

MTP

Arshdeep Singh

2019csm1001

Supervisor : Dr. Shashi Shekhar Jha

Topic

Multi-UAVs (Unmanned Aerial vehicles) Safe Navigation through Multi Agent Reinforcement Learning

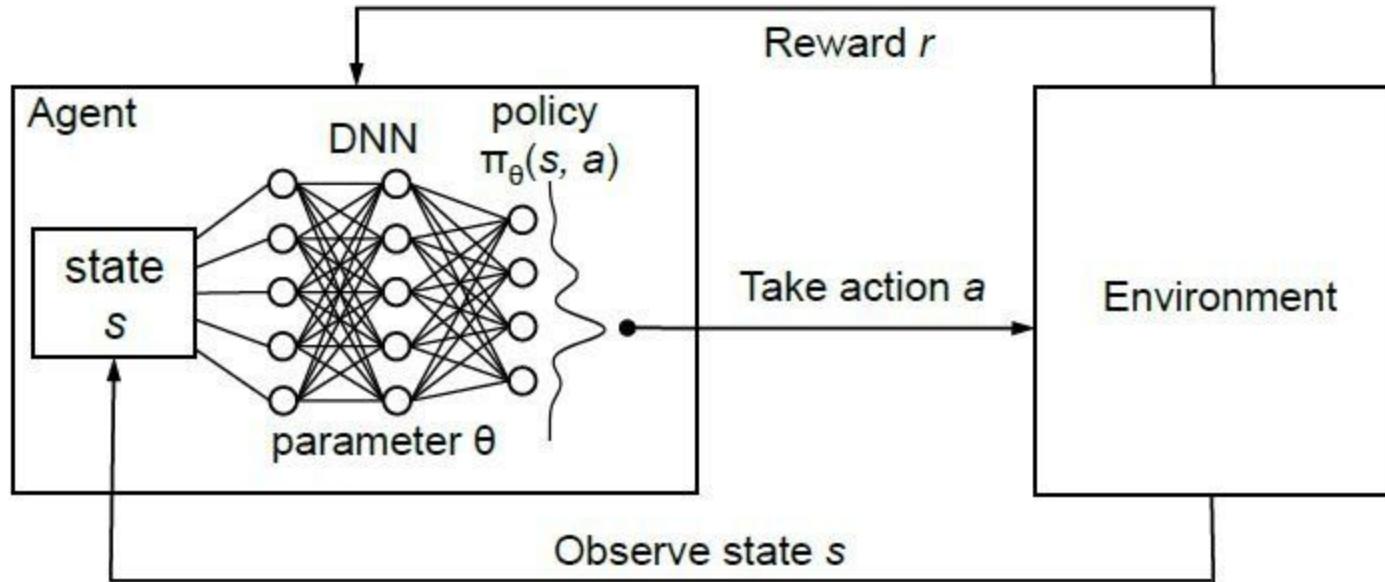
Multi-UAV Applications

1. Mapping and Surveying
2. Crop Monitoring
3. Emergency Response

Formulation

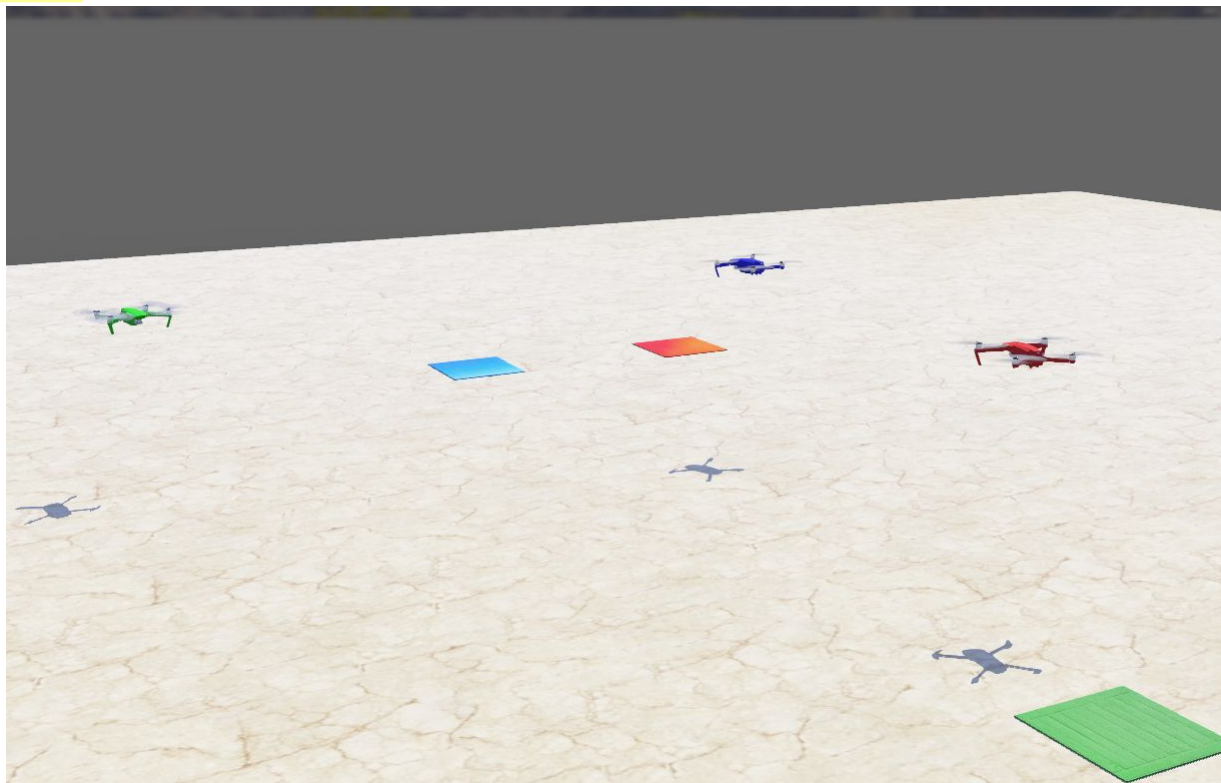
1. Multi Agent Reinforcement Learning is when we are considering various intelligent agents interacting with an environment
2. Multi Agent Systems is still an under-explored area of reinforcement learning
3. Multi Agent scenario
 - a. Cooperative
 - b. Competitive
 - c. Mixture of both

RL world



RL feedback loop[1]

Continued



Formulation

1. Actor-Critic Methods
2. Each agent is modelled as an Actor
 - a. Updates its policy as suggested by the Critic
3. The Critic estimates the values of the state
 - a. It guides the actor to prefer good actions over bad ones

MADDPG (Lowe *et al.*)[2]

