# DSDF: An approach to handle stochastic agents in collaborative multi-agent reinforcement learning

Satheesh K. Perepu and Kaushik Dey
Ericsson Research (Artificial Intelligence)
Chennai, Tamil Nadi
India 600096
`perepu.satheesh.kumar@ericsson.com, deykaushik@ericsson.com`

## Abstract

Multi-Agent reinforcement learning has received lot of attention in recent years and have applications in many different areas. Existing methods involving Centralized Training and Decentralized execution, attempts to train the agents towards learning a pattern of coordinated actions to arrive at optimal joint policy. However if some agents are stochastic to varying degrees of stochasticity, the above methods often fail to converge and provides poor coordination among agents. In this paper we show how this stochasticity of agents, which could be a result of malfunction or aging of robots, can add to the uncertainty in coordination and there contribute to unsatisfactory global coordination. In this case, the deterministic agents have to understand the behavior and limitations of the stochastic agents while arriving at optimal joint policy. Our solution, DSDF which tunes the discounted factor for the agents according to uncertainty and use the values to update the utility networks of individual agents. DSDF also helps in imparting an extent of reliability in coordination thereby granting stochastic agents tasks which are immediate and of shorter trajectory with deterministic ones taking the tasks which involve longer planning. Such an method enables joint co-ordinations of agents some of which may be partially performing and thereby can reduce or delay the investment of agent/robot replacement in many circumstances. Results on benchmark environment for different scenarios shows the efficacy of the proposed approach when compared with existing approaches.

## 1 Introduction

Multi-agent reinforcement learning (MARL) has been applied to wide variety of application which involve participation of collaborative agents such as Traffic management [2], power distribution [15], fleet management [12] etc. There are different ways in to achieve this collaboration. Some set of algorithms focus on learning a centralized policies [4, 10] while some on decentralized policies [24]. To improve the performance of the decentralized policies, some works assumed centralized training while learning these policies [22, 17, 13].

The first multi-agent methods were based on updation of Q-values arranged in form of table [9]. This approach is not scalable with increase in number of states and/or actions as one need to store all these values in memory. With the advent of deep networks, people started modelling value function using deep Q networks [11, 14, 3, 6, 7, 20]. In literature, people call these approaches under the name of deep Q learning.

Deep Reinforcement learning has been effective to deal with challenges of Multi-Agent coordination. The simplest technique was to forego a centralized coordination and let the agents learn independent value functions known as independent Q-learning (IQL) [22]. However the same does not handle non-stationarity introduced by actions executed by other agents. Although there exists some variants of IQl like [23], they miss some vital information on how other agents are performing. With introduction of deep Q-learning in Multi-Agent problems [21] complex environment representations were possible but the problems of non-stationarity still persisted. Subsequently, with introduction of centralized experience replay [4, 10] the problem of coordination started to be addressed using centralized approach. Further, some recent work involving COMA [5] which uses a centralized critic and actors with counterfactuals to address non stationarity

and credit assignment has proved to effectively address Multi-agent coordination. However this requires on-policy learning and thus more number of samples to learn optimal joint policy. Also, this method is not easily scalable with number of agents.

Sunehag et. al. [19] proposes a value decomposition method (VDN) which updates the Q-networks of individual agents based on sum of all the agents value function using centralized training and decentralized execution. However the same may not consider the extra state information of the environment and also it cannot be applied to all the general MARL problems and particularly where joint Q function is not a linear function of individual Q functions. To address this problem, Rashid et. al. [17] proposed a QMIX method which lies between extremes of VDN and COMA. The proposed approach uses a mixing network which mixes the individual agents value function through a mixing network which is then used to obtain $Q_{tot}$. Further the agents' value function are trained based on the $Q_{tot}$ and the mixing network is trained by conditioning the same on the state of the environment. In this way we can add the state of the environment while training of the individual agents and also the agents are trained based on other agents performance. However, the issue with this approach is that the value function of individual agents should monotonically increase w.r.t the $Q_{tot}$. To handle the stricter non-monotonic assumption, Mahajan et. al. [13] proposed a method known as MAVEN to train decentralized policies for agents condition their behavior on the shared latent variable controlled by a hierarchical policy.

All these above methods generally assume the agents behave exactly in the way as policy instructed it. However in many cases, the agents can behave in entirely stochastic way i.e. they execute the actions different that of the actions given by the policy. The degree of stochasticity can be different from different agents and over the time it can become constant. This is a common phenomena in industries where we learn policies using agents which when they get old and undergo wear and tear may not always be able to follow the strict demands of the policy. To explain it better let us consider the following example.

