Open in app

# Adrian Chow

About

# Multi-Agent Control using Deep Reinforcement Learning

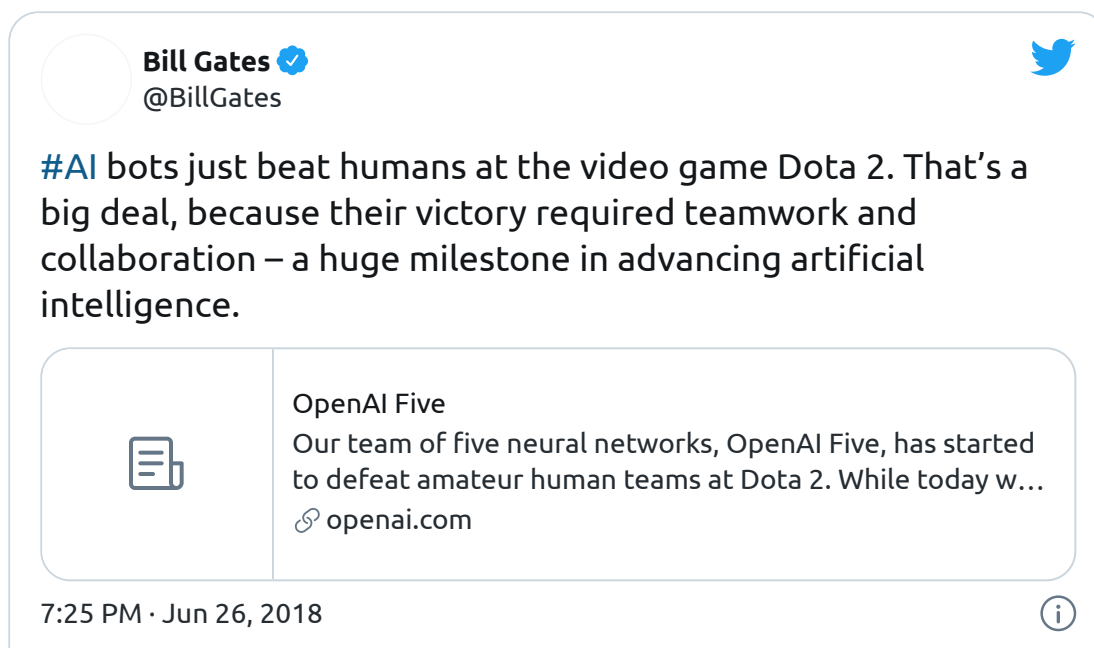A   Adrian Chow  2 days ago · 4 min read

One difficult task in the Deep Reinforcement Learning space is having multiple agents interact and learn. The reality is that the world is a multi-agent environment in which intelligence is developed by interacting with multiple agents. A major breakthrough in this regard was when researchers at DeepMind developed an AI engine, alphago zero, which learned to play Go by training in a multi-agent environment. The agent developed such a high level of competency it was able to beat LeeSedol, a professional Go player. Extensions of this finding led to OpenAI Five which was able to conquer a complex game such as Dota 2. Clearly, DRL in a multi-agent environment has expanded the capabilities of Markov learning.
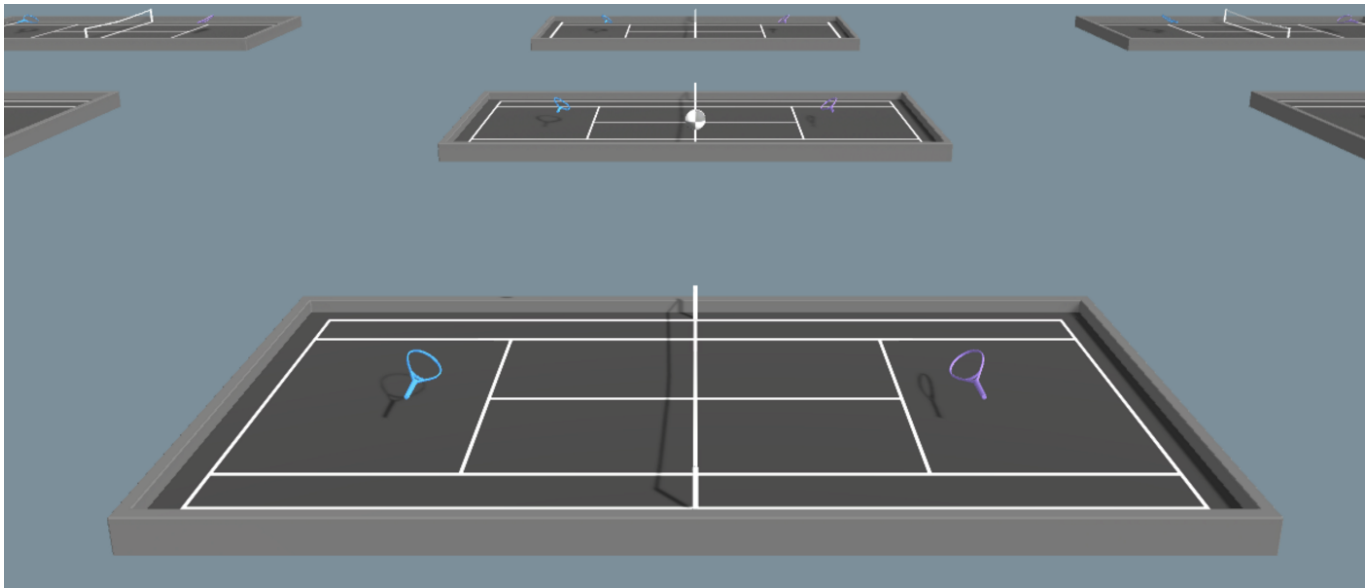
Source: OpenAI



I have taken the challenge to myself to implement a multi-environment agent using my pre-existing knowledge of RL to solve the Tennis environment on the Unity ML-Agents. The goal of this project is to teach two agents to play a low-level version of tennis. By applying deep reinforcement learning to this environment, two separate agents compete in order to define their individual policies. This method allows competition to essential become collaboration as agents compete to find the optimal policy. The following project has been an implementation based on OpenAi's **Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments** research paper.

## The Environment

Source: Unity-Technologies

## RL Problem Specifications

- Goal of Agent: keep the ball in play

- Rewards : `+0.1` every time agent hits the ball over the net, `-0.1` every time agent lets a ball hit the ground or hits the ball out of bounds

