

Project 2: Tennis

Report

The game is episodic with state space of 8 variables and action space of size 2. The solution was implemented as a DDPG actor-critic agent architecture. The two agents use a shared replay buffer while training, however they learn separately, each with its own actor and critic local and target networks .

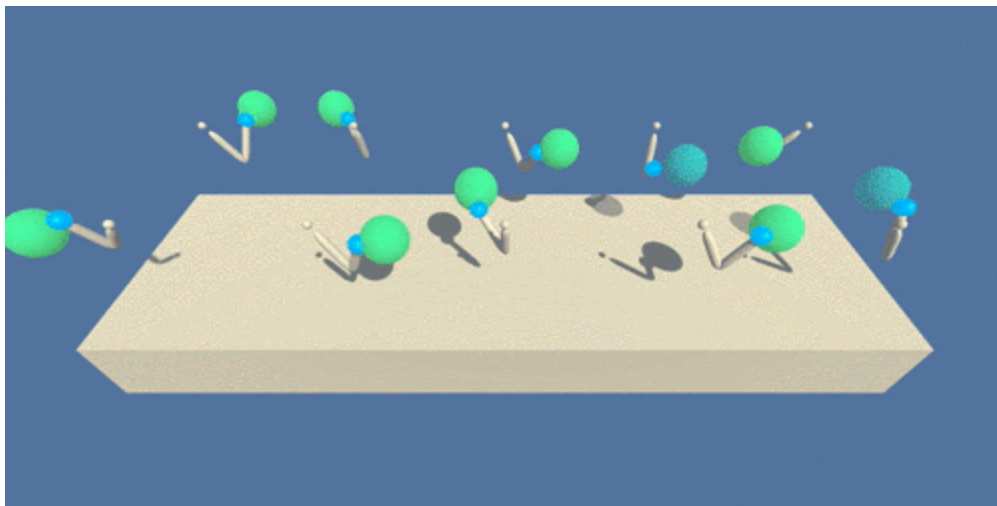
A $1e6$ length buffer was defined to collect the experience of both agents while learning. This buffer is then periodically sampled to get a 128 long batch which is then used by both agents for learning.

Agent structure:

The agent actor is a simple 2 fully connected (linear+relu) layers of size 256+128 in addition to the output layer with a size equal to the action space size (2).

The critic is identical to the structure of the actor except that the actor uses tanh for the output layer to get actions within the +1/-1 range.

Both actor & critic learning rates were set to $1e-3$



Learning:

Given an experience tuple both agent goes through the following steps to update its actor & critic:

Local networks update:

Critic update:

- Compute the **Q-expected** passing the states and actions to the critic local network
- Compute the next actions for the next states using the actor target network without noise addition (deterministic policy). Then compute the **Q-target** passing the next states and previously computed next actions to the critic target network
- Compute the loss function from the **Q-expected/Q-target**, set the gradient clip normalization and train the critic.

Actor update:

- Compute the actions for the experience states. Here we don't use the actions from the experience tuple because they were noise-added to ensure stochastic exploring interaction policy and we need to learn a deterministic policy.
- Compute the actor loss function from the states and the previously computed deterministic actions and train the actor

Target networks update:

Actor & Critic target networks are updated through soft update where the local network parameters are tau-weighted and then added to (1-tau)-weighted target parameters to obtain the new updated parameters. Tau was set to $6e-2$

The agent has the following capabilities:

1. Act: given the current state of the environment, it runs the state through the available agent nn model and selects the action according to the epsilon-greedy algorithm.
2. Step: given a complete experience, it saves it in the replay memory and if it's the fifth experience since the last learning it calls the learning routine.
3. Learn: given a sampled batch of size 128 experiences and the learning discount factor, it computes the error of the current agent model (which is computed as the `mmean_square_error` between the agent current q-values of the given states to the

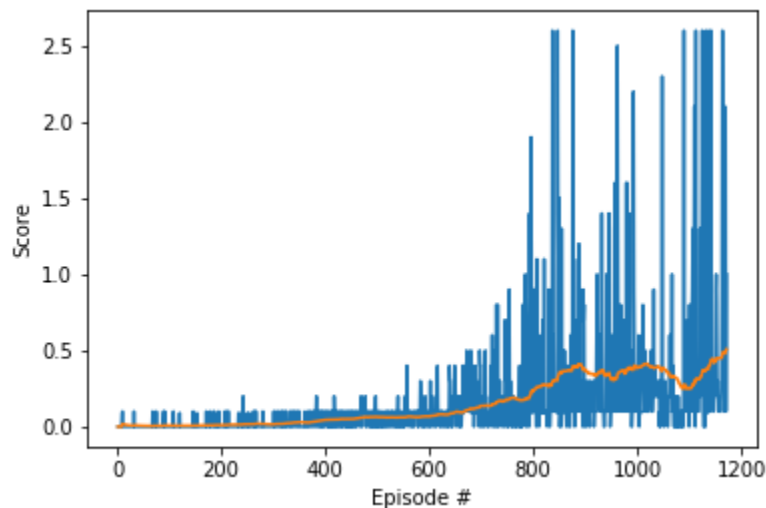
target q-values of the same states computed from the target nn) and propagates it back to update the agent model.

4. Soft_update: given the target-nn model update factor and the current agent-nn model, it updates the target-nn model in a way to ensure smooth slower change in the target-nn model so as to improve the whole DQN stability.

Performance:

The two agents were set to play for (1200) episodes and learn throughout. Finally the agents learnt models are saved periodically every 100 episodes.

Below is the plot showing the learning performance. The average score exceeds +0.5 after 1171 episodes.



Future Ideas for improving theAgent performance:

The solution here implemented has used two ddpg agent learning separately. One interesting approach to try is the MADDPG where both agents share a central critic while maintaining separate actors (with their local and target networks). This approach might be of better performance because the shared critic will have more information about the environment than a separate critic, since the two agents' observations are not identical (the info about both rackets will be seen only by a shared critic).