Let us assume there are robots in a warehouse performing some task. At the start of the experiment we have all the robots performing up to expectations and assume we train the agents to learn a collaborative policy to perform a task. However over the time robots can go under wear and tear and need not perform according to expectations. For example, at the start of the experiment applying an action "forward" to motors move the robot forward with a speed of 16 KM/HR. However, over the time motors can undergo wear and tear and motors may not perform well. Now, with the same action of "forward" results in moving some robots forward with a speeds varying between 8 KM/HR to 16 KM/HR while some may still perform per expectations. As one can see the joint policy learned for the previous case will no longer perform well for this case. It should be noted that only some of the agents behave in stochastic way as it will depend on the factors such as extent of usage, defects during manufacturing etc. One solution is to replace those stochastic agents but it is not recommended as it involves some cost with replacement.

Hence we need to retrain the policy from scratch when some of the agents are behaving in this stochastic way. Also, the changes in the agent won't be drastic and stays at the same degree of stochasticity for some period of time. Hence it can be concluded that the behavior of agents will be same for both training and execution phase and the policy learnt during training will work for execution phase also for a reasonable amount of time.

Intuitively, the success of coordination will depend on the deterministic (and good) agents understanding the limitations of stochastic agents, their behavior and tune the respective decentralized execution policies such that reward of overall joint policy is maximized. In this work, we achieve this coordination by using different discounted factors for different agents so that the respective dependencies of the value function on the future values is effected appropriately. For highly stochastic agents, we will use lower discounted factor so that their value function will depend on short-term actions and not so much on the long-term actions. For deterministic agents we will choose higher discounted factor to make them more depend on future values and let them plan effectively towards long term strategies.

In this paper we came with two methods to obtain the discounted factor for all the agents. One is the iterative penalization method which will penalize the discounted factor for the agent for every stochastic action taken. Another is to use a Deep Stochastic Discounted Factor (DSDF) method which will predict the discounted factor based on local observations and state of the environment. More details on both the methods are proposed in Section 3. For the sake of comparison, we compared the results with QMIX [17] and IQL [22] and corresponding plots are analyzed in Section 3.

## 2 Background

In this work, we assume a fully cooperative multi-agent task which is described with a decentralized partially observable Markov decision process (Dec-POMDP) which is defined with a tuple $G = < S, \mathbf{U}, P, r, Z, O, N, \gamma >$ where $s \in S$ describes the state of the environment [1]. At each step in time, each agent $i$ out of these $N$ agents will take an action $a_i$ and for all the agents the join action is represented by $U$. Due to the applied actions, the system will transit to $P(s'|s, \mathbf{U}) : \mathbf{S} \times \mathbf{U} \times \mathbf{S} \longrightarrow [0, 1]$. All the agents share common reward function $r(s, \mathbf{U} : \mathbf{S} \times \mathbf{U} \in \mathcal{R}$ and $\gamma : N \times 1 \longrightarrow [0, 1]$ is the discounted factor chosen for N agents.

Here we consider the environment is partially observable in which agent draws individual observation $z \in \mathbf{Z}$ according to a observation space $O(s, a) : \mathbf{S} \times \mathbf{A} \longrightarrow \mathbf{Z}$. Each agent can have their own observation history $\tau_i \in \tau : (\mathbf{Z} \times \mathbf{U}$ which influences the underlying stochastic policy $\pi^i(u^i|\tau^i)$. The joint policy $\pi$ has a join action-value function $Q^\pi(s_t, u_t) = E_{s_{t+1:\infty}, u_{t:\infty}}[R_t|S_t, u_t]$, where $R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$.

## 3 Proposed method

In this section, we describe the proposed approach to handle the stochastic agents in MARL approach. As described earlier, the proposed approach results in learning the most optimal joint policy of all agents even when some agents are stochastic. An important assumption here is "All the agents are deterministic while starting the experiment. However, after sometime the agents may become stochastic and we need to retrain the policy from scratch so that we address these stochastic agents. During the execution phase also we assume the same agents are stochastic with same degree of stochasticity. This is required to prevent the usage of continuous learning during execution."

The proposed approach is explained with QMIX [17] as collaborative mechanism between agents since it is considered as one of state of art collaborative mechanism. Although it can be extended to any other collaborative mechanisms which satisfies centralized training and decentralized execution.

In QMIX we compute the utility function for agent $i$ out of all $N$ agents as $Q_i(\tau_i, a_i)$, where $\tau_i$ is the history of observation space and action for agent $i$. The agent utility values $Q_i(\tau_i, a_i)$ are obtained by sending the observation spaces $(o_i, a_i)$ through a network with parameters $\theta^i$, where $i$ is the agent index out of $N$ agents. The obtained utility values $Q_i(\mathbf{o}_i, a_i) \ \forall \ i = 1, \cdots, N$ are mixed using a mixing network to obtain combined mixing value $Q_{\text{tot}}(\tau, \mathbf{U})$. Now, each agent $i$ network $\theta^i$ is updated by total value function $Q_{\text{tot}}$ instead of local utility function $Q_i(\mathbf{o}_i, a_i)$. The advantage here is that each agent network will have information on reward obtained and also indirectly other agent performance. In this way, we can incorporate other agents information without actually collecting them and able to arrive at joint policy.

