

Washington University in St. Louis

Washington University Open Scholarship

Electrical and Systems Engineering
Undergraduate Research

Electrical & Systems Engineering

12-6-2023

Comparative Study of Reinforcement Learning Algorithms: Deep Q-Networks, Deep Deterministic Policy Gradients and Proximal Policy Optimization

Haoyi Wang

Washington University in St. Louis, haoyi@wustl.edu

Follow this and additional works at: https://openscholarship.wustl.edu/eseundergraduate_research

Recommended Citation

Wang, Haoyi, "Comparative Study of Reinforcement Learning Algorithms: Deep Q-Networks, Deep Deterministic Policy Gradients and Proximal Policy Optimization" (2023). *Electrical and Systems Engineering Undergraduate Research*. 16.

https://openscholarship.wustl.edu/eseundergraduate_research/16

This Article is brought to you for free and open access by the Electrical & Systems Engineering at Washington University Open Scholarship. It has been accepted for inclusion in Electrical and Systems Engineering Undergraduate Research by an authorized administrator of Washington University Open Scholarship. For more information, please contact digital@wumail.wustl.edu.

Comparative Study of Reinforcement Learning Algorithms: Deep Q-Networks, Deep Deterministic Policy Gradients and Proximal Policy Optimization

Haoyi Wang
haoyi@wustl.edu
Department of Electrical & System Engineering, Washington University in St. Louis, St. Louis, MO, USA
Research Supervisor: Bruno Sinopoli

Introduction

Background

Artificial Intelligence (AI) has dramatically advanced through Reinforcement Learning (RL), particularly with the development of algorithms like Deep Q-Networks (DQN), Deep Deterministic Policy Gradients (DDPG), and Proximal Policy Optimization (PPO). These algorithms have been pivotal in diverse RL challenges.

Purpose of the Project

This study compares DQN, DDPG, and PPO, aiming to clarify their mechanics, efficiencies, and practical applications. We assess their theoretical foundations, performance in standard benchmarks, and adaptability to various environments. The goal is to offer insights for choosing the right algorithms for different environments.

Theoretical Background

Deep Q-Networks (DQN)

- Given state information, output the best policy based on the Q-value function, which estimates the expected cumulative rewards for each action in each state.
- Integrates deep neural networks with Q-learning.
- Experience Replay method is introduced.

```

Algorithm 1 Deep Q-learning with Experience Replay
Initialize replay memory D to capacity N
Initialize action-value function Q with random weights
for episode = 1, M do
  Initialize sequence s1 = {x1} and preprocessed sequenced φ1 = φ(s1)
  for t = 1, T do
    With probability ε select a random action at
    otherwise select at = max_a Q*(φ(st), a; θ)
    Execute action at in emulator and observe reward rt and image xt+1
    Set st+1 = st, at, xt+1 and preprocess φt+1 = φ(st+1)
    Store transition (φt, at, rt, φt+1) in D
    Sample random minibatch of transitions (φj, aj, rj, φj+1) from D
    Set yj = { rj for terminal φj+1
              rj + γ max_a' Q(φj+1, a'; θ) for non-terminal φj+1
    Perform a gradient descent step on (yj - Q(φj, aj; θ))^2 according to equation 3
  end for
end for
    
```

Deep Deterministic Policy Gradients (DDPG)

- Given state and action information, output the best policy based on the estimation of Action value and Q-value.
- Actor-critic approach: Actor learns the policy; Critic learns the Q-value function.

```

Algorithm 1 DDPG algorithm
Randomly initialize critic network Q(s, a|θ^Q) and actor μ(s|θ^μ) with weights θ^Q and θ^μ.
Initialize target network Q' and μ' with weights θ'^Q ← θ^Q, θ'^μ ← θ^μ
Initialize replay buffer R
for episode = 1, M do
  Initialize a random process N for action exploration
  Receive initial observation state s1
  for t = 1, T do
    Select action at = μ(st|θ^μ) + Nt according to the current policy and exploration noise
    Execute action at and observe reward rt and observe new state st+1
    Store transition (st, at, rt, st+1) in R
    Sample a random minibatch of N transitions (si, ai, ri, si+1) from R
    Set yi = ri + γ Q'(si+1, μ'(si+1|θ'^μ))|θ'^Q
    Update critic by minimizing the loss: L = 1/N ∑ (yi - Q(si, ai|θ^Q))^2
    Update the actor policy using the sampled policy gradient:
      ∇_θ^μ J ≈ 1/N ∑ ∇_a Q(s, a|θ^Q)|_{s=si, a=μ(si)} ∇_θ^μ μ(s|θ^μ)|_{si}
    Update the target networks:
      θ'^Q ← τθ^Q + (1 - τ)θ'^Q
      θ'^μ ← τθ^μ + (1 - τ)θ'^μ
  end for
end for
    
```

Proximal Policy Optimization (PPO)

- Given current state, output action probabilities or specific actions and a value estimate for those actions.
- Use the Clipping mechanism to prevent large policy updates.