- Action Space — Continuous, 4-Dimension Vector Space from [-1, +1]

- State Space — Continuous, 8-Dimension, 3 Stacked Observations

- Solving Condition: Average score of +0.5 over 100 consecutive episodes.

The observation space consists of 8 variables corresponding to the position and velocity of the ball and racket. Each agent receives its own, local observation. Two continuous actions are available, corresponding to the movement toward (or away from) the net, and jumping.

## Learning Algorithm

The algorithm is a multi-agent variation of the standard DDPG algorithm that I have implemented in the past. If you are unfamiliar with the Deep Deterministic Policy Gradient algorithm, you can check out my other medium post. Essentially this article is

an extension of the previous DDPG I have implemented. Here is the basis of the algorithm:

---

**Algorithm 1** DDPG algorithm

---

Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights $\theta^Q$ and $\theta^\mu$.
Initialize target network $Q'$ and $\mu'$ with weights $\theta^{Q'} \leftarrow \theta^Q$, $\theta^{\mu'} \leftarrow \theta^\mu$
Initialize replay buffer $R$
**for** episode = 1, M **do**
   Initialize a random process $\mathcal{N}$ for action exploration
   Receive initial observation state $s_1$
   **for** t = 1, T **do**
      Select action $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise
      Execute action $a_t$ and observe reward $r_t$ and observe new state $s_{t+1}$
      Store transition $(s_t, a_t, r_t, s_{t+1})$ in $R$
      Sample a random minibatch of $N$ transitions $(s_i, a_i, r_i, s_{i+1})$ from $R$
      Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$
      Update critic by minimizing the loss: $L = \frac{1}{N}\sum_i(y_i - Q(s_i, a_i|\theta^Q))^2$
      Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N}\sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu}\mu(s|\theta^\mu)|_{s_i}$$

