For $|\mathcal{K}| = 1$, we have that

$$\begin{aligned}
|\mathbb{E}_{\mu_1}(f) - \mathbb{E}_{\nu_1}(f)| &= \left| \mathbb{E}_{\mu_1} \left(f - \frac{M+m}{2} \right) - \mathbb{E}_{\nu_1} \left(f - \frac{M+m}{2} \right) \right| \\
&= \frac{M-m}{2} \left| \mathbb{E}_{\mu_1} \left(\frac{2f}{M-m} - \frac{M+m}{M-m} \right) - \mathbb{E}_{\nu_1} \left(\frac{2f}{M-m} - \frac{M+m}{M-m} \right) \right|
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{M-m}{2} \max_{|h| \leq 1} |\mathbb{E}_{\mu_1}(h) - \mathbb{E}_{\nu_1}(h)| \\
&= TV(\mu_1, \nu_1) \delta_1(f).
\end{aligned}$$

As induction assumption, assume that Lemma 11 holds for $|\mathcal{K}| - 1$. Then:

$$\begin{aligned}
|\mathbb{E}_{z \sim \mu} f(z) - \mathbb{E}_{z \sim \nu} f(z)| &= |\mathbb{E}_{z_1 \sim \mu_1} \mathbb{E}_{z_{2:n} \sim \mu_{2:n}} f(z) - \mathbb{E}_{z_1 \sim \nu_1} \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z)| \\
&\leq |\mathbb{E}_{z_1 \sim \mu_1} \mathbb{E}_{z_{2:n} \sim \mu_{2:n}} f(z) - \mathbb{E}_{z_1 \sim \mu_1} \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z)| \\
&\quad + |\mathbb{E}_{z_1 \sim \mu_1} \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z) - \mathbb{E}_{z_1 \sim \nu_1} \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z)| \\
&\leq \mathbb{E}_{z_1 \sim \mu_1} |\mathbb{E}_{z_{2:n} \sim \mu_{2:n}} f(z) - \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z)| \\
&\quad + \left| \mathbb{E}_{z_1 \sim \mu_1} \tilde{f}(z_1) - \mathbb{E}_{z_1 \sim \nu_1} \tilde{f}(z_1) \right|.
\end{aligned}$$

where $\tilde{f}(z_1) = \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z)$. By induction assumption:

$$|\mathbb{E}_{z_{2:n} \sim \mu_{2:n}} f(z) - \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z)| \leq \sum_{k \neq 1 \in \mathcal{K}} TV(\mu_k, \nu_k) \delta_k(f(z_1, \cdot)) \leq \sum_{k \neq 1 \in \mathcal{K}} TV(\mu_k, \nu_k) \delta_k(f).$$

Since

$$\begin{aligned}
\delta_1(\tilde{f}) &= \sup_{z_1, z'_1} |\mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z_1, z_{2:n}) - \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} f(z'_1, z_{2:n})| \\
&\leq \sup_{z_1, z'_1} \mathbb{E}_{z_{2:n} \sim \nu_{2:n}} |f(z_1, z_{2:n}) - f(z'_1, z_{2:n})| \\
&\leq \sup_{z_1, z'_1} \sup_{z_{2:n}} |f(z_1, z_{2:n}) - f(z'_1, z_{2:n})| = \delta_1(f),
\end{aligned}$$

we have

$$\begin{aligned}
|\mathbb{E}_{z \sim \mu} f(z) - \mathbb{E}_{z \sim \nu} f(z)| &\leq \mathbb{E}_{z_1 \sim \mu_1} \sum_{k \neq 1 \in \mathcal{K}} TV(\mu_k, \nu_k) \delta_k(f) + TV(\mu_1, \nu_1) \delta_1(f) \\
&\leq \sum_{k \in \mathcal{K}} TV(\mu_k, \nu_k) \delta_k(f),
\end{aligned}$$

which concludes the induction. ■

Lemma 12 Consider a Markov Chain with state $z \in \mathcal{Z}$, where $\mathcal{Z} = \prod_{k \in \mathcal{K}} \mathcal{Z}_k$ and \mathcal{K} is defined as in Section 2. Suppose its transition probability factorizes as

$$P(z(t+1)|z(t)) = \prod_{k \in \mathcal{K}} P_k(z_k(t+1)|z(t)).$$

Let $C \in \mathbb{R}^{K \times K}$ be a matrix whose elements respect the condition

$$C_{ij} \geq \sup_{z_j, z_{-j}, z'_j} TV(P_i(\cdot|z_j, z_{-j}), P_i(\cdot|z'_j, z_{-j})).$$

If $\sum_{j \in \mathcal{K}} e^{\beta d(j,k)} C_{kj} \leq \rho$, then, $\forall \mathcal{J} \subseteq \mathcal{K}$,

$$\sup_{z_j, z_{-j}, z'_j} TV(P_i(\cdot|z_{\mathcal{J}}, z_{-\mathcal{J}}), P_i(\cdot|z'_{\mathcal{J}}, z_{-\mathcal{J}})) \leq \sum_{j \in \mathcal{J}} C_{ij}$$

and

$$\sup_{z_{\mathcal{J}}, z_{-\mathcal{J}}, z'_{\mathcal{J}}} TV(P_i(\cdot|z_{\mathcal{J}}, z_{-\mathcal{J}}), P_i(\cdot|z'_{\mathcal{J}}, z_{-\mathcal{J}})) \leq \rho e^{-\beta d(\mathcal{J}, i)},$$

where $d(\mathcal{J}, i) = \min_{j \in \mathcal{J}} d(j, i)$.

Proof: We prove the first claim of Lemma 12 by induction. The first claim clearly holds if $|\mathcal{J}| = 1$. As induction assumption, assume that the first claim holds for a generic \mathcal{J} . Then, it holds for $\mathcal{J}' = \mathcal{J} + \{k\}$:

$$\begin{aligned} & \sup_{z_j, z_{-j}, z'_j} TV(P_i(\cdot|z_{\mathcal{J}'}, z_{-\mathcal{J}'}) , P_i(\cdot|z'_{\mathcal{J}'}, z_{-\mathcal{J}'})) = \sup_{\substack{A \subseteq \mathcal{Z}_i \\ z_j, z_{-j}, z'_j}} |P_i(A|z_{\mathcal{J}'}, z_{-\mathcal{J}'}) - P_i(A|z'_{\mathcal{J}'}, z_{-\mathcal{J}'})| \\ & \leq \sup_{\substack{A \subseteq \mathcal{Z}_i \\ z_j, z_{-j}, z'_j}} |P_i(A|z_{\mathcal{J}'}, z_{-\mathcal{J}'}) - P_i(A|z'_{\mathcal{J}}, z_{-\mathcal{J}})| + \sup_{\substack{A \subseteq \mathcal{Z}_i \\ z_j, z_{-j}, z'_j}} |P_i(A|z'_{\mathcal{J}}, z_{-\mathcal{J}}) - P_i(A|z'_{\mathcal{J}'}, z_{-\mathcal{J}'})| \\ & \leq \sum_{j \in \mathcal{J}} C_{ij} + C_{ik} = \sum_{j \in \mathcal{J}'} C_{ij}. \end{aligned}$$

The second claim follows immediately, since

$$e^{\beta d(\mathcal{J}, i)} \sum_{j \in \mathcal{J}} C_{ij} \leq \sum_{j \in \mathcal{J}} e^{\beta d(j, i)} C_{ij} \leq \sum_{j \in \mathcal{K}} e^{\beta d(j, i)} C_{ij} \leq \rho,$$

and

$$\sum_{j \in \mathcal{J}} C_{ij} \leq \rho e^{-\beta d(\mathcal{J}, i)}.$$

■

Lemma 13 Consider the setting of Lemma 12. For a generic value of r , denote by d_t and \tilde{d}_t the distribution of $z(t)$ with starting state, respectively, $z = (z_{N_k^r}, z_{N_{-k}^r})$ and $\tilde{z} = (z_{N_k^r}, \tilde{z}_{N_{-k}^r})$. Then, if $\sum_{j \in \mathcal{K}} e^{\beta d(j,k)} C_{kj} \leq \rho$, we have that $TV(d_{t,k}, \tilde{d}_{t,k}) \leq \rho^t e^{-\beta r}$, $\forall k \in \mathcal{K}$.

Proof: We prove Lemma 13 by induction. The case where $t = 1$ follows from Lemma 12. As induction assumption, assume that Lemma 13 holds for t . Then,

$$\begin{aligned}
& \left| \mathbb{E}_{s \sim d_{t+1,k}(s)} \mathbf{1}_A(s) - \mathbb{E}_{s \sim \tilde{d}_{t+1,k}} \mathbf{1}_A(s) \right| \\
&= \left| \mathbb{E}_{z \sim d_t} E_{s \sim P_k(\cdot|z)} \mathbf{1}_A(s) - \mathbb{E}_{z \sim \tilde{d}_t} E_{s \sim P_k(\cdot|z)} \mathbf{1}_A(s) \right| \\
&\leq \sum_{j \in \mathcal{K}} TV(d_{t,j}, \tilde{d}_{t,j}) \delta_j(E_{s \sim P_k(\cdot|z)} \mathbf{1}_A(s)) \\
&\leq \sum_{j \in \mathcal{K}} TV(d_{t,j}, \tilde{d}_{t,j}) C_{kj} \\
&\leq \sum_{j \in \mathcal{K}} \rho^t e^{-\beta(r-d(j,k))} C_{kj} \\
&= \rho^t e^{-\beta r} \sum_{j \in \mathcal{K}} e^{\beta d(j,k)} C_{kj} \leq \rho^{t+1} e^{-\beta r},
\end{aligned}$$

where we used Lemma 11 in the first inequality. ■

Lemma 14 Consider the setting of Lemma 12. Let $P^t(z'|z) = P(z(t) = z' | z(0) = z)$ and

$$\delta_i P_k^t = \sup_{z_i, z_{-i}, z'_i} TV(P_k^t(\cdot | z_i, z_{-i}), P_k^t(\cdot | z'_i, z_{-i})).$$

If $\sum_{j \in \mathcal{K}} e^{\beta d(j,k)} C_{kj} \leq \rho$, then $\forall k \in \mathcal{K}$

$$\sum_{i \in \mathcal{K}} e^{\beta d(k,i)} \delta_i P_k^t \leq \rho^t.$$

Proof: We prove Lemma 14 by induction. The claim holds for $t = 1$:

$$\sum_{i \in \mathcal{K}} e^{\beta d(k,i)} \delta_i P_k^t = \sum_{i \in \mathcal{K}} e^{\beta d(k,i)} C_{k,i} \leq \rho$$

As induction assumption, we assume that the claim holds for t . Then, using Lemma 11,

$$\begin{aligned}
\delta_i P_k^{t+1} &= \sup_{\substack{A \subseteq \mathcal{S}_k \\ z_i, z_{-i}, z'_i}} \left| \mathbb{E}_{s \sim P_k^{t+1}(\cdot | z_i, z_{-i})} \mathbf{1}_A(s) - \mathbb{E}_{s \sim P_k^{t+1}(\cdot | z'_i, z_{-i})} \mathbf{1}_A(s) \right| \\
&= \sup_{\substack{A \subseteq \mathcal{S}_k \\ z_i, z_{-i}, z'_i}} \left| \mathbb{E}_{x \sim P^t(\cdot | z_i, z_{-i})} \mathbb{E}_{s \sim P_k(\cdot | x)} \mathbf{1}_A(s) - \mathbb{E}_{x \sim P^t(\cdot | z'_i, z_{-i})} \mathbb{E}_{s \sim P_k(\cdot | x)} \mathbf{1}_A(s) \right| \\
&\leq \sup_{z_i, z_{-i}, z'_i} \sum_{j \in \mathcal{K}} TV(P_j^t(\cdot | z_i, z_{-i}), P_j^t(\cdot | z'_i, z_{-i})) \delta_j (E_{s \sim P_k(\cdot | \cdot)} \mathbf{1}_A(s)) \\
&\leq \sum_{j \in \mathcal{K}} \delta_i P_j^t C_{kj}
\end{aligned}$$

and, using the inverse triangle inequality,

$$\begin{aligned}
\sum_{i \in \mathcal{K}} e^{\beta d(k,i)} \delta_i P_k^{t+1} &\leq \sum_{i \in \mathcal{K}} e^{\beta d(k,i)} \sum_{j \in \mathcal{K}} \delta_i P_j^t C_{kj} \\
&\leq \sum_{j \in \mathcal{K}} e^{\beta d(k,j)} C_{kj} \sum_{i \in \mathcal{K}} e^{\beta(d(k,i) - d(k,j))} \delta_i P_j^t \\
&\leq \sum_{j \in \mathcal{K}} e^{\beta d(k,j)} C_{kj} \sum_{i \in \mathcal{K}} e^{\beta d(j,i)} \delta_i P_j^t \leq \rho^{t+1},
\end{aligned}$$

which concludes the induction. ■

B.2 Main Result

Proof:[of Proposition 6] The following holds for every $k \in \mathcal{K}$. Let $s, \tilde{s} \in \mathcal{S}$, $a, \tilde{a} \in \mathcal{A}$ be such that $s_{N_k^r} = \tilde{s}_{N_k^r}$ and $a_{N_k^r} = \tilde{a}_{N_k^r}$. Notice that

$$\begin{aligned}
&|Q_k^\pi(s, a) - Q_k^\pi(\tilde{s}, \tilde{a})| \\
&\leq \sum_{t=0}^{\infty} \gamma^t \left| \mathbb{E} [r_k(s_k(t), a_k(t)) | \pi, s(0) = s, a(0) = a] - \mathbb{E} [r_k(s_k(t), a_k(t)) | \pi, s(0) = \tilde{s}, a(0) = \tilde{a}] \right| \\
&\leq \sum_{t=1}^{\infty} \gamma^t \left| \mathbb{E} [r_k(s_k(t), a_k(t)) | \pi, s(0) = s, a(0) = a] - \mathbb{E} [r_k(s_k(t), a_k(t)) | \pi, s(0) = \tilde{s}, a(0) = \tilde{a}] \right| \\
&\leq \sum_{t=1}^{\infty} \gamma^t TV(d_{t,k}, \tilde{d}_{t,k}),
\end{aligned}$$

where $d_{t,k}(s_k, a_k)$ and $\tilde{d}_{t,k}(s_k, a_k)$ are the distributions of s_k, a_k at time t with starting point (s, a) and (\tilde{s}, \tilde{a}) , respectively. We use Lemma 13 to bound $TV(d_{t,k}, \tilde{d}_{t,k})$. The structure of our MDP implies that:

$$P(s(t+1), a(t+1)|s(t), a(t)) = \prod_{k \in \mathcal{K}} \pi^k(a_k(t+1)|s(t+1)) P_k(s_k(t+1)|s(t), a(t)).$$

Let C be defined as in Assumption 1 and note that

$$\begin{aligned} C_{kj} &= \sup_{s_j, s_{-j}, a_j, a_{-j}, s'_j, a'_j} TV(P_k(\cdot | s_j, s_{-j}, a_j, a_{-j}), P_k(\cdot | s'_j, s_{-j}, a'_j, a_{-j})) \\ &\geq \sup_{s_j, s_{-j}, a_j, a_{-j}, s'_j, a'_j} TV(P_k(\cdot, \cdot | s_j, s_{-j}, a_j, a_{-j}), P_k(\cdot, \cdot | s'_j, s_{-j}, a'_j, a_{-j})). \end{aligned}$$

Then, if Assumption 1 holds, the requirements of Lemma 13 are satisfied. Therefore, $TV(d_{t,k}, \tilde{d}_{t,k}) \leq \rho^t e^{-\beta r}$ and

$$|Q_k^\pi(s, a) - Q_k^\pi(\tilde{s}, \tilde{a})| \leq \sum_{t=1}^{\infty} \gamma^t TV(d_{t,k}, \tilde{d}_{t,k}) \leq e^{-\beta r} \sum_{t=1}^{\infty} \gamma^t \rho^t = \frac{\gamma \rho e^{-\beta r}}{1 - \gamma \rho}.$$

We use this result to prove the exponential decay property for the value function. Let

$$\delta_j Q_k^\pi(s, a) = \sup_{s_j, s_{-j}, a_j, a_{-j}, s'_j, a'_j} |Q_k^\pi(s_j, s_{-j}, a_j, a_{-j}) - Q_k^\pi(s'_j, s_{-j}, a'_j, a_{-j})|.$$

Using Lemma 11 and Assumption 2, we have that

$$\begin{aligned} |V_k^\pi(s) - V_k^\pi(\tilde{s})| &= |\mathbb{E}_{a \sim \pi(\cdot | s)} Q_k^\pi(s, a) - \mathbb{E}_{a \sim \pi(\cdot | \tilde{s})} Q_k^\pi(\tilde{s}, a)| \\ &\leq |\mathbb{E}_{a \sim \pi(\cdot | s)} Q_k^\pi(s, a) - \mathbb{E}_{a \sim \pi(\cdot | \tilde{s})} Q_k^\pi(s, a)| + |\mathbb{E}_{a \sim \pi(\cdot | \tilde{s})} Q_k^\pi(s, a) - \mathbb{E}_{a \sim \pi(\cdot | \tilde{s})} Q_k^\pi(\tilde{s}, a)| \\ &\leq \sum_{i \in \mathcal{K}} TV(\pi_i(\cdot | s), \pi_i(\cdot | \tilde{s})) \delta_i Q_k^\pi(s, a) + \frac{\gamma \rho e^{-\beta r}}{1 - \gamma \rho} \\ &\leq \xi e^{-\beta r} \sum_{i \in \mathcal{K}} e^{-\beta d(i,k)} \delta_i Q_k^\pi(s, a) + \frac{\gamma \rho e^{-\beta r}}{1 - \gamma \rho}. \end{aligned}$$

We have already shown that the MDP satisfies the condition of Lemma 14. Using it we obtain

$$\sum_{i \in \mathcal{K}} e^{\beta d(k,i)} \delta_i (Q_k^\pi(s, \cdot)) \leq \sum_{t=1}^{\infty} \gamma^t \sum_{i \in \mathcal{K}} e^{\beta d(k,i)} \delta_i P_k^t \leq \sum_{t=1}^{\infty} \gamma^t \rho^t = \frac{\gamma \rho}{1 - \gamma \rho},$$

where

$$\delta_i P_k^t = \sup_{s_j, s_{-j}, a_j, a_{-j}, s'_j, a'_j} TV(P_k^t(\cdot, \cdot | s_j, s_{-j}, a_j, a_{-j}), P_k^t(\cdot, \cdot | s'_j, s_{-j}, a'_j, a_{-j})).$$

Then, we have that

$$|V_k^\pi(s) - V_k^\pi(\tilde{s})| \leq \frac{\gamma\rho(1+\xi)e^{-\beta r}}{1-\gamma\rho}.$$

■

C Decay for the Optimal Policy

Lemma 15 *Assume that the exponential decay property holds for the Q -function with parameters (c, ψ) . Then the exponential decay property holds also for the value function associated with the optimal policy, with parameters $(3c, \psi)$.*

Proof: The following holds for every $k \in \mathcal{K}$. Let $s, \tilde{s} \in \mathcal{S}$ be such that $s_{N_k^r} = \tilde{s}_{N_k^r}$ and let $a_{N_{-k}^r} \in \mathcal{A}_{N_{-k}^r}$.

$$\begin{aligned} |V_k^*(s) - V_k^*(\tilde{s})| &= \left| \mathbb{E}_{a \sim \pi^*(\cdot|s)} Q_k^*(s, a) - \mathbb{E}_{a \sim \pi^*(\cdot|\tilde{s})} Q_k^*(\tilde{s}, a) \right| = \left| \max_a Q_k^*(s, a) - \max_a Q_k^*(\tilde{s}, a) \right| \\ &= \left| \max_a Q_k^*(s, a) - \max_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r}) + \max_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r}) \right. \\ &\quad \left. - \max_a Q_k^*(\tilde{s}, a) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) + \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) \right| \\ &= \left| \max_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r}) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) \right| \\ &\quad + \left| \max_a Q_k^*(\tilde{s}, a) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) \right| + \left| \max_a Q_k^*(s, a) - \max_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r}) \right| \\ &\leq \left| \max_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r}) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) \right| + 2c\psi^{r+1}. \end{aligned}$$

Let $a'_{N_k^r} \in \operatorname{argmax}_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r})$, then

$$\begin{aligned} &Q_k^*(s, a'_{N_k^r}, a_{N_{-k}^r}) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) \\ &\leq Q_k^*(\tilde{s}, a'_{N_k^r}, a_{N_{-k}^r}) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) + c\psi^{r+1} \end{aligned}$$

$$\leq \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) - \max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) + c\psi^{r+1} = c\psi^{r+1}.$$

The same holds for $\max_{a_{N_k^r}} Q_k^*(\tilde{s}, a_{N_k^r}, a_{N_{-k}^r}) - \max_{a_{N_k^r}} Q_k^*(s, a_{N_k^r}, a_{N_{-k}^r})$. The lemma follows immediately. \blacksquare

We make an assumption on the minimum influence that the action of an agent has on its expected future rewards. Assumption 16 and Proposition 17 hold for any $k \in \mathcal{K}$.

Assumption 16 Let $\mathcal{A}^* = \operatorname{argmax}_{a \in \mathcal{A}} Q_k^*(s, a)$, $\forall s \in \mathcal{S}$. Assume that, if \tilde{a} is such that $\tilde{a}_k \notin \mathcal{A}_k^*$, then

$$|Q_k^*(s, a) - Q_k^*(s, \tilde{a})| \geq R$$

Proposition 17 Assume that the exponential decay property holds for the Q -function with parameters (c, ψ) and that Assumption 16 holds. Let $s, \tilde{s} \in \mathcal{S}$ be such that $s_{N_k^r} = \tilde{s}_{N_k^r}$. Let $\mathcal{A}^* = \operatorname{argmax}_{a \in \mathcal{A}} Q_k^*(s, a)$ and $\tilde{a} \in \operatorname{argmax}_{a \in \mathcal{A}} Q_k^*(\tilde{s}, a)$. If $r > \log_\psi R/4c$, then $\tilde{a}_k \in \mathcal{A}_k^*$.

Proof: We prove this by contradiction. Lemma 15 shows that, $\forall r > 0$, if $s, \tilde{s} \in \mathcal{S}$ are such that $s_{N_k^r} = \tilde{s}_{N_k^r}$,

$$|V_k^*(s) - V_k^*(\tilde{s})| = \left| \max_a Q_k^*(s, a) - \max_a Q_k^*(\tilde{s}, a) \right| \leq 3c\psi^{r+1}.$$

Let $a \in \operatorname{argmax}_{a \in \mathcal{A}} Q_k^*(s, a)$ and $\tilde{a} \in \operatorname{argmax}_{a \in \mathcal{A}} Q_k^*(\tilde{s}, a)$. Let $\mathcal{A}^* = \operatorname{argmax}_{a \in \mathcal{A}} Q_k^*(s, a)$ and assume that $\tilde{a}_k \notin \mathcal{A}_k^*$. Then

$$\begin{aligned} |Q_k^*(s, a) - Q_k^*(\tilde{s}, \tilde{a})| &= |Q_k^*(s, a) - Q_k^*(s, \tilde{a}) + Q_k^*(s, \tilde{a}) - Q_k^*(\tilde{s}, \tilde{a})| \\ &\geq ||Q_k^*(s, a) - Q_k^*(s, \tilde{a})| - |Q_k^*(s, \tilde{a}) - Q_k^*(\tilde{s}, \tilde{a})|| \\ &\geq |Q_k^*(s, a) - Q_k^*(s, \tilde{a})| - c\psi^{r+1} \geq R - c\psi^{r+1} \end{aligned}$$

where we used Assumption 16 in the last passage. Then, due to Lemma 15, if $r > \log_\psi R/4c$ we have a contradiction. \blacksquare

Proposition 17 shows that the optimal policy of an agent is not influenced by distant agents. Assumption 2 ensures that the policy class we consider respects this condition.

D Independent Agents

By completely independent agents we mean agents whose transition dynamics are independent and whose policy is defined, for $s \in \mathcal{S}, a \in \mathcal{A}$, as:

$$\pi_\theta(a|s) = \prod_{k \in \mathcal{K}} \pi_{\theta_k}(a_k|s_k) = \prod_k \frac{e^{f_{\theta_k}(s_k, a_k)}}{\sum_{a' \in \mathcal{A}_k} e^{f_{\theta_k}(s_k, a')}}.$$

In accordance to previous assumptions, we also assume that $\pi_{\theta_k}(a_k|s_k)$ is δ -smooth.

Proof:[of Remark 8] Firstly we show that the two applications of the algorithm coincide. Let

$$L(w) = \mathbb{E}_{s \sim d_{\bar{v}}, a \sim \pi_{\theta}(\cdot|s)} \left[(A^{\pi_{\theta}}(s, a) - w \cdot \nabla_{\theta} \log \pi_{\theta}(a|s))^2 \right].$$

In Agarwal et al. [2020], the NPG update is

$$w_{\star} \in \underset{w}{\operatorname{argmin}} L(w).$$

We now show that the gradient of the loss function is the same as the one in the single agent setting. This is enough to show that the two algorithms coincide, since every other operation coincides. The only exception is the projection step, problem that can be side-stepped by projecting each single component of the gradient vector instead of the whole vector. For each agent k we have that

$$\begin{aligned} \nabla_{\theta_k} L(w) &= \mathbb{E}_{s \sim d_{\bar{v}}, a \sim \pi_{\theta}(\cdot|s)} \left[(A^{\pi_{\theta}}(s, a) - w \cdot \nabla_{\theta} \log \pi_{\theta}(a|s)) \nabla_{\theta_k} \log \pi_{\theta}(a|s) \right] \\ &= \sum_{j \in \mathcal{K}} \mathbb{E}_{s \sim d_{\bar{v}}, a \sim \pi_{\theta}(\cdot|s)} \left[\left(\frac{1}{K} A_j^{\pi_{\theta}}(s_j, a_j) - w_j \cdot \nabla_{\theta_j} \log \pi_{\theta}(a_j|s_j) \right) \nabla_{\theta_k} \log \pi_{\theta}(a|s) \right] \\ &= \mathbb{E}_{s_k \sim d_{\nu, k}^{\pi}, a_k \sim \pi_{\theta_k}^k(\cdot|s_k)} \left[\left(\frac{1}{K} A_k^{\pi_{\theta}}(s_k, a_k) - w_k \cdot \nabla_{\theta_k} \log \pi_{\theta}(a_k|s_k) \right) \nabla_{\theta_k} \log \pi_{\theta}(a_k|s_k) \right], \end{aligned}$$

which corresponds to the gradient for the single agent setting.

With regards to guarantees, since the problem is decoupled, for any agent k we have the same result as in Agarwal et al. 2020:

$$\mathbb{E} \left[\min_{t < T} \{V_k^{\pi^*}(\rho) - V_k^{(t)}(\rho)\} \right] \leq \frac{W}{1 - \gamma} \sqrt{\frac{2\delta \log |\mathcal{A}_k|}{T}} + \sqrt{\frac{\kappa \varepsilon_{\text{stat}, k}}{(1 - \gamma)^3}} + \frac{\sqrt{\varepsilon_{\text{bias}, k}}}{1 - \gamma},$$

where we assumed that

$$\mathbb{E} \left[L_k(w_k^{(t)}, \theta_k^{(t)}, d_k^{(t)}) - L(w_{\star, k}^{(t)}, \theta_k^{(t)}, d_k^{(t)}) | \theta_k^{(t)} \right] \leq \varepsilon_{\text{stat}, k},$$

$$\mathbb{E} \left[L(w_{\star, k}^{(t)}, \theta_k^{(t)}, d_k^{\star}) \right] \leq \varepsilon_{\text{bias}, k},$$

where $d_k^{(t)} = d_{\nu}^{\pi_k^{(t)}}$, $d_k^{\star} = d_{\nu}^{\pi_k^{\star}}$ and

$$L_k(w, \theta_k, d) = \mathbb{E}_{s_k, a_k \sim d} \left[\left(A_k^{\pi_{\theta_k}}(s_k, a_k) - w \cdot \nabla_{\theta_k} \log \pi_{\theta_k}(a_k|s_k) \right)^2 \right],$$

$$w_{\star, k}^{(t)} \in \underset{\|w\|_2 \leq W}{\operatorname{argmin}} L_k(w, \theta_k^{(t)}, d_k^{(t)}).$$

The same result holds for the whole network:

$$\begin{aligned} \mathbb{E} \left[\min_{t < T} \{V^{\pi^*}(\rho) - V^{(t)}(\rho)\} \right] &= \mathbb{E} \left[\min_{t < T} \left\{ \frac{1}{K} \sum_{k \in \mathcal{K}} \left(V_k^{\pi^*}(\rho) - V_k^{(t)}(\rho) \right) \right\} \right] \\ &\leq \mathbb{E} \left[\min_{t < T} \max_{k \in \mathcal{K}} \left\{ V_k^{\pi^*}(\rho) - V_k^{(t)}(\rho) \right\} \right]. \end{aligned}$$

■

E Proof of Theorem 9

We follow the proof in Agarwal et al. [2020] modifying it where necessary. We start by proving a modified NPG Regret Lemma.

Lemma 18 *Consider the setting of Theorem 9, then we have that:*

$$\mathbb{E} \left[\min_{t < T} \{V^{\pi^*}(\rho) - V^{(t)}(\rho)\} \right] \leq \frac{W}{1 - \gamma} \sqrt{\frac{2\delta \max_{k \in \mathcal{K}} \log |\mathcal{A}_k|}{T}} + \mathbb{E} \left[\frac{1}{T(1 - \gamma)} \sum_{t=0}^{T-1} err_t \right],$$

where

$$err_t = \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[A^{(t)}(s, a) - \frac{1}{K} \nabla_{\theta} \log \pi^{(t)}(a|s) \cdot w^{(t)} \right].$$

Proof: We assume $\log \pi_{\theta}(a|s)$ is a δ -smooth function of θ . By smoothness we have:

$$\begin{aligned} \log \frac{\pi^{(t+1)}(a|s)}{\pi^{(t)}(a|s)} &\geq \nabla_{\theta} \log \pi^{(t)}(a|s) \cdot (\theta^{(t+1)} - \theta^{(t)}) - \frac{\delta}{2} \|\theta^{(t+1)} - \theta^{(t)}\|_2^2 \\ &= \eta \nabla_{\theta} \log \pi^{(t)}(a|s) \cdot w^{(t)} - \eta^2 \frac{\delta}{2} \|w^{(t)}\|_2^2. \end{aligned}$$

Then

$$\begin{aligned} \frac{1}{K} \mathbb{E}_{s \sim d_p^*} (\text{KL}(\pi_s^* || \pi_s^{(t)}) - \text{KL}(\pi_s^* || \pi_s^{(t+1)})) &= \frac{1}{K} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\log \frac{\pi^{(t+1)}(a|s)}{\pi^{(t)}(a|s)} \right] \\ &\geq \frac{\eta}{K} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\nabla_{\theta} \log \pi^{(t)}(a|s) \cdot w^{(t)} \right] - \eta^2 \frac{\delta}{2K} \|w^{(t)}\|_2^2 \end{aligned}$$

$$\begin{aligned}
&= \eta \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} [A^{(t)}(s, a)] + \eta \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\frac{1}{K} \nabla_{\theta} \log \pi^{(t)}(a|s) \cdot w^{(t)} - A^{(t)}(s, a) \right] \\
&\quad - \eta^2 \frac{\delta}{2K} \|w^{(t)}\|_2^2 \\
&= (1 - \gamma) \eta (V^{\pi^*}(\rho) - V^{(t)}(\rho)) - \eta^2 \frac{\delta}{2K} \|w^{(t)}\|_2^2 - \eta \text{err}_t,
\end{aligned}$$

where

$$\text{err}_t = \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[A^{(t)}(s, a) - \frac{1}{K} \nabla_{\theta} \log \pi^{(t)}(a|s) \cdot w^{(t)} \right].$$

Rearranging

$$V^{\pi^*}(\rho) - V^{(t)}(\rho) \leq \frac{1}{1 - \gamma} \left(\frac{1}{\eta K} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} (\text{KL}(\pi_s^* || \pi_s^{(t)}) - \text{KL}(\pi_s^* || \pi_s^{(t+1)})) + \frac{\eta \delta}{2} W^2 + \text{err}_t \right)$$

and summing over t

$$\begin{aligned}
V^{\pi^*}(\rho) - \frac{1}{T} \sum_{t=0}^{T-1} V^{(t)}(\rho) &\leq \frac{1}{\eta K T (1 - \gamma)} \sum_{t=0}^{T-1} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} (\text{KL}(\pi_s^* || \pi_s^{(t)}) - \text{KL}(\pi_s^* || \pi_s^{(t+1)})) \\
&\quad + \frac{1}{T(1 - \gamma)} \sum_{t=0}^{T-1} \left(\frac{\eta \delta}{2} W^2 + \text{err}_t \right) \\
&\leq \frac{\mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \text{KL}(\pi_s^* || \pi_s^{(0)})}{\eta K T (1 - \gamma)} + \frac{\eta \delta W^2}{2(1 - \gamma)} + \frac{1}{T(1 - \gamma)} \sum_{t=0}^{T-1} \text{err}_t \\
&\leq \frac{\log |\mathcal{A}|}{\eta K T (1 - \gamma)} + \frac{\eta \delta W^2}{2(1 - \gamma)} + \frac{1}{T(1 - \gamma)} \sum_{t=0}^{T-1} \text{err}_t.
\end{aligned}$$

Optimizing for η , we obtain the lemma. ■

We can now prove Theorem 9

Proof:[of Theorem 9] After using the NPG Regret Lemma we want to bound the err_t term. We do so by dividing it in 3 parts. Let $(w^{(t)})_{t=1, \dots, T}$ be an update sequence such that, for

each t, k , $(w_k^{(t)})_{N_{-k}^r} = 0$, then

$$\begin{aligned}
\text{err}_t &= \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[A^{(t)}(s, a) - \frac{1}{K} \nabla_{\theta} \log \pi^{(t)}(a|s) \cdot w^{(t)} \right] \\
&= \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[A_k^{(t)}(s, a) - \nabla_{\theta_k} \log \pi^{(t)}(a|s) \cdot w_k^{(t)} \right] \\
&= \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[A_k^{(t)}(s, a) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w_k^{(t)})_{N_k^r} \right] \\
&= \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w_{k,\star}^{(t)})_{N_k^r} \right] \\
&\quad + \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[A_k^{(t)}(s, a) - \tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) \right] \\
&\quad + \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot ((w_{k,\star}^{(t)})_{N_k^r} - (w_k^{(t)})_{N_k^r}) \right],
\end{aligned}$$

where $\forall k \in \mathcal{K}$

$$w_{k,\star}^{(t)} \in \underset{\substack{\|w\|_2 \leq W \\ w_{N_{-k}^r} = 0}}{\text{argmin}} \mathbb{E}_{s \sim d^{(t)}, a \sim \pi^{(t)}(\cdot|s)} \left[\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a_k | s_{N_k^r}) \cdot w \right]^2.$$

We now analyse each term separately.

Term 1

$$\begin{aligned}
&\frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w_{k,\star}^{(t)})_{N_k^r} \right] \\
&\leq \frac{1}{K} \sum_{k \in \mathcal{K}} \sqrt{\mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w_{k,\star}^{(t)})_{N_k^r} \right]^2} \leq \sqrt{\varepsilon_{\text{bias}}}
\end{aligned}$$

Term 2 Firstly, we have that $\forall k \in \mathcal{K}$

$$\tilde{Q}_k^\pi(s_{N_k^r}, a_{N_k^r}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_k(s_k(t), a_k(t)) \mid \pi_{N_k^r}, s_{N_k^r}(0) = s_{N_k^r}, a_{N_k^r}(0) = a_{N_k^r} \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_k(s_k(t), a_k(t)) \middle| \pi_{N_k^r}, s_{N_k^r}(0) = s_{N_k^r}, a_{N_k^r}(0) = a_{N_k^r}, \right. \right. \\
&\quad \left. \left. s_{N_{-k}^r}(0), a_{N_{-k}^r}(0) \right] \right] \\
&= \mathbb{E} \left[Q_k^\pi \left(s_{N_k^r}, a_{N_k^r}, s_{N_{-k}^r}(0), a_{N_{-k}^r}(0) \right) \right],
\end{aligned}$$

for some distribution of $s(0)_{N_{-k}^r}$ and $a(0)_{N_{-k}^r}$. Similarly $\forall k \in \mathcal{K}$

$$V_k^\pi(s) - \tilde{V}_k^\pi(s_{N_k^r}) = \mathbb{E} \left[V_k^\pi(s) - V_k^\pi \left(s_{N_k^r}, s(0)_{N_{-k}^r} \right) \right].$$

So

$$A_k^\pi(s, a) - \tilde{A}_k^\pi(s_{N_k^r}, a_{N_k^r}) = Q_k^\pi(s, a) - \tilde{Q}_k^\pi(s_{N_k^r}, a_{N_k^r}) - V_k^\pi(s) + \tilde{V}_k^\pi(s_{N_k^r}),$$

$$\begin{aligned}
V_k^\pi(s) - \tilde{V}_k^\pi(s_{N_k^r}) &= \mathbb{E} \left[V_k^\pi(s) - V_k^\pi \left(s_{N_k^r}, s(0)_{N_{-k}^r} \right) \right] \\
&\leq \mathbb{E} [c' \xi^{r+1}] = c' \xi^{r+1}
\end{aligned}$$

and

$$\begin{aligned}
Q_k^\pi(s, a) - \tilde{Q}_k^\pi(s_{N_k^r}, a_{N_k^r}) &= \mathbb{E} \left[Q_k^\pi(s, a) - Q_k^\pi \left(s_{N_k^r}, a_{N_k^r}, s(0)_{N_{-k}^r}, a(0)_{N_{-k}^r} \right) \right] \\
&\leq \mathbb{E} [c\rho^{r+1}] = c\rho^{r+1}.
\end{aligned}$$

Then $\forall k \in \mathcal{K}$

$$\left| A_k^\pi(s, a) - \tilde{A}_k^\pi(s_{N_k^r}, a_{N_k^r}) \right| \leq c\rho^{r+1} + c' \xi^{r+1}.$$

Term 3 In this paragraph, we denote, $\forall k \in \mathcal{K}$, $w_{k,\star}^{(t)} = (w_{k,\star}^{(t)})_{N_k^r}$ and $w_k^{(t)} = (w_k^{(t)})_{N_k^r}$ for clarity of exposition. Remember that

$$\Sigma_{\nu,k}^\theta = \mathbb{E}_{s,a \sim \nu} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi_\theta(a_k | s_{N_k^r}) (\nabla_{(\theta_k)_{N_k^r}} \log \pi_\theta(a_k | s_{N_k^r}))^\top \right]$$

and that we assume, $\forall k$,

$$\sup_{w \in \mathbb{R}^d} \frac{w^\top \Sigma_{d^*,k}^{(t)} w}{w^\top \Sigma_{\nu,k}^{(t)} w} \leq \kappa'.$$

Then $\forall k \in \mathcal{K}$

$$\begin{aligned}
& \mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w_{k,\star}^{(t)} - w_k^{(t)}) \right] \\
& \leq \sqrt{(w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{d^*,k}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)})} \\
& = \sqrt{\frac{(w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{d^*,k}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)})}{(w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{\nu,k}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)})} (w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{\nu,k}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)})} \\
& \leq \sqrt{\kappa' (w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{\nu,k}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)})} \\
& \text{[using that } (1 - \gamma)\nu \leq d_\nu^{\pi^{(t)}} \text{]} \\
& \leq \sqrt{\frac{\kappa'}{1 - \gamma} (w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{d^{(t)}}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)})}.
\end{aligned}$$

Since $w_{k,\star}^{(t)}$ optimizes $\tilde{L}_k(w, \theta^{(t)}, d^{(t)})$ the first order optimality condition implies that, $\forall w, \forall k \in \mathcal{K}$,

$$(w - w_{k,\star}^{(t)}) \cdot \nabla \tilde{L}_k(w_{k,\star}^{(t)}, \theta^{(t)}, d^{(t)}) \geq 0.$$

So $\forall w, \forall k \in \mathcal{K}$,

$$\begin{aligned}
& \tilde{L}_k(w, \theta^{(t)}, d^{(t)}) - \tilde{L}_k(w_{k,\star}^{(t)}, \theta^{(t)}, d^{(t)}) \\
& = \mathbb{E}_{s \sim d^{(t)}, a \sim \pi^{(t)}(\cdot|s)} \left[\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot w \right. \\
& \quad \left. + \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot w_{k,\star}^{(t)} - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot w_{k,\star}^{(t)} \right]^2 \\
& \quad - \tilde{L}_k(w_{k,\star}^{(t)}, \theta^{(t)}, d^{(t)}) \\
& = \mathbb{E}_{s \sim d^{(t)}, a \sim \pi^{(t)}(\cdot|s)} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot w_{k,\star}^{(t)} - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot w \right]^2 \\
& \quad + 2(w - w_{k,\star}^{(t)}) \mathbb{E}_{s \sim d^{(t)}, a \sim \pi^{(t)}(\cdot|s)} \left[\left(\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot w_{k,\star}^{(t)} \right) \right. \\
& \quad \left. \cdot \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \right] \\
& = (w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{d^{(t)}}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)}) + (w - w_{k,\star}^{(t)}) \cdot \nabla \tilde{L}_k(w_{k,\star}^{(t)}, \theta^{(t)}, d^{(t)})
\end{aligned}$$

$$\geq (w_{k,\star}^{(t)} - w_k^{(t)})^\top \Sigma_{d^{(t)}}^{(t)} (w_{k,\star}^{(t)} - w_k^{(t)}).$$

Therefore

$$\begin{aligned} & \mathbb{E} \left[\mathbb{E}_{s \sim d^*, a \sim \pi^*(\cdot|s)} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w_{k,\star}^{(t)} - w_k^{(t)}) \right] \right] \\ & \leq \sqrt{\mathbb{E} \left[\frac{\kappa'}{1-\gamma} (\tilde{L}_k(w, \theta^{(t)}, d^{(t)}) - \tilde{L}_k(w_{k,\star}^{(t)}, \theta^{(t)}, d^{(t)})) \right]} \\ & = \sqrt{\mathbb{E} \left[\frac{\kappa'}{1-\gamma} \mathbb{E} \left[(\tilde{L}_k(w, \theta^{(t)}, d^{(t)}) - \tilde{L}_k(w_{k,\star}^{(t)}, \theta^{(t)}, d^{(t)})) | \theta^{(t)} \right] \right]} \\ & \leq \sqrt{\frac{\kappa' \varepsilon_{\text{stat}}}{1-\gamma}}, \end{aligned}$$

which completes the proof. ■

F Statistical Error

Assume access to a global clock, then Algorithm 2 is a sampler of $d^{(t)}$ and an unbiased sampler of $\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r})$, for every $k \in \mathcal{K}$.

Proposition 19 *Consider the setting of Theorem 9 and assume access to the sampler in Algorithm 2. Assume that $\left\| \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \right\|_2 \leq B$. Then, solving the minimization problem in (8) with stochastic projected gradient descent for N steps and step size $\alpha = W/(8B(BW + \frac{1}{1-\gamma})\sqrt{N})$ gives*

$$\varepsilon_{\text{stat}} \leq \frac{8BW(BW + \frac{1}{1-\gamma})}{\sqrt{N}}.$$

Proof: The proposition follows from a result on stochastic projected gradient descent [Shalev-Shwartz and Ben-David, 2014]. Consider the minimization problem $\min_{x \in C} f(x)$ for a convex function f , then performing the update

$$x_{t+1} = P_C(x_t - \alpha v_t),$$

where $C = \{x : \|x\|_2 \leq W\}$ for some $W \geq 0$, $P_C(\cdot)$ is the projection on C , v_t is such that $\mathbb{E}[v_t|x_t] = \nabla f(x_t)$ and $\|v_t\| \leq \rho$, for N steps and with $\alpha = \sqrt{\frac{W^2}{\rho^2 N}}$, gives

$$\mathbb{E} \left[f \left(\frac{1}{N} \sum_{t=1}^N x_t \right) \right] - f(x^*) \leq \frac{W\rho}{\sqrt{N}}.$$

Noticing that $\tilde{A}_k^{(t)}(s, a) \leq 2/(1 - \gamma)$ and that the sampled gradient is therefore bounded by $8B(BW + \frac{1}{1-\gamma})$ gives the proposition. \blacksquare

The minimization problem in (8) can therefore be solved by each agent k through stochastic projected gradient descent, which only depends on the states and actions of N_k^r and on the parameters $(\theta_k)_{N_k^r}$ and with a computational cost that does not depend on K .

Algorithm 2: Sampler

Input: Starting state-action distribution ν

For each agent $k \in \mathcal{K}$ set $\hat{Q}_k = 0, \hat{V}_k = 0$.

$\forall k \in \mathcal{K}$, sample $s_k(0), a_k(0) \sim \nu_k$ and start at state $s_k(0)$;

(d_ν^r sampling) **At every time-step** $h \geq 0$:

with probability $\gamma, \forall k \in \mathcal{K}$ execute $a_k(h)$, transition to $s_k(h + 1)$ and sample

$a_k(h + 1) \sim \pi_k(\cdot | s_{N_k^r}(h + 1))$;

else accept $(s(h), a(h))$ as the sample and exit the loop, each agent only saves

$(s_{N_k^r}(h), a_{N_k^r}(h))$.

Set SampleQ=True with probability 1/2;

if SampleQ=True **then**

$\forall k \in \mathcal{K}$ execute $a_k(h)$ and then, for every time-step $h' > h$ with termination probability γ , transition to $s_k(h')$, sample $a_k(h' + 1) \sim \pi_k(\cdot | s_{N_k^r}(h' + 1))$ and set

$\hat{Q}_k = \hat{Q}_k + r(s_k(h'), a_k(h'))$;

end

else

$\forall k \in \mathcal{K}$ sample $a_k(h) \sim \pi_k(\cdot | s_{N_k^r}(h))$ and execute $a_k(h)$ and then, for every time-step $h' > h$ with termination probability γ , transition to $s_k(h')$, sample

$a_k(h' + 1) \sim \pi_k(\cdot | s_{N_k^r}(h' + 1))$ and set $\hat{V}_k = \hat{V}_k + r(s_k(h'), a_k(h'))$;

end

Result: $(s(h), a(h))$ and $\left(\hat{A}_k(s_{N_k^r}(h), a_{N_k^r}(h)) = 2 \left(\hat{Q}_k - \hat{V}_k \right) \right)_{k \in \mathcal{K}}$.

G Bias analysis

Let \tilde{L}_k^r and $w_{k,\star}$ be defined as in Section 4. For every $k \in \mathcal{K}$, let

$$L_k(w, \theta, \nu) = \mathbb{E}_{s,a \sim \nu} \left[(A_k^{\pi_\theta}(s, a) - \nabla_{\theta_k} \log \pi_{\theta_k}(a_k | s) \cdot w)^2 \right].$$

Let $\varepsilon_{\text{bias},r}$ be the smallest possible localized bias term associated to the range parameter r , i.e.

$$\varepsilon_{\text{bias},r} = \max_{k \in \mathcal{K}} \tilde{L}_k^r(w_{k,\star}, \theta^{(t)}, d^\star),$$

which is the same assumption as in Theorem 9, but with an equality instead of an inequality. Let $\varepsilon_{\text{bias}}$ be the smallest possible non-localized bias term associated with an infinite range parameter, i.e. the case where we make no approximation, namely,

$$\varepsilon_{\text{bias}} = \max_{k \in \mathcal{K}} L_k(w'_{k,\star}, \theta^{(t)}, d^\star)$$

where

$$w'_{k,\star} \in \underset{\|w\|_2 \leq W}{\operatorname{argmin}} L_k(w, \theta^{(t)}, d^{(t)}).$$

Proposition 20 *Assume that*

$$\|\nabla_{\theta_k} \log \pi^{(t)}(a|s)\|_2 \leq B$$

and that

$$\left\| \nabla_{(\theta_k)_{N_k^r}} \pi_\theta(a_k|s) \right\|_2 \leq \omega_r.$$

Then

$$\varepsilon_{\text{bias},r} \leq \varepsilon_{\text{bias}} + e_r + 2 \left(\frac{2}{1-\gamma} + WB \right) \sqrt{\frac{2\kappa' e_r}{1-\gamma} + \frac{2\kappa' e_r}{1-\gamma}},$$

where

$$e_r = \left(\frac{4}{1-\gamma} + 2WB \right) W\omega_r + \left(\frac{4}{1-\gamma} + 2WB \right) (c\psi^{r+1} + c'\phi^{r+1}).$$

The second assumption can be respected by choosing a policy belonging to the policy class described in Appendix A, as ω_r can be controlled by setting the parameters $\{\alpha_r\}$ small enough. In particular, it is possible to set it to be exponentially small in r . Therefore, the proposition shows that the localized bias is at most the bias associated with an infinite range parameter plus a quantity that goes to 0 exponentially fast in r . We present the proof of this result below.

Proof:[of Proposition 20] To prove the proposition we need two intermediate results. For any w, θ and ν , we have that

$$\begin{aligned} & \left| \tilde{L}_k^r((w)_{N_k^r}, \theta, \nu) - L_k(w, \theta, \nu) \right| \\ &= \mathbb{E}_{s, a \sim \nu} \left[\left(\tilde{A}_k^{\pi_\theta}(s_{N_k^r}, a_{N_k^r}) \right)^2 - \left(A_k^{\pi_\theta}(s, a) \right)^2 \right] \\ & \quad + \mathbb{E}_{s, a \sim \nu} \left[\left(\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w)_{N_k^r} \right)^2 - \left(\nabla_{\theta_k} \log \pi_{\theta_k}(a_k|s) \cdot w \right)^2 \right] \end{aligned}$$

$$\begin{aligned}
& + 2\mathbb{E}_{s,a\sim\nu} \left[A_k^{\pi_\theta}(s, a) \nabla_{\theta_k} \log \pi_{\theta_k}(a_k|s) \cdot w - \tilde{A}_k^{\pi_\theta}(s_{N_k^r}, a_{N_k^r}) \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s) \cdot (w)_{N_k^r} \right] \\
& \leq \frac{4}{1-\gamma} (c\psi^{r+1} + c'\phi^{r+1}) + 2W^2 B\omega_r + 2WB (c\psi^{r+1} + c'\phi^{r+1}) + \frac{4}{1-\gamma} W\omega_r \\
& = e_r.
\end{aligned}$$

Following the same passages in the proof of Theorem 9 in Appendix E, we have that

$$\begin{aligned}
& \max_{k \in \mathcal{K}} \mathbb{E}_{s,a \sim d^*} \left[\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s_{N_k^r})^\top \left((w'_{k,\star})_{N_k^r} - w_{k,\star} \right) \right] \\
& \leq \sqrt{\frac{\kappa'}{1-\gamma} \left(\tilde{L}_k^r((w'_{k,\star})_{N_k^r}, \theta^{(t)}, d^{(t)}) - \tilde{L}_k^r(w_{k,\star}, \theta^{(t)}, d^{(t)}) \right)} \\
& \text{[using the first result]} \\
& \leq \sqrt{\frac{\kappa'}{1-\gamma} \left| L_k(w'_{k,\star}, \theta^{(t)}, d^{(t)}) - \tilde{L}_k^r(w_{k,\star}, \theta^{(t)}, d^{(t)}) \right| + \frac{\kappa' e_r}{1-\gamma}} \\
& = \sqrt{\frac{\kappa'}{1-\gamma} \left| \min_{\|w\|_2 \leq W} L_k(w, \theta^{(t)}, d^{(t)}) - \min_{\|w\|_2 \leq W} \tilde{L}_k^r((w)_{N_k^r}, \theta^{(t)}, d^{(t)}) \right| + \frac{\kappa' e_r}{1-\gamma}} \\
& \leq \sqrt{\frac{2\kappa' e_r}{1-\gamma}}.
\end{aligned}$$

We can now prove the proposition.

$$\begin{aligned}
\varepsilon_{\text{bias},r} & = \max_{k \in \mathcal{K}} \mathbb{E}_{s,a \sim d^*} \left[\left(\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s_{N_k^r}) \cdot w_{k,\star} \right)^2 \right] \\
& = \max_{k \in \mathcal{K}} \mathbb{E}_{s,a \sim d^*} \left[\left(\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s_{N_k^r}) \cdot (w'_{k,\star})_{N_k^r} \right)^2 \right. \\
& \quad + 2 \left(\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s_{N_k^r}) \cdot (w'_{k,\star})_{N_k^r} \right) \\
& \quad \cdot \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s_{N_k^r})^\top \left((w'_{k,\star})_{N_k^r} - w_{k,\star} \right) \\
& \quad \left. + \left(\nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a|s_{N_k^r})^\top \left((w'_{k,\star})_{N_k^r} - w_{k,\star} \right) \right)^2 \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \max_{k \in \mathcal{K}} \mathbb{E}_{s, a \sim d^*} \left[\left(\tilde{A}_k^{(t)}(s_{N_k^r}, a_{N_k^r}) - \nabla_{(\theta_k)_{N_k^r}} \log \pi^{(t)}(a | s_{N_k^r}) \cdot (w'_{k, \star})_{N_k^r} \right)^2 \right] \\
&\quad + 2 \left(\frac{2}{1-\gamma} + WB \right) \sqrt{\frac{2\kappa' e_r}{1-\gamma}} + \frac{2\kappa' e_r}{1-\gamma} \\
&= \max_{k \in \mathcal{K}} \tilde{L}_k^r((w'_{k, \star})_{N_k^r}, \theta^{(t)}, d^*) + 2 \left(\frac{2}{1-\gamma} + WB \right) \sqrt{\frac{2\kappa' e_r}{1-\gamma}} + \frac{2\kappa' e_r}{1-\gamma} \\
&\leq \varepsilon_{\text{bias}} + e_r + 2 \left(\frac{2}{1-\gamma} + WB \right) \sqrt{\frac{2\kappa' e_r}{1-\gamma}} + \frac{2\kappa' e_r}{1-\gamma}.
\end{aligned}$$

■