

Introduction to Reinforcement Learning for Physicists

Frank Cichos

2024-06-17

Introduction

Reinforcement Learning (RL) is a subfield of machine learning where an agent learns to make decisions by interacting with an environment. The agent receives feedback in the form of rewards and aims to maximize the cumulative reward over time. This document provides an overview of RL concepts and their applications in physics.

Basic Concepts

Agent and Environment

In RL, the agent interacts with the environment by taking actions. The environment responds by providing new states and rewards. The goal of the agent is to learn a policy that maximizes the cumulative reward.

- **State** (s): A representation of the current situation of the environment.
- **Action** (a): A decision made by the agent that affects the state.
- **Reward** (r): Feedback from the environment indicating the success of an action.

Policy and Value Functions

- **Policy** (π): A strategy used by the agent to decide actions based on the current state. It can be deterministic or stochastic.
- **Value Function**: Measures the expected cumulative reward from a given state or state-action pair.
 - **State-Value Function** ($V(s)$): Expected reward starting from state s .

- **Action-Value Function** ($Q(s, a)$): Expected reward starting from state s and taking action a .

Markov Decision Processes (MDPs)

Definition

A Markov Decision Process (MDP) is a mathematical framework for modeling decision-making problems where outcomes are partly random and partly under the control of a decision-maker. An MDP is defined by a tuple $\langle S, A, P, R, \gamma \rangle$ where:

- S : A finite set of states.
- A : A finite set of actions.
- P : State transition probability matrix, $P(s'|s, a)$, which represents the probability of transitioning to state s' from state s under action a .
- R : Reward function, $R(s, a)$, which gives the expected reward received after transitioning from state s to state s' under action a .
- γ : Discount factor, $0 \leq \gamma \leq 1$, which represents the importance of future rewards.

Markov Property

The Markov property states that the future state depends only on the current state and action, not on the sequence of events that preceded it. Formally, this is expressed as:

$$P(s_{t+1}|s_t, a_t) = P(s_{t+1}|s_1, a_1, s_2, a_2, \dots, s_t, a_t)$$

Derivation of the Markov Property

To derive the Markov property, consider a stochastic process (X_t) with the Markov property. For any times $t_1 < t_2 < \dots < t_n$ and any states x_1, x_2, \dots, x_n , the Markov property can be written as:

$$P(X_{t_{n+1}} = x_{n+1} | X_{t_n} = x_n, X_{t_{n-1}} = x_{n-1}, \dots, X_{t_1} = x_1) = P(X_{t_{n+1}} = x_{n+1} | X_{t_n} = x_n)$$

This means that the conditional probability of the future state $X_{t_{n+1}}$ depends only on the present state X_{t_n} and not on the past states $X_{t_{n-1}}, \dots, X_{t_1}$.

Bellman Equation

Bellman Expectation Equation

The Bellman Expectation Equation provides a recursive decomposition of the value function. For a given policy π , the state-value function $V^\pi(s)$ is defined as:

$$V^\pi(s) = \mathbb{E}_\pi [R_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s]$$

This equation states that the value of a state under a policy π is the expected reward for taking an action a in state s , plus the discounted value of the next state.

Derivation of the Bellman Expectation Equation

To derive the Bellman Expectation Equation, consider the definition of the state-value function $V^\pi(s)$:

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid s_t = s \right]$$

We can decompose the sum into the immediate reward and the future rewards:

$$V^\pi(s) = \mathbb{E}_\pi \left[R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \mid s_t = s \right]$$

Notice that the term inside the sum is the value function at the next state s_{t+1} :

$$V^\pi(s) = \mathbb{E}_\pi [R_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s]$$

This is the Bellman Expectation Equation.

Bellman Optimality Equation

The Bellman Optimality Equation defines the value of a state under the optimal policy π^* . The state-value function $V^*(s)$ is given by:

$$V^*(s) = \max_a \mathbb{E} [R_{t+1} + \gamma V^*(s_{t+1}) \mid s_t = s, a_t = a]$$

Similarly, the action-value function $Q^*(s, a)$ is defined as:

$$Q^*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') \mid s_t = s, a_t = a \right]$$

Derivation of the Bellman Optimality Equation

To derive the Bellman Optimality Equation, consider the definition of the optimal state-value function $V^*(s)$:

$$V^*(s) = \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid s_t = s \right]$$

We can decompose the sum into the immediate reward and the future rewards:

$$V^*(s) = \max_{\pi} \mathbb{E}_{\pi} \left[R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \mid s_t = s \right]$$

Notice that the term inside the sum is the value function at the next state s_{t+1} :

$$V^*(s) = \max_a \mathbb{E} [R_{t+1} + \gamma V^*(s_{t+1}) \mid s_t = s, a_t = a]$$

This is the Bellman Optimality Equation.