      Update the target networks:
$$\theta^{Q'} \leftarrow \tau\theta^Q + (1-\tau)\theta^{Q'}$$
$$\theta^{\mu'} \leftarrow \tau\theta^\mu + (1-\tau)\theta^{\mu'}$$

   **end for**
**end for**

---

DDPG Algorithm (Source: Here)

Some constraints/assumptions to make when implementing our multi-agent policy include:

- The learned policies can only use local information (i.e. their own observations) at execution time

- Do not assume a differentiable model of the environment dynamics

- Do not assume any particular structure on the communication method between agents

*Multi-Agent DDPG Initialization:*

```
1   class MADDPG:
2       """ Class for training multiple agents in the multi-agent environment"""
3
4       def __init__(self, state_size, action_size, num_agents, seed):
5
6           super(MADDPG, self).__init__()
7           self.state_size = state_size
8           self.action_size = action_size
9           self.num_agents = num_agents
10          self.seed = random.seed(seed)
11
12          # Initialize agents in the multi-agent environment
13          self.DDPGs = [DDPG(state_size, action_size, num_agents, seed) for i in range(num
14
15          # Replay Buffer (shared by all agents)
16          self.memory = ReplayBuffer(action_size, num_agents, BUFFER_SIZE, BATCH_SIZE, see
```

**MAPPG_init.py** hosted with ❤ by **GitHub**                                                      **view raw**

- Notice the number of agents is passed in as a parameter for the MADDPG. In the case of the Tennis Environment, there are two agents.

- Individual agents with their own set of actors and critics are created based on the number of agents, see line 13.

- The replay buffer will be the same as before, except instead this time an experience will contain the states and actions for both agents in a given timestep.

