

Learning Complex Group Behaviours in a Multi-Agent Competitive Environment

Muhammed Kazım Sanlav Afshan Nabi

1. Introduction

The aim of this project is to use reinforcement learning algorithms to explore predator-prey dynamics, such as those seen when predator fish and prey fish interact in nature. In order to do so, we have adapted OpenAI’s multiagent-particle-envs (Lowe et al.) to simulate two competing groups of fish- a predator group (has a single predator) and a prey group (2 or 10 prey). We used 2 different Q-learning algorithms (Q-Network with Experience Replay & Q-Network with Experience Replay and Fixed Q-Target) and a Monte-Carlo Policy Gradient algorithm to train these predator and prey agents. The predator is rewarded based on the number of times they collide with (“eat”) prey. The prey are rewarded negatively for each collision with the prey (“every time they are eaten”). Using such reward functions, the prey group is expected to learn defense strategies and the predator is expected to learn attack strategies.

2. Methods

Reinforcement learning is a branch of machine learning along with supervised and unsupervised learning. The aim of reinforcement learning is for an agent to interact with its environment. By interaction with the environment, the agent tries to learn the sequence of actions to take in particular states that lead to the maximization of a reward signal.

In addition to an environment and an agent, there are three main elements of a reinforcement learning system. A *policy* defines the way an agent behaves at a given time by mapping states of the environment to actions to be taken when in those states. A *reward signal* defines the goal of the reinforcement learning agent. At each step, the environment outputs a scalar reward to the agent. The aim of the agent is the maximize the long term reward it receives. A *value function* of a particular state is the total amount of reward the agent can expect to accumulate in the future starting from that particular state.

In the multi-agent setting, more than one agent is present and each agent learns independently of all other agents.

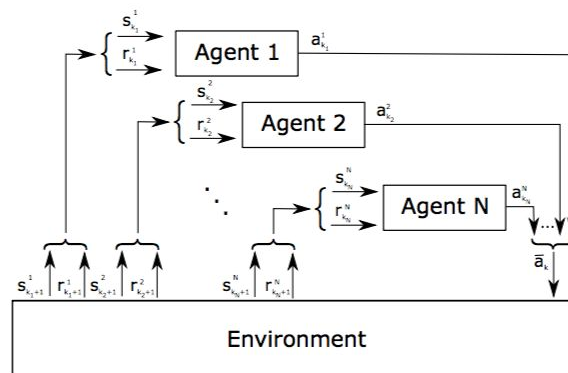


Figure 1: A multiagent environment (Chincoli et. al).

Each agent interacts with environment at discrete time steps from $t = 0, 1, 2, \dots$. At each time step, the each agent receives a state $S_t \in S$ and based on that, selects an action $A_t \in A$. As a consequence of taking the action, each agent receives a scalar reward $R_{t+1} \in R$ and a new state S_{t+1} . In order to perform well in the long-run, the agent needs to take into account the immediate reward as well as the future rewards it expects to get. The agent tries to choose action A_t in order to maximize the expected discounted return:

$$G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

where, γ is a parameter, $0 \leq \gamma \leq 1$, called the *discount rate*. The discount rate determines the present value of future rewards. If $\gamma = 0$, the agent is concerned only with maximizing the immediate rewards. If $\gamma = 1$, future rewards are taken into consideration equally. Returns at successive steps are recursively related to each other:

$$G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} \dots = R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) = R_{t+1} + \gamma G_{t+1}$$

A *stochastic policy* is a mapping from states to the probabilities of selecting each possible action. The *action-value function* $q_{\pi}(s, a)$ for policy π is the value of taking an action a in state s under policy π . It is the expected return when agent starts in state s , takes action a and subsequently follows policy π :

$$q_{\pi}(s, a) := E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

The action value functions can also be written recursively (*Bellman Expectation Equation* for q_{π}):

$$q_{\pi}(s, a) := E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, a_{t+1}) | S_t = s, A_t = a]$$

Solving a reinforcement learning problem means finding an *optimal policy* that achieves maximum reward. The action-value function for the optimal policy is called the *optimal action-value function* q_* :

$$q_*(s, a) := \max_{\pi} q_{\pi}(s, a)$$

The *Bellman Optimality Equation* for $q_*(s, a)$ can be written as:

$$q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a]$$

Once we know $q_*(s, a)$, the optimal policy can be followed just by choosing the action that maximizes $q_*(s, a)$ at each step.

2.1 Q-Learning

Q-learning solves the Bellman Optimality Equation iteratively. In Q-learning, a learned action-value function Q is used to directly approximate q_* , the optimal action-value function. In each step of each episode, Q is updated as:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') - Q(S_t, A_t)]$$

In case the problem is too large to learn all action values in all states separately, it is useful to learn a parameterized action value function $Q(s, a; \theta_t)$. The standard Q-learning update for the parameters after taking action A_t in state S_t and observing immediate reward R_{t+1} and resulting state S_{t+1} is:

$$\theta_{t+1} = \theta_t + \alpha(Y_t^Q - Q(S_t, A_t; \theta_t)) \nabla_{\theta_t} Q(S_t, A_t; \theta_t)$$

where α is the scalar step size and the target Y_t^Q is defined as:

$$Y_t^Q := R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t)$$

A Deep Q-Network (DQN) is a neural network that for a given state s outputs a vector of action values $Q(s, \cdot, \theta)$ where θ are the parameters of the network. For an n -dimensional state space and an m -dimensional action space, the DQN is a function from R^n to R^m . In this project, two algorithms using Q-Networks were implemented.

2.1.1 Q-learning with experience replay

A memory bank with a fixed size is initialized and the agent's experiences at each time step $e_t = [S_t, A_t, R_{t+1}, S_{t+1}]$ are stored in this memory bank $D = e_1, e_2, \dots$ over many episodes. This memory bank is called the *replay memory*. During learning, mini batch updates are applied to samples drawn at random from the replay memory.

Algorithm 1: Q-learning with experience replay

```

Initialize a replay memory to capacity N
Initialize action-value function Q with random weights  $\theta$ 
For episode 1:M do:
    Initialize environment to get initial state  $s_1$ 
    For step 1:T do:
        With probability  $\epsilon$  select random action  $a_t$ 
        Otherwise select  $a_t = \operatorname{argmax}_a Q(s_t, a; \theta)$ 
        Execute action  $a_t$  and observe reward  $r_{t+1}, s_{t+1}$ 
        Store experience  $s_t, a_t, r_{t+1}, s_{t+1}$  in replay memory
        Set  $s_{t+1} \leftarrow s_t$ 
        Sample a batch  $s_j, a_j, r_{j+1}, s_{j+1}$  from replay memory
        Set target  $Y_j = r_j + \gamma \max_{a'} Q(s_{j+1}, a'; \theta)$ 
        Perform gradient descent step on  $(Y_j - Q(s_j, a_j; \theta))^2$ 
    End For
End For

```

2.1.2 Q-learning with experience replay and a fixed Q-target network

The target network is a neural network with the same architecture as the Q-network. The target network parameters θ^- are the same as the Q-network parameters except that the target network weights are updated after τ steps by copying the weights from the Q-network. At other times, the target network weights remain unchanged. The target network used by the DQN is:

$$Y_t^{DQN} := R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t^-)$$

DQNs successfully learn the action value function q_* corresponding to the optimal policy by minimizing the MSE loss between the Q-network and Q-learning targets. In this case, Q-learning targets are calculated with respect to the old fixed parameters θ^- :

$$L(\theta) = E_{s,a,r,s_{t+1} \sim D} [(Y_t^{DQN} - Q(s, a; \theta_t))^2]$$

Algorithm 2: Q-learning with experience replay & fixed Q-target

```

Initialize a replay memory to capacity N
Initialize action-value function Q with random weights  $\theta$ 
Initialize target action-value function  $\hat{Q}$  with random weights  $\hat{\theta}$ 
For episode 1:M do:
    Initialize environment to get initial state  $s_1$ 
    For step 1:T do:
        With probability  $\epsilon$  select random action  $a_t$ 
        Otherwise select  $a_t = \operatorname{argmax}_a Q(s_t, a; \theta)$ 
        Execute action  $a_t$  and observe reward  $r_{t+1}, s_{t+1}$ 
        Store experience  $s_t, a_t, r_{t+1}, s_{t+1}$  in replay memory
        Set  $s_{t+1} \leftarrow s_t$ 
        Sample a batch  $s_j, a_j, r_{j+1}, s_{j+1}$  from replay memory
        Set target  $Y_j = r_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a'; \hat{\theta})$ 
        Perform gradient descent step on  $(Y_j - Q(s_j, a_j; \theta))^2$ 
        Every C step update  $\hat{\theta} = \theta$ 
    End For
End For

```

2.2 Monte-Carlo Policy Gradient

In policy-gradient methods, instead of learning an action-value function, a parametrized policy will be learned that can choose actions without using a value function. Policy will be parametrized with a vector, so taking an action a in time t given that environment is in state s and have parameters θ is:

$$\pi(a|s, \theta) = \Pr\{A_t = a \mid S_t = s, \theta_t = \theta\}$$

Parameters will be learned using the gradient of an performance measure $J(\theta)$. As we want to maximize our agents performance, we will use *gradient ascent*.

