

A Comprehensive Review of Deep Reinforcement Learning

Abstract

Deep Reinforcement Learning (DRL) has emerged as a powerful approach for solving complex sequential decision-making problems. This review paper provides a detailed analysis and synthesis of the major advancements in DRL, focusing on key algorithms, methodologies, and applications. The reviewed works span topics including sample efficiency, stability, generalization, and real-world applicability of DRL methods. We critically analyze the contributions of each paper, highlight their strengths and limitations, and identify future research directions. This comprehensive review aims to serve as a valuable resource for researchers and practitioners in the field of DRL.

1 Introduction

Deep Reinforcement Learning (DRL) combines reinforcement learning with deep learning to enable agents to learn policies directly from high-dimensional inputs, such as images. This combination has led to significant breakthroughs in various domains, from playing Atari games to controlling robotic systems. Despite these successes, DRL faces several challenges, such as sample inefficiency, stability issues, and difficulty in generalizing to new tasks. This paper reviews recent advancements in DRL, providing a critical analysis of key contributions and identifying promising research directions.

2 Main Body

2.1 Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor

Haarnoja et al. (?) propose Soft Actor-Critic (SAC), an off-policy actor-critic algorithm based on the maximum entropy framework. SAC aims to improve both sample efficiency and stability by maximizing expected reward while encouraging exploration through entropy maximization. The algorithm achieves state-of-the-art performance on various continuous control tasks and demonstrates robustness across different random seeds. Despite its advantages, SAC’s reliance on the entropy term introduces additional hyperparameters, which may complicate tuning in some environments.

2.2 Asynchronous Methods for Deep Reinforcement Learning

Mnih et al. (?) introduce asynchronous methods for DRL, particularly Asynchronous Advantage Actor-Critic (A3C). By employing parallel actor-learners, A3C stabilizes training and reduces the reliance on GPU hardware, enabling efficient training on multi-core CPUs. A3C achieves impressive results on the Atari domain and continuous control tasks, significantly reducing training time compared to previous methods. However, the asynchronous nature of the algorithm may lead to non-deterministic behavior, complicating reproducibility.

2.3 Human-level Control Through Deep Reinforcement Learning

Mnih et al. (?) demonstrate the power of DRL with their Deep Q-Network (DQN) algorithm, which learns control policies directly from raw pixel inputs. DQN outperforms previous methods on multiple Atari games, achieving human-level performance on several tasks. This work highlights the potential of DRL for learning complex behaviors from high-dimensional data. Nonetheless, DQN suffers from overestimation bias and instability, which later works such as Double DQN (?) address.

2.4 Continuous Control with Deep Reinforcement Learning

Lillicrap et al. (2017) extend the ideas of DQN to continuous action spaces with the Deep Deterministic Policy Gradient (DDPG) algorithm. DDPG combines the actor-critic approach with deterministic policy gradients, enabling efficient learning in high-dimensional continuous environments. The algorithm demonstrates robust performance on various simulated physics tasks. However, DDPG is sensitive to hyperparameter settings and requires careful tuning to achieve optimal performance.

2.5 Deep Reinforcement Learning with Double Q-Learning

Van Hasselt et al. (2016) address the overestimation bias in Q-learning with the Double Q-learning algorithm. By decoupling the action selection and evaluation steps, Double Q-learning provides more accurate value estimates, leading to improved performance on the Atari domain. This work highlights the importance of addressing overestimation bias in DRL algorithms to achieve stable and reliable learning.

2.6 Dueling Network Architectures for Deep Reinforcement Learning

Wang et al. (2016) propose a dueling network architecture for DRL, which separately estimates the state value function and the state-dependent action advantage function. This architecture improves policy evaluation by generalizing learning across actions, leading to better performance on tasks with many similar-valued actions. The dueling architecture outperforms standard architectures on the Atari domain, demonstrating its effectiveness in enhancing DRL performance.

2.7 Deep Reinforcement Learning from Human Preferences

Christiano et al. (2017) explore the integration of human preferences into DRL, enabling agents to learn complex tasks without explicit reward functions. By leveraging human feedback, the proposed method effectively solves challenging tasks with minimal human oversight. This approach reduces the cost

of human involvement and demonstrates the potential of combining human intuition with DRL. However, scaling this method to more complex tasks may require additional advancements in human-computer interaction.

2.8 Rainbow: Combining Improvements in Deep Reinforcement Learning

Hessel et al. (2018) combine several enhancements to the DQN algorithm, including Double Q-learning, prioritized experience replay, and dueling network architectures, into a single algorithm called Rainbow. Rainbow achieves state-of-the-art performance on the Atari benchmark, demonstrating the complementary nature of these improvements. This work underscores the importance of integrating multiple advancements to achieve robust and efficient DRL algorithms.

3 Conclusion

Deep Reinforcement Learning has made significant strides in recent years, addressing various challenges and achieving impressive results across numerous domains. This review paper has highlighted key contributions, including Soft Actor-Critic, Asynchronous Methods, Deep Q-Networks, and advancements in continuous control. Despite these successes, several challenges remain, such as sample inefficiency, stability, and generalization. Future research should focus on addressing these challenges, exploring new algorithms, and extending DRL applications to more complex and real-world scenarios.

References