- Note the DDPG follows the same class implemented in the previous project. See below:

```
1   class DDPG:
2       """ Base Class for an Agent in the multi-agent environment"""
3
4       def __init__(self,
5                    state_size,
6                    action_size,
7                    num_agents,
8                    seed,
9                    hidden_in_actor=200,
10                   hidden_out_actor=150,
11                   hidden_in_critic=200,
```

```python
12                      hidden_out_critic=150,
13                  lr_actor=LR_ACTOR,
14                  lr_critic=LR_CRITIC):
15
16          super(DDPG, self).__init__()
17
18          self.state_size = state_size
19          self.action_size = action_size
20          self.num_agents = num_agents
21          self.seed = random.seed(seed)
22
23          self.actor = Actor(state_size,
24                             action_size,
25                             seed,
26                             hidden_in_actor,
27                             hidden_out_actor).to(device)
28
29          self.target_actor = Actor(state_size,
30                                    action_size,
31                                    seed,
32                                    hidden_in_actor,
33                                    hidden_out_actor).to(device)
34
35          self.critic = Critic(state_size,
36                               action_size,
37                               num_agents,
38                               seed,
39                               hidden_in_critic,
40                               hidden_out_critic).to(device)
41
42          self.target_critic = Critic(state_size,
43                                      action_size,
44                                      num_agents,
45                                      seed,
46                                      hidden_in_critic,
47                                      hidden_out_critic).to(device)
48
49          # Ornstein Uhlenbeck Noise for Action Space Exploration
50          self.noise = OrnsteinUhlenbeckNoise(action_size, seed)
51
52          # Initialize targets same as original networks
53          self.copy_weights(self.critic, self.target_critic)
54          self.copy_weights(self.actor, self.target_actor)
55
56          # Actor and Critic Adam Optimizers
```

```python
57          self.actor_optimizer = optim.Adam(self.actor.parameters(), lr=lr_actor)
58          self.critic_optimizer = optim.Adam(self.critic.parameters(), lr=lr_critic)
59
60
61      def act(self, state, add_noise=True):
62          """Returns actions for given state as per current policy."""
63
64          state = torch.from_numpy(state).float().to(device)
65          self.actor.eval()
66          with torch.no_grad():
67              action = self.actor(state).cpu().data.numpy()
68          self.actor.train()
69          if add_noise:
70              action += self.noise.noise()
71          return np.clip(action, -1, 1)
72
73
74      def step(self, experiences, gamma):
75          self.learn(experiences, gamma)
76
77
78      def learn(self, experiences, gamma):
79          """Update policy and value parameters using given batch of experience tuples.
80          Q_targets = r + y * target_critic(next_state, target_actor(next_state))
81          where:
82              target_actor(state) -> action
83              target_critic(state, action) -> Q-value
84          Params
85          ======
86              experiences (Tuple[torch.Tensor]): tuple of (s, a, r, s', a2, s2' done) tup
87              gamma (float): discount factor
88          """
89
90          states, actions, rewards, next_states, dones = experiences
91
92          next_states_tensor = torch.cat(next_states, dim=1).to(device)
93          states_tensor = torch.cat(states, dim=1).to(device)
94          actions_tensor = torch.cat(actions, dim=1).to(device)
95
96          # --------------------------- update critic --------------------------- #
97          # Get predicted next-state actions and Q values from target models
98
99          next_actions = [self.actor(state) for state in states]
100         next_actions_tensor = torch.cat(next_actions, dim=1).to(device)
```

```
101             Q_targets_next = self.target_critic(next_states_tensor, next_actions_tensor)
102             # Compute Q targets for current states (y_i)
103             Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))
104             # Compute critic loss
105             Q_expected = self.critic(states_tensor, actions_tensor)
106             critic_loss = F.mse_loss(Q_expected, Q_targets)
107             # Minimize the loss
108             self.critic_optimizer.zero_grad()
109             critic_loss.backward()
110             self.critic_optimizer.step()
111
112             # --------------------------- update actor --------------------------- #
113             # Compute actor loss
114             actions_pred = [self.actor(state) for state in states]
115             actions_pred_tensor = torch.cat(actions_pred, dim=1).to(device)
116             actor_loss = -self.critic(states_tensor, actions_pred_tensor).mean()
117             # Minimize the loss, thererby maximizing the reward
118             self.actor_optimizer.zero_grad()
119             actor_loss.backward()
120             self.actor_optimizer.step()
121
122             # ---------------------- update target networks ---------------------- #
123             self.soft_update(self.critic, self.target_critic, TAU)
124             self.soft_update(self.actor, self.target_actor, TAU)
125
126
127     def soft_update(self, local_model, target_model, tau):
128         """Soft update model parameters.
129         θ_target = τ*θ_local + (1 - τ)*θ_target
130         Params
131         ======
132             local_model: PyTorch model (weights will be copied from)
133             target_model: PyTorch model (weights will be copied to)
134             tau (float): interpolation parameter
135         """
136         for target_param, local_param in zip(target_model.parameters(), local_model.par
137             target_param.data.copy_(tau*local_param.data + (1.0-tau)*target_param.data)
138
139
140     def reset(self):
141         self.noise.reset()
142
143
144     def copy_weights(self, source, target):
145         """Copies the weights from the source to the target"""
```

```
146              for target_param, source_param in zip(target.parameters(), source.parameters())
147                  target_param.data.copy_(source_param.data)
```

ddpg.py hosted with ❤ by **GitHub**                                                              **view raw**

DDPG, <u>Source</u>

- DDPG is still using then Ornstein Uhlenbeck Noise for action space exploration.