$\theta_t = \theta_t + \alpha \nabla J(\theta_t)$, where $\nabla J(\theta_t)$ is a stochastic estimate whose expectation approximates the gradient of the performance measure, $\Delta_{\theta} J(\theta)$.

When action space is discrete and not too big (as in our case), common way for parameterizing is to form parameterized numerical preferences $h(s, a, \theta)$ for each state-action pair. In each state, action with the highest preference will be selected, using for example exponential soft-max distribution.

The action preferences $h(s, \alpha, \theta)$ can be parameterized with a artificial neural network, where θ is vector of all the weights in NN. An advantage of using soft-max and policy gradient method is, actions will be chosen stochastically with arbitrary probabilities, which sometimes a better policy than a deterministic way.

It is also known that, action probabilities change smoothly with policy gradient whereas it does not hold for the ε -greedy selection.

In monte-carlo policy gradient algorithm, each increment is proportional to the product of a return G_t and a vector $\nabla J(\theta_t)$, the gradient of the probability of taking the action actually taken divided by the probability of taking that action. $\nabla \ln \pi(A_t | S_t, \theta_t)$ is used for $\nabla J(\theta_t)$ as $\nabla \ln x = \nabla x/x$. This is a result of Policy Gradient Theorem.

Algorithm 3: Monte-Carlo Policy Gradient

```

Initialize policy parametrization with random weights  $\theta$ 
For episode 1:M do:
    Generate an episode following  $\pi_\theta$ 
    For step 1:T do:
        Set  $G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$ 
        Set  $\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta_t)$ 
    End For
End For

```

3. Simulation Setup

Environment: In this project, we modified the multi-agent environment developed by OpenAI Gym to obtain a new predator-prey environment which we call “Swarm”. In this predator-prey environment, a red predator chases green prey which are slightly faster and smaller in size than it. Multiple predators and prey can exist in the same environment and the goal of each agent is to maximize its own reward. This environment defines the state and action space for each agent in the environment. At each step, each agent performs an action and the environment returns its reward and the next state.

States: The agent's state consists of agent's position and velocity and the relative velocity and position of all there agents.

Rewards: Red predator receives a reward of +10 if they are able to collide with ('eat') any green prey. Green prey receive a negative reward of -10 if they collide with ('are eaten by') any predator. Also, all agents are rewarded negatively for moving out of the screen. This is done so that agents remain in the boundaries of the screen.

Action space: The discrete action space for each agent is the set of four elements which represent whether the agent should move left/right/up/down.

Neural Network Architecture and hyper-parameters: In algorithms 1 & 2, the Q-network for each individual agent has 2 hidden layers. Each hidden layer has twice the number of neurons as the input layer. For each agent these networks take as input the state of the agent and return the estimated Q-value for all possible actions. During training, mini-batches of size 32 are sampled from the replay memory of 2000 states, actions, rewards and next states tuples. We trained our agents with 2 different learning rates : 0.0001 & 0.001. In the beginning each episode consists of 100 steps. After every 500 episodes, the step size is increased by 50. We implemented this so that as agents learn more, they have more time to explore more of the state space, which can help agents learn longer. Both networks are implemented using Keras. Despite the fact that we coded the same neural network architecture from scratch ourselves(appendix), we preferred to use neural networks implemented in Keras for training for getting better performance. We used cloud computing to train our agents with different hyperparameters in parallel. Finally, we saved the statistics (loss, reward, collisions for each agent) for each episode. We also saved the weights of the neural network every 50 episodes to be able to recover any of them later.

For the policy gradient algorithm, the architecture of the network is the same as for the Q-network. In this case the output of the neural network is a soft-max layer that returns the probabilities for selecting each action. This was implemented in TensorFlow. The best performance (least loss) was obtained by using a learning rate

4. Performance

4.1 Algorithm 1: DQN + Experience Replay 2-Prey vs. 1-Predator; Trained for 16 hours

In first row of **Figure 2** (below), note that epsilon decreases as the number of episodes increase. This means that in the beginning, agents explore the environment and as the number of episodes increases, they learn and exploit what they have learnt rather than exploring the environment. Rows 2,3,4 and of Figure 2 show the number of collisions between predator (agent 0) and prey (agent 1 & agent 2) for each episode. Rows 5,6,7 show the rewards obtained by agents for each episode.

Note that in the beginning, there are no collisions between predator and prey. However, the reward obtained by all agents is negative. This is because in the beginning, agents leave the screen. At approx. episode 1000, predator begins to obtain positive reward; this means that predator learns to collide with/ "eat" the prey. Note also that near episode 4000, the number of collisions becomes close to 0 again. Simultaneously, the rewards for all agents are almost 0. This suggests that all the agents remain within the screen, but there are no collisions between predator and prey. In other words, the prey have learnt to evade the predator. Furthermore, this pattern of having a number of episodes where predator learns to catch prey followed by a number of episodes where prey learn to evade predator continues. For instance, near episodes 6000, 8000, 12000, 14000, 18000 the number of collisions between the predator and prey is large; predator learns strategies to catch prey. But in between these regions with high collisions, the number of collisions is small suggesting that the prey learn strategies to evade predator. Some other points of interest are as follows. Near episode 12000, the number of collisions between predator and prey 1 are high while there are no collisions between predator and prey 2. This suggests that while prey 2 has learnt a strategy to evade predator, prey 1 has not. This is interesting because in the beginning there was no difference between the two prey agents. In the same environment, the same prey have learnt differently. Another point of interest is near episode 18000: the predator seems to have learnt a great strategy to catch both prey as the number of

collisions between each predator and prey increases upto 300. This means that in one episode the predator “eats” upto 600 prey!

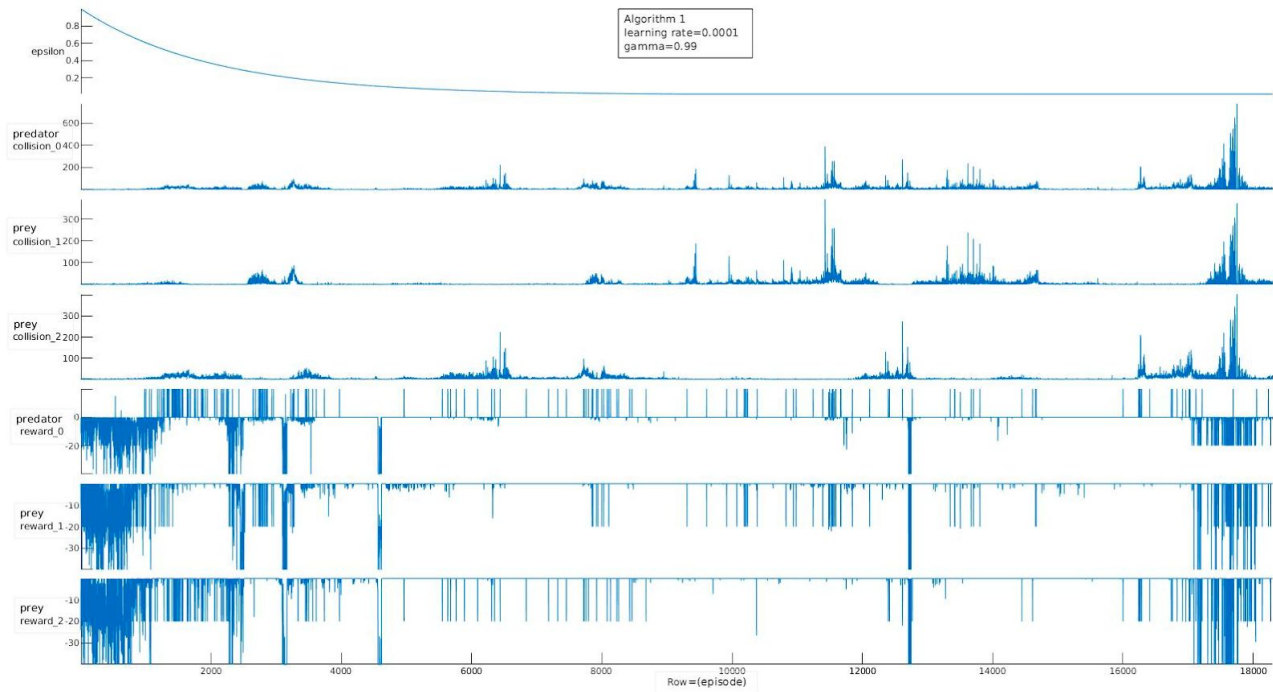


Figure 2 : Algorithm 1 *Epsilon, Collision & Reward versus Episode* are plotted for each agent. For this model, the learning rate = 0.0001 & gamma = 0.99

4.2 Algorithm 2: DQN + Experience Replay + Fixed Q-Target

4.2.1 2-Prey vs. 1-Predator; Trained for 16 hours

The first row of **Figure 3** shows that epsilon decreases with increase in number of episodes meaning that in the beginning agents explore the environment and as time passes and learning occurs, the agents exploit what they have learnt. Rows 2,3,4 and of Figure 3 show the number of collisions between predator (agent 0) and prey (agent 1 & agent 2) for each episode. Rows 5,6,7 show the rewards obtained by agents for each episode.