In normal circumstances where all the agents are deterministic the agents together learn a joint policy since all the agents will know how other agents will work. However in the case where some agents are stochastic, the deterministic agents have to adjust to other agents stochastic behavior and tune their own behavior to arrive at good joint policy. This can be achieved by tuning the value function of stochastic agents to depend only on the current values and not on the future values. This can be achieved through the tuning the value of discounted factor for each agent.

However, in existing approach it assumes all agents are deterministic and we use a single discounted factor. However when some of the agents are stochastic using a single discounted factor can often results in poor update of network. For example, if the agent is highly stochastic, then it will take random actions with high probability which can be different from the action chosen by the policy. Now, if we use higher discounted factor here to update the network, it will result in poor update as the future actions can be totally different from planned actions (even in the case of exploration) with high probability and the same would upset the planned coordination of all the agents. Now, ideally we need to use smaller value of discounted factor for the higher stochastic agents and vice-versa to restrict the value function of those stochastic agents to short-term value rather than the future long-term value.

Hence, in this work we propose to compute the discounted factor for each agent based on the current observation and also the global state of the environment. In this work we propose two methods to compute the optimal $\gamma_i$ value for agent $i$. (i) Iterative penalization method where we penalize the discounted factor of each agent for every stochastic action taken and (ii) Fully connected network to compute the optimal
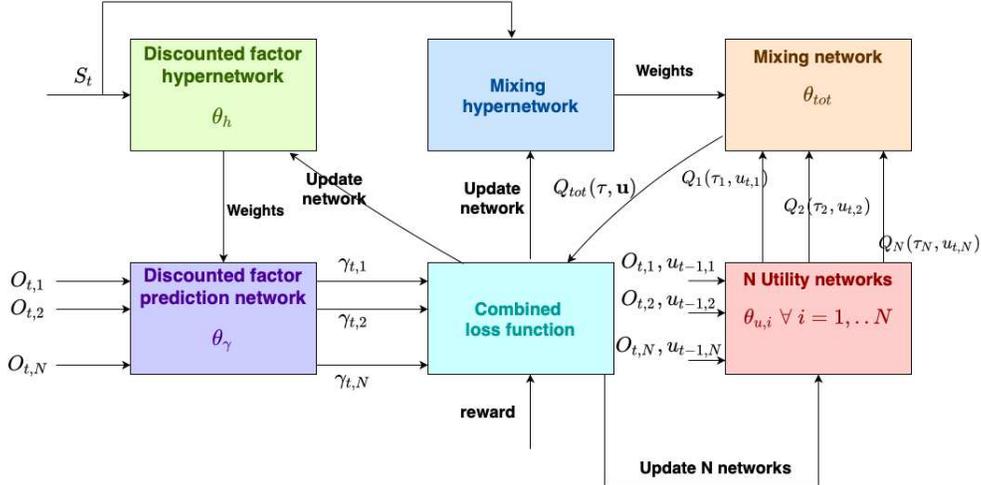
Figure 1: Proposed approach to compute discounted factor

discounted factors for each agent. For this we train a fully connected neural network which will output the discounted factor values for each agent $i$. In both the above cases, the computed discounted factor is inputted to the individual utility networks and use the computed discounted factor to compute the $y$, where $y = r + Q_{tot}(s', \mathbf{u}, \theta^-, \gamma)$ $\theta^-$ is the target network. The previous equation can be rewritten as $y = r + g(s', \gamma_i Q_i(o_i, a_i), \mathbf{u}, \theta_{tot})$, where $\theta_{tot}$ is the parameters of mixing network. Now, we update both the utility networks, mixing network (hypernetwork to be exact) using the different discounted factors for different agents. In this way, we can handle stochastic agents.

## 3.1 Proposed methods to calculate optimal $\gamma$ for all the agents

In this work, we propose two methods to compute the optimal $\gamma$ for all the agents. Below we explain them in detail.

### 3.1.1 Iterative penalization Method

In this method, we assume the discounted factor is 1 for all the agents during the starting of the experiment i.e. during retraining from scratch. If the action executed by the agent $i$ is different from that of action given by the policy we will penalize the discounted factor for the agent $i$ by a factor $P$. At every one time step, with every mismatch we will decrease the discounted factor with the factor $P$. Now, we will use the latest discounted factor at the time step to compute the utility function. Finally, we use the regular QMIX approach to update the utility networks and mixing networks.