- The replay buffer is shared between agents, that is why it is stored in the MADDPG class.

- Since there are multiple states, actions per experience they are concatenated into single tensors, see lines 100 and 115.

*Multi-Agent DDPG Methods:*

```
1    class MADDPG:
2        """ Class for training multiple agents in the multi-agent environment"""
3
4        def __init__(self, state_size, action_size, num_agents, seed):
5
6            super(MADDPG, self).__init__()
7            self.state_size = state_size
8            self.action_size = action_size
9            self.num_agents = num_agents
10           self.seed = random.seed(seed)
11
12           # Initialize agents in the multi-agent environment
13           self.DDPGs = [DDPG(state_size, action_size, num_agents, seed) for i in range(num
14
15           # Replay Buffer (shared by all agents)
16           self.memory = ReplayBuffer(action_size, num_agents, BUFFER_SIZE, BATCH_SIZE, see
17
18
19       def act(self, states, add_noise=True):
20           """Returns actions for each agent."""
21
22           return [agent.act(state, add_noise) for agent, state in zip(self.DDPGs, states)]
23
24
25       def step(self, states, actions, rewards, next_states, dones):
26
```

```
27              self.memory.add(states, actions, rewards, next_states, dones)

28

29          for agent in self.DDPGs:
30              if len(self.memory) > BATCH_SIZE:
31                  experiences = self.memory.sample()
32                  agent.step(experiences, GAMMA)
```

**maddpg_methods.py** hosted with ❤ by **GitHub**                              **view raw**

MADDPG, <u>Source</u>

- Per the time step for MADDPG, each DDPG is updated normally (randomly sampling from replay buffer) to allow each agent to learn independently from each other.

*Multi-Agent DDPG Training:*

- Before we get into the training aspect here are some of the hyperparameters that have worked for me in this project. Notice that I have increased the batch_size, due to the increase in learning capacity required.

```
1   BUFFER_SIZE = int(1e5)  # replay buffer size
2   BATCH_SIZE = 250         # minibatch size
3   GAMMA = 0.99            # discount factor
4   TAU = 1e-3             # for soft update of target parameters
5   LR_ACTOR = 1e-4         # learning rate of the actor
6   LR_CRITIC = 1e-3        # learning rate of the critic
```