One interesting thing to note about Algorithm 2 is that there are no alternating sets of episodes where predator learns to catch prey and prey learn to evade predator. In fact, after approx episode 1000, when all agents learn to not leave the screen, we note that there do not seem to many collisions between predator and prey. However, near episode 15000, the number of collisions between predator and prey 2 increases upto a maximum of 800. This suggests that predator learns to catch only prey 2. Again, this is similar to what we observed in Figure 2; in the same environment, agents with similar characteristics learn differently. Close to episode 30000, we observe that all agents obtain highly negative rewards despite the absence of any collisions. This suggests that all agents once again exit the screen. This observation is also interesting since the agents had previously learnt not to exit the screen. This suggests that we might have “over-trained” the model.

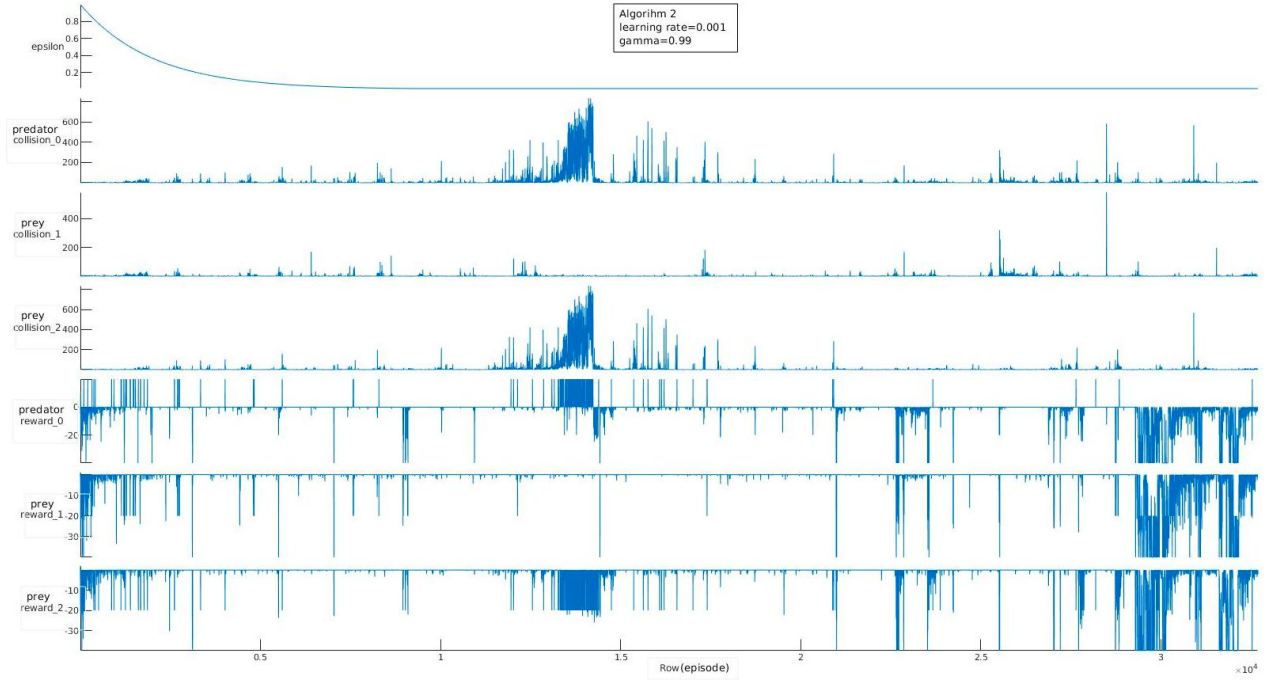
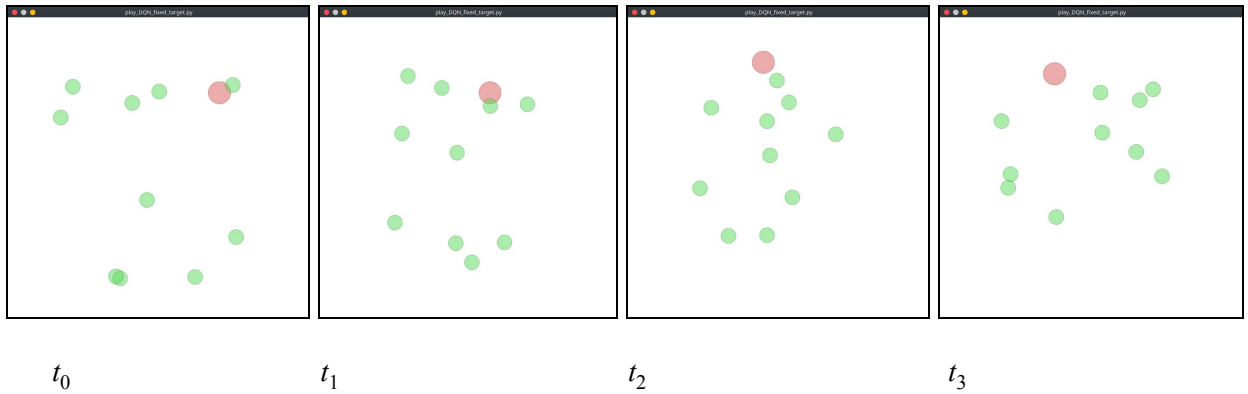


Figure 3 : Algorithm 2 Epsilon, Collision & Reward versus Episode are plotted for each agent. For this model, the learning rate = 0.001 & gamma = 0.99

4.2.1 10-Prey vs. 1-Predator; Trained for 20 hours

Figure 4 shows eight time-steps of a 10-vs-1 environment. The single red agent is the predator and ten green agents are the prey. At the first time step t_0 , agents are placed randomly on the screen. At t_1 , we notice the red predator colliding with (“eating”) a green prey. Over time, we notice that the green prey group learns a defense strategy: the prey distribute themselves at equal distance from the predator. This seems to prevent the predator from catching any of the prey by causing the predator to oscillate about its position. The predator oscillates because the prey have learnt to maintain equal distance from the predator, this prevents the predator from choosing a single prey to chase. This shows how the prey learn a cooperation strategy without being explicitly coded.



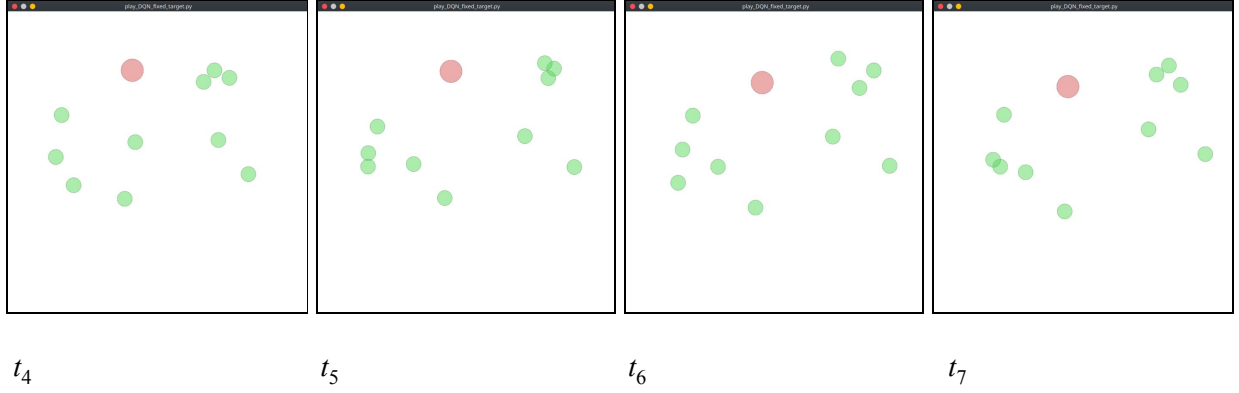


Figure 4: From Left-Right and Top-Bottom, we see 8 time-steps of an episode of 10vs1 game and how agents behave during these time steps.