1. Multi Agent Deep Deterministic Policy Gradient
2. Centralized Critic network
3. Decentralized execution (multiple independent actors)
4. Deterministic Policy function

Architecture

- An actor network for using observations for *deterministic* actions
- An identical target actor network for training stability
- A critic network that uses *joint* states action pairs to estimate Q-values

Issues

- No obstacle avoidance : static, dynamic
- No safety guarantees
- Large state space configuration
- Techniques :
 - Safety Controllers[3]
 - Control Barrier functions[4]

Safe Reinforcement Learning

1. Learning optimal policies while satisfying Safety Constraints
2. Safe RL will help in transitioning the control to physical systems
3. A major challenge in using reinforcement learning is safety - control policies learned using reinforcement learning typically do not provide any safety guarantees.

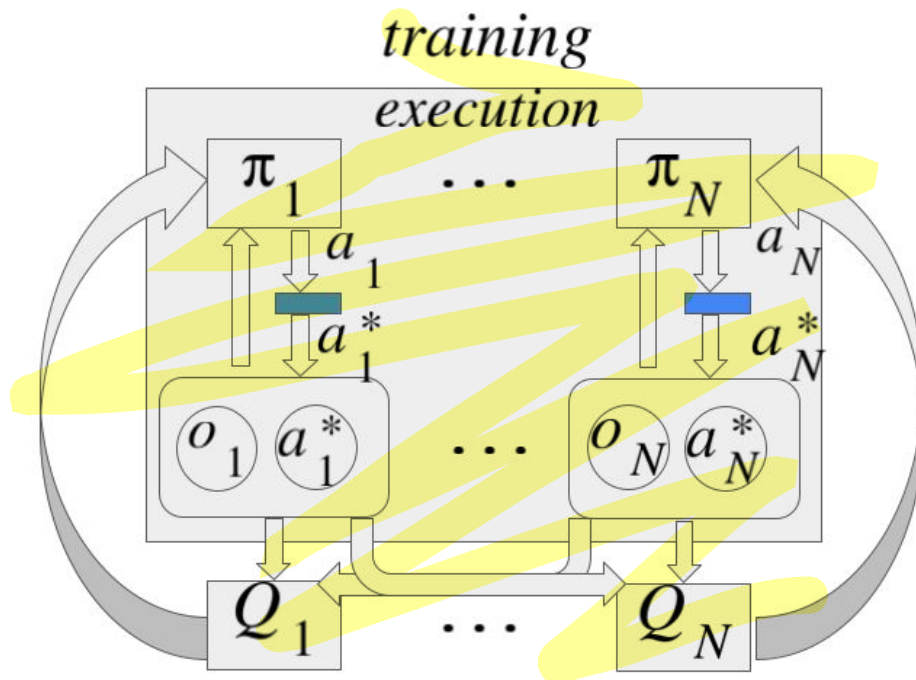
Safety Layer (Dalal *et al.*)[3]