**hyperparameters.py** hosted with ❤ by **GitHub**                              **view raw**

Hyperparameters

- The critic and actor networks have stayed the same for the most part. One minor change is the fact that the critic now accounts for the states of both agents when calculating the return. This minor change allows for the critic to adapt to a broader state space and converge better.

```
1   def hidden_init(layer):
2       fan_in = layer.weight.data.size()[0]
3       lim = 1. / np.sqrt(fan_in)
4       return (-lim, lim)

5

6   class Actor(nn.Module):
```

```python
 6   class Actor(nn.Module):
 7       """Actor (Policy) Model."""
 8
 9       def __init__(self, state_size, action_size, seed, fc1_units=200, fc2_units=150):
10           """Initialize parameters and build model.
11           Params
12           ======
13               state_size (int): Dimension of each state
14               action_size (int): Dimension of each action
15               seed (int): Random seed
16               fc1_units (int): Number of nodes in first hidden layer
17               fc2_units (int): Number of nodes in second hidden layer
18           """
19           super(Actor, self).__init__()
20           self.seed = torch.manual_seed(seed)
21           self.fc1 = nn.Linear(state_size, fc1_units)
22           self.fc2 = nn.Linear(fc1_units, fc2_units)
23           self.fc3 = nn.Linear(fc2_units, action_size)
24           self.reset_parameters()
25
26       def reset_parameters(self):
27           self.fc1.weight.data.uniform_(*hidden_init(self.fc1))
28           self.fc2.weight.data.uniform_(*hidden_init(self.fc2))
29           self.fc3.weight.data.uniform_(-3e-3, 3e-3)
30
31       def forward(self, state):
32           """Build an actor (policy) network that maps states -> actions."""
33           x = F.relu(self.fc1(state))
34           x = F.relu(self.fc2(x))
35           return F.torch.tanh(self.fc3(x))
36
37
38   class Critic(nn.Module):
39       """Critic (Value) Model."""
40
41       def __init__(self, state_size, action_size, num_agents, seed, fcs1_units=200, fc2_un
42           """Initialize parameters and build model.
43           Params
44           ======
45               state_size (int): Dimension of each state
46               action_size (int): Dimension of each action
47               seed (int): Random seed
48               fcs1_units (int): Number of nodes in the first hidden layer
49               fc2_units (int): Number of nodes in the second hidden layer
50           """
```

```python
51              super(Critic, self).__init__()
52              self.seed = torch.manual_seed(seed)
53              self.num_agents = num_agents
54              self.fcs1 = nn.Linear((state_size+action_size) * num_agents, fcs1_units)
55              self.fc2 = nn.Linear(fcs1_units, fc2_units)
56              self.fc3 = nn.Linear(fc2_units, 1)
57              self.reset_parameters()
58
59          def reset_parameters(self):
60              self.fcs1.weight.data.uniform_(*hidden_init(self.fcs1))
61              self.fc2.weight.data.uniform_(*hidden_init(self.fc2))
62              self.fc3.weight.data.uniform_(-3e-3, 3e-3)
63
64          def forward(self, state, action):
65              """Build a critic (value) network that maps (state, action) pairs -> Q-values."""
66              xs = torch.cat((state, action), dim=1)
67              x = F.relu(self.fcs1(xs))
68              x = F.relu(self.fc2(x))
69              return self.fc3(x)
```

networks_maddpg.py hosted with ❤ by **GitHub**                                                                    **view raw**