4.3 Algorithm 3: Monte-Carlo Policy Gradient

Figure 5 shows the collisions, rewards & loss. Agent 0 is the predator and agent 1 and agent 2 are the prey. One point of interest about row 1,2,3 of Figure 5 is that the predator seems to be collide with a single prey (prey 1) most of the time. This observation also explains the rewards obtained by the agent (rows 4,5,6). The predator consistently obtains positive reward for colliding with the prey, while prey 1 consistently obtains negative reward for being “eaten” by the predator. Reward for prey 2 is more negative in the beginning and afterwards, it is rarely caught by the predator. This suggests that prey 2 has learnt a strategy to evade the predator, while prey 1 has not. Loss functions (rows 7,8,9) for each of the three agents approach 0 as the number of episodes increases. A difference between algorithm 3 and algorithm 1 is that in the former, there are no alternating sets of episodes where the predator or the prey learn one-by-one. Rather, learning seems to occur constantly. The predator and prey 2 seem to be able to maximize their rewards easily, but prey 1 does not.

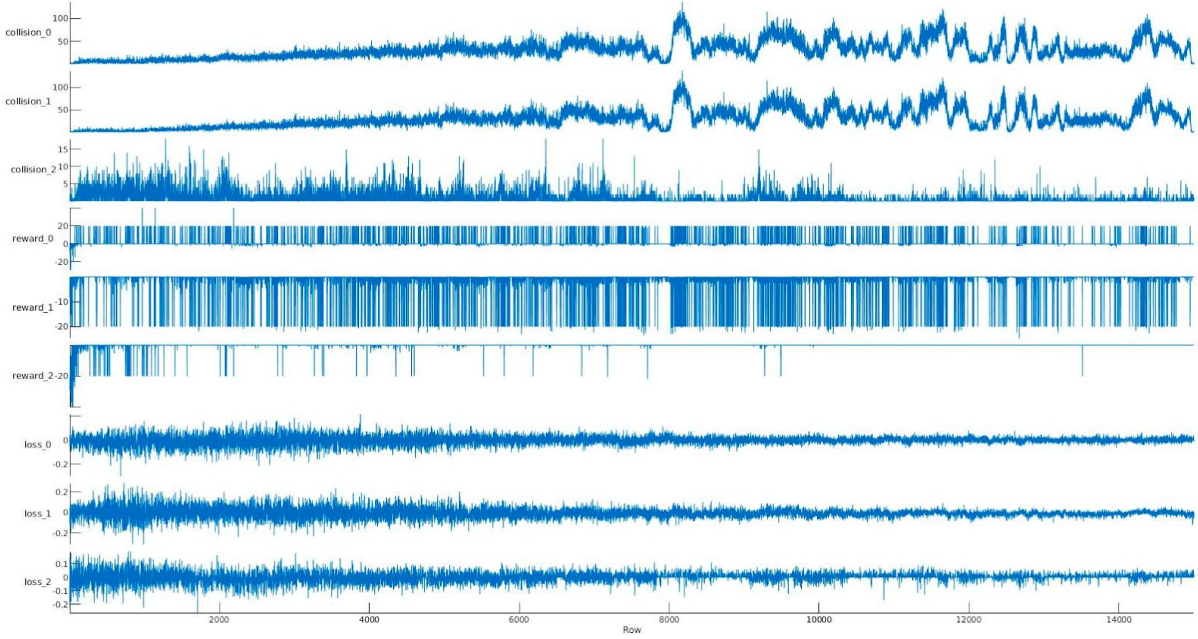


Figure 5: Algorithm 3 Collision, Reward & Loss versus Episode are plotted for each agent. Agent 0: predator; agent 1 & agent 2: prey.

Discussion

We were able to train our agents in the multi-agent environment using three different reinforcement learning algorithms. From the figures provided we note that though learning occurs in all algorithms, it seems that algorithm 3 Monte-Carlo Policy Gradient seem to perform best in terms of the predator consistently being able to “eat” prey. However, in this algorithm, as seen from Figure 5, prey 1 does not seem to learn at all and consistently obtains negative rewards. We observe such an issue in Algorithm 2 as well, (Figure 3) where prey 1 seems to learn evasion better than prey 2. The difference though is that in Algorithm 2, such an effect is not seen throughout, instead it is only observed before episode 15000. Thus, in terms of the all agents learn to increase the reward they collect over time in both the models.

We note that the prey agents develop different strategies over time to evade the predator. One such strategy is shown in Figure 4, where we the green prey agents learn to distribute themselves equally from the predator. This seems to confuse the predator and causes it to oscillate about its position unsure of which prey to chase. We also tried different number of predators and prey (1-predator-vs-2-prey & 1-predator-vs-10 prey) and in each case, found that different strategies were learnt by agents. These will be shown in the demo, in the form of videos.

In conclusion, we note that using reinforcement learning algorithms and simple reward functions for agents, it is possible for the prey agents to discover cooperative strategies that prevent the predator from “eating” them. It is also interesting to note that though the prey agents are the same in the beginning, as training proceeds, the same prey agents learn different strategies. While some are able to successfully evade predators, others are not. Such a concept is similar to the concept of “fitness” observed in real biological systems, where some individuals are able to survive better than others. Extensions of this work could incorporate more agents to study the different strategies developed as the number of individuals in the group increases. Moreover, instead of a 2-D environment, a 3-D

environment might be used where the strategies developed by agents might resemble those observed in real predator-prey interactions.

Changes from Phase 2 report

Figure format changed for clarity. Inclusion of only the best figures obtained using different learning rate and gamma parameters. (Previously all figures for all models were included in the appendix). Addition of 3rd algorithm and related figure.

Response to Feedback for Phase 2 Report

1. Policy gradient algorithm implemented.
2. Report written in a clearer fashion. Only the best plots included. Explanation for each plot was present in previous report as well (except for plots in appendix, which have been removed for clarity.)
3. Each algorithm has its own separate section. Each section includes relevant plots. Explanation of why results achieved and what they correspond to added.

References

- Chincoli, Michele, and Antonio Liotta. "Self-learning power control in wireless sensor networks." *Sensors* 18.2 (2018): 375.
- Lowe, Ryan, et al. "Multi-agent actor-critic for mixed cooperative-competitive environments." *Advances in Neural Information Processing Systems*. 2017.
- Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." *arXiv preprint arXiv:1312.5602* (2013).
- Matiisen, Tanel. "Demystifying Deep Reinforcement Learning." 2015
- Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Van Hasselt, Hado, Arthur Guez, and David Silver. "Deep Reinforcement Learning with Double Q-Learning." *AAAI*. Vol. 2. 2016.

Appendix

Code

Swarm.py

(Swarm.py was modified for our needs from the OpenAI's multi-agent gym environment)

```
import numpy as np
from multiagent.core import World, Agent, Landmark
from multiagent.scenario import BaseScenario
class Scenario(BaseScenario):
    def make_world(self):
        world = World()
        # set any world properties first
        world.dim_c = 0
        num_good_agents = 2
        num_adversaries = 1
        num_agents = num_adversaries + num_good_agents
        num_landmarks = 0
        # add agents
        world.agents = [Agent() for i in range(num_agents)]
        for i, agent in enumerate(world.agents):
            agent.name = 'agent %d' % i
            agent.collide = True
            agent.silent = True
            agent.adversary = True if i < num_adversaries else False
            agent.size = 0.075 if agent.adversary else 0.05
            agent.accel = 3.0 if agent.adversary else 4.0
            #agent.accel = 20.0 if agent.adversary else 25.0
            agent.max_speed = 1.0 if agent.adversary else 1.3
            #! agent.max_speed = 1.0/5 if agent.adversary else 1.3/5
        # add landmarks
        world.landmarks = [Landmark() for i in range(num_landmarks)]
        for i, landmark in enumerate(world.landmarks):
            landmark.name = 'landmark %d' % i
            landmark.collide = True
            landmark.movable = False
            landmark.size = 0.2
            landmark.boundary = False
        # make initial conditions
        self.reset_world(world)
        return world
    def reset_world(self, world):
        # random properties for agents
        for i, agent in enumerate(world.agents):
            agent.color = np.array([0.35, 0.85, 0.35]) if not agent.adversary else np.array([0.85, 0.35, 0.35])
        # random properties for landmarks
        for i, landmark in enumerate(world.landmarks):
            landmark.color = np.array([0.25, 0.25, 0.25])
        # set random initial states
        for agent in world.agents:
```

```

agent.state.p_pos = np.random.uniform(-1, +1, world.dim_p)
agent.state.p_vel = np.zeros(world.dim_p)
agent.state.c = np.zeros(world.dim_c)
for i, landmark in enumerate(world.landmarks):
    if not landmark.boundary:
        landmark.state.p_pos = np.random.uniform(-0.9, +0.9, world.dim_p)
        landmark.state.p_vel = np.zeros(world.dim_p)
def benchmark_data(self, agent, world):
    # returns number of collisions for adversary agent
    if agent.adversary:
        collisions = 0
        for ga in self.good_agents(world):
            if self.is_collision(ga, agent):
                collisions += 1
        return collisions
    if not agent.adversary:
        collisions = 0
        for adv in self.adversaries(world):
            if self.is_collision(adv, agent):
                collisions += 1
        return collisions
    else:
        return -1
def is_collision(self, agent1, agent2):
    delta_pos = agent1.state.p_pos - agent2.state.p_pos
    dist = np.sqrt(np.sum(np.square(delta_pos)))
    dist_min = agent1.size + agent2.size
    return True if dist < dist_min else False
# return all agents that are not adversaries
def good_agents(self, world):
    return [agent for agent in world.agents if not agent.adversary]
# return all adversarial agents
def adversaries(self, world):
    return [agent for agent in world.agents if agent.adversary]
def reward(self, agent, world):
    # Agents are rewarded based on minimum agent distance to each landmark
    main_reward = self.adversary_reward(agent, world) if agent.adversary else self.agent_reward(agent, world)
    return main_reward
def bound(self, x):
    # agents are penalized for exiting the screen, so that they can be caught by the adversaries
    if x < 0.9:
        return 0
    if x < 1.0:
        return (x - 0.9) * 10
    return min(np.exp(2 * x - 2), 10)
def agent_reward(self, agent, world):
    # Agents are negatively rewarded if caught by adversaries
    rew = 0
    shape = False
    adversaries = self.adversaries(world)
    if shape: # reward can optionally be shaped (increased reward for increased distance from adversary)
        for adv in adversaries:
            rew += 0.1 * np.sqrt(np.sum(np.square(agent.state.p_pos - adv.state.p_pos)))
    if agent.collide:

```

```

    for a in adversaries:
        if self.is_collision(a, agent):
            rew -= 10
    #! punishment
    for p in range(world.dim_p):
        x = abs(agent.state.p_pos[p])
        rew -= self.bound(x)
    return rew
def adversary_reward(self, agent, world):
    # Adversaries are rewarded for collisions with agents
    rew = 0
    shape = False #!False
    agents = self.good_agents(world)
    adversaries = self.adversaries(world)
    if shape: # reward can optionally be shaped (decreased reward for increased distance from agents)
        for adv in adversaries:
            rew -= 0.1 * min([np.sqrt(np.sum(np.square(a.state.p_pos - adv.state.p_pos))) for a in agents])
    if agent.collide:
        for ag in agents:
            for adv in adversaries:
                if self.is_collision(ag, adv):
                    rew += 10
    #! punishment
    for p in range(world.dim_p):
        x = abs(agent.state.p_pos[p])
        rew -= self.bound(x)
    return rew
def observation(self, agent, world):
    # get positions of all entities in this agent's reference frame
    entity_pos = []
    for entity in world.landmarks:
        if not entity.boundary:
            entity_pos.append(entity.state.p_pos - agent.state.p_pos)
    # communication of all other agents
    comm = []
    other_pos = []
    other_vel = []
    for other in world.agents:
        if other is agent: continue
        comm.append(other.state.c)
        other_pos.append(other.state.p_pos - agent.state.p_pos)
        # if not other.adversary:
        other_vel.append(other.state.p_vel)
    return np.concatenate([agent.state.p_vel] + [agent.state.p_pos] + entity_pos + other_pos + other_vel)

```

play_DQN.py

```

import random
import gym
import make_env_
import numpy as np
import csv
from collections import deque
from keras.models import Sequential
from keras.layers import Dense

```

```

from keras.optimizers import Adam
import os # for creating directories
#^ Set parameters
env = make_env._make_env('swarm', benchmark=True)
num_of_agents = env.n
state_size = (2+2+2*(num_of_agents-1)*2) # [agent's velocity(2d vector) + agent's position(2d vector) +
# other agent's relative position((n-1)*2d vector) +
# other agent's relative velocity((n-1)*2d vector)]
# in 3 agent case it is 2+2+2*2+2*2=12
action_size = 4 # discrete action space [up,down,left,right]
batch_size = 32 # used for batch gradient descent update
testing = True # render or not, expodation vs. exploration
n_episodes = 100000 if not testing else 5 # number of simulations
n_steps = 100 if not testing else 500 # number of steps
load_episode = 15100
output_dir = 'model_output/swarm/DQN_2v1'
if testing:
    import pyautogui
#^ Define agent
class DQNAgent:
    def __init__(self, state_size, action_size):
        self.state_size = state_size # defined above
        self.action_size = action_size # defined above
        self.memory = deque(maxlen=2000) # double-ended queue; removes the oldest element each time that you
add a new element.
        self.gamma = 0.95 # discount rate
        self.epsilon = 1.0 if not testing else 0.1 # exploration rate: how much to act randomly; more initially than later
due to epsilon decay
        self.epsilon_decay = (1-0.001) # exponential decay rate for exploration prob
        self.epsilon_min = 0.01 # minimum amount of random exploration permitted
        self.learning_rate = 0.001 # learning rate of NN
        self.model = self._build_model()
    def _build_model(self):
        # neural net for approximating Q-value function:  $Q^*(s,a) \sim Q(s,a;W)$ 
        model = Sequential() #fully connected NN
        model.add(Dense(state_size*2, input_dim=self.state_size, activation='relu')) # 1st hidden layer
        model.add(Dense(state_size*2, activation='relu')) # 2nd hidden layer
        model.add(Dense(self.action_size, activation='linear')) # 4 actions, so 4 output neurons
        model.compile(loss='mse', optimizer=Adam(lr=self.learning_rate))
        return model
    def remember(self, state, action, reward, next_state, done):
        self.memory.append((state, action, reward, next_state, done)) # list of previous experiences, enabling
re-training later
    def act(self, state):
        if np.random.rand() <= self.epsilon: # take random action with epsilon probability
            onehot_action = np.zeros(action_size+1)
            onehot_action[random.randint(1,4)] = 1
            return onehot_action
        act_values = self.model.predict(state) # predict reward value based on current state
        # print(act_values)
        act_index = np.argmax(act_values[0]) # pick the action with highest value
        onehot_action = np.zeros(action_size+1)
        onehot_action[act_index+1] = 1
        return onehot_action

```



```

def replay(self, batch_size): # method that trains NN with experiences sampled from memory
    minibatch = random.sample(self.memory, batch_size) # sample a minibatch from memory
    for state, action, reward, next_state, done in minibatch: # extract data for each minibatch sample
        target = reward # if done then target = reward
        state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
        next_state = np.reshape(next_state, [1, state_size]) #! reshape the state for DQN model
        if not done: # if not done, then predict future discounted reward
            target = (reward + self.gamma * # (target) = reward + (discount rate gamma) *
                      np.amax(self.model.predict(next_state))) # (maximum target Q based on future action a')
        target_f = self.model.predict(state) # approximately map current state to future discounted reward
        target_f[0][np.argmax(action)-1] = target
        history = self.model.fit(state, target_f, epochs=1, verbose=0)
        # single epoch of training with x=state, y=target_f
    if self.epsilon > self.epsilon_min:
        self.epsilon *= self.epsilon_decay
    return history

def load(self, name):
    self.model.load_weights(name)

def save(self, name):
    self.model.save_weights(name)

#^ Interact with environment
agents = [ DQNAgent(state_size, action_size) for agent in range(num_of_agents) ] # initialise agents
#! create model output folders
for i, agent in enumerate(agents):
    if not os.path.exists(output_dir + "/weights/agent{}".format(i)):
        os.makedirs(output_dir + "/weights/agent{}".format(i))
#! load weights if exist
for i, agent in enumerate(agents):
    file_name = (output_dir + "/weights/agent{}".format(i) + "weights_" + '{:04d}'.format(load_episode) + ".hdf5")
    if os.path.isfile(file_name):
        print("Loading weights to use for agent {}".format(i))
        agent.load(file_name)
#! statistics
collision_ = ['collision_{}'.format(i) for i in range(num_of_agents)]
loss_ = ['loss_{}'.format(i) for i in range(num_of_agents)]
reward_ = ['reward_{}'.format(i) for i in range(num_of_agents)]
statistics = ['episode', 'epsilon'] + collision_ + reward_ + loss_
if not testing:
    with open(output_dir + '/statistics.csv', 'a') as csvFile:
        writer = csv.writer(csvFile)
        writer.writerow(statistics)
    csvFile.close()
for episode in range(1, n_episodes+1): # iterate over new episodes of the game
    if(episode % 500 == 0): n_steps+=50
    #^ for statistics
    statistics_row=[]
    collisions = [0]*num_of_agents
    rewards = [0]*num_of_agents
    losses = [0]*num_of_agents
    states = env.reset() # reset states at start of each new episode of the game
    for step in range(1, n_steps+1): # for every step
        if (testing):
            env.render()
            # if (step % 4 == 0):