1. Modifications to original policy
2. Learning immediate constraint functions
3. After Safety Layer is trained, safe actions a^* will be used by Agents(UAVs):

$$a_i^* = \arg \min_{a|i} \frac{1}{2} \| a_i - \mu_{i|\theta}(o) \|^2$$

$$c_k(o_i, a) + g(o, w)^T a_i \leq C_k, \forall k \in K$$

Safe-MADDPG



Safe-MADDPG

STEP 1:

```
1 while  $epochs \leq TOTAL\ EPOCHS$  do
2   while  $step \leq TOTAL\ SAMPLE\ STEPS$  do
3     Reset Environment(state =  $o_i, \forall i \in N$ );
4     Calculate constraint values  $c_i$ ;
5     Execute actions ( $a_1, \dots, a_N$ ) and observe reward
       $r_i$  and new obs  $o'_i$  for each agent  $\forall i \in N$ ;
6     Calculate constraint values  $c_i^{next}$ ;
7     Store ( $a_i, o_i, c_i, c_i^{next}$ ) in replay buffer  $D$ ;
8   end
```

Safety Network update:

```
9   while update via mini-batch do
10     Sample a mini-batch for each agent  $i \in N$  :
      ( $a_i, o_i, c_i, c_i^{next}$ );
11     Observe safety layer :  $g_i(o_i, w_i)$ ;
12     Predict next constraint values;
13      $c_i^{next.predicted} = c_i + g_i(o_i, w_i)a_i$ ;
14     Calculate loss and update safety layer weights;
15      $L^k = \| c_i^{next} - c_i^{next.predicted} \|^2, \forall k \in K$ 
16   end
17 end
```

Safe-MADDPG

STEP 2:

```

18 while episode number  $\leq$  TOTAL EPISODES do
19   Initialize exploration noise;
20   Reset environment (state =  $o_i$ ,  $\forall i \in N$ );
21   while timestep( $t$ )  $\leq$  TOTAL EPISODE LENGTH do
22     For each agent  $i$ , select action
23        $a_i = \mu_{i|\theta}(o_i) + \mathcal{N}_t$  w.r.t. the current policy
24       and exploration;
25     calculate lagrange multipliers :
26        $\lambda_i^* = \frac{g(o_i, w_i)^T \mu_{i|\theta}(o_i) + c_i(o_i) - C_i}{g(o_i, w_i)^T g(o_i, w_i)}$ ;
27     correct each agent's action as :
28        $a_i^* = \mu_{i|\theta}(s) - \lambda_i^* g(o_i, w_i)$ ;
29     Execute corrected actions ( $a_1^*, \dots, a_N^*$ ) and
30       observe reward  $r_i$  and new obs  $o'_i$  for agent  $i$ ,
31        $\forall i \in N$ ;
32     Store ( $o_i, a_i^*, r_i, o'_i$ ) in replay buffer  $D$ ;
33      $o_i \leftarrow o'_i$ ;

```

```

28 while agent  $i = 1$  to  $N$  do
29   Sampling of a randomized minibatch of  $M$ 
30   samples ( $o^j, a^j, r^j, o'^j$ ) from  $D$ ;
31   set  $y_i^j =$ 
32      $r_i^j + \gamma Q_i^{\mu'}(o_i'^j, a_1^j, \dots, a_N^j)|_{a_i' = \mu_{i|\theta}(o_i^j)}$ ;
33   Update of the critic by minimization of the
34   loss  $L(\theta_c) =$ 
35      $\frac{1}{M} \sum_j (y_i^j - Q_i^\mu(o_i^j, a_1^j, \dots, a_N^j))^2$ ;
36   Update of the actor by the use of sampled
37   policy gradient;
38    $\nabla_{\theta_i} J \approx \frac{1}{M} \sum_j \nabla_{\theta_i} \mu_i(o_i^j)$ 
39    $\nabla_{a_i} Q_i^\mu(x_i^j, a_1^j, \dots, a_N^j)|_{a_i = \mu_i(o_i^j)}$ ;
40 end
41 Update of the target network parameters for
42 each agent  $i$ ;
43  $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$ ;
44 end

```

Environments

1. Multiple UAVs and their Target locations
2. Quadcopter: DJI Mavic
3. Multi-Agent Particle Environments(MAPE) by OpenAI (2D)
4. Webots Simulation by Cyberbotics (3D)
 - a. Real world physics engine
 - b. Deepbots[3]: Deepbots is a simple framework which is used as middleware between the Webots Simulator and Reinforcement Learning algorithm
 - c. Emitter-Receiver scheme
 - d. Combined robot-supervisor controller scheme

Continued..