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

```

Algorithm 1 PPO, Actor-Critic Style
for iteration=1, 2, ... do
  for actor=1, 2, ... N do
    Run policy π_θ^old in environment for T timesteps
    Compute advantage estimates A1, ..., AT
  end for
  Optimize surrogate L wrt θ, with K epochs and minibatch size M ≤ NT
  θ^old ← θ
end for
    
```

Analysis & Results

Experiment Setup

The study utilizes classic control environments from OpenAI Gym, specifically the Mountain Car (Discrete Action Space), Mountain Car Continuous (Continuous Action Space), and Cart Pole (Discrete Action Space).

The experiment is implemented by code projects running locally and is divided into two parts:

- Mountain Car for Q-learning and Mountain Car Continuous for DDPG. Compare the performance by measuring the time it takes for the car to reach the top of the hill in the simulation.
- Cart Pole for DQN and PPO. Compare the performance of the two algorithms' rewards and the holding time for the inverted pendulum in the simulation.

Results

Both Q-learning and DDPG successfully solved the related problems. Since the environments used are different, we cannot directly compare the rewards of the algorithms.

Q-learning trained for 48 seconds for 5000 episodes, while DDPG trained for 7 minutes and 29 seconds for 200 episodes.

In terms of performance in simulations of the trained cars, DDPG outperforms Q-learning, using only 1 second compared to the 5 seconds required by Q-learning.

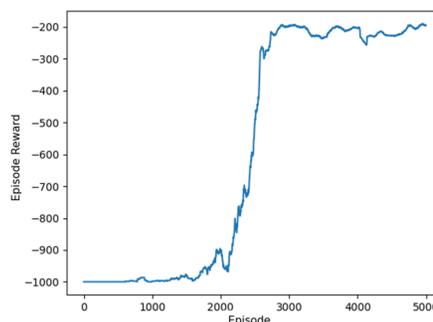


Figure 1. Training Reward for Mountain Car Q-learning

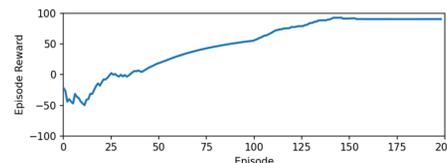


Figure 2. Training Reward for Mountain Car Continuous DDPG

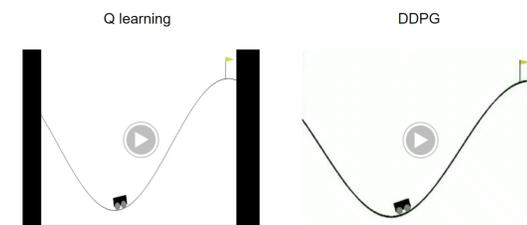


Figure 3: Simulation videos of Q-learning and DDPG

By comparing the reward graphs, we can conclude that PPO outperforms DQN.

DQN trained for 3 hours and 5 minutes, while PPO trained for 7 minutes and 34 seconds.

In terms of performance in simulations of the trained pendulums, PPO outperforms DQN by stabilizing the inverted pendulum for a longer time, e.g., at the 120th episode, PPO holds for 54 seconds while DQN holds for only 14 seconds.

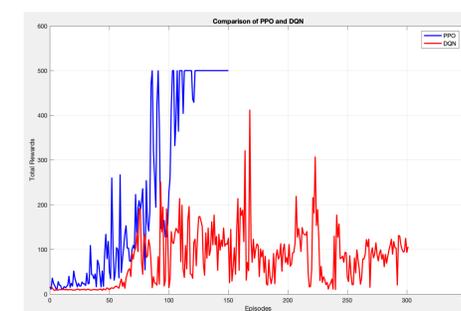


Figure 4. Training Reward for CartPole DQN and PPO

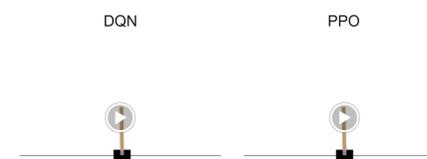


Figure 5: Simulation videos of DQN and PPO

Conclusion

Our analysis has demonstrated that while DQN excels in discrete action spaces, DDPG is more suited for continuous control tasks, and PPO shows consistent performance offering versatility and ease of tuning. This research not only aids in informed algorithm selection but also sets a foundation for future exploration in this area of study.

Related Publications

- [1] V. Mnih et al. "Playing atari with deep reinforcement learning". In: arXiv preprint arXiv:1312.5602 (2013).
- [2] T. P. Lillicrap et al. "Continuous control with deep reinforcement learning". In: arXiv preprint arXiv:1509.02971 (2015).
- [3] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. "Proximal policy optimization algorithms". In: arXiv preprint arXiv:1707.06347 (2017).

Acknowledgment

Carmel Fisco, Haoyu Yin and Chengyu Li