```

```

# # Take screenshot
# pic = pyautogui.screenshot()
# # Save the image
# pic.save(output_dir+'screenshots/Screenshot_{}.png'.format(step))
all_actions=[]
for state,agent in zip(states,agents):
    state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
    action_i = agent.act(state)
    all_actions.append(action_i)
next_states, rewards, dones, infos = env.step(all_actions) # take a step (update all agents)
## collision, reward statistics
for i in range(num_of_agents):
    collisions[i] += (infos['collision'][i])
    rewards[i] += (rewards[i])
for state in next_states:
    state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
for i,agent in enumerate(agents):
    agent.remember(states[i], all_actions[i], rewards[i], next_states[i], dones[i])
    # remember the previous timestep's state, actions, reward vs.
states = next_states # update the states
print("episode: {}/{}, collisions: {}, epsilon: {:.2}".format(episode, n_episodes, collisions[0], agent.epsilon))
for i,agent in enumerate(agents):
    if len(agent.memory) > batch_size:
        history = agent.replay(batch_size) # train the agent by replaying the experiences of the episode
        losses[i] += history.history['loss'][0]
#! episode, epsilon, collisions, rewards, losses statistics written
statistics_row.append(episode)
statistics_row.append(agents[0].epsilon)
statistics_row += (collisions)
statistics_row += (rewards)
statistics_row += (losses)
if not testing:
    with open(output_dir + '/statistics.csv', 'a') as csvFile:
        writer = csv.writer(csvFile)
        writer.writerow(statistics_row)
    csvFile.close()
#! save weights
if not testing:
    if episode % 50 == 0:
        for i,agent in enumerate(agents):
            agent.save(output_dir + "/weights/agent{}/".format(i) + "weights_" + '{:04d}'.format(episode) + ".hdf5")

```

play_DQN_fixed_target.py

```

import random
import gym
import make_env_
import numpy as np
import csv
from collections import deque
from keras.models import Sequential
from keras.layers import Dense
from keras.optimizers import Adam
import os # for creating directories

```

```

#^ Set parameters
env = make_env_.make_env('swarm',benchmark=True)
num_of_agents = env.n
state_size = (2+2+2*(num_of_agents-1)*2) # [agent's velocity(2d vector) + agent's position(2d vector) +
# other agent's relative position((n-1)*2d vector) +
# other agent's relative velocity((n-1)*2d vector)]
# in 3 agent case it is 2+2+2*2+2*2=12
action_size = 4 # discrete action space [up,down,left,right]
batch_size = 32 # used for batch gradient descent update
testing = True # render or not, expodation vs. exploration
n_episodes = 100000 if not testing else 1 # number of simulations
n_steps = 100 if not testing else 32 # number of steps
load_episode = 5250
updating_target_freq = 50 # rate C, reset W' <- W
output_dir = 'model_output/swarm/DQQ_fixed_target_10v1'
if testing:
    import pyautogui

```

```

#^ Define agent
class DQNAgent:
    def __init__(self, state_size, action_size):
        self.state_size = state_size # defined above
        self.action_size = action_size # defined above
        self.memory = deque(maxlen=2000) # double-ended queue; removes the oldest element each time that you
add a new element.
        self.gamma = 0.95 # discount rate
        self.epsilon = 1.0 if not testing else 0.1 # exploration rate: how much to act randomly; more initially than later
due to epsilon decay
        self.epsilon_decay = (1-0.0005) # exponential decay rate for exploration prob
        self.epsilon_min = 0.01 # minimum amount of random exploration permitted
        self.learning_rate = 0.0005 # learning rate of NN
        self.evaluation_model = self._build_model()
        self.target_model = self._build_model()
    def _build_model(self):
        # neural net for approximating Q-value function: Q*(s,a) ~ Q(s,a;W)
        model = Sequential() #fully connected NN
        model.add(Dense(state_size*2, input_dim=self.state_size, activation='relu')) # 1st hidden layer
        model.add(Dense(state_size*2, activation='relu')) # 2nd hidden layer
        model.add(Dense(self.action_size, activation='linear')) # 4 actions, so 4 output neurons
        model.compile(loss='mse',optimizer=Adam(lr=self.learning_rate))
        return model
    def update_target_weights(self):
        self.target_model.set_weights(self.evaluation_model.get_weights())
    def remember(self, state, action, reward, next_state, done):
        self.memory.append((state, action, reward, next_state, done)) # list of previous experiences, enabling
re-training later
    def act(self, state):
        if np.random.rand() <= self.epsilon: # take random action with epsilon probability
            onehot_action = np.zeros(action_size+1)
            onehot_action[random.randint(1,4)] = 1
            return onehot_action
        act_values = self.evaluation_model.predict(state) # predict reward value based on current state
        # print(act_values)

```

```

act_index = np.argmax(act_values[0]) # pick the action with highest value
onehot_action = np.zeros(action_size+1)
onehot_action[act_index+1] = 1
return onehot_action
def replay(self, batch_size): # method that trains NN with experiences sampled from memory
    minibatch = random.sample(self.memory, batch_size) # sample a minibatch from memory
    for state, action, reward, next_state, done in minibatch: # extract data for each minibatch sample
        target = reward # if done then target = reward
        state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
        next_state = np.reshape(next_state, [1, state_size]) #! reshape the state for DQN model
        if not done: # if not done, then predict future discounted reward
            target = (reward + self.gamma * # (target) = reward + (discount rate gamma) *
                    np.amax(self.target_model.predict(next_state))) # (maximum target Q based on future action a')
            target_f = self.evaluation_model.predict(state) # approximately map current state to future discounted
reward
        target_f[0][np.argmax(action)-1] = target
        history = self.evaluation_model.fit(state, target_f, epochs=1, verbose=0)
        # single epoch of training with x=state
        if self.epsilon > self.epsilon_min:
            self.epsilon *= self.epsilon_decay
        return history
def load(self, name):
    self.evaluation_model.load_weights(name)
    self.update_target_weights()
def save(self, name):
    self.evaluation_model.save_weights(name)
#^ Interact with environment
agents = [ DQNAgent(state_size, action_size) for agent in range(num_of_agents) ] # initialise agents
#! create model output folders
for i,agent in enumerate(agents):
    if not os.path.exists(output_dir + "/weights/agent{}".format(i)):
        os.makedirs(output_dir + "/weights/agent{}".format(i))
#! load weights if exist
for i,agent in enumerate(agents):
    file_name = (output_dir + "/weights/agent{}".format(i) + "weights_" + '{:04d}'.format(load_episode) + ".hdf5")
    if os.path.isfile(file_name):
        print("Loading of model weights to use for agent {}".format(i))
        agent.load(file_name)
#! statistics

collision_ = ['collision_{}'.format(i) for i in range(num_of_agents)]
loss_ = ['loss_{}'.format(i) for i in range(num_of_agents)]
reward_ = ['reward_{}'.format(i) for i in range(num_of_agents)]
statistics = ['episode','epsilon']+collision_+reward_+loss_
if not testing:
    with open(output_dir + '/statistics.csv', 'a') as csvFile:
        writer = csv.writer(csvFile)
        writer.writerow(statistics)
    csvFile.close()

for episode in range(1,n_episodes+1): # iterate over new episodes of the game
    if(episode % 500 == 0):

```

```

n_steps+=50
updating_target_freq+=25



---



#^ for statistics
statistics_row=[]
collisions = [0]*num_of_agents
rewards = [0]*num_of_agents
losses = [0]*num_of_agents



---



states = env.reset() # reset states at start of each new episode of the gam
for step in range(1,n_steps+1): # for every step



---



#! reset target model weights
if(step % updating_target_freq == 0):
    for agent in agents:
        agent.update_target_weights()
if (testing):
    env.render()
if (step % 4 == 0):
    # Take screenshot
    pic = pyautogui.screenshot()
    # Save the image
    pic.save(output_dir+'/screenshots/Screenshot_{}.png'.format(step))



---



# if(episode > 100 and episode < 110): env.render();
# if(episode > 500 and episode < 510): env.render();
# if(episode > 950 and episode < 1000): env.render();



---



all_actions=[]
for state,agent in zip(states,agents):
    state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
    action_i = agent.act(state)
    all_actions.append(action_i)
next_states, rewards, dones, infos = env.step(all_actions) # take a step (update all agents)
#* collision,reward statistics
for i in range(num_of_agents):
    collisions[i] += (infos['collision'][i])
    rewards[i] += (rewards[i])
for state in next_states:
    state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
for i,agent in enumerate(agents):
    agent.remember(states[i], all_actions[i], rewards[i], next_states[i], dones[i])
    # remember the previous timestep's state, actions, reward vs.
    states = next_states # update the states
print("episode: {}/ {}, collisions: {}, epsilon: {:.2}".format(episode, n_episodes, collisions[0], agent.epsilon))
for i,agent in enumerate(agents):
    if len(agent.memory) > batch_size:
        history = agent.replay(batch_size) # train the agent by replaying the experiences of the episode
        losses[i] += history.history['loss'][0]

```

```

## episode,epsilon,collisions,rewards,losses statistics written
statistics_row.append(episode)
statistics_row.append(agents[0].epsilon)
statistics_row += (collisions)
statistics_row += (rewards)
statistics_row += (losses)
if not testing:
    with open(output_dir + '/statistics.csv', 'a') as csvFile:
        writer = csv.writer(csvFile)
        writer.writerow(statistics_row)
    csvFile.close()
## save weights
if not testing:
    if episode % 50 == 0:
        for i,agent in enumerate(agents):
            file_name = (output_dir + "/weights/agent{}/".format(i) + "weights_" + '{:04d}'.format(episode) +
".hdf5")
            agent.save(file_name)

```

DPG.py

```

import numpy as np
import tensorflow as tf
np.random.seed(1)
tf.set_random_seed(1)
class PolicyGradientAgent:
    def __init__(self, state_size, action_size, learning_rate=0.01, gamma=0.95):
        self.state_size = state_size
        self.action_size = action_size
        self.learning_rate = learning_rate
        self.gamma = gamma
        self.observations = []
        self.actions = []
        self.rewards = []
        self._build_model()
        self.sess = tf.Session()
        self.sess.run(tf.global_variables_initializer())
        self.saver = tf.train.Saver(max_to_keep=100000)
    def _build_model(self):
        tf.reset_default_graph()
        with tf.name_scope('inputs'):
            # placeholders
            self.tf_obs = tf.placeholder(
                tf.float32, [None, self.state_size], name="observations")
            self.tf_acts = tf.placeholder(
                tf.int32, [None, ], name="actions_indexes")
            self.tf_vt = tf.placeholder(
                tf.float32, [None, ], name="actions_values")
        # layer1
        layer = tf.layers.dense(
            inputs=self.tf_obs,
            units=self.state_size*2,
            activation=tf.nn.tanh, # tanh activation
            kernel_initializer=tf.random_normal_initializer(

```

```

        mean=0, stddev=0.3),
        bias_initializer=tf.constant_initializer(0.1),
        name='layer1'
    )
    # layer2
    layer = tf.layers.dense(
        inputs=layer,
        units=self.state_size*2,
        activation=tf.nn.tanh, # tanh activation
        kernel_initializer=tf.random_normal_initializer(
            mean=0, stddev=0.3),
        bias_initializer=tf.constant_initializer(0.1),
        name='layer2'
    )
    # layer3
    all_act = tf.layers.dense(
        inputs=layer,
        units=self.action_size,
        activation=None,
        kernel_initializer=tf.random_normal_initializer(
            mean=0, stddev=0.3),
        bias_initializer=tf.constant_initializer(0.1),
        name='layer3'
    )
    # use softmax to convert to probability
    self.all_act_prob = tf.nn.softmax(all_act, name='act_prob')
    with tf.name_scope('loss'):
        # maximizing total reward (log_p * R) is equal to minimizing
        # -(log_p * R), tensorflow has only have minimizing(loss)
        neg_log_prob = tf.nn.sparse_softmax_cross_entropy_with_logits(
            logits=all_act, labels=self.tf_acts)
        # this is negative log of the chosen action
        # reward guided loss
        loss = tf.reduce_mean(neg_log_prob * self.tf_vt)
        self.loss = loss
    with tf.name_scope('train'):
        self.train_op = tf.train.AdamOptimizer(
            self.learning_rate).minimize(loss)
def act(self, observation):
    """
    Choose actions with respect to their probabilities
    """
    # runs one "step" of TensorFlow computation
    prob_weights = self.sess.run(self.all_act_prob, feed_dict={
        self.tf_obs: observation[np.newaxis, :]})
    action = np.random.choice(
        range(prob_weights.shape[1]), p=prob_weights.ravel())
    return action
def remember(self, state, action, reward):
    """
    Add state,action,reward to the memory
    """
    self.observations.append(state)
    self.actions.append(action)

```

```

        self.rewards.append(reward)
def learn(self):
    """
    Training of the PG agent
    """
    discounted_normalized_rewards = self._discount_and_normalize_rewards()
    _, loss = self.sess.run((self.train_op, self.loss),
                            feed_dict={
                                # shape=[None, n_obs]
                                self.tf_obs: np.vstack(self.observations),
                                # shape=[None, ]
                                self.tf_acts: np.array(self.actions),
                                # shape=[None, ]
                                self.tf_vt: discounted_normalized_rewards,
                            })
    self.observations = []
    self.actions = []
    self.rewards = []
    return discounted_normalized_rewards, loss
def _discount_and_normalize_rewards(self):
    """
    discount and normalize the reward of the episode
    """
    discounted_rewards = np.zeros_like(self.rewards, dtype=np.float64)
    running_add = 0
    for t in reversed(range(0, len(self.rewards))):
        running_add = running_add * self.gamma + self.rewards[t]
        discounted_rewards[t] = running_add
    # normalize episode rewards
    discounted_rewards -= np.mean(discounted_rewards, dtype=np.float64)
    discounted_rewards /= (np.std(discounted_rewards,
                                   dtype=np.float64)+1e-6)
    return discounted_rewards
def load(self, path):
    self.saver.restore(self.sess, path)
def save(self, path):
    self.saver.save(self.sess, path)

```

play_DPG.py

```

import random
import gym
import make_env_
import numpy as np
import csv
# from collections import deque
import tensorflow as tf
import os # for creating directories
from DPG import PolicyGradientAgent # PG agent
np.random.seed(1)
tf.set_random_seed(1)
# ^ Set parameters
env = make_env_.make_env('swarm', benchmark=True)
num_of_agents = env.n

```



```

state_size = (2+2+2*(num_of_agents-1)*2)
# [agent's velocity(2d vector) + agent's position(2d vector) +
# other agent's relative position((n-1)*2d vector) +
# other agent's relative velocity((n-1)*2d vector)) ]
# in 3 agent case it is 2+2+2*2+2*2=12
action_size = 4 # discrete action space [up,down,left,right]
testing = False # render or not, expodation vs. exploration
render = False
n_episodes = 2000 if not testing else 7 # number of simulations
n_steps = 200 if not testing else 300 # number of steps
load_episode = 0
output_dir = 'model_output/swarm/DPG'
# ^ Interact with environment
agents = [PolicyGradientAgent(state_size, action_size)
          for agent in range(num_of_agents)] # initialize agents
#! create model output folders
for i, agent in enumerate(agents):
    if not os.path.exists(output_dir + "/weights/agent{}".format(i)):
        os.makedirs(output_dir + "/weights/agent{}".format(i))
#! load weights if exist
for i, agent in enumerate(agents):
    file_name = (output_dir + "/weights/agent{}".format(i) +
                 "weights_" + '{:04d}'.format(load_episode))
    try:
        agent.load(file_name)
        print("Loaded weights to use for agent {}".format(i))
    except:
        print("No weights to use for agent {}".format(i))
    finally:
        pass
#! statistics
collision_ = ['collision_{}'.format(i) for i in range(num_of_agents)]
loss_ = ['loss_{}'.format(i) for i in range(num_of_agents)]
reward_ = ['reward_{}'.format(i) for i in range(num_of_agents)]
statistics = ['episode']+collision_+reward_+loss_
if not testing:
    with open(output_dir + '/statistics.csv', 'a') as csvFile:
        writer = csv.writer(csvFile)
        writer.writerow(statistics)
    csvFile.close()
for episode in range(1, n_episodes+1): # iterate over new episodes of the game
    # if(episode % 500 == 0):
    #     n_steps += 50
    # ^ for statistics
    statistics_row = []
    collisions = [0]*num_of_agents
    rewards_ = [0]*num_of_agents
    losses = [0]*num_of_agents
    states = env.reset() # reset states at start of each new episode of the game
    for step in range(1, n_steps+1): # for every step
        if (render):
            env.render()
            all_actions = []
            all_actions_index = []