- Observation Space(UAVs and Targets):

Position of each entity: $P_i : [p_i^x, p_i^y, p_i^z]$

Vel. of each UAV: $V_i : [l_i^x, l_i^y, l_i^z, v_i^x, v_i^y, v_i^z]$

Concatenated Tuple $\{P^i, V^i, P_{Target}^i, \bigcup_{j \in N; j \neq i} P^j\}$

- Action Space:

For each UAV: $a_i : \{pitch_i, yaw_i\}$

- Reward Structure:

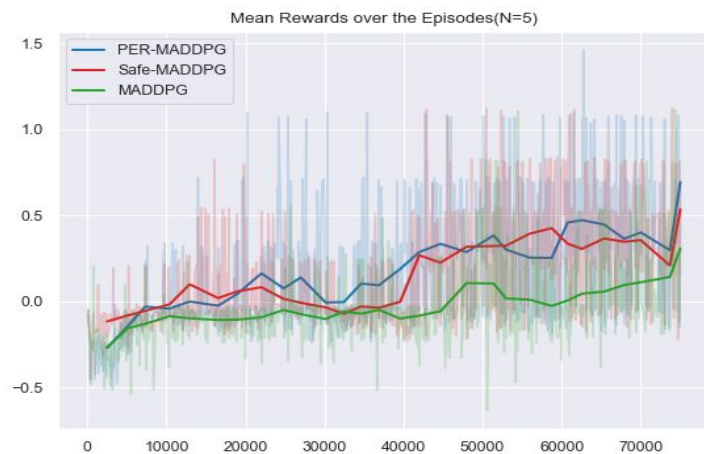
Distance, Angular and Collision Penalty

| Reward Value | Conditions |
|--------------|--|
| -4 | if the angle between agent and its target is > 1.5 Radians |
| -2 | if angle < 1.5 Radians and $curr_dist > prev_dist$ |
| -2 | if angle < 0.1 Radians and $curr_dist > prev_dist$ |
| +5 | if angle < 1.5 Radians and $curr_dist < prev_dist$ |
| +10 | if angle < 0.1 Radians and $curr_dist < prev_dist$ |
| -10 | if there is a collision |

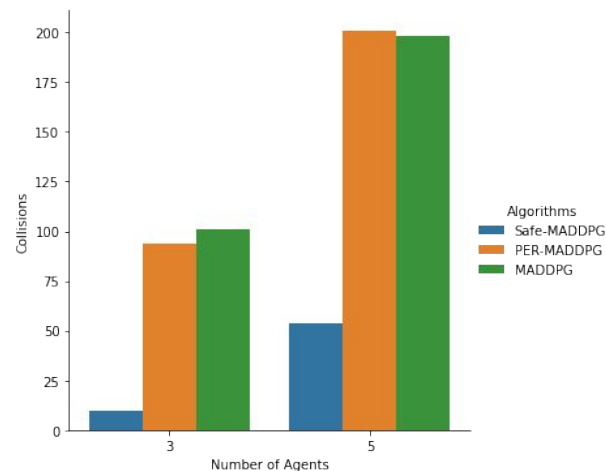
Results of 2D MAPE environment

Evaluation

- Mean Rewards

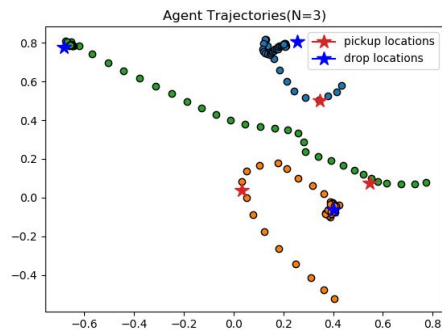


- Testing final Policies

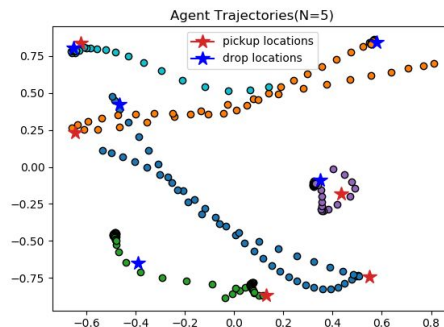


Trajectory Plotting(N=3) and (N=5)