Algorithms

Policy Gradient Methods

Policy gradient methods directly optimize the policy by adjusting the parameters of the policy network to maximize the expected reward. These methods are effective in high-dimensional action spaces.

Details of Policy Gradient Methods

1. **Objective Function:** The goal is to maximize the expected return $J(\theta)$ where θ are the policy parameters.
2. **Gradient Ascent:** The policy parameters are updated using gradient ascent:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

3. **Policy Gradient Theorem:** The gradient of the expected return can be expressed as:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi}(s, a)]$$

Definition of $J(\theta)$

Formally, $J(\theta)$ is defined as the expected return when following the policy π_{θ} :

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)]$$

where: τ denotes a trajectory, which is a sequence of states, actions, and rewards: $\tau = (s_0, a_0, r_1, s_1, a_1, r_2, \dots)$. $R(\tau)$ is the cumulative reward for the trajectory τ , often defined as the sum of rewards along the trajectory: $R(\tau) = \sum_{t=0}^T \gamma^t r_t$, where γ is the discount factor. π_{θ} is the policy parameterized by θ , which defines the probability distribution over actions given states.

Objective of Policy Gradient Methods

The primary objective of policy gradient methods is to find the optimal policy parameters θ^* that maximize the expected cumulative reward $J(\theta)$:

$$\theta^* = \arg \max_{\theta} J(\theta)$$

Gradient of $J(\theta)$

To optimize $J(\theta)$, policy gradient methods use the gradient of the expected reward with respect to the policy parameters, $\nabla_{\theta} J(\theta)$. This gradient provides the direction in which the policy parameters should be adjusted to increase the expected reward. The policy gradient theorem provides a practical form of this gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) R(\tau) \right]$$

This expression can be estimated using samples of trajectories generated by following the policy π_{θ} .

```
# Example of a simple policy gradient algorithm in Python
import numpy as np

def policy_gradient():
    # Initialize policy parameters
    theta = np.random.rand()
    for episode in range(1000):
        # Generate an episode
        states, actions, rewards = generate_episode(theta)
        # Update policy parameters
        theta += alpha * np.sum(rewards) * np.sum(actions)
    return theta
```

Value-Based Methods

Value-based methods, such as Q-learning, focus on estimating the value functions. The agent selects actions that maximize the estimated value.

Details of Value-Based Methods

1. **Q-Learning Update Rule:** The Q-value is updated using the following rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

2. **Exploration-Exploitation Trade-off:** An ϵ -greedy policy is often used to balance exploration and exploitation.
3. **Convergence:** Q-learning is guaranteed to converge to the optimal policy under certain conditions, such as sufficient exploration and a decaying learning rate.

```
# Example of Q-learning algorithm in Python
import numpy as np

def q_learning():
    # Initialize Q-table
    Q = np.zeros((state_space, action_space))
    for episode in range(1000):
        state = env.reset()
        done = False
        while not done:
            action = select_action(Q, state)
            next_state, reward, done = env.step(action)
```

```

        Q[state, action] += alpha * (reward + gamma * np.max(Q[next_state])) - Q[state, action]
        state = next_state
    return Q

```

Actor-Critic Methods

Actor-critic methods combine policy gradient and value-based methods. The actor updates the policy, while the critic evaluates the action by estimating the value function.

Details of Actor-Critic Methods

1. **Actor:** The policy function, which selects actions based on the current state.
2. **Critic:** The value function, which evaluates the actions taken by the actor.
3. **TD Error:** The temporal difference (TD) error is used to update both the actor and the critic:

$$\delta = r + \gamma V(s') - V(s)$$

4. **Update Rules:**

- Actor: $\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} \log \pi_{\theta}(a|s)$
- Critic: $w \leftarrow w + \beta \delta \nabla_w V_w(s)$

Tip

In reinforcement learning, the temporal difference (TD) error, denoted as δ , is a crucial concept used in various algorithms, including TD learning, Q-learning, and actor-critic methods. The TD error measures the difference between the predicted value of a state and the actual observed value after taking an action. This document explains the TD error and its significance in reinforcement learning.