```

```

for state, agent in zip(states, agents):
    # state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
    act_index = agent.act(state)
    all_actions_index.append(act_index)
    onehot_action = np.zeros(action_size+1)
    onehot_action[act_index+1] = 1
    all_actions.append(onehot_action)
next_states, rewards, dones, infos = env.step(
    all_actions) # take a step (update all agents)
# * collision, reward statistics
for i in range(num_of_agents):
    collisions[i] += (infos['collision'][i])
    rewards_[i] += (rewards[i])
# for state in next_states:
# state = np.reshape(state, [1, state_size]) #! reshape the state for DQN model
for i, agent in enumerate(agents):
    agent.remember(states[i], all_actions_index[i], rewards[i])
    # remember the previous timestep's state, actions, reward vs.
states = next_states # update the states
# End of the episode
for i, agent in enumerate(agents):
    rewards_sum = sum(agent.rewards)
    if 'running_reward' not in globals():
        running_reward = rewards_sum
    else:
        running_reward = running_reward * \
            agent.gamma + rewards_sum * (1-agent.gamma)
    value, loss = agent.learn()
    losses[i] = loss
print("\n episode: {}/ {}, collisions: {}, \
rewards: {:.2f}|{:.2f}|{:.2f}, \
losses: {:.2f}|{:.2f}|{:.2f}".format(episode,
                                    n_episodes,
                                    collisions[0],
                                    rewards_[0],
                                    rewards_[1],
                                    rewards_[2],
                                    losses[0],
                                    losses[1],
                                    losses[2]))
#! episode, collisions, rewards, losses statistics written
statistics_row.append(episode)
statistics_row += (collisions)
statistics_row += (rewards_)
statistics_row += (losses)
if not testing:
    with open(output_dir + '/statistics.csv', 'a') as csvFile:
        writer = csv.writer(csvFile)
        writer.writerow(statistics_row)
    csvFile.close()
#! save weights
if not testing:
    if episode % 50 == 0:
        for i, agent in enumerate(agents):

```

```
agent.save(output_dir + "/weights/agent{}/".format(i) +
           "weights_" + '{:04d}'.format(episode))
```

Neural Network Implementation

```
class _init_py(object):
    """description of class"""
    """

Using some libraries to iterate over the data
"""

from typing import Iterator, NamedTuple # inspired by PEP style coding
import numpy as np
from neuralnet.tensor import Tensor # neuralnet is the folder that contains
# elements of NN in separate files
Batch = NamedTuple("Batch", [("inputs", Tensor), ("targets", Tensor)])
#giving an ordered dict mapping field names to types
class DataIterator:
    def __call__(self, inputs: Tensor, targets: Tensor) -> Iterator[Batch]:
        raise NotImplementedError #base class
class BatchIterator(DataIterator):
    def __init__(self, batch_size: int = 32, shuffle: bool = True) -> None:
        self.batch_size = batch_size
        self.shuffle = shuffle
    def __call__(self, inputs: Tensor, targets: Tensor) -> Iterator[Batch]:
        starts = np.arange(0, len(inputs), self.batch_size)
        if self.shuffle:
            np.random.shuffle(starts)
        for start in starts:
            end = start + self.batch_size # get batch size index array
            batch_inputs = inputs[start:end] # get input from batch index
            batch_targets = targets[start:end] # get target from same index
            yield Batch(batch_inputs, batch_targets) # returns one batch every call, like range() function
    """

NN consists of some layers and activation function
between them
"""

from typing import Dict, Callable
import numpy as np
from neuralnet.tensor import Tensor #ndarray

class Layer: #base class
    def __init__(self) -> None:
        self.params: Dict[str, Tensor] = {} # typing library functions enables to declare types
        self.grads: Dict[str, Tensor] = {}
    def forward(self, inputs: Tensor) -> Tensor:
        """
        forward the input signal
        """
        raise NotImplementedError
    def backward(self, grad: Tensor) -> Tensor:
        """
        backpropagate the error signal
        """
        raise NotImplementedError
```

```

class Linear(Layer):
    """
    computes output =  $XW+b$ 
    """
    def __init__(self, input_size: int, output_size: int) -> None:
        super().__init__()
        self.params["w"] = np.random.randn(input_size, output_size)
        self.params["b"] = np.random.randn(output_size)
    def forward(self, inputs: Tensor) -> Tensor:
        """
        return =  $XW + b$ 
        """
        self.inputs = inputs
        return inputs @ self.params["w"] + self.params["b"]
    def backward(self, grad: Tensor) -> Tensor:
        """
        using chain rule
        """
        self.grads["b"] = np.sum(grad, axis=0)
        self.grads["w"] = self.inputs.T @ grad
        return grad @ self.params["w"].T
F = Callable[[Tensor], Tensor]
class Activation(Layer):
    """
    in an activation layer a function
    is applied to its inputs elementwise
    """
    def __init__(self, f: F, f_prime: F) -> None:
        super().__init__()
        self.f = f
        self.f_prime = f_prime #derivative
    def forward(self, inputs: Tensor) -> Tensor:
        self.inputs = inputs
        return self.f(inputs)
    def backward(self, grad: Tensor) -> Tensor:
        return self.f_prime(self.inputs) * grad
def tanh(x: Tensor) -> Tensor:
    return np.tanh(x)
def tanh_prime(x: Tensor) -> Tensor: #derivative
    y = tanh(x)
    return 1 - y ** 2
#Todo: implement relu, sigmoid
class Tanh(Activation):
    def __init__(self):
        super().__init__(tanh, tanh_prime)
import numpy as np
from neuralnet.tensor import Tensor
class Loss:
    def loss(self, predicted: Tensor, actual: Tensor) -> float:
        raise NotImplementedError
    def grad(self, predicted: Tensor, actual: Tensor) -> Tensor:
        raise NotImplementedError
class MSE(Loss):
    def loss(self, predicted: Tensor, actual: Tensor) -> float:

```

```

    return np.sum((predicted - actual) ** 2)
def grad(self, predicted: Tensor, actual: Tensor) -> Tensor: #derivative of loss at output layer
    return 2 * (predicted - actual)
from typing import Sequence, Iterator, Tuple
from neuralnet.tensor import Tensor
from neuralnet.layers import Layer
class NeuralNet:
    def __init__(self, layers: Sequence[Layer]) -> None:
        self.layers = layers
    def forward(self, inputs: Tensor) -> Tensor:
        for layer in self.layers:
            inputs = layer.forward(inputs)
        return inputs
    def backward(self, grad: Tensor) -> Tensor:
        for layer in reversed(self.layers):
            grad = layer.backward(grad)
        return grad
    def params_and_grads(self) -> Iterator[Tuple[Tensor, Tensor]]:
        for layer in self.layers:
            for name, param in layer.params.items():
                grad = layer.grads[name]
                yield param, grad # yields parameters and partial gradients of
                # each layer
from neuralnet.nn import NeuralNet
class Optimizer:
    def step(self, net: NeuralNet) -> None:
        raise NotImplementedError
class SGD(Optimizer):
    def __init__(self, lr: float = 0.01) -> None:
        self.lr = lr
    def step(self, net: NeuralNet) -> None:
        for param, grad in net.params_and_grads():
            param -= self.lr * grad
"""
using numpy's nd array as tensor
"""
from numpy import ndarray as Tensor
"""
training the NN
"""
import matplotlib.pyplot as plt
from neuralnet.tensor import Tensor
from neuralnet.nn import NeuralNet
from neuralnet.loss import Loss, MSE
from neuralnet.optim import Optimizer, SGD
from neuralnet.data import DataIterator, BatchIterator
def train(net: NeuralNet,
          inputs: Tensor,
          targets: Tensor,
          num_epochs: int = 5000,
          iterator: DataIterator = BatchIterator(),
          loss: Loss = MSE(),
          optimizer: Optimizer = SGD()) -> None:
    losses = []

```

```

for epoch in range(num_epochs):
    epoch_loss = 0.0
    for batch in iterator(inputs, targets):
        predicted = net.forward(batch.inputs)
        epoch_loss += loss.loss(predicted, batch.targets)
        grad = loss.grad(predicted, batch.targets)
        net.backward(grad)
        optimizer.step(net)
    losses.append(epoch_loss)
    if (epoch % 100 == 0):
        print('epoch num:', epoch, 'epoch_loss:', epoch_loss)
plt.plot(losses)
plt.title('History of Error')
plt.xlabel('Iteration')
plt.ylabel('MSE'); plt.show()

```

Q1.py

```

#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Dec 5 18:24:20 2018
@author: kazim
Coding a neural nets with 2 hidden layers to fit
on the data in PS3 Q1 (tried for regression)
"""

#%% imports
import numpy as np
import scipy.io as sio
from neuralnet.train import train
from neuralnet.nn import NeuralNet
from neuralnet.layers import Linear, Tanh
from neuralnet.optim import SGD

#%% read data
data = sio.loadmat('perceptrontest.mat')['D']
X = data[:, :-1]
Y = data[:, -1][:, np.newaxis]

#%% Create NN
NN = NeuralNet([
    Linear(input_size=11, output_size=20), Tanh(),
    Linear(input_size=20, output_size=20), Tanh(),
    Linear(input_size=20, output_size=1)
])

#%% Train NN
train(NN, X, Y, num_epochs=500, optimizer=SGD(lr=0.001))

```

epoch num: 0 epoch_loss: 4365.586284819617
epoch num: 100 epoch_loss: 75.22132375934444
epoch num: 200 epoch_loss: 52.50037650515755
epoch num: 300 epoch_loss: 43.16964098990116
epoch num: 400 epoch_loss: 38.781758082490704