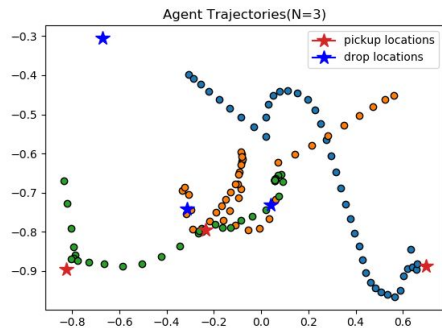
Max. reward : 56.19



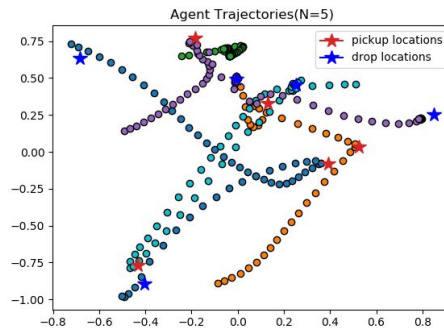
Max. reward : 88.02



Min. reward : 11.02

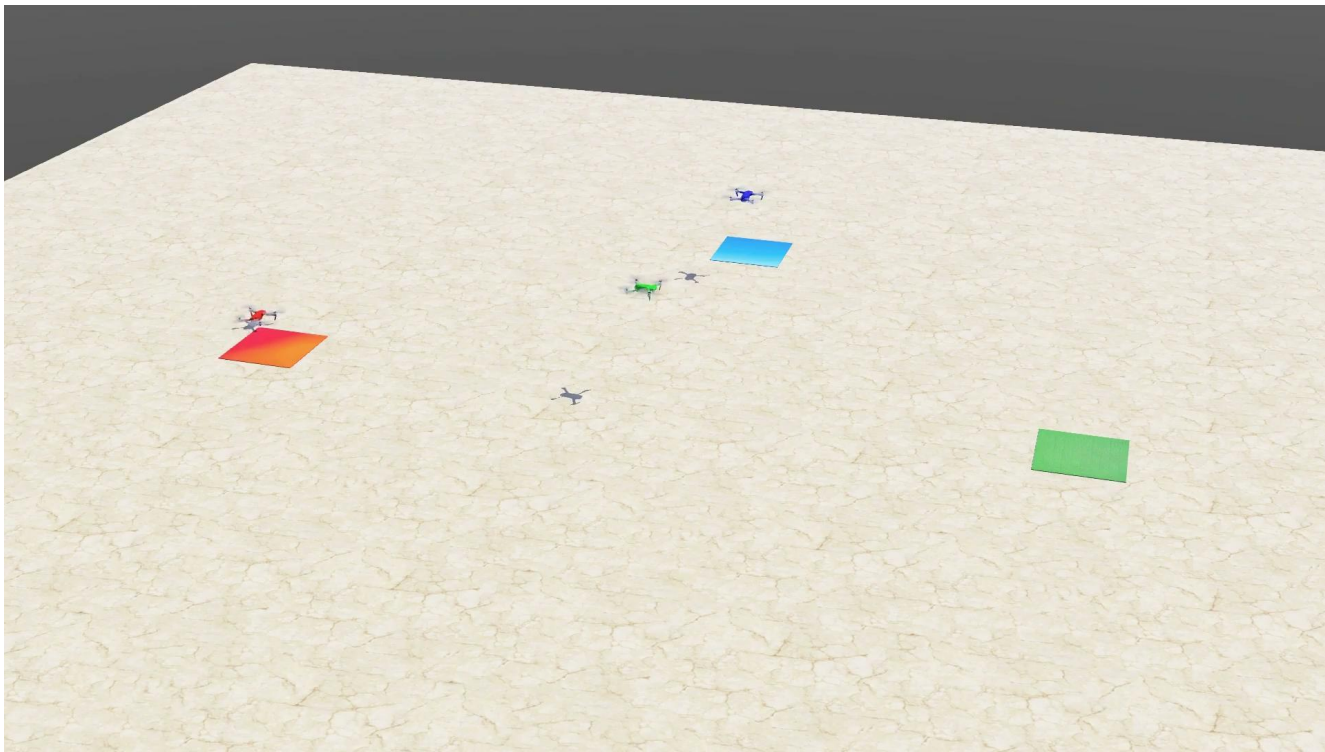


Min. reward : 24.11



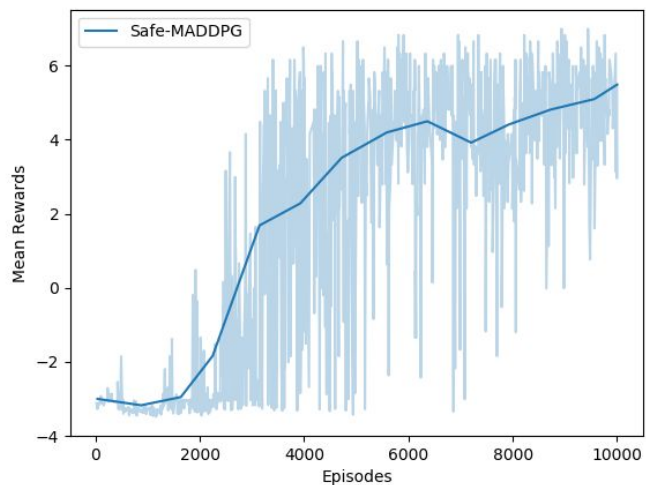
Results of 3D Webots environment

3-Agent Scenario

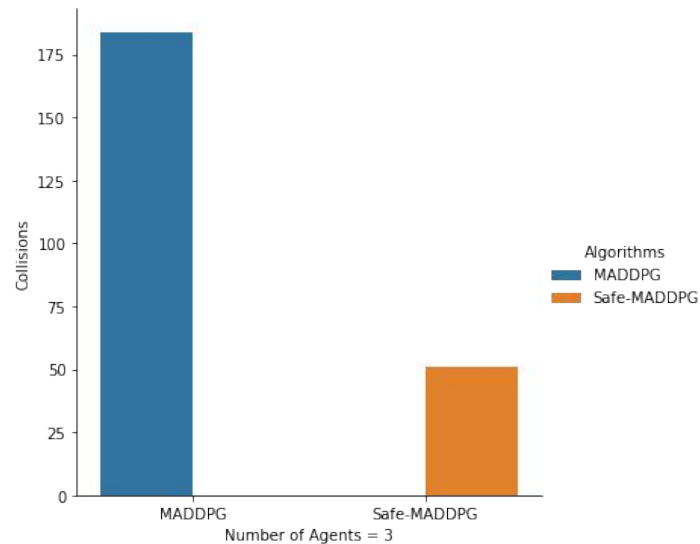


Evaluation

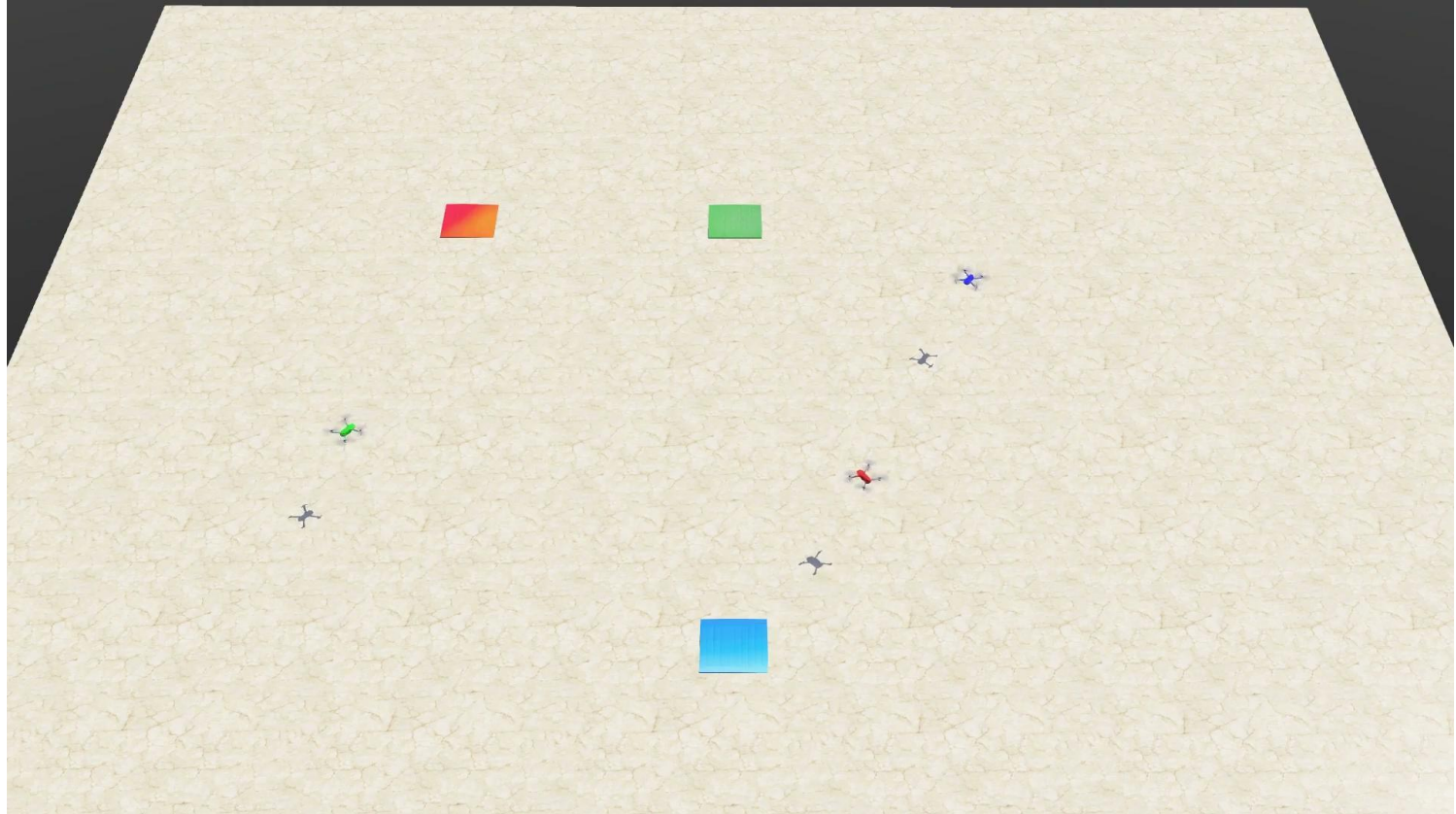
- Mean Rewards



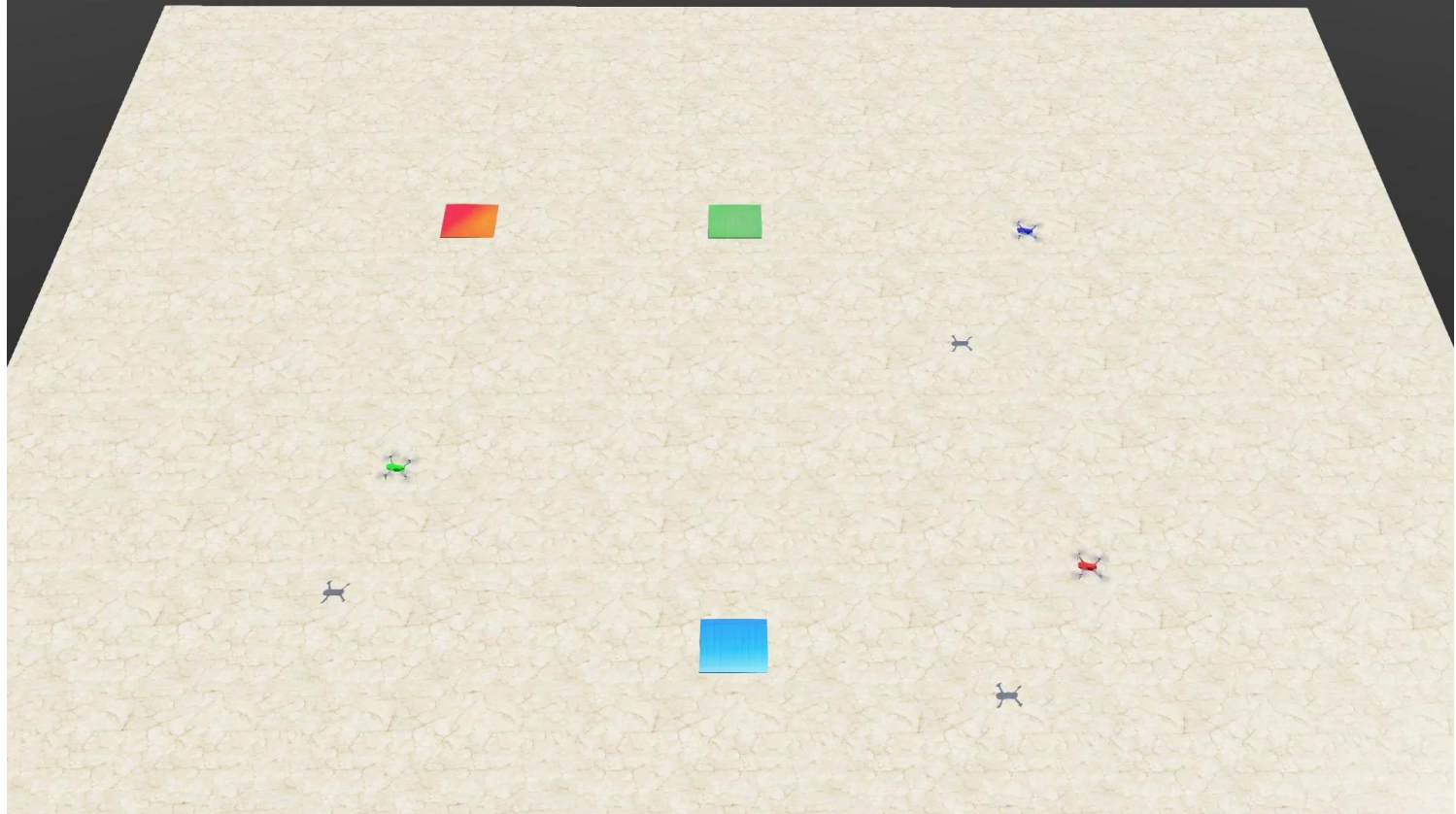