- The training loop is the same as usual

```python
1
2    def maddpg_train(agent,
3                     env,
4                     brain_name,
5                     n_agents,
6                     n_episodes=10000,
7                     max_t=10000,
8                     print_every=100,
9                     win_condition=0.5):
10
11       """
12
13       Tennis using MADDPG.
14
15       Params
16       ======
17           n_episodes (int): maximum number of training episodes
18           max_t (int): maximum number of timesteps per episode
19           print_every (int): how many episodes before printing scores
```

```
20            n_agents (int): how many agents are in the environment
21
22        """
23
24        scores_deque = deque(maxlen=100)
25        scores = []
26        average_scores_list = []
27
28        for i_episode in tqdm(range(1, n_episodes+1)):
29
30            env_info = env.reset(train_mode=True)[brain_name] # reset the environment
31            states = env_info.vector_observations          # get the current state (for e
32            score = np.zeros(n_agents)                     # initialize the score (for each
33
34            for t in range(max_t):
35                actions = agent.act(states)               # consult agent for actions
36                env_info = env.step(actions)[brain_name]   # take a step in the env
37                next_states = env_info.vector_observations # get next state (for each age
38                rewards = env_info.rewards                 # get reward (for each agent)
39                dones = env_info.local_done                # see if episode finished
40                agent.step(states, actions, rewards, next_states, dones)  # take a learning
41                score += rewards                           # update the score (for each a
42                states = next_states                       # roll over states to next tim
43                if np.any(dones):                          # exit loop if episode finishe
44                    break
45
46            score_max = np.max(score)
47            scores.append(score_max)
48            scores_deque.append(score_max)
49            average_score = np.mean(scores_deque)
50            average_scores_list.append(average_score)
51
52            print('\rEpisode {}\tAverage Score: {:.3f}'.format(i_episode, np.mean(scores_deq
53
54            if i_episode % 100 == 0:
55                print('\rEpisode {}\tAverage score: {:.3f}'.format(i_episode , average_score
56
57            # Winning condition + save model parameters
58            if average_score >= win_condition:
59                print("\rSolved in episode: {} \tAverage score: {:.3f}".format(i_episode , a
60                break
61
62        execution_info = {'last_score': scores.pop(),
63                          'solved_in': i_episode,
64                          'last_100_avg': average_score}
```
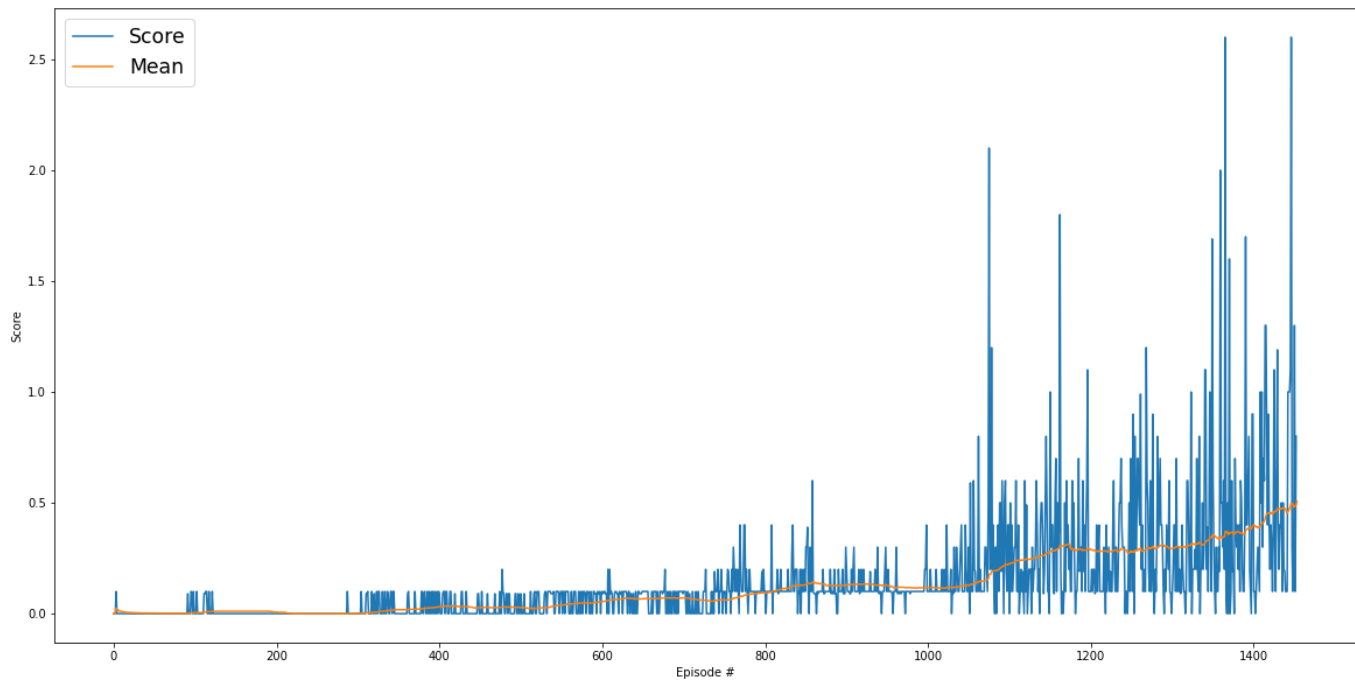
```
65
66          return scores, scores_deque, execution_info
```

- The results from training were approximately on-point with the baseline training requirements presented for this environment. I was able to solve the environment in 1454 episodes.



Training Score-Episodes Plot. Source

*Possible Improvements*

- Modify the critic and actor to account for the agent's belief of the opponent's next move. If the agent can understand the opponent's policy it may be able to perform better?

- Update the agents selectively or stochastically.

- Decrease action exploration as the training episodes increase.