The penalization factor $P$ should decrease with time steps just like we do in exploration factor in $\epsilon$-greedy method. The choice of optimal value for $P$ can be itself posed as optimization problem which is out of scope of the paper. However, the method proposed in Section 3.1.2 does not require any hyperparameter selection.

### 3.1.2 Deep Stochastic Discounted Factor (DSDF) Method

In this case, we proposed a method to compute the optimal $\gamma_i$, $i = 1, \cdots, N$ using a trained network. The proposed method is shown the Figure 1.

The idea is to utilize the agent local observations $o_{t,i}$ and global state $s_t$ at time instant $t$ of the environment to compute the discounted factor for individual agents. The reason behind is explained below:

Since we assume the underlying collaborative mechanism works for deterministic agents, the only issue behind poor results is due to the stochasticity of agents. Each agent have their perspective of the global state in form of local observations. Hence local observations will have an idea on the stochasticity of the agent and we need to estimate the discounted factor depending on the local observations and global state of the environment.

4

Returning to the main discussion, we use the local observations of agents to estimate the discounted factors using a fully connected network. However, we require additional information about the global state to the process. Since the local observations and global state are in different scale, we cannot combine together in same network. Also, since we require the discounted factor values between 0 to 1 from the network, we cannot train this network directly. Hence we use the concept of hypernetwork described in [8].

The local observations of all $N$ agents are sent to the network $\theta_\gamma$ which will compute the discounted factor values $\gamma_i, \quad i = 1, \cdots, N$. We utilize the hypernetwork $\theta_h$ to arrive at the network weights $\theta_\gamma$. The training process to update the $\theta_h$ is explained below.

As one can see the $\theta_h$, $\theta_u$ (utility networks) and $\theta_m$ are interlinked with each other. The general loss function of QMIX is

$$\mathcal{L}(\theta) = \sum_{t=1}^{B} \left( y^t - Q_{tot}(\tau, \mathbf{u}, s : \theta) \right) \tag{1}$$

where $\theta$ is the set of parameters of $N$ utility agents and mixing hypernetwork. Now, if we expand $y^t$

$$y^t = r + \gamma \max_{u'} Q_{tot}(\tau', \mathbf{u}', s' : \theta^-) \tag{2}$$

Now, instead of using single $\gamma$ value we will take the $\gamma$ inside the equation and we use our predicted network to compute the discounted factor.

$$y^t = r + \max_{u'} g(\mathbf{u}', s', \gamma_i Q_{u,i}, \gamma_{tot}) \tag{3}$$

Here $g(.)$ is the mixing network architecture which is parametrized by $\theta_{tot}$, $Q_{u,i}$ is the individual utility function of agent $i$. Now, we will replace the $\gamma_i$ with output of the network $\theta_\gamma$. The replaced equation is

$$y^t = r + \max_{u'} g(\mathbf{u}', s', f_\gamma(o_1^t, \cdots, o_N^t, \theta_\gamma) Q_{u,i}, \theta_{tot})$$
$$= r + \max_{u'} g(\mathbf{u}', s', f_\gamma(o_1^t, \cdots, o_N^t, f_h(\theta_h, s')) Q_{u,i}, \theta_{tot}) \tag{4}$$

where $f_\gamma$ is the discounted factor hyper network which is parametrized by $\theta_\gamma$ and $f_h$ is the hypernetwork function which is parametrized by $\theta_h$.

Replacing the value of $y_i^{tot}$ from (4) with that of (1) we obtain the loss function as

$$\mathcal{L}(\theta, \theta_h) = \sum_{t=1}^{B} \left( r^t + \max_{u'} g(\mathbf{u}', s', f_\gamma(o_1^t, \cdots, o_N^t, f_h(\theta_h, s')) Q_{u,i}, \theta) - Q_{tot}(\tau, \mathbf{u}, s : \theta) \right) \tag{5}$$

There are two unknown parameters in the above equation (i) $\theta$, (ii) $\theta_h$. Since the parameters are interdependent on each other i.e. $\theta_h$ on $\theta$ and vice-versa, we need to solve them in an iterative fashion. For every $B$ samples, first we will update the hyper network $\theta_h$ for fixed $\theta$ and then update $\theta$ for the computed $\theta_h$. So at each step we update the both $\theta_h$ and $\theta$ iteratively.

A important point to note here is that $\theta$ needs to be updated for every batch of samples as each batch will have new information on the environment. However, since the degree of stochasticity of the agents are assumed to be constant the $\theta_\gamma$ network will converge after some iterations and hence we don't want to update the $\theta_\gamma$ after some batches. Hence, it may be noted that we are adding only one training step to the existing mechanism for some batches. It can be inferred that we are adding only a small complexity and hence computational complexity would be closer to the existing mechanism.

The proposed DSDF method to estimate the discounted factor with existing QMIX method is explained in algorithm 1. The proposed DSDF has the following two advantages when compared with iterative penalization method.

1. DSDF method can dynamically change the discounted factor whereas the simpler method of penalizing it will only decrease it. Also, for the deterministic agents, the discounted factor is always 1 and hence it may degrade those agents performance for accidental deviation.

2. DSDF method calculates the discounted factor of an agent relative to other agents. Hence it can consider complex interactions between agents while deciding discounted factor for the agents.

**Algorithm 1** DSDF method to estimate discounted factor with QMIX

---

**Require:** Initialize parameter vector $\theta_h$, hypernetwork parameters and $\theta$ (agents utility networks, maxing network, hyper network), Learning rate $\leftarrow \alpha_\gamma$ and $\alpha_\theta$, $\mathcal{B} \leftarrow \{\}$

**Require:** step $= 0$, $\theta^- = \theta$

  **while** step $<$ step$_{max}$ **do**

    $t = 0$, $s_0 =$ Initial state

    **while** $t \neq$ terminal and $t <$ episode limit **do**

      **for** each agent $i$ **do**

        $\tau_t^i = \tau_{t-1}^i \cup \{(o_t, u_{t-1}\}$

        $\epsilon =$ epsilon-schedule (step)

$$u_a^t = \begin{cases} \underset{u_t^i}{\text{argmax}} \ \ Q(\tau_t^i, u_t^i) & \text{with probability } 1 - \epsilon \\ \text{randint}(1, |U|) & \text{with probability } \epsilon \end{cases}$$

      **end for**

      $s_{t+1} = p(s_{t+1}|s_t, \mathbf{u}_t)$

      $\mathcal{B} = \mathcal{B} \cup \{(s_t, \mathbf{u}_t, r_t, s_{t+1}\}$

      $t = t + 1, step = step + 1$

    **end while**

    **if** $|\mathcal{B}| >$ batch-size **then**

      b $\leftarrow$ random batch of episodes from $\mathcal{B}$

      **if** $\theta_h$ not converged **then**

        Update $\theta_h = \theta_h - \alpha_\gamma \nabla_{\theta_h}(\Delta Q_{tot})^2$

      **end if**

      Update $\theta_\gamma = f_\gamma(O, \theta_h)$, where $O$ is the set of observations for all agents in the sampled batch.

      Update $Q_{tot}$ using the latest updated $\theta_\gamma$.

      Update $\theta = \theta - \alpha_\theta \nabla_\theta(\Delta Q_{tot})^2$

    **end if**

    **if** update-interval steps have passed **then**

      $\theta^- = \theta$

    **end if**

  **end while**

---

Table 1: Type of agents used in simulation

| Agent Index | Type of Agent | Degree of stochasticity |
|:-:|:-:|:-:|
| 1 | Deterministic | NA |
| 2 | Stochastic | 0.2 |
| 3 | Deterministic | NA |
| 4 | Deterministic | NA |
| 5 | Stochastic | 0.4 |
| 6 | Stochastic | 0.6 |

# 4 Results and Discussion

We performed experiments on the modified lbforaging [16] using the method discussed in Section 3.1.2 and results are presented in Section 4.2. For the sake of stochasticity, we added a component which will decide whether to execute the action given by the policy or to take random action.

## 4.1 Environment

This environment contains agents and food resources randomly placed in a grid world. The agents navigate in the grid world and collect food resources by cooperating with other agents. We modified it by adding a couple of additions to the environment.

1. We added two additional actions 'Carry on' and 'Leave' to the agents. The 'Carry on' action enables the agent to store the food resources which they can subsequently leave in another time step and/or state for consumption of another agent. The 'Leave' action enables to agent to drop the resources which they consumed.

2. Each agent is given a target of the resources they need to consume. For this, we modify the reward function by adding targets to the them. Our eventual goal is to ensure all agents reach their targets. This means if some of the agents consumed in excess they must realize and give up the extra resources for benefit of other agents.

We placed 100 food resources in the $30 \times 30$ grid and we chosen 6 agents to eat these resources. The episode is terminated either when there are no food resources available in the grid or number of time steps reached 500. Here out of 6 agents, three agents 3 are deterministic (1,3,4) and 3 are stochastic (2,5,6). The level of stochasticity $\beta$ means the will perform actions given by the policy with $1 - \beta$ probability and will perform random actions with $\beta$ probability. The type of agents along with degree of stochasticity are shown in Table 1 We chosen the targets for the individual agents for two scenarios (i) When there are correct amount of food resources to achieve targets and (ii) When there are not enough resources to achieve the targets.

## 4.2 Results on Environment

In this section, we performed experiments on the above two cases with both the proposed DSDF method, iterative penalization method, QMIX method and IQL method. For this we utilized pymarl library used in [18]. The individual targets for respective agents are $20, 20, 30, 60, 30$ and $60$. For both cases, we restrict sum of all food levels available in the grid to a number $T$ such that total sum of targets for individual agents is equal to $T$ in Case 1 and $> T$ in Case 2. **An important point to be noted that the value of $T$ is not known to any agent during training or in execution.**

   **Case 1 Enough Resources**: Here we chosen sum of targets to be $T = 220$. Now, the agents have to collaborate within themselves to reach their respective targets.

   The predictions of the discounted factor using DSDF method for all the time steps is shown in the Figure 2a. As you can see the agents discounted factor changes whenever we update the discounted factor network. From the plot, we can observe the values got almost saturated after some update instants which shows our hypernetwork is converging and hence after this time we don't want to update the discounted factor hyper network.
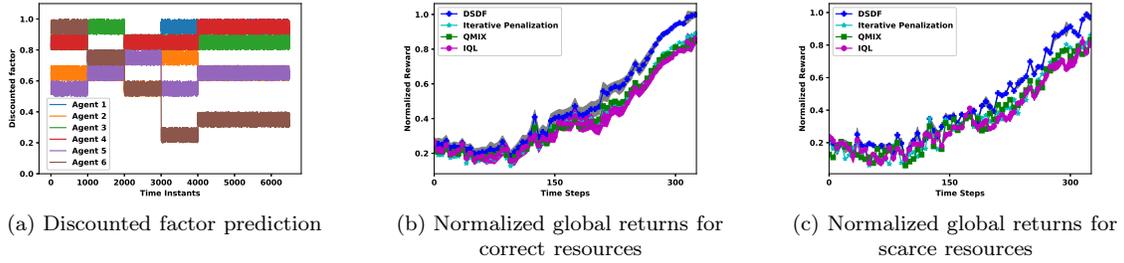
(a) Discounted factor prediction     (b) Normalized global returns for correct resources     (c) Normalized global returns for scarce resources

Figure 2

Also, one can observe the discounted factor values for deterministic agents are higher ($>= 0.9$) which suggests we can make the respective utility functions depend more on the future values. This is quite evident as these agents have to look more into future and decide current actions. On the other hand stochastic agents should choose lower discounted factor so that they won't take future values into consideration. From Figure 2a, one can observe that the agents with higher degree of stochasticity i.e. they have less probability of executing actions given by the policy should have lesser discounted factor. This is agreeing with our assumption that the higher stochastic agents require lower discounted factor and vice-versa.

Regarding the agent's performance, we gave the agents performance in last 10 training episodes and 10 execution episodes and results for individual agents are shown in Figure 3. In all the plots, first 10 indices correspond to latest 10 training episodes and next indices correspond to 10 execution episodes. From the plots, it can be concluded that deterministic agents except agent 4 trained using the proposed method reached their targets in all the execution episodes. On the other hand, the deterministic agents trained with vanilla QMIX and IQL have also reached their respective targets. The performance of the stochastic agents trained using the proposed DSDF method as well as iterative penalization method is well better than the vanilla QMIX and IQL. All the stochastic agents with smaller targets reached their targets while the agent with higher target settled closer to the target value. For the case of deterministic agents, the agents trained using proposed DSDF method either performed well or shown comparable performance when compared with IQL and QMIX. On the other hand, the agents trained with iterative penalization method shown equal performance with QMIX and IQl. This is due to the fact that value of discounted factor is 1 for deterministic agents, which means the agents depend more on the future and because of this it may spoil the agents performance.

The mean reward obtained for every time step with 95% confidence interval during evaluation is shown in Figure 2b. From the plot, it is evident that the proposed DSDF method resulted in higher average reward when compared to iterative penalization method and other methods. The agents trained with iterative penalization method also fared comparably well when compared with QMIX and IQL methods.

**Case 2: Scarce Resources**: In this case, the sum of targets $T$ is chosen to be 200 i.e. we do not have enough resources to satisfy all the agents. In this case, we are skipping the discounted factor plot as it resembles to the one obtained in Figure 2a.

Since we do not have enough resources the deterministic agents have to achieve all the targets and have to go beyond their targets since the dependency of the global reward on deterministic agents have to be made higher than stochastic agents. The agents performance for both the training and execution episodes are shown in Figure 4. From the plots, one can observe all the deterministic agents trained with DSDF method reached their target and settled beyond their targets in the execution episodes. On other hand, agents trained with QMIX and IQL methods settled below the proposed method showing the efficacy of the proposed method. The proposed iterative penalization method also gave good results for stochastic agents with higher degree of stochasticity.

Figure 2c shows the mean reward obtained with 95% confidence interval in this scarce resources case with the agents trained with different methods. From the plot it is evident that the DSDF method resulted in good mean reward when compared with existing methods and also the iterative penalization method. In addition, the agents trained with iterative penalization method also fare comparably well when compared with existing methods.

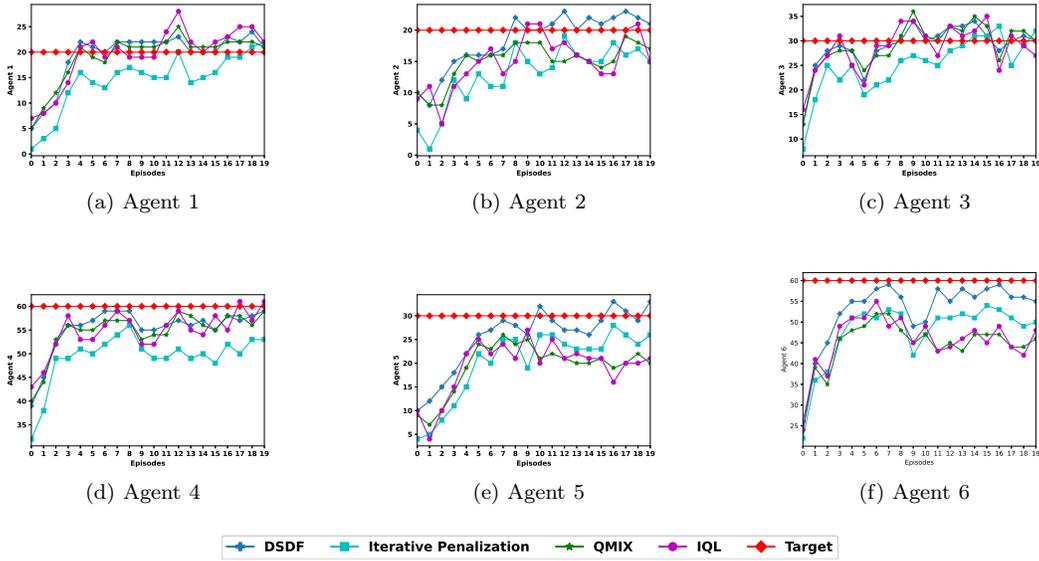|   |   |   |
|---|---|---|
| (a) Agent 1 | (b) Agent 2 | (c) Agent 3 |
| (d) Agent 4 | (e) Agent 5 | (f) Agent 6 |

Figure 3: Comparison of different methods for the case of correct resources

From the both scenarios, it can be concluded that the DSDF method results in better performance than the agents trained with QMIX and IQL method even when they are stochastic agents. The proposed iterative penalization method also gave good results for stochastic agents with higher degree of stochasticity. For the case of deterministic agents, the performance of the iterative penalization is greater than or equal to the existing methods.

## 5 Conclusion

In this paper, we propose a novel method DSDF to handle a mix of deterministic and stochastic agents which together learns a collaborative joint policy. Along with the above, an Iterative Penalization method is proposed, which though has lesser gains than DSDF. Our proposed methods can be combined with any state of art MARL algorithms without much impact to existing computation complexity.

The proposed method is tested on the lbforaging environment in two different scenarios and demonstrated clear improvement in results when compared with QMIX and IQL methods both for individual and global returns. Future directions include extending the technique to situations where environment might also be stochastic in regions ( like oil patches on floor). To our knowledge, this is the first time, such an complicated interaction involving stochastic and deterministic agents have been explored.

## References

[1] Daniel S Bernstein, Robert Givan, Neil Immerman, and Shlomo Zilberstein. The complexity of decentralized control of markov decision processes. *Mathematics of operations research*, 27(4):819–840, 2002.

[2] Tianshu Chu, Jie Wang, Lara Codecà, and Zhaojian Li. Multi-agent deep reinforcement learning for large-scale traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 21(3):1086–1095, 2019.

[3] Maxim Egorov. Multi-agent deep reinforcement learning. *CS231n: convolutional neural networks for visual recognition*, pages 1–8, 2016.

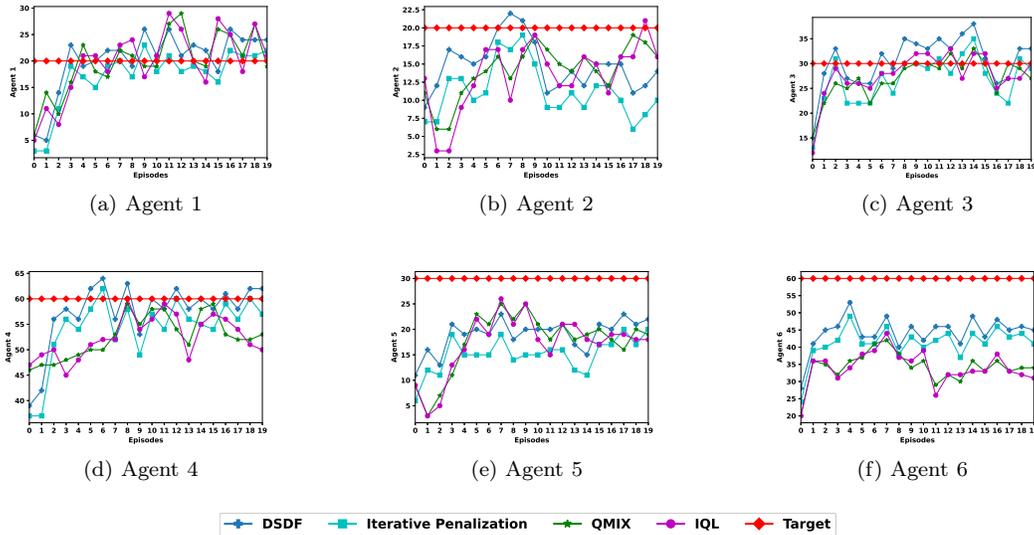(a) Agent 1  (b) Agent 2  (c) Agent 3

(d) Agent 4  (e) Agent 5  (f) Agent 6

Figure 4: Comparison of different methods for the case of scarce resources

[4] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[5] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[6] Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *International Conference on Autonomous Agents and Multiagent Systems*, pages 66–83. Springer, 2017.

[7] Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *International Conference on Autonomous Agents and Multiagent Systems*, pages 66–83. Springer, 2017.

[8] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.

[9] Junling Hu and Michael P. Wellman. Multiagent reinforcement learning: Theoretical framework and an algorithm. In *Proceedings of the Fifteenth International Conference on Machine Learning*, ICML '98, page 242–250, San Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc.

[10] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, DJ Strouse, Joel Z Leibo, and Nando De Freitas. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. In *International Conference on Machine Learning*, pages 3040–3049. PMLR, 2019.

[11] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

[12] Kaixiang Lin, Renyu Zhao, Zhe Xu, and Jiayu Zhou. Efficient large-scale fleet management via multi-agent deep reinforcement learning. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1774–1783, 2018.

[13] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent variational exploration. *arXiv preprint arXiv:1910.07483*, 2019.

[14] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.

[15] Yasar Sinan Nasir and Dongning Guo. Multi-agent deep reinforcement learning for dynamic power allocation in wireless networks. *IEEE Journal on Selected Areas in Communications*, 37(10):2239–2250, 2019.

[16] Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V Albrecht. Comparative evaluation of multi-agent deep reinforcement learning algorithms. *arXiv preprint arXiv:2006.07869*, 2020.

[17] 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. PMLR, 2018.

[18] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philiph H. S. Torr, Jakob Foerster, and Shimon Whiteson. The StarCraft Multi-Agent Challenge. *CoRR*, abs/1902.04043, 2019.

[19] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.

[20] Ardi Tampuu, Tambet Matiisen, Dorian Kodelja, Ilya Kuzovkin, Kristjan Korjus, Juhan Aru, Jaan Aru, and Raul Vicente. Multiagent cooperation and competition with deep reinforcement learning. *PloS one*, 12(4):e0172395, 2017.

[21] Ardi Tampuu, Tambet Matiisen, Dorian Kodelja, Ilya Kuzovkin, Kristjan Korjus, Juhan Aru, Jaan Aru, and Raul Vicente. Multiagent cooperation and competition with deep reinforcement learning. *PLOS ONE*, 12(4):1–15, 04 2017.

[22] Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *In Proceedings of the Tenth International Conference on Machine Learning*, pages 330–337. Morgan Kaufmann, 1993.

[23] Xing Xu, Rongpeng Li, Zhifeng Zhao, and Honggang Zhang. Stigmergic independent reinforcement learning for multiagent collaboration. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

[24] 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. PMLR, 2018.

This figure "Block_Diagram_NIPS.png" is available in "png" format from:

http://arxiv.org/ps/2109.06609v1

This figure "Discounted_Factor.png" is available in "png" format from:

http://arxiv.org/ps/2109.06609v1

This figure "Discounted_Factor_Correct.png" is available in "png" format from:

http://arxiv.org/ps/2109.06609v1

This figure "LB_Foraging.jpg" is available in "jpg" format from:

http://arxiv.org/ps/2109.06609v1

This figure "legend.png" is available in "png" format from:

http://arxiv.org/ps/2109.06609v1