- Testing Final Policies



Behaviour-MADDPG

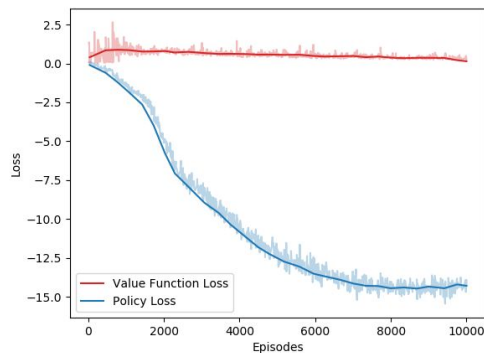


Behaviour-SafeMADDPG

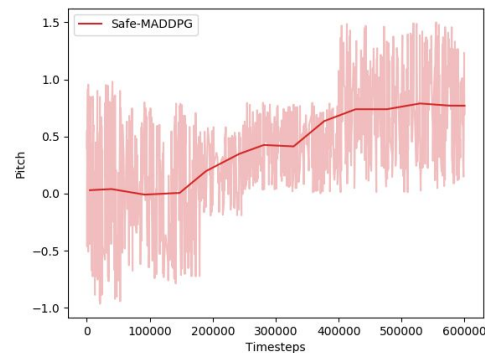


Evaluation Continued..

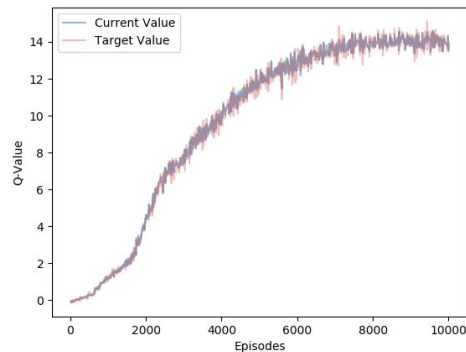
Value & Policy Loss



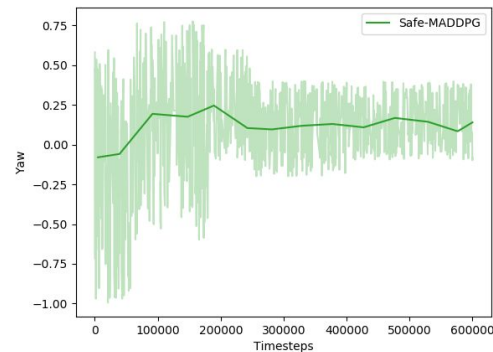
Pitch variation



Current - Target Q-Value



Yaw variation



Reviewer Comments

- It is not very clear what is observations, what is state space and what is action space. Why if an agent is taking an action a_i based on the observation o_i , then state of every other agent is changing. What is θ_i ?
- On page 3, the student mentions that "the bias is also corrected in the loss function during the training". What kind of bias are we talking about?
- In safe reinforcement learning, what is C_k ? how easy or difficult in general it is to specify the constraints?
- If the constraints are already enforced, what is the use of the safety layer?
- In algorithm 1, where is L^k used in rest of the algorithm?

Continued..

- I understand that the "safety" aspect of the MADDPG the central novelty claim of the paper. This "safety" aspect should be clearly defined and illustrated right in the beginning sections of the paper.

References

1. https://www.researchgate.net/figure/Deep-Reinforcement-Learning-13_fig2_35109551
2. <https://arxiv.org/abs/1706.02275>
3. <https://arxiv.org/abs/1801.08757>
4. <https://arxiv.org/abs/1903.08792>