

1) Introduction

- The goal is to get an average score of +0.5 over 100 consecutive episodes.

2) Environment

- Environment yields 2 (potentially different) scores. The maximum of these 2 scores is taken

3) Learning Algorithm

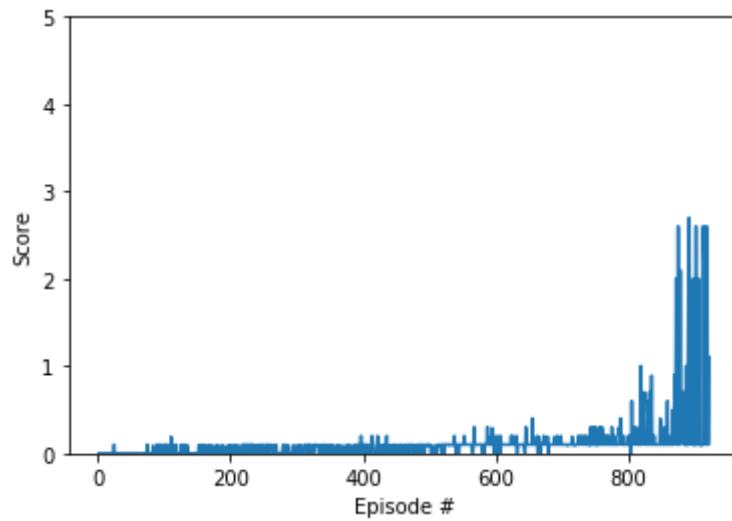
- Independent two DDPG agents

4) Parameters

- Actor: 2 Hidden fully-connected layers(64, 64 units)
- Critic: 2 Hidden fully-connected layers(64, 64 units)
- Actor learning rate: $1e-3$
- Critic learning rate: $1e-3$
- Gamma: 0.99
- Soft update(τ): $6e-2$
- Memory size: $1e6$
- Batch size: 128
- Optimizer: Adam

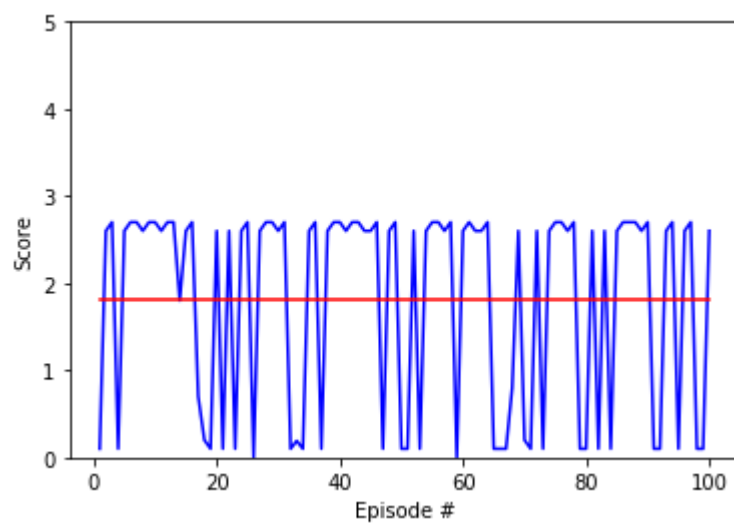
5) Training

- Training episodes: 923



6) Results

- Average score: 1.81



7) Future work

- To improve learning method, Multi-Agent Deterministic Policy Gradients(MADDPG) will be implemented