Definition of TD Error

The TD error δ is defined as:

$$\delta = r + \gamma V(s') - V(s)$$

where: - r is the reward received after transitioning from state s to state s' . - γ is the discount factor, which represents the importance of future rewards. - $V(s)$ is the value of the current state s . - $V(s')$ is the value of the next state s' .

Explanation of TD Error

Components of TD Error

1. **Immediate Reward (r):** This is the reward received immediately after taking an action in state s and transitioning to state s' .
2. **Discounted Future Value ($\gamma V(s')$):** This term represents the discounted value of the next state s' . The discount factor γ (where $0 \leq \gamma \leq 1$) determines how much future rewards are valued compared to immediate rewards.
3. **Current State Value ($V(s)$):** This is the predicted value of the current state s before taking the action.

Interpretation of TD Error

The TD error δ represents the difference between the expected value of the current state and the observed value after taking an action. Specifically: - If $\delta > 0$, the observed value is higher than expected, indicating that the action taken was better than predicted. - If $\delta < 0$, the observed value is lower than expected, indicating that the action taken was worse than predicted. - If $\delta = 0$, the observed value matches the expected value, indicating that the prediction was accurate.

Significance of TD Error

The TD error is used to update the value function and policy in various reinforcement learning algorithms. It provides a signal for learning and improving the agent's predictions and actions.

TD Learning

In TD learning, the value function is updated using the TD error:

$$V(s) \leftarrow V(s) + \alpha \delta$$

where α is the learning rate. This update rule adjusts the value of the current state s based on the TD error, gradually improving the accuracy of the value function.

Q-Learning

In Q-learning, the Q-value is updated using the TD error:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

Here, the TD error is used to adjust the Q-value of the state-action pair (s, a) , guiding the agent towards actions that maximize future rewards.

Actor-Critic Methods

In actor-critic methods, the TD error is used to update both the policy (actor) and the value function (critic):

- **Critic Update:** The value function is updated using the TD error:

$$V(s) \leftarrow V(s) + \alpha \delta$$

- **Actor Update:** The policy parameters are updated using the TD error:

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} \log \pi_{\theta}(a|s)$$

The TD error provides a feedback signal that helps the actor improve the policy by favoring actions that lead to higher rewards.

Example: TD Error in Action

Consider an agent navigating a grid world. The agent receives a reward of $r = 1$ for reaching a goal state. The discount factor is $\gamma = 0.9$. Suppose the current state value is $V(s) = 0.5$ and the next state value is $V(s') = 0.8$.

The TD error is calculated as:

$$\delta = r + \gamma V(s') - V(s) = 1 + 0.9 \cdot 0.8 - 0.5 = 1 + 0.72 - 0.5 = 1.22$$

Since $\delta > 0$, the observed value is higher than expected, indicating that the action taken was better than predicted. The value function and policy will be updated accordingly to reflect this positive outcome.

```
# Example of an actor-critic algorithm in Python
import numpy as np

def actor_critic():
    # Initialize policy and value parameters
    theta = np.random.rand()
    w = np.random.rand()
    for episode in range(1000):
        state = env.reset()
        done = False
        while not done:
            action = select_action(theta, state)
            next_state, reward, done = env.step(action)
            delta = reward + gamma * value(w, next_state) - value(w, state)
```

```
theta += alpha * delta * grad_log_policy(theta, state, action)
w += beta * delta * grad_value(w, state)
state = next_state
return theta, w
```

On-Policy vs. Off-Policy Methods

On-Policy Methods

On-policy methods evaluate and improve the same policy that is used to make decisions. Examples include SARSA and Policy Gradient methods.

Details of On-Policy Methods

1. **SARSA**: An on-policy TD control algorithm that updates the Q-value based on the action actually taken by the current policy.

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$$

2. **Policy Gradient**: Directly optimizes the policy by following the gradient of the expected return.

Off-Policy Methods

Off-policy methods evaluate and improve a policy different from the one used to generate the data. Examples include Q-learning and Deep Q-Networks (DQN).

Details of Off-Policy Methods

1. **Q-Learning**: An off-policy TD control algorithm that updates the Q-value based on the maximum reward of the next state.

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

2. **Experience Replay**: Stores past experiences and samples them randomly to break the correlation between consecutive experiences, improving learning stability.

💡 Tip

In reinforcement learning (RL), on-policy and off-policy methods represent two fundamental approaches to how an agent learns from its interactions with the environment. These approaches significantly influence the agent's learning strategy, convergence properties, and efficiency. