

## A. Convergence Analysis

### A.1. Preliminaries

Let  $\mathfrak{J} = \phi \circ \mathbf{J} : \mathbb{R}^{N \times K} \rightarrow \mathbb{R}$  denote the *Social Welfare Objective* (SWO), as the composition of the social welfare function  $\phi$  with the vectorial objective  $\mathbf{J}$ , where  $N$  is the number of agents and  $K$  denotes the dimension of the parameter  $\theta_i$  for the policy network of an agent  $i \in [N]$ .

Under a gradient descent optimization scheme, the direction of the update should follow the gradient at iteration  $k$ , i.e.,

$$\mathbf{g}^{k,*} := \nabla \mathfrak{J}(\boldsymbol{\theta}^k) = \nabla \phi(\mathbf{J}(\boldsymbol{\theta}^k))^\top \cdot \nabla \mathbf{J}(\boldsymbol{\theta}^k), \quad (10)$$

where  $\nabla \mathbf{J}$  denotes the Jacobian of  $\mathbf{J}$ . While the subscript on  $\boldsymbol{\theta}$  refers to the different agents, superscript  $k$  refers to the successive iterations of policy; it may be omitted below to lighten the notations.

Adapting the Policy Gradient Theorem (Sutton & Barto, 1998) to this fair-multi-agent context, leads us to the following gradient direction for agent  $i$ :

**Lemma A.1.**

$$\nabla_{\theta_i} \mathfrak{J}(\boldsymbol{\theta}^k) = \beta_k \mathbb{E}_{\boldsymbol{\theta}} \left[ \mathbf{A}^{k,SWF}(\mathbf{s}, \mathbf{a}) \cdot \nabla_{\theta_i} \log \pi_{\theta_i}^k(a_i | o_i) \right], \quad (11)$$

where  $\mathbf{A}^{k,SWF} := \nabla \phi(\mathbf{J}(\boldsymbol{\theta}^k))^\top \cdot \mathbf{A}^k : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , and  $\beta_k$  is the average length of an episode while following the policy  $\pi^k$  in the episodic case or  $\beta_k = 1$  in the infinite horizon case.<sup>4</sup>

*Proof.* Let us look at  $\mathbf{J}_{(j)}(\boldsymbol{\theta}) := \mathbf{J}_{I_j}(\boldsymbol{\theta})$ , the vectorial objective of agent  $i$ , and roughly follow the classical policy gradient theorem's proof. We immediately have:

$$\nabla_i \mathbf{V}_{(j)}^\pi(\mathbf{s}) = \sum_{\mathbf{a} \in \mathcal{A}} \left( \nabla_i \pi_{\boldsymbol{\theta}}(\mathbf{a} | \mathbf{s}) \mathbf{Q}_{(j)}^\pi(\mathbf{s}, \mathbf{a}) + \pi_{\boldsymbol{\theta}}(\mathbf{a} | \mathbf{s}) \nabla_i \mathbf{Q}_{(j)}^\pi(\mathbf{s}, \mathbf{a}) \right)$$

where  $\nabla_i$ ,  $\mathbf{Q}_{(j)}$  and  $\mathbf{V}_{(j)}$  stands for  $\nabla_{\theta_i}$ ,  $\mathbf{Q}_{I_j}$  and  $\mathbf{V}_{I_j}$  respectively. For the second term, we can develop  $\mathbf{Q}_{(j)}^\pi$  as  $\sum_{s', r} P(s', r | s, a)(r + \mathbf{V}_{(j)}^\pi(s'))$ . Since nor the reward nor the transition probability depends on  $\theta$ , it leads to, after simplifications:

$$\nabla_i \mathbf{V}_{(j)}^\pi(\mathbf{s}) = \psi_i(\mathbf{s}) + \sum_{\mathbf{a} \in \mathcal{A}} \pi_{\boldsymbol{\theta}}(\mathbf{a} | \mathbf{s}) \sum_{s' \in \mathcal{S}} P(s' | \mathbf{s}, \mathbf{a}) \nabla_i \mathbf{V}_{(j)}^\pi(s') = \psi_i(\mathbf{s}) + \sum_{s' \in \mathcal{S}} \rho^\pi(\mathbf{s} \rightarrow s', 1) \nabla_i \mathbf{V}_{(j)}^\pi(s'), \quad (12)$$

where  $\psi_{i,j}(\mathbf{s}) := \sum_{\mathbf{a} \in \mathcal{A}} \nabla_i \pi_{\boldsymbol{\theta}}(\mathbf{a} | \mathbf{s}) \mathbf{Q}_{(j)}^\pi(\mathbf{s}, \mathbf{a})$ , and  $\rho^\pi(\mathbf{s} \rightarrow s', k)$  is the probability of transitioning from state  $\mathbf{s}$  to state  $s'$  while following the policy  $\pi$  after  $k$  steps. Notably  $\rho^\pi(\mathbf{s} \rightarrow s', 1) = \sum_{\mathbf{a} \in \mathcal{A}} \pi_{\boldsymbol{\theta}}(\mathbf{a} | \mathbf{s}) P(s' | \mathbf{s}, \mathbf{a})$ . The equation (12) enables us to smoothly unroll a recursive process:

$$\begin{aligned} \nabla_i \mathbf{V}_{(j)}^\pi(\mathbf{s}) &= \psi_{i,j}(\mathbf{s}) + \sum_{s' \in \mathcal{S}} \rho^\pi(\mathbf{s} \rightarrow s', 1) \left[ \psi_{i,j}(s') + \sum_{s'' \in \mathcal{S}} \rho^\pi(s' \rightarrow s'', 1) \nabla_i \mathbf{V}_{(j)}^\pi(s'') \right] \\ &= \psi_{i,j}(\mathbf{s}) + \sum_{s' \in \mathcal{S}} \rho^\pi(\mathbf{s} \rightarrow s', 1) \psi_{i,j}(s') + \sum_{s' \in \mathcal{S}} \sum_{s'' \in \mathcal{S}} \rho^\pi(\mathbf{s} \rightarrow s'', 2) \nabla_i \mathbf{V}_{(j)}^\pi(s'') \\ &= \dots \text{ (iterating) } \dots \\ &= \sum_{s' \in \mathcal{S}} \sum_{k=0}^{\infty} \rho^\pi(\mathbf{s} \rightarrow s', k) \psi_{i,j}(s'), \end{aligned}$$

Coming back to the derivative of the vectorial objective of agent  $j$ :

$$\nabla_i \mathbf{J}_{(j)}(\boldsymbol{\theta}) = \sum_{\mathbf{s} \in \mathcal{S}} \chi^\pi(\mathbf{s}) \psi_{i,j}(\mathbf{s}) = \beta^\pi \left( \sum_{\mathbf{s} \in \mathcal{S}} d^\pi(\mathbf{s}) \psi_{i,j}(\mathbf{s}) \right),$$

<sup>4</sup>In the infinite horizon case, we can follow a nearly-identical proof, upon the assumption the underlying MDP is ergodic.

where  $\chi^\pi(\mathbf{s}) := \sum_{k=0}^{\infty} \rho^\pi(\mathbf{s}_0 \rightarrow \mathbf{s}, k)$ ,  $d^\pi(\mathbf{s}) = \frac{\chi^\pi(\mathbf{s})}{\sum_{\mathbf{s}' \in \mathcal{S}} \chi^\pi(\mathbf{s}' )}$  is the stationary distribution and  $\beta^\pi := \sum_{\mathbf{s} \in \mathcal{S}} \chi^\pi(\mathbf{s})$ . Since  $\pi_{\theta}(\mathbf{a}|\mathbf{s}) = \sum_{\mathbf{o}} \Omega(\mathbf{o}|\mathbf{s}) \prod_i \pi_{\theta_i}(a_i|o_i)$ :

$$\begin{aligned} \nabla_i \mathbf{J}_{(j)}(\boldsymbol{\theta}) &= \beta^\pi \sum_{\mathbf{s} \in \mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} d^\pi(\mathbf{s}) \mathbf{Q}_{(j)}^\pi(\mathbf{s}, \mathbf{a}) \nabla_i \pi_{\theta}(\mathbf{a}|\mathbf{s}) = \beta^\pi \sum_{\mathbf{s} \in \mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} d^\pi(\mathbf{s}) \pi_{\theta}(\mathbf{a}|\mathbf{s}) \mathbf{Q}_{(j)}^\pi(\mathbf{s}, \mathbf{a}) \nabla_i \log \pi_{\theta}(\mathbf{a}|\mathbf{s}) \\ &= \beta^\pi \mathbb{E}_{\substack{\mathbf{s} \sim d^\pi(\bullet) \\ \mathbf{a} \sim \pi(\bullet|\mathbf{s}) \\ \mathbf{o} \sim \Omega(\bullet|\mathbf{s})}} \left[ \mathbf{Q}_{(j)}^\pi(\mathbf{s}, \mathbf{a}) \nabla_{\theta_i} \log \pi_{\theta_i}(a_i|o_i) \right] \end{aligned}$$

These computations may be equivalently rewritten with the advantage function:

$$\nabla_i \mathbf{J}(\boldsymbol{\theta}) = \beta^\pi \mathbb{E} \left[ \mathbf{A}^\pi(\mathbf{s}, \mathbf{a}) \nabla_{\theta_i} \log \pi_{\theta_i}(a_i|o_i) \right]$$

Hence, returning to the SWO  $\mathfrak{J}$ , it boils down to the Lemma's claim:

$$\nabla_i \mathfrak{J}(\boldsymbol{\theta}) = \beta^\pi \nabla \phi(\mathbf{J}(\boldsymbol{\theta}))^T \cdot \mathbb{E} \left[ \mathbf{A}^\pi(\mathbf{s}, \mathbf{a}) \nabla_{\theta_i} \log \pi_{\theta_i}(a_i|o_i) \right] = \beta^\pi \mathbb{E} \left[ \mathbf{A}^{\pi, \text{SWF}}(\mathbf{s}, \mathbf{a}) \nabla_{\theta_i} \log \pi_{\theta_i}(a_i|o_i) \right]$$

□

As mentioned previously, since we adopt a fully decentralised framework with partial observability and since the advantage function is not known, we instead approximate our gradient update according to the following direction  $\tilde{\mathbf{g}}_i^k$  for an agent  $i$ :

$$\tilde{\mathbf{g}}_i^k := \mathbb{E}_{\theta} \left[ \hat{\mathbf{A}}^{k, \text{SWF}}(\mathbf{o}, \mathbf{a}) \cdot \nabla_{\theta_i} \log \pi_{\theta_i}(a_i|o_i) \right],$$

where  $\hat{\mathbf{A}}^k = (\hat{\mathbf{A}}_i^k)_{i \in [N]}$  aggregates the estimate of the local advantage at step  $k$  for each agent  $i$ ,  $\hat{\mathbf{A}}_i^k : \mathcal{O}_i \times \mathcal{A}_i \rightarrow \mathbb{R}^{|\mathcal{I}_i|}$ , and  $\hat{\mathbf{A}}^{k, \text{SWF}}(\mathbf{o}, \mathbf{a}) = \nabla \phi(\hat{\mathbf{J}}(\boldsymbol{\theta}^k))^\top \cdot \hat{\mathbf{A}}^k(\mathbf{o}, \mathbf{a})$ .

Let  $\hat{\mathbf{g}}^k \in \mathbb{R}^{N \times K}$  denote the estimated gradient at step  $k$ , estimated from sampling a minibatch of transitions of size  $m$ . The update at step  $k$  of this stochastic gradient ascent can be written as:

$$\boldsymbol{\theta}^{k+1} \leftarrow \boldsymbol{\theta}^k + \alpha_k \hat{\mathbf{g}}^k \quad (13)$$

Comparing the ideal and estimated gradients by introducing the error term  $\boldsymbol{\epsilon}^k$ , using (11):

$$\hat{\mathbf{g}}^k = \bar{\mathbf{g}}^k + \boldsymbol{\epsilon}^k = \beta_k \bar{\mathbf{g}}^{k,*} + \boldsymbol{\epsilon}^k, \quad \text{where} \quad \bar{\mathbf{g}}_i^k := \mathbb{E}_{\theta} \left[ \mathbf{A}^{k, \text{SWF}}(\mathbf{s}, \mathbf{a}) \cdot \nabla_{\theta_i} \log \pi_{\theta_i}(a_i|o_i) \right]. \quad (14)$$

Let us decompose this error term into two parts:

$$\boldsymbol{\epsilon}^k = \boldsymbol{\nu}^k + \boldsymbol{\eta}^k, \quad \text{where} \quad \boldsymbol{\nu}^k := \hat{\mathbf{g}}^k - \tilde{\mathbf{g}}^k \quad \text{and} \quad \boldsymbol{\eta}^k := \tilde{\mathbf{g}}^k - \bar{\mathbf{g}}^k \quad (15)$$

While the first term  $\boldsymbol{\nu}^k$  tends towards zero when  $m$  tends towards infinity, some assumptions would be made below to control the second part of this error  $\boldsymbol{\eta}^k$ .

## A.2. Assumptions

Here we detail and discuss the assumptions made in our theoretical analysis. The first part of the assumptions regarding the social welfare function  $\phi$  has already been justified in the context of fairness considerations in Section 3.2. We simply assume in addition that the norm of the gradient of  $\phi$  is bounded by  $M_\phi \in \mathbb{R}^+$  on the image of  $\mathbf{J}$ .

**Assumption 1.**  $\phi$  is concave and non-decreasing in each argument. It is (sub-)differentiable<sup>5</sup> and the norm of its gradient is assumed to be bounded by  $M_\phi \in \mathbb{R}^+$  on the image of  $\mathbf{J}$ . This implies that  $\phi$  is  $L_\phi$ -smooth on the image of  $\mathbf{J}$ , with  $L_\phi \leq M_\phi$ .

Adopting the boundedness of the reward function is a common and legitimate hypothesis in the literature of policy gradient or actor-critic algorithms. If  $B_R$  is uniformly bounding the reward, in the discounted case,  $\frac{B_R}{1-\gamma}$  would uniformly bound both the value or action-value function, and  $\frac{2B_R}{1-\gamma}$  would bound the advantage function. From these considerations, it seems legitimate to assume the estimated advantage function to be bounded:

<sup>5</sup>We could adapt the proof to the case of sub-differentiability of  $\phi$  and  $\mathbf{J}$ , yet their sub-gradient should satisfy Lipschitz-inequalities.

**Assumption 2.** The estimated advantage function  $\hat{A}^k$  is uniformly bounded by  $B_A$ , for all  $k$ .

Here, we adopt classic assumptions concerning the vectorial objective function  $\mathbf{J}$ :

**Assumption 3.** The function  $\mathbf{J} : \mathbb{R}^{N \times K} \rightarrow \mathbb{R}^D$  is concave and Lipschitz-smooth with constant  $L_J$ , i.e., its gradient is  $L_J$ -Lipschitz continuous. These two hypotheses imply respectively:

$$\begin{aligned} \forall i \in [D], \quad \mathbf{J}_i(\boldsymbol{\theta}') - \mathbf{J}_i(\boldsymbol{\theta}) &\leq \langle \nabla \mathbf{J}_i(\boldsymbol{\theta}), \boldsymbol{\theta}' - \boldsymbol{\theta} \rangle \\ \|\nabla \mathbf{J}(\boldsymbol{\theta}') - \nabla \mathbf{J}(\boldsymbol{\theta})\| &\leq L_J \|\boldsymbol{\theta}' - \boldsymbol{\theta}\| \end{aligned}$$

This assumption notably covers the case where a linear approximation scheme is used. By Cauchy-Schwartz, integrating the last equation, and using the concavity hypothesis:

$$\mathbf{J}(\boldsymbol{\theta}') - \mathbf{J}(\boldsymbol{\theta}) \geq \nabla \mathbf{J}(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) - \frac{L_J}{2} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|^2 \quad (16)$$

The previous assumptions imply:

**Lemma A.2.** (i) Under Assumptions (1, 3), the SWO  $\mathfrak{J}$  is concave.

(ii) Under Assumptions (1, 3), with the additional hypothesis that the gradient of the vectorial objective  $\nabla \mathbf{J}$  is bounded by  $M_J$ , then the SWO  $\mathfrak{J}$  is  $L_{\mathfrak{J}}$ -smooth, with:

$$L_{\mathfrak{J}} \leq L_\phi M_J^2 + M_\phi L_J.$$

*Proof.* (i) From Assumptions 1 and 3, both  $\phi$  and  $\mathbf{J}$  are concave, and  $\phi$  is non-decreasing in each argument. It is then straightforward to prove that the fair objective  $\mathfrak{J} = \phi \circ \mathbf{J}$  is concave.

(ii) For the last assertion, we can show:

$$\begin{aligned} \|\nabla \mathfrak{J}(\boldsymbol{\theta}') - \nabla \mathfrak{J}(\boldsymbol{\theta})\| &= \|\nabla \phi(\mathbf{J}(\boldsymbol{\theta}'))^\top \nabla \mathbf{J}(\boldsymbol{\theta}') - \nabla \phi(\mathbf{J}(\boldsymbol{\theta}))^\top \nabla \mathbf{J}(\boldsymbol{\theta})\| \\ &\leq \|\nabla \phi(\mathbf{J}(\boldsymbol{\theta}'))^\top \nabla \mathbf{J}(\boldsymbol{\theta}') - \nabla \phi(\mathbf{J}(\boldsymbol{\theta}))^\top \nabla \mathbf{J}(\boldsymbol{\theta}')\| + \|\nabla \phi(\mathbf{J}(\boldsymbol{\theta}))^\top \nabla \mathbf{J}(\boldsymbol{\theta}') - \nabla \phi(\mathbf{J}(\boldsymbol{\theta}))^\top \nabla \mathbf{J}(\boldsymbol{\theta})\| \\ &\leq L_\phi \|\mathbf{J}(\boldsymbol{\theta}') - \mathbf{J}(\boldsymbol{\theta})\| M_J + L_J M_\phi \|\boldsymbol{\theta}' - \boldsymbol{\theta}\| \leq (L_\phi M_J^2 + L_J M_\phi) \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|. \end{aligned}$$

□

Let  $\Theta^*$  denote the set of maximizers of  $\mathfrak{J}$ , i.e.,  $\Theta^* := \operatorname{argmax}\{\mathfrak{J}(\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \mathbb{R}^{N \times K}\}$ . The concavity of  $\mathfrak{J}$  implies that the set  $\Theta^*$  is convex. We furthermore assume:

**Assumption 4.** The optimal set  $\Theta^*$  is non empty and there exists a radius bound  $R_0$  such that:

$$\max_{\boldsymbol{\theta}^* \in \Theta^*} \max_{\boldsymbol{\theta}} \{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \mid \mathfrak{J}(\boldsymbol{\theta}) \geq \mathfrak{J}(\boldsymbol{\theta}^*)\} \leq R_0,$$

where  $\boldsymbol{\theta}^0$  is assumed to be the initial parameter value.<sup>6</sup>

A common condition on the learning rate in convex optimization is to assume it is smaller than the inverse of the Lipschitz constant. Here, we use a slight variant of this condition:

**Assumption 5.** Given Lipschitz constant  $L_J$  from Assumption 3, and gradient bound  $M_\phi$  from Assumption 1, the learning rate satisfies:

$$\alpha_k \leq \frac{1}{L_J M_\phi \beta_k},$$

where  $\beta_k$  is the average length of an episode while following the policy  $\pi^k$  in the episodic case or  $\beta_k = 1$  in the infinite horizon case.

<sup>6</sup>However, if such hypothesis is not satisfied for  $\boldsymbol{\theta}^0$ , we may assume it is satisfied for another  $\boldsymbol{\theta}^{i_0}$ ; it alters the inequality formulation in our theorem —as we need to shift the indices by  $i_0$ —, but not the convergence result itself.

Adopting classic assumptions regarding the policy, in the policy gradient or actor-critic literature of convergence analysis:<sup>7</sup>

**Assumption 6.** (i) The policy  $\pi_\theta$  is differentiable with respect to  $\theta$ , and  $\nabla \log \pi_\theta(\mathbf{a}|s)$ , referred to as the score function, exists. Moreover, there exists a bound  $M_\pi \in \mathbb{R}^+$  which uniformly bounds  $\|\nabla \log \pi_\theta\|$  for all  $\theta$ :

$$\sup_{\theta, s, \mathbf{a}} \|\nabla \log \pi_\theta(\mathbf{a}|s)\| = M_\pi$$

(ii) The variance of the estimated gradient is uniformly bounded:

$$\sup_{\theta, s, \mathbf{a}} \text{Var}(\hat{\mathbf{g}}^\theta(s, \mathbf{a})) \leq \sigma^\#$$

In our context of fully decentralized and partially observable multi-agent RL, we have to consider the error made by the local estimates of the advantage function, as mentioned in (7). Within such an approximation lies several sources of error that we will consider separately below.

First, we need to approximate the objective  $J$ . Since we want to focus our analysis on the convergence of the policy<sup>8</sup>, we make the following assumptions:

**Assumption 7.** The approximation error of  $J$  is bounded in probability:<sup>9</sup>

$$\forall k > 0, \quad \mathbb{P}\left(\left\|\hat{J}(\theta^k) - J(\theta^k)\right\| \leq \Xi_J\right) = 1$$

Then, we have to control the approximation error of the local advantage function, which stems from several sources (such as the finiteness of the sample set on which it would be estimated, and the expressiveness of the class of functions estimators used for  $\hat{Q}$  and  $\hat{V}$ ):

**Assumption 8.** The approximation error of the local advantage  $A$  is uniformly bounded in probability:

$$\forall k > 0, \quad \mathbb{P}\left(\left\|\hat{A}^k(\mathbf{o}, \mathbf{a}) - A^k(\mathbf{o}, \mathbf{a})\right\|_\infty \leq \Xi_A\right) = 1$$

Finally, the mismatch between the local and global advantage functions has to be taken into consideration:

**Assumption 9.** The error stemming from the use of the decentralized local advantage is uniformly bounded in probability:

$$\forall k > 0, \forall i \in [N], \forall j \in [N], j \neq i, \quad \mathbb{P}\left(\left\|\mathbf{A}_{j/i}^k(o_i, a_i)\right\|_\infty \leq \Xi_L\right) = 1$$

Or, less restrictively, after application of the SWF  $\phi$ , we can assume instead:

$$\forall k > 0, \forall i \in [N], \quad \mathbb{P}\left(\left\|\mathbf{A}_{\bullet/i}^{(-i),k,SWF}(o_i, a_i)\right\|_\infty \leq \Xi_L\right) = 1, \quad (17)$$

with  $\mathbf{A}_{j/i}^k(o_i, a_i) := \mathbb{E}_{\substack{o_{-i} \sim \mathbb{P}(\cdot|o_i), \\ \mathbf{a}_{-i} \sim \pi^{-i}(\cdot|o_{-i})}} \left[ \mathbf{A}_j^k(\{o_i, \mathbf{o}_{-i}\}, \{a_i, \mathbf{a}_{-i}\}) \right]$  and  $\left(\mathbf{A}_{\bullet/i}^{(-i),k}(o_i, a_i)\right)_j := \begin{cases} \mathbf{A}_{j/i}^k(o_i, a_i) & j \neq i \\ 0 & j = i \end{cases}$ .

Here, the notations  $\mathbf{o}_{-i}$  resp.  $\mathbf{a}_{-i}$  designate the set  $\{o_j\}_{j \neq i}$  resp.  $\{a_j\}_{j \neq i}$ . The notation  $\mathbf{A}_{j/i}$  stands for the advantage of the actions of agent  $i$  relative to another agent  $j$ . As made clear in the proof of Lemma A.3, when comparing  $A(s, \mathbf{a})$  and  $A(\mathbf{o}, \mathbf{a})$ , these inter-agent advantages  $\mathbf{A}_{\bullet/i}^{(-i)}$  appear. This assumption tackles the discrepancy between each agent's local advantage and each agent's global advantage, which would take into consideration the benefit of the actions of other agents.

<sup>7</sup>Assumption 5.1 in (Xu et al., 2020), Assumption 3.1 in (Zhang et al., 2020), Assumption 4.1 in (Papini et al., 2018), Assumption 1 in (Kumar et al., 2019), Assumption 4.5 in (Qiu et al., 2021), etc.

<sup>8</sup>In Actor-Critic schemes, it is common to adopt a two-timescale analysis, in which we usually need to ensure that the actor update moves slowly while the critic update chases the slowly-moving target defined by the actor by enforcing a condition on the ratio of their learning rates (cf. (Konda & Tsitsiklis, 2000)). Here, we do not focus on bounding the error made from the advantage estimation in probability, in function of the size of the mini-batch  $m$ , but the analysis could be extended in later work.

<sup>9</sup>If we adopt GGF as SWF, this hypothesis could be relaxed, as we would only care that the estimate of  $J$  preserves the order of the different  $J_i(\theta^k)$  for each agent, not the actual value.

Looking at (17), let us remark that the benefits (positive or negative) of one agent action, may, upon application of the SWF  $\phi$ , average out in some situations. Therefore, ignoring them may not have too much impact.

Note that the conditional probability arising in the above expected value satisfies, with the previous notations:

$$P(\mathbf{o}_{-i} | o_i) = \sum_{\mathbf{s}} d^\pi(\mathbf{s}) P(\mathbf{o}_{-i} | o_i, \mathbf{s}) = \sum_{\mathbf{s}} \frac{d^\pi(\mathbf{s}) \Omega(\mathbf{o} | \mathbf{s})}{P(o_i | \mathbf{s})} = \sum_{\mathbf{s}} \frac{d^\pi(\mathbf{s}) \Omega(\mathbf{o} | \mathbf{s})}{\Omega(o_i | \mathbf{s})}$$

These assumptions enable to bound the error term  $\boldsymbol{\eta}^k$  as defined in (15):

**Lemma A.3.** *Under assumptions 1,2,6,7,8, 9,  $\boldsymbol{\eta}^k$  is bounded in probability:*

$$\forall k \quad P(\|\boldsymbol{\eta}^k\| \leq \Xi) = 1, \quad \text{with} \quad \Xi = (L_\phi B_A \Xi_J + M_\phi \Xi_A + M_\phi \Xi_L) M_\pi,$$

with the bounds as previously defined.

*Proof.* Splitting into three terms, from the definition of  $\mathbf{A}^{SWF}$ :

$$\begin{aligned} \|\boldsymbol{\eta}_i^k\| &= \|\tilde{\mathbf{g}}_i^k - \bar{\mathbf{g}}_i^k\| = \left\| \mathbb{E}_\theta \left[ \left( \hat{\mathbf{A}}^{k,SWF}(\mathbf{o}, \mathbf{a}) - \mathbf{A}^{k,SWF}(\mathbf{s}, \mathbf{a}) \right) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] \right\| \\ &\leq \left\| \mathbb{E}_\theta \left[ \left( \nabla \phi(\hat{\mathbf{J}}(\boldsymbol{\theta}^k)) - \nabla \phi(\mathbf{J}(\boldsymbol{\theta}^k)) \right)^\top \cdot \hat{\mathbf{A}}^k(\mathbf{o}, \mathbf{a}) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] \right\| \\ &\quad + \left\| \mathbb{E}_\theta \left[ \nabla \phi(\mathbf{J}(\boldsymbol{\theta}^k))^\top \cdot (\hat{\mathbf{A}}^k(\mathbf{o}, \mathbf{a}) - \mathbf{A}^k(\mathbf{o}, \mathbf{a})) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] \right\| \\ &\quad + \left\| \mathbb{E}_\theta \left[ \nabla \phi(\mathbf{J}(\boldsymbol{\theta}^k))^\top \cdot (\mathbf{A}^k(\mathbf{o}, \mathbf{a}) - \mathbf{A}^k(\mathbf{s}, \mathbf{a})) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] \right\| \end{aligned} \quad (18)$$

Since  $\|\mathbb{E}(X)\| \leq \mathbb{E}(\|X\|)$ :

$$\begin{aligned} \|\boldsymbol{\eta}_i^k\| &\leq \left( L_\phi B_A \|\mathbf{J}(\boldsymbol{\theta}^k) - \mathbf{J}(\boldsymbol{\theta})\| + M_\phi \Xi_A + M_\phi \Xi_L \right) \mathbb{E}_\theta \left[ \|\nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i)\| \right] \\ &\leq \left( L_\phi B_A \Xi_J + M_\phi \Xi_A + M_\phi \Xi_L \right) M_\pi, \end{aligned} \quad (19)$$

with the previously defined bounds; for the first inequality we used the fact that  $\phi$  is  $L_\phi$ -smooth, and Assumptions 2 resp. 8 resp. 9 for the different terms; for the second inequality, we used Assumption 7 and the boundedness of  $\|\nabla \ln \pi\|$  from Assumption 6.

Let us justify the computations for the third term of (18). On the one hand, from the definition of  $\mathbf{A}_{j/i}$ , after some simplifications, we arrive at:

$$\mathbb{E}_\theta \left[ \mathbf{A}_j^k(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] = \mathbb{E}_\theta \left[ \mathbf{A}_{j/i}^k(o_i, a_i) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right]$$

On the other hand, the other terms may be simplified to:

$$\mathbb{E}_\theta \left[ \mathbf{A}_j^k(o_j, a_j) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] = \begin{cases} 0 & j \neq i \\ \mathbb{E}_\theta \left[ \mathbf{A}_i^k(o_i, a_i) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] & j = i \end{cases}$$

Therefore, the third term of the right side in (18), boils down to:

$$\left\| \mathbb{E}_\theta \left[ \nabla \phi(\mathbf{J}(\boldsymbol{\theta}^k))^\top \cdot (\mathbf{A}_{\bullet/i}^{(-i),k}(o_i, a_i)) \nabla_{\boldsymbol{\theta}_i} \log \pi_{\boldsymbol{\theta}_i}^k(a_i | o_i) \right] \right\|,$$

where  $\mathbf{A}_{\bullet/i}^{(-i),k}(o_i, a_i)$  denotes the vector whose  $j$ -th component is:  $\left( \mathbf{A}_{\bullet/i}^{(-i),k}(o_i, a_i) \right)_j := \begin{cases} \mathbf{A}_{j/i}^k(o_i, a_i) & j \neq i \\ 0 & j = i \end{cases}$ .

Assumption 9 enable us to conclude. □

**Locally concave Extension** To extend our analysis to the locally concave scenario, two previous assumptions (Assumptions 3, 4) need to be reformulated as follows:

**Assumption (3').** *There exists a neighborhood  $\mathcal{U} \subset \mathbb{R}^{N \times K}$  containing the initial parameter value  $\theta^0$ , such that:*

- *The function  $\mathbf{J} : \mathbb{R}^{N \times K} \rightarrow \mathbb{R}^D$  is concave on  $\mathcal{U}$ . It is Lipschitz-smooth on  $\mathcal{U}$  with constant  $L_J$ , i.e., its gradient is  $L_J$ -Lipschitz continuous on  $\mathcal{U}$ .*
- *The set of optimal values in  $\mathcal{U}$   $\Theta^{\mathcal{U},*}$  is non empty. There exists a radius bound  $R_0^{\mathcal{U}}$  such that:*

$$\max_{\theta^* \in \Theta^{\mathcal{U},*}} \max_{\theta \in \mathcal{U}} \{\|\theta - \theta^*\| \mid \mathfrak{J}(\theta) \geq \mathfrak{J}(\theta^0)\} \leq R_0^{\mathcal{U}}.$$

This scenario covers the more general case where a neural network is used for the approximation. As previously, together with Assumption (1), (i) implies the local concavity of  $\mathfrak{J}$ :

**Lemma A.4.** *Under Assumptions (1, 3'), the SWO  $\mathfrak{J}$  is locally concave.*

### A.3. Suboptimality Analysis

**Lemma A.5.** *Under Assumptions (1, 3, 5) the SWO between two consecutive update satisfies:*

$$\forall k > 0, \quad \mathfrak{J}(\theta^{k+1}) - \mathfrak{J}(\theta^k) \geq \gamma_k \|\nabla \mathfrak{J}(\theta^k)\|^2 - \delta_k \|\epsilon^k\|^2, \quad \text{with } \gamma_k = \frac{\alpha_k \beta_k}{2}, \text{ and } \delta_k = \frac{\alpha_k}{2\beta_k}. \quad (20)$$

*Proof.* Since  $\phi$  is assumed concave (Assumption 1):

$$\mathfrak{J}(\theta^{k+1}) - \mathfrak{J}(\theta^k) \geq \nabla \phi(\mathbf{J}(\theta^k))^\top (\mathbf{J}(\theta^{k+1}) - \mathbf{J}(\theta^k)). \quad (21)$$

Since  $\mathbf{J}$  is assumed  $L_J$ -smooth and concave (Assumption 3), the following inequality holds:

$$\mathbf{J}(\theta^{k+1}) - \mathbf{J}(\theta^k) \geq \nabla \mathbf{J}(\theta^k)^\top (\theta^{k+1} - \theta^k) - \frac{L_J}{2} \|\theta^{k+1} - \theta^k\|^2 \quad (22)$$

Since  $\phi$  is assumed non decreasing in each component (Assumption 1), we can easily combine (21) and (22) to deduce:

$$\mathfrak{J}(\theta^{k+1}) - \mathfrak{J}(\theta^k) \geq \nabla \phi(\mathbf{J}(\theta^k))^\top (\nabla \mathbf{J}(\theta^k)) (\theta^{k+1} - \theta^k)^\top - \|\phi(\mathbf{J}(\theta^k))\|_1 \frac{L_J}{2} \|\theta^{k+1} - \theta^k\|^2 \quad (23)$$

Using Assumption 1 to bound the gradient of  $\phi$ , and given the definition of  $\mathbf{g}^{k,*}$  (11), and the update rule (13) the above inequality turns into:

$$\mathfrak{J}(\theta^{k+1}) - \mathfrak{J}(\theta^k) \geq \alpha_k \langle \mathbf{g}^{k,*}, \hat{\mathbf{g}}^k \rangle - \frac{L_J M_\phi}{2} \alpha_k^2 \|\hat{\mathbf{g}}^k\|^2 \quad (24)$$

Rewriting (24) with the error in (14), and using the condition on the learning rate (Assumption 5) would provide the result, since:

$$\begin{aligned} \alpha_k \langle \mathbf{g}^{k,*}, \hat{\mathbf{g}}^k \rangle - \frac{L_J M_\phi}{2} \alpha_k^2 \|\hat{\mathbf{g}}^k\|^2 &\geq \alpha_k \langle \mathbf{g}^{k,*}, \hat{\mathbf{g}}^k \rangle - \delta_k \|\hat{\mathbf{g}}^k\|^2 \\ &= (\alpha_k \beta_k - \delta_k (\beta_k)^2) \langle \mathbf{g}^{k,*}, \mathbf{g}^{k,*} \rangle + (\alpha_k - 2\delta_k \beta_k) \langle \mathbf{g}^{k,*}, \epsilon^k \rangle - \delta_k \|\epsilon^k\|^2 \\ &= \gamma_k \langle \mathbf{g}^{k,*}, \mathbf{g}^{k,*} \rangle - \delta_k \|\epsilon^k\|^2, \end{aligned}$$

where we set  $\delta_k = \frac{\alpha_k}{2\beta_k}$ ,  $\gamma_k = \frac{\alpha_k \beta_k}{2}$ . The first inequality comes from Assumption 5, as it implies  $\frac{L_J M_\phi}{2} \alpha_k^2 \leq \delta_k$ . As introduced previously, let us remind the constant  $\beta_k$  is the one stemming from the Policy Gradient Theorem, and corresponds in the episodic case to the average length of an episode.  $\square$

Before stating the main theorem, let us invoke a simple general lemma needed for evaluating the rate of convergence, in plain form and in probability.

**Lemma A.6.** (i) If a sequence  $u_k$  satisfies an equation of the form:

$$\forall k \leq k_0, \quad u_k - u_{k+1} \geq \mu u_k^2, \quad \text{with } \mu \in \mathbb{R}^+, u_0 > 0, \quad (25)$$

Then, it has a sublinear rate of convergence:

$$\forall k \leq k_0, \quad u_k \leq \frac{u_0}{(\mu u_0)k + 1}$$

(ii) If a sequence  $u_k$  satisfies:

$$\forall k \leq k_0, \quad \mathbb{P}(u_k - u_{k+1} \geq \mu u_k^2) > 1 - \iota, \quad \text{with } \mu \in \mathbb{R}^+, u_0 > 0. \quad (26)$$

Then, it has a sublinear rate of convergence with a certain probability:

$$\forall k \leq k_0, \quad \mathbb{P}\left(u_k \leq \frac{u_0}{(\mu u_0)k + 1}\right) \geq (1 - \iota)^k$$

*Proof.* (i) First, let us remark that the assumption (25) implies that  $(u_k)$  is a decreasing sequence. If the sequence reaches 0, the claim becomes trivial, so we would exclude this case in the proof, and focus on  $k$  such that  $u_k > 0$ . Exploiting the assumption in (25):

$$\frac{1}{u_k} - \frac{1}{u_{k-1}} = \frac{u_{k-1} - u_k}{u_{k-1}u_k} \geq \mu \frac{u_{k-1}}{u_k} \geq \mu. \quad (27)$$

By summing these inequalities, the following inequality holds:

$$\forall k \leq k_0, \quad \frac{1}{u_k} \geq \frac{1}{u_0} + k\mu \geq \frac{1 + k\mu u_0}{u_0}$$

It entails the claim (i) of the Lemma.

(ii) Transposing (27) to the probabilistic scenario, the assumption (26) provides the following inequality:

$$\mathbb{P}\left(\frac{1}{u_k} - \frac{1}{u_{k-1}} \geq \mu\right) > 1 - \iota$$

By considering the intersection of these events, we deduce:

$$\mathbb{P}\left(\frac{1}{u_k} - \frac{1}{u_0} \geq k\mu\right) \geq \mathbb{P}\left(\bigcap_{l=0}^k \left(\frac{1}{u_l} - \frac{1}{u_{l-1}} \geq \mu\right)\right) \geq (1 - \iota)^k$$

We can easily conclude. □

Let  $\beta^b$  denote a lower bound for the average length of an episode, and  $\alpha^b$  (resp.  $\alpha^\#$ ) a lower (resp. upper bound) on the learning rate.

**Theorem A.7.** (i) Under Assumptions (1, 2, 3, 4, 5, 6, 7, 8, 9), and assuming the size of the minibatch  $m$  tends towards infinity, the social welfare objective  $\mathfrak{J}(\theta^k)$  converges in probability and with a sub-linear convergence rate within a radius of convergence  $\tilde{\tau}$  of the optimal value where:

$$\tilde{\tau} = \frac{\Xi(2R_0 + \alpha^\#)}{2\beta^b}, \quad \text{with } \Xi := (L_\phi B_A \Xi_J + M_\phi \Xi_A + M_\phi \Xi_L) M_\pi.$$

More precisely, the suboptimality after  $k$  iterations may be bounded:

$$\lim_{m \rightarrow \infty} \mathbb{P}\left(\mathfrak{J}^* - \mathfrak{J}(\theta^k) \leq \max\left(\tilde{\tau}, \frac{j_0}{\kappa j_0 k + 1}\right)\right) = 1$$

$$\text{with } j_0 := \mathfrak{J}^* - \mathfrak{J}(\theta^0), \kappa := \frac{v-1}{v} \frac{\alpha^b \beta^b}{2R_0^2}, \text{ and } v := \left(1 + \frac{\alpha^\#}{2R_0}\right)^2$$

- (ii) Under Assumptions (1, 2, 3, 4, 5, 6, 7, 8, 9), the social welfare objective  $\mathfrak{J}(\boldsymbol{\theta}^k)$  converges, with a certain probability, within the radius of convergence  $\tilde{\tau}$  of the optimal value.  
More precisely, the suboptimality after  $k$  iterations may be bounded:

$$\forall \epsilon > 0, \forall k > 0, \quad \mathbb{P} \left( \mathfrak{J}^* - \mathfrak{J}(\boldsymbol{\theta}^k) \leq \max \left( \tilde{\tau}_\epsilon, \frac{j_0}{\kappa j_0 k + 1} \right) \right) \geq 1 - k \iota_{\epsilon, m} + o(\iota_{\epsilon, m}), \quad (28)$$

$$\text{with } \tilde{\tau}_\epsilon := \frac{(\Xi + \epsilon) (2R_0 + \alpha^\#)}{2\beta^b}.$$

- (iii) Similarly, under Assumptions (1, 2, 3', 5, 6, 7, 8, 9), if  $\mathbf{J}$  is simply locally concave, the SWO converges in probability and with a sublinear rate of convergence within the radius  $\tilde{\tau}$  of a local optimum, under the limit where  $m$  tends toward infinity. Under these same assumptions, in the general case, we can, as in (28) bound the suboptimality relative to a local optima with a certain probability.

This theorem implies that we can bound the number of iterations required to have  $\|\mathfrak{J}(\boldsymbol{\theta}^k)\| < \epsilon$  in probability:

**Corollary A.8.** Under Assumptions (1, 2, 3, 4, 5, 6, 7, 8, 9), with the previous notations:

$$\forall \epsilon \text{ s.t. } \tilde{\tau} < \epsilon < j_0, \quad \forall k \geq k_\epsilon, \quad \lim_{m \rightarrow \infty} \mathbb{P}(\mathfrak{J}^* - \mathfrak{J}(\boldsymbol{\theta}^k) < \epsilon) = 1, \quad \text{with } k_\epsilon := \frac{v}{v-1} \frac{2R_0^2}{\alpha^b \beta^b} \left( \frac{1}{\epsilon} - \frac{1}{j_0} \right).$$

Assuming that the terms  $\Xi_J, \Xi_A, \Xi_L$  stemming from Assumptions 7, 8, 9 may be as small as desired, the convergence radius would tend towards 0, under the limit  $m$  tends to infinity. However, in case of interdependent agents, particularly in cooperative or competitive settings, it is not reasonable to neglect the term stemming from Assumption 9. In contexts where the interdependence between agents is tamer, these different advantages (of an agent depending on other agents actions) could average themselves out in (17).

*Proof of Corollary A.8.* Directly following from the assertion (i) of Theorem A.7, since it is straightforward to check that  $\forall k \geq k_\epsilon, \frac{j_0}{\kappa j_0 k + 1} \leq \epsilon$ .  $\square$

*Proof of Theorem A.7.* As previously stated in Lemma A.2, under Assumptions (1, 3),  $\mathfrak{J}$  is concave. From the concavity of  $\mathfrak{J}$ , applying Cauchy–Schwartz inequality:

$$\mathfrak{J}^* - \mathfrak{J}(\boldsymbol{\theta}^k) \leq \|\nabla \mathfrak{J}(\boldsymbol{\theta}^k)\| \|\boldsymbol{\theta}^* - \boldsymbol{\theta}^k\| \leq R_0 \|\nabla \mathfrak{J}(\boldsymbol{\theta}^k)\|, \quad (29)$$

where the last inequality stems from Assumption 4. Indeed, Assumption 4 is satisfied for  $\boldsymbol{\theta}^0$  and it is easy to prove by recursion, that it holds for all the following  $\boldsymbol{\theta}^k$ 's; given (30), we can prove recursively  $\mathfrak{J}(\boldsymbol{\theta}^k) \geq \mathfrak{J}(\boldsymbol{\theta}^0)$ .

Combining (29) with Lemma A.5:

$$\mathfrak{J}(\boldsymbol{\theta}^{k+1}) - \mathfrak{J}(\boldsymbol{\theta}^k) \geq \kappa_k (\mathfrak{J}(\boldsymbol{\theta}^k) - \mathfrak{J}^*)^2 - \mathcal{E}_k = \kappa_k \left[ (\mathfrak{J}(\boldsymbol{\theta}^k) - \mathfrak{J}^*)^2 - \mathfrak{d}_k^2 \right] \quad \text{where} \quad \begin{cases} \kappa_k &= \frac{\alpha_k \beta_k}{2R_0^2} \\ \mathcal{E}_k &= \frac{\alpha_k}{2\beta_k} \|\epsilon^k\|^2 \\ \mathfrak{d}_k &= \sqrt{\frac{\mathcal{E}_k}{\kappa_k}} = \frac{\|\epsilon^k\| R_0}{\beta_k} \end{cases} . \quad (30)$$

Rewriting (30), with  $j_k = \mathfrak{J}^* - \mathfrak{J}(\boldsymbol{\theta}^k)$ :

$$j_k - j_{k+1} \geq \kappa_k (j_k^2 - \mathfrak{d}_k^2), \quad (31)$$

Below, we first proceed to the convergence analysis while assuming the error term  $\|\epsilon^k\|$  is bounded by  $\tilde{\Xi}$ , in order to familiarise the reader with the proof's argument. In a second step, we refine this analysis, in order to work in probability.



**Prior Analysis:** In this prior analysis, we assume the error term is bounded by  $\tilde{\Xi}$ . It enables us to define an upper bound for  $\mathfrak{d}_k$ :

$$\mathfrak{d}_k \leq \mathfrak{r} := \frac{\tilde{\Xi} R_0}{\beta^b} \quad (32)$$

From (30), we deduce that the sequence  $(\mathfrak{J}(\boldsymbol{\theta}^{k+1}))_k$  is increasing at least until it reaches a distance  $\mathfrak{r}$  from the optimal value  $\mathfrak{J}^*$ . Moreover, by setting  $\mathfrak{r}^v := \sqrt{v}\mathfrak{r}$ , combining (31) and (32), we deduce:

$$\forall v > 1, \forall k, \quad j_k \geq \mathfrak{r}^v \Rightarrow j_k - j_{k+1} \geq \kappa_k \left( j_k^2 - \frac{(\mathfrak{r}^v)^2}{v} \right) \geq \frac{v-1}{v} \kappa_k j_k^2 \geq \kappa^v j_k^2, \quad \text{with } \kappa^v := \frac{v-1}{v} \frac{\alpha^b \beta^b}{2R_0^2}. \quad (33)$$

From there, applying Lemma A.6 would enable us to roughly conclude that  $j_k$  converges with a sublinear rate of convergence towards the optimal value, at least until it reaches  $\mathfrak{r}$ . More precisely, fixing  $k^v$  such that  $j_k > \mathfrak{r}^v$  for all  $k < k^v$ , Lemma A.6 asserts:

$$\forall k \leq k^v, \quad j_k \leq \frac{j_0}{\kappa^v j_0 k + 1} \quad (34)$$

Since this is valid for all  $v > 1$ , and  $\lim_{v \rightarrow 1} \mathfrak{r}^v = \mathfrak{r}$ , it suggests the sequence  $j_k$  converges towards the radius  $\mathfrak{r}$ . However, this previous radius of convergence  $\mathfrak{r}$  is not guaranteed to be stable. Indeed, once reached the radius  $\mathfrak{r}$ , the right side of (31) would be negative, and we can not conclude the sequence  $(j_k)$  is still decreasing. However, we can rethink a stable radius of convergence. Let us choose  $k_1$ —assuming it exists—such that  $j_{k_1}$  lies within the radius  $\mathfrak{r}$ , but  $j_{k_1-1}$  does not. Then, rewriting (31):

$$j_{k_1} - j_{k_1-1} \leq \kappa_{k_1-1} (\mathfrak{d}_{k_1-1}^2 - j_{k_1-1}^2) \quad (35)$$

Therefore, from the previous definitions and assumptions:

$$\|j_{k_1}\| \leq \|j_{k_1-1}\| + \kappa_{k_1-1} \|\mathfrak{d}_{k_1-1}^2 - j_{k_1-1}^2\| \leq \mathfrak{r} + \kappa_{k_1-1} \mathfrak{d}_{k_1-1}^2 = \mathfrak{r} + \mathcal{E}_{k_1-1} \leq \frac{\tilde{\Xi}}{\beta^b} R_0 + \frac{\alpha_{k_1-1}}{2\beta_{k_1-1}} \tilde{\Xi} \leq \tilde{\mathfrak{r}}, \quad (36)$$

where  $\alpha^\#$  is an upper bound over the learning rate, and with our redefined radius of convergence:

$$\tilde{\mathfrak{r}} := \frac{\tilde{\Xi}}{\beta^b} \left( R_0 + \frac{\alpha^\#}{2} \right).$$

By (36),  $j_{k_1}$  lies within this radius. We claim that all the following values  $(j_k)_{k > k_1}$  stay within the radius  $\tilde{\mathfrak{r}}$ . Indeed, by (31), we can deduce that if  $j_{k_1} > \mathfrak{r}$ , the successive values  $(j_k)_{k > k_1}$  will be again decreasing at least until reaching the radius  $\mathfrak{r}$ . Therefore, even when some values in the sequence  $(j_k)$  ventures out of the original radius of convergence  $\mathfrak{r}$ , they will stay within the new radius of convergence  $\tilde{\mathfrak{r}}$ .

**Analysis in Probability** In order to finalise the theorem proof, we will have to adapt the above arguments in probability. First, let us examine the error term, which, from the decomposition (14), satisfies:

$$\|\boldsymbol{\epsilon}^k\| \leq \|\boldsymbol{\nu}^k\| + \|\boldsymbol{\eta}^k\|$$

A consequence of the Chebyshev inequality, upon assumption 6 bounding the variance of the estimated gradient:

$$\forall \epsilon > 0, \quad \mathbb{P}(\|\boldsymbol{\nu}^k\| < \epsilon) = \mathbb{P}(\|\hat{\mathbf{g}}^k - \tilde{\mathbf{g}}^k\| < \epsilon) \geq 1 - \frac{(\sigma^\#)^2}{\epsilon^2 m},$$

where  $m$  is the size of the minibatch on which the gradient  $\hat{\mathbf{g}}^k$  is estimated. Meanwhile, as demonstrated in Lemma A.3,  $\|\boldsymbol{\eta}^k\|$  may be bounded by  $\Xi$  in probability, under Assumptions 1, 2, 6, 7, 8, 9. Hence:

$$\forall \epsilon > 0, \quad \mathbb{P}(\|\boldsymbol{\epsilon}^k\| < \Xi_\epsilon) \geq 1 - \nu_{\epsilon, m}, \quad \text{with } \nu_{\epsilon, m} := \frac{(\sigma^\#)^2}{\epsilon^2 m}, \quad \text{and } \Xi_\epsilon := \Xi + \epsilon \quad (37)$$

In particular, upon taking the limit over  $m$ :

$$\forall \epsilon > 0, \quad \lim_{m \rightarrow \infty} \mathbb{P}(\|\boldsymbol{\epsilon}^k\| < \Xi_\epsilon) = 1. \quad (38)$$

Instead of the (32) from the prior analysis, we get the following probabilistic form:

$$\forall \epsilon > 0, \quad \mathbb{P}(\mathfrak{d}_k \leq \tau_\epsilon) \geq 1 - \iota_{\epsilon, m}, \quad \text{where } \tau_\epsilon := \frac{\Xi_\epsilon R_0}{\beta^b}. \quad (39)$$

By setting  $\tau_\epsilon^v := \sqrt{v}\tau_\epsilon$ , it induces the following —since  $j_k^2 \geq v\mathfrak{d}_k^2 \Rightarrow j_k^2 - \mathfrak{d}_k^2 \geq j_k^2 \frac{v-1}{v}$ :

$$\forall \epsilon > 0, \forall v > 1, \forall k, \quad \mathbb{P}\left(j_k^2 - \mathfrak{d}_k^2 \geq \frac{v-1}{v} j_k^2 \mid j_k \geq \tau_\epsilon^v\right) \geq \mathbb{P}(j_k^2 \geq v\mathfrak{d}_k^2 \mid j_k \geq \tau_\epsilon^v) \geq \mathbb{P}(\mathfrak{d}_k \leq \tau_\epsilon) \geq 1 - \iota_{\epsilon, m}.$$

With such an inequality in hand, we deduce from (31):

$$\forall \epsilon > 0, \forall v > 1, \quad \mathbb{P}(j_k - j_{k+1} \geq \kappa^v j_k^2 \mid j_k \geq \tau_\epsilon^v) \geq 1 - \iota_{\epsilon, m}, \quad \text{where } \kappa^v := \frac{v-1}{v} \frac{\alpha^b \beta^b}{2R_0^2} \quad (40)$$

Notably, passing to the limit over  $v$  in (40), since  $\tau_\epsilon = \lim_{v \rightarrow 1} \tau_\epsilon^v$ :

$$\mathbb{P}(j_k \geq j_{k+1} \mid j_k \geq \tau_\epsilon) \geq 1 - \iota_{\epsilon, m} \quad (41)$$

The sequence  $j_k$  is therefore decreasing until it reaches  $\tau_\epsilon$ , with a certain probability. At the limit where  $m \rightarrow \infty$ ,  $j_k$  is almost surely decreasing, until reaching the radius  $\tau$ , since we can tie  $\epsilon$  to  $m$  such that  $\lim_{m \rightarrow \infty, \epsilon \rightarrow \infty} \iota_{\epsilon, m} = 0$ , e.g. with  $\epsilon = m^{-\frac{1}{3}}$ . Applying Lemma A.6 to the last inequality would enable us to roughly conclude that  $j_k$  converges with a sublinear rate a convergence with a certain (bounded) probability until it reaches a radius  $\tau_\epsilon$ .

More precisely, let us fix  $k_\epsilon^v$  such that  $\forall k < k_\epsilon^v, j_k > \tau_\epsilon^v$ . Then, similarly to (34), Lemma A.6 leads to:

$$\forall k \leq k_\epsilon^v, \quad \mathbb{P}\left(j_k \leq \frac{j_0}{j_0 \kappa^v k + 1}\right) \geq (1 - \iota_{\epsilon, m})^k = 1 - k\iota_{\epsilon, m} + o(\iota_{\epsilon, m}) \quad (42)$$

As this is valid for all  $v > 1$  (although it changes the factor  $\kappa^v$  of the convergence rate), it signifies the sequence  $j_k$  is decreasing until  $k$  and converges towards  $\tau_\epsilon$  with a certain probability.

However, as highlighted previously in (i), this radius  $\tau_\epsilon$  is not guaranteed to be stable. Let us fix  $k_{0, \epsilon}$  the first  $k$  such that  $j_k \leq \tau_\epsilon$  and  $k_{1, \epsilon}$  the first  $k > k_{0, \epsilon}$ , such that  $j_k > \tau_\epsilon$  again —assuming they exists. Similarly to (36), from (35):

$$j_{k_{1, \epsilon}} \leq j_{k_{1, \epsilon} - 1} + \kappa_{k_{1, \epsilon} - 1} \left\| \mathfrak{d}_{k_{1, \epsilon} - 1}^2 - j_{k_{1, \epsilon} - 1}^2 \right\| \leq \tau_\epsilon + \kappa_{k_{1, \epsilon} - 1} \mathfrak{d}_{k_{1, \epsilon} - 1}^2 = \tau_\epsilon + \mathcal{E}_{k_{1, \epsilon} - 1}, \quad (43)$$

From the assumption on the error term and definition of  $\mathcal{E}_k$ , with the previous notations:

$$\forall k, \quad \mathbb{P}\left(\mathcal{E}_k < \frac{\alpha^\#}{2\beta^b} \Xi_\epsilon\right) > 1 - \iota_{\epsilon, m} \quad (44)$$

Combining (43) and (44), and introducing the new convergence radius, it implies:

$$\mathbb{P}(j_{k_{1, \epsilon}} \leq \tilde{\tau}_\epsilon) \geq 1 - \iota_{\epsilon, m}, \quad \text{where } \tilde{\tau}_\epsilon := \frac{\Xi_\epsilon}{\beta^b} \left( R_0 + \frac{\alpha^\#}{2} \right).$$

By (40) and the initial assumption on  $k_{1, \epsilon}$ , the sequence  $j_k$  is now decreasing with probability  $(1 - \iota_{\epsilon, m})$ . Thereupon, the following value  $j_{k_{1, \epsilon} + 1}$  satisfies:

$$\mathbb{P}(j_{k_{1, \epsilon} + 1} \leq \tilde{\tau}_\epsilon) \geq \mathbb{P}(j_{k_{1, \epsilon}} \leq \tilde{\tau}_\epsilon) \mathbb{P}(j_{k_{1, \epsilon} + 1} \leq j_{k_{1, \epsilon}}) \geq (1 - \iota_{\epsilon, m})^2$$

Iterating such argument, would lead us to the following bound:

$$\forall k > k_{1, \epsilon}, \quad \mathbb{P}(j_k \leq \tilde{\tau}_\epsilon) \geq (1 - \iota_{\epsilon, m})^{k - k_{1, \epsilon}} \geq (1 - \iota_{\epsilon, m})^k. \quad (45)$$

This argument may be repeated for successive values of the sequence  $j_k$ , if such behaviors happens. Gathering the previous results, (42) and (45), leads us to the following claim:

$$\forall \epsilon > 0, \forall k > 0, \quad \mathbb{P}\left[j_k \leq \max\left(\tilde{\tau}_\epsilon, \frac{j_0}{\kappa^v j_0 k + 1}\right)\right] \geq (1 - \nu_{\epsilon, m})^k = 1 - k\nu_{\epsilon, m} + o(\nu_{\epsilon, m}), \quad (46)$$

under the condition of  $v$  satisfying:

$$\tau_\epsilon^v \leq \tilde{\tau}_\epsilon \iff \sqrt{v}R_0 \leq \left(R_0 + \frac{\alpha^b}{2}\right) \iff v \leq \left(1 + \frac{\alpha^b}{2R_0}\right)^2.$$

For the theorem, we adopted  $v = \left(1 + \frac{\alpha^b}{2R_0}\right)^2$ .

**Limit case:** As previously mentioned, by taking the limit  $m \rightarrow \infty$ , and  $\epsilon \rightarrow 0$ , such that  $\epsilon^2 m \rightarrow \infty$  (e.g.  $\epsilon = m^{-\frac{1}{3}}$ ),  $\nu_{\epsilon, m}$  would tend towards 0. Hence, (46) reveals that  $j_k$  converges in probability within the radius  $\tilde{\tau}$  with a sublinear rate of convergence, as claimed.

**Locally Concave Case:** The case of assuming  $J$  simply locally concave would lead to a similar convergence result, under the corresponding set of assumptions, yet for a local optimum on the neighborhood  $\mathcal{U}$ . Since the proof is almost identical to the previous case, it is left to the reader. □

## B. Additional Experimental Results

We present here the additional experimental results that were not included in the main paper. Through those experiments, we try to answer the following questions:

- (A) How do SOTO and FEN perform with unequal access to resources in a controlled domain?
- (B) Can SOTO cope with a different number of users and agents?
- (C) What if Basic could also learn from  $\hat{A}_i^{\text{IND}}$  then  $\hat{A}^{\text{SWF}}$  gradually with  $\beta$ ?
- (D) What if we use the same architecture proposed in FEN to optimize a SWF?
- (E) Where does the bad performance of FEN in some domains come from?
- (F) Does SOTO still outperform FEN if it is also allowed to share scores to neighbors?
- (G) Is Basic capable enough to find fair solutions in small problems?

**Question (A)** In those experiments, we evaluate how our method performs when there is an unequal access to resources. We modified the Matthew Effect environment to restrict one of the 3 resources to be only collectible by the first two agents. They have an additional observation to determine whenever the resource is restricted. Other agents neither observe nor collect those restricted resources.

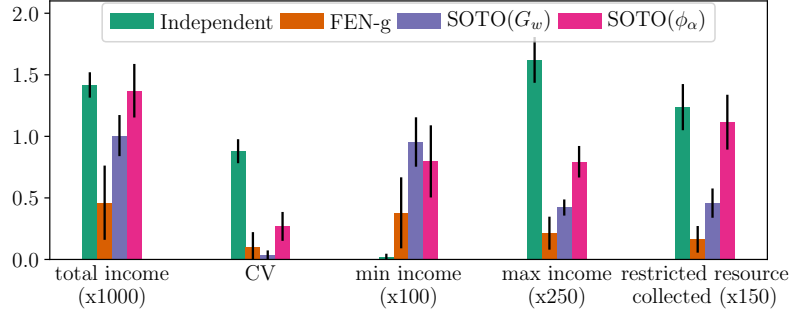


Figure 10: Comparison of Independent, FEN, and SOTO on Matthew Effect with unequal access to resources in the CLDE scenario. Higher total income is better, lower CV is better, higher min income is better, higher max income is better and higher restricted resource collected is better.

Compared to Matthew Effect without resource restriction (Figure 17), in Figure 10, as expected, Independent achieves a lower total income since it is more difficult to collect resources for the majority of the agents. FEN is Pareto-dominated by SOTO( $G_w$ ). Because of the unequal access to resources, the CV criterion has less importance in this experience: the first two agents should collect the maximum of restricted resources. This behavior is achieved by SOTO( $\phi_\alpha$ ), which collected a high number of restricted resources, which also Pareto-dominates FEN if we ignore the CV criterion. This shows that our method can cope with unequal access to resources whereas the first two agents of FEN do not collect more resources than the others.

**Question (B)** To answer (B), we turn to a more complex SUMO environment where we split the objectives per road instead of per intersection. The objective is to minimize the waiting time of the two roads of an intersection, hence each agent observes two users ( $2N = D$ ). Accordingly, the output size of the critics in the SOTO architecture is two.

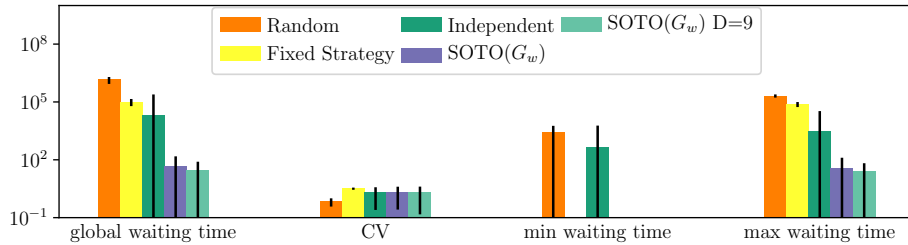


Figure 11: SUMO simulation with  $D=18$  in the CLDE scenario. Lower global waiting time is better, lower CV is better, lower min waiting time is better and lower max waiting time is better.

In Figure 11, we can see that our method copes with a different number of users and agents. However, this artificial splitting does not allow for a fairer solution than the one when  $D = 9$ .

**Question (C)** To verify that the SOTO architecture is necessary and the benefits do not come only from training with  $\beta$ , we performed another control experiment where Basic is also trained with the two different objectives:  $\hat{A}_i^{IND}$  and then  $\hat{A}^{SWF}$ .

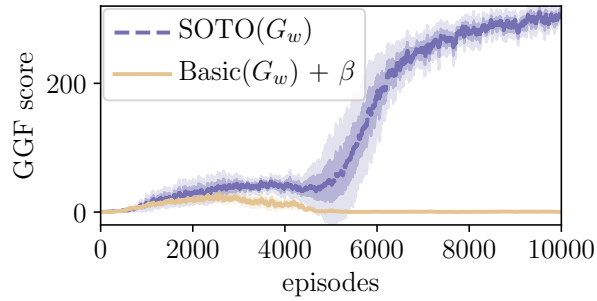


Figure 12: Comparison between SOTO and Basic updated by the two objectives with  $\beta$  on Matthew Effect.

Figure 12 can be analyzed in two parts. In the first stage, when  $\hat{A}_i^{IND}$  is mostly used to update the network, both approaches achieve an almost similar performance. However, the more  $\hat{A}_i^{SWF}$  is used, the more catastrophic forgetting intervenes. This demonstrates the importance of having two separate networks for each of the objectives.

**Question (D)** Here we want to verify if we could use the FEN architecture instead of the SOTO one. To do so, we extend FEN such that the meta-controller is rewarded by the current GGF value (FEN-g GGF). We also extended FEN such that the meta-controller is updated with  $\hat{A}_i^{SWF}$  (FEN-g  $A^{SWF}$ ).

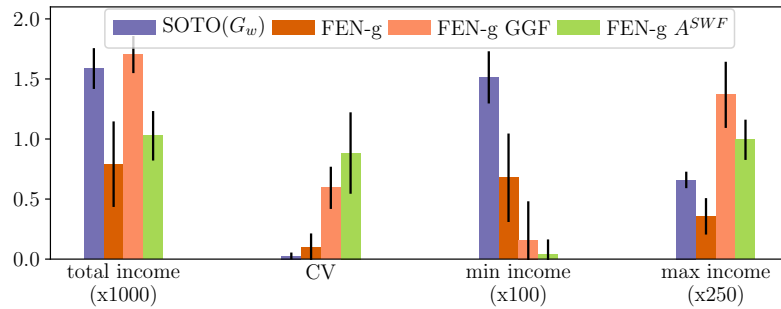


Figure 13: Comparison between SOTO, FEN without gossip, and FEN without gossip when it optimizes a SWF on the Matthew Effect. Higher total income is better, lower CV is better, higher min income is better, higher max income is better and higher restricted resource collected is better.

Figure 13 shows that both baselines are able to learn a good behavior (in terms of total income). However, none of them learn to achieve the better CV than FEN-g and SOTO. From this experiment, we conclude again that the SOTO architecture is essential to optimize a SWF without a centralized critic.

**Question (E)** Since FEN was not able to learn a good policy on the SUMO domain, we also tried to perform a pretraining of the first sub-policy as suggested by the authors for Plant Manufacturing. For the first 400 episodes, the first sub-policy is trained without interventions of the meta-controller.

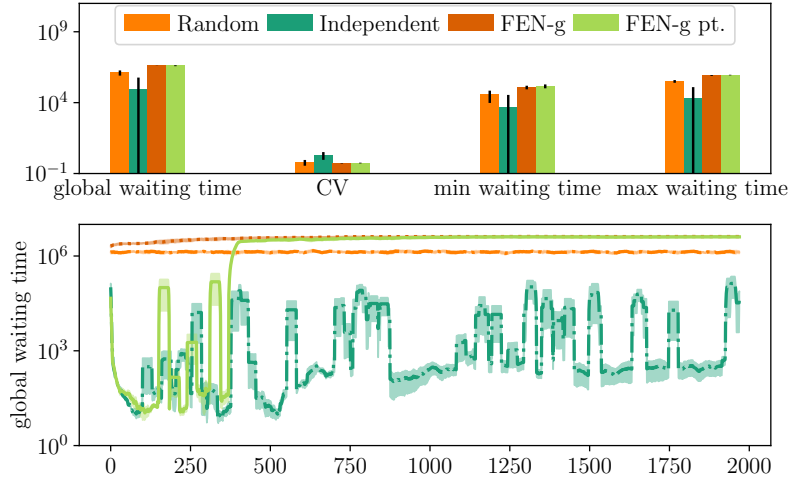


Figure 14: Comparison between FEN without gossip and FEN without gossip with pretraining (without gossip). Lower global waiting time is better, lower CV is better, lower min waiting time is better and lower max waiting time is better.

Figure 14 shows that during the first 400 episodes, the performance of FEN and Independent are indeed similar. However, when the meta-controller starts to switch other sub-policies the divergence begins. Note that when the pretraining is over, the first sub-policy is not updated anymore. So, if the performance drops it means that the meta-controller selects less and less often the first sub-policy. This is probably due to the fact that several agents select a wrong sub-policy at the same time, so the average performance is greatly reduced. Therefore, the other agents who chose the first sub-policy will also choose bad sub-policies for the next time to reduce the CV leading to a vicious circle until the worst possible average is reached. We observed this phenomenon in several domains.

**Question (F)** To verify that our better performances are not due to the fact that our networks have more input information, we provided FEN with the same additional inputs for its neural networks (i.e. the  $J(\theta)$  of the neighbors).

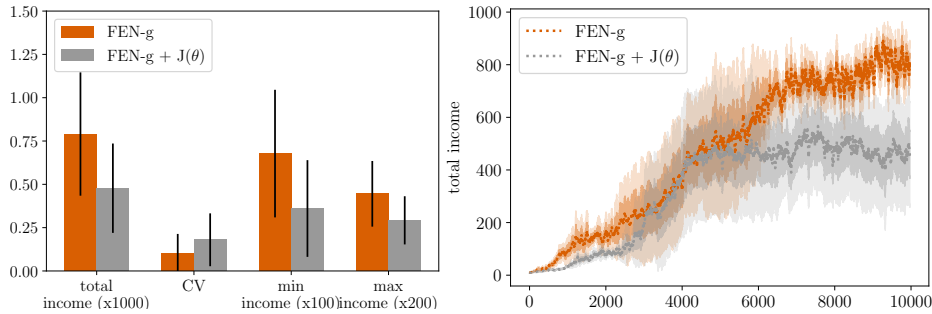


Figure 15: Comparison between FEN and FEN with more communications (without gossip).

From Figure 15, we can confirm that the additional communications are not helping FEN in Matthew Effect.

**Question (G)** In order to answer (G), we discuss the experimental results of Independent, FEN, Basic and SOTO on Job Scheduling in the CLDE scenario. Figure 16 illustrates the performances in terms of resource utilization, CV, minimum and maximum of utilities. As expected, Independent performs worse as it has the highest CV. Our algorithms achieve lower CV than FEN which indicates Basic and SOTO methods are able to find more fairer solution than FEN. Among all the algorithms, Basic( $G_w$ ) performs the best as it has the lowest CV. From these results, we can conclude that in simple domains Basic can find optimal solution without our proposed neural network architecture while in difficult task such as in Matthew Effect 17 and Plant Manufacturing 18 it performs worse.

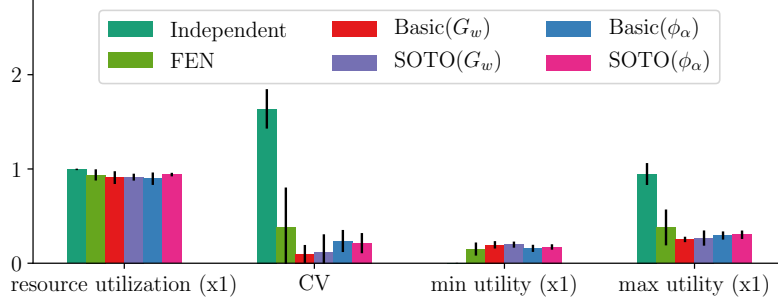


Figure 16: Comparison of Independent, FEN, Basic and SOTO on Job Scheduling in the CLDE scenario.

## C. Complete Experimental Results Per Domain

In this part, we provide the results on the Job Scheduling, Matthew Effect and Plant Manufacturing environments (Table 1, Figure 17 and Figure 18). We also present an additional plot to visualize the solutions on Matthew Effect in the FD scenario (Figure 4).

### C.1. Job Scheduling

Table 1: Job scheduling benchmark. The first part refers to the CLDE scenario. The second part refers to the FD scenario.

Method	Resource utilization	CV	min utility	max utility
Independent	1.00 $\pm$ 0.00	1.64 $\pm$ 0.19	0.00 $\pm$ 0.00	0.95 $\pm$ 0.11
FEN without gossip	0.93 $\pm$ 0.07	0.37 $\pm$ 0.42	0.15 $\pm$ 0.07	0.37 $\pm$ 0.19
Basic( $G_w$ )	0.91 $\pm$ 0.07	0.10 $\pm$ 0.10	0.20 $\pm$ 0.04	0.25 $\pm$ 0.03
SOTO( $G_w$ )	0.91 $\pm$ 0.04	0.12 $\pm$ 0.19	0.20 $\pm$ 0.03	0.27 $\pm$ 0.08
Basic( $\phi_\alpha$ )	0.90 $\pm$ 0.07	0.24 $\pm$ 0.12	0.16 $\pm$ 0.04	0.29 $\pm$ 0.04
SOTO( $\phi_\alpha$ )	0.94 $\pm$ 0.02	0.23 $\pm$ 0.10	0.16 $\pm$ 0.03	0.31 $\pm$ 0.04
Independent	1.00 $\pm$ 0.00	1.64 $\pm$ 0.19	0.00 $\pm$ 0.00	0.95 $\pm$ 0.11
FEN	0.72 $\pm$ 0.24	0.84 $\pm$ 0.38	0.04 $\pm$ 0.03	0.44 $\pm$ 0.26
FD Basic( $G_w$ )	0.79 $\pm$ 0.22	1.12 $\pm$ 0.46	0.02 $\pm$ 0.03	0.57 $\pm$ 0.28
FD SOTO( $G_w$ )	0.87 $\pm$ 0.12	1.07 $\pm$ 0.42	0.02 $\pm$ 0.04	0.60 $\pm$ 0.23
FD Basic( $\phi_\alpha$ )	0.71 $\pm$ 0.31	1.36 $\pm$ 0.42	0.01 $\pm$ 0.02	0.61 $\pm$ 0.34
FD SOTO( $\phi_\alpha$ )	0.94 $\pm$ 0.11	1.61 $\pm$ 0.26	0.00 $\pm$ 0.00	0.89 $\pm$ 0.18

### C.2. Matthew Effect

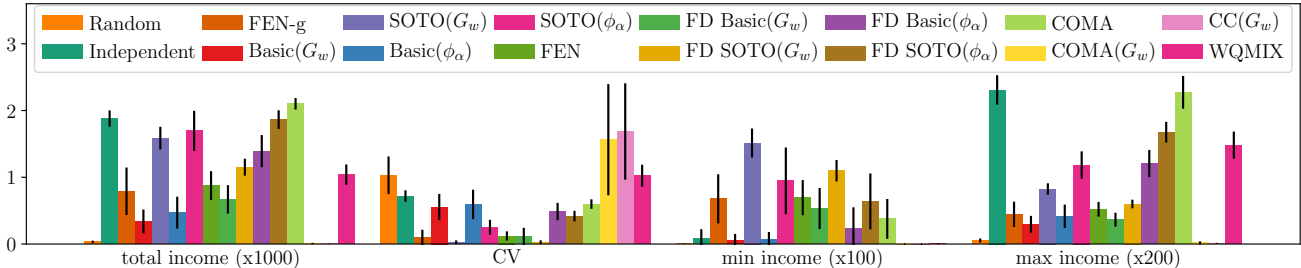


Figure 17: Comparison of different methods on Matthew Effect.

C.3. Plant Manufacturing

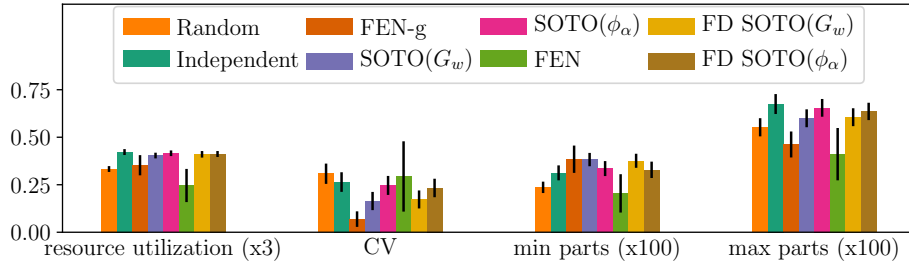


Figure 18: Comparison of different methods on Plant Manufacturing.

In Plant Manufacturing, Figure 18, in both scenarios, SOTO( $G_w$ ) and FEN are able to build the maximum number of products (equivalent to the minimum number of parts for each agent).

D. More details on experimentation

Note that compared to the results presented in (Jiang & Lu, 2019), our Independent method works almost twice better because we normalize the inputs, used Generalized Advantage Estimation (Schulman et al., 2016) instead of pure Monte-Carlo, and added bias units to the neural networks. All the presented methods benefit from those improvements.

In both scenarios (CLDE and FD), FEN is not restricted to send messages to neighbors for its gossip algorithm. To emphasize decentralization, all the agents are learning independent parameters without weights sharing (except for the COMA and WQMIX baselines).

The source code is available at <https://gitlab.com/AAAL/DFRL> with the configuration files used to generate our experimental results. The hyperparameters of the algorithms were optimized by grid search using Lightweight HyperParameter Optimizer (LHPO)<sup>10</sup>, an open source library used to run asynchronous distributed experiments (Zimmer, 2018).

The complete SOTO architecture (including critics) is provided in Figure 19.

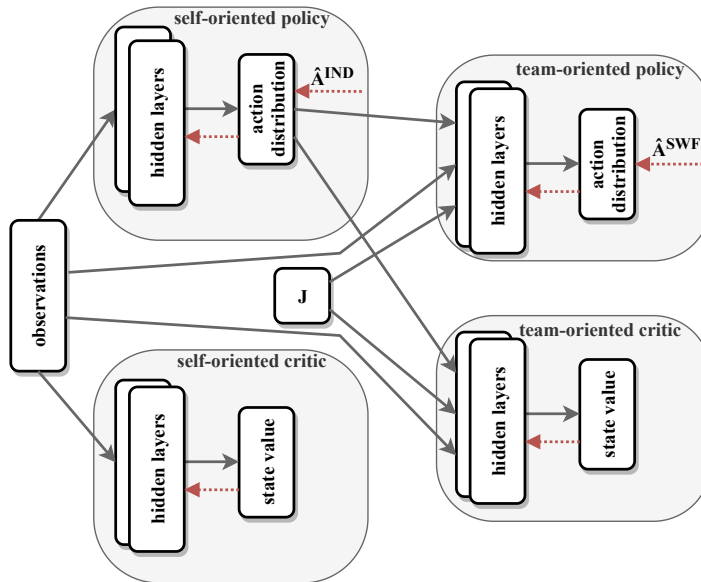


Figure 19: The complete SOTO architecture.

<sup>10</sup><https://github.com/matthieu637/lhpo>



### D.1. Hyperparameters

For all the experiments, we use PPO. The hidden layers of the neural networks are composed of two ReLU layers with 256 units. We used the ADAM (Kingma & Ba, 2015) optimization method with  $2.5 \times 10^{-4}$  as learning rate for the actor and  $1 \times 10^{-3}$  for the critic. The exploration bonus and the clipping ratio of the importance sampling in PPO are set to 0.03 and 0.1 respectively. In both scenarios, the FEN gossip algorithm uses  $g = 10$  gossip rounds and  $\tilde{k} = 3$  agents receiving the message. In every environment, we used Generalized Advantage Estimation with  $\lambda = 0.97$  except in the Job Scheduling environment where we kept pure Monte-Carlo estimation to reproduce the state-of-the-art results. The discount factor  $\gamma$  is 0.98 for Plant Manufacturing, Matthew effect and Job Scheduling, 0.99 for SUMO and 0.95 for Iroko. Batch sizes are respectively 25, 50, 50, 50 and 128 in Job Scheduling, Matthew Effect, Plant Manufacturing, SUMO and Iroko. In continuous action space (Iroko), we tuned the standard deviation of the Gaussian distribution for each instance of algorithms instead of learning it (so it is independent of the state). The best standard deviations for Independent, FEN-g, SOTO( $G_w$ ) and FD SOTO( $G_w$ ) are 0.5, 0.5, 0.05, and 0.1 respectively.

### D.2. Hardware

We performed the experiments on different computers equipped with 2 x Intel Xeon CPU E5-2678 v3 or Intel Core i7-8700 for SUMO and Iroko and 2 x Intel Xeon E5-2630 v3 or 2 x Intel Xeon Gold 6130 for the other domains.

Table 2: Comparison of the average computation times in hours.

Method	Job Scheduling	Matthew Effect	Plant Manufacturing	SUMO (D=9)	Iroko
Independent	9	27	8	20	5
FEN-g	13	37	64	60	6
Basic( $G_w$ )	11	30			
SOTO( $G_w$ )	11	34	12	14	7
Basic( $\phi_\alpha$ )	11	23			
SOTO( $\phi_\alpha$ )	14	29	11	27	
FEN	14	34	80+		
FD Basic( $G_w$ )	11	26			
FD SOTO( $G_w$ )	13	27	12	14	7
FD Basic( $\phi_\alpha$ )	14	26			
FD SOTO( $\phi_\alpha$ )	16	31	11		
Random		4	8	45	4
Fixed Strategy				45	4
COMA	13	23			
COMA( $G_w$ )	15	29			
CC( $G_w$ )	13	24			
WQMIX	24	80+		70	

In Table 2, the large computation time of FEN in Plant Manufacturing is due to the fact that the episode can last 10000 steps if all the 800 gems have not been collected.

### D.3. Environments

**Job Scheduling** In Job Scheduling, a permanent and unique resource is placed on a grid with 4 agents, which they must learn to share it. In Job Scheduling, neighbors are limited to being one block away from the agent (the number of neighbors may vary over time). The different neighborhoods correspond to the grid observations.

**Plant Manufacturing** In Plant Manufacturing, 3 types of gems are randomly placed in a grid along with 5 agents. Once a gem is collected, it reappears at a random location with a random type. Each agent needs a specific combination of the gem types to build a part. Once each agent have a part, a product is built. However, the agents are not aware of the complete definition of that reward function, they are rewarded for each gem and part, not for the full product.

In Plant Manufacturing, neighbors are limited to being two blocks away from the agent (the number of neighbors may vary

over time). The different neighborhoods correspond to the grid observations.

The reported parameter of the Plant Manufacturing between the original paper and the code provided by Jiang & Lu (2019) differed<sup>11</sup>. Thus, we run all the algorithms with the following parameters: 800 gems are available for an episode, the reward bonus for collecting a gem is 0.01, the size of a minibatch is 50, and the type of resource is encoded by a one-hot vector.

**Traffic Light Control** We evaluate our approach in classic traffic light control problems. In such problems the goal is usually to minimize the accumulated waiting time of all vehicles over all agents. We consider the accumulated waiting time for each agent at each intersection separately. We use Simulation of Urban MObility (SUMO)<sup>12</sup> to simulate a 3x3 intersection grid that has total of 9 agents where each agent controls the traffic light phase of one intersection. Each intersection has 4 roads and a total of 8 lanes with different numbers of vehicles. A traffic light phase controls the traffic and specifies which lanes have the green light. In our setting, we assume that each intersection have 4 phases. Depending on the type of intersections, different types of phase pattern can be considered.

The state space is composed of the current traffic light phase and for each lane, its queue and density of cars stopped at the intersection. An action corresponds to choose the current traffic light phase. The environment simulates intersections for 5000 seconds, with an episode length of 500 decision steps. At each time step, new vehicles enter into the intersection with a fixed destination.

The minibatch size used is 50 and  $\gamma = 0.99$ . The neighborhood definition is based on the position on the grid (the number of neighbors vary from 2 to 4).

The reward function is defined as the total waiting time at the current intersection subtracted from the accumulated waiting time at the last step. This definition of reward function motivates agent to minimize the accumulated waiting time which is our objective in this domain. We have performed experiments in both cases where the reward is the accumulated waiting time for each agent ( $D = N$ ) and where each agent’s reward is further split into two components corresponds to two roads of the network ( $D = 2N$ ).

Note that in this simulation if the agents’ decisions are too bad, a situation of complete blockage can occur. If this is the case, then regardless of the actions taken by the agents for the rest of the episode, traffic will remain blocked everywhere.

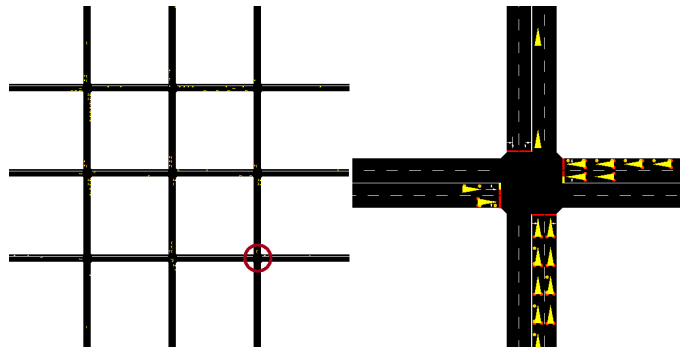


Figure 20: SUMO simulation with nine intersections in a grid. A zoom on one intersection is depicted at the bottom left. Cars can arrive from any direction and go to any of the 3 remaining ones.

**Data Center Control** In the Data Center Traffic Control problem, controllers manage the computer network that is shared by certain number of hosts. For the network topology, a fat-tree topology is considered (see Figure 21). In this topology, there are 16 hosts that are connected at the bottom with 20 switches. Each switch is connected with 4 ports, which results in a total of 80 queues in the network. We used Mininet<sup>13</sup> to simulate the network and Gobin<sup>14</sup> to generate network traffic. In classic data center control problems, the goal is usually to maximize the sum of host bandwidths in order to avoid the network congestion. However, we instead aim at maximizing the bandwidth of each host/agent in order to ensure fairness.

<sup>11</sup> Authors confirmed it in private communications. They also mentioned that a pretraining of the first sub-policy was required in this environment.

<sup>12</sup><https://github.com/eclipse/sumo>

<sup>13</sup><https://github.com/mininet/mininet>

<sup>14</sup><https://github.com/udhos/goben>

## Learning Fair Policies in Decentralized Cooperative Multi-Agent Reinforcement Learning

The minibatch size used is 128 and the  $\gamma = 0.95$ . We consider the neighbors as the 3 closest hosts which can be linked with a maximum of two switches. The total number of agents are 16, where each agent is dedicated to each host.

The global state is a  $n \times m$  matrix, where  $n$  is the number of ports in each switch and  $m$  is the number of network features collected by Goban. The continuous action corresponds to the allowed bandwidth for each host.

The  $D$ -dimensional reward vector is defined as follows:

$$\mathbf{r}_{\mathbf{a},s} = \mathbf{a} - 2 * \mathbf{a} * \max_i q_i(s)$$

where  $\mathbf{a}$  is the vector action that represents the bandwidth allocation and  $q_i(s)$  represents the  $i$ -th queue length. The reward is a vector, adapted from the (Ruffy et al., 2019), whose components are bandwidths per host penalized by the maximum of queue lengths. In the original environment, the average of  $\mathbf{a}$  is taken to define a scalar reward.

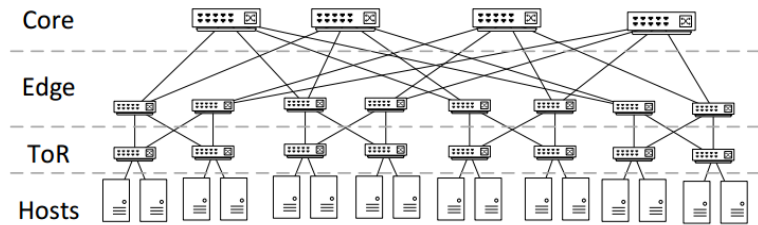


Figure 21: Network with a fat-tree topology from (Ruffy et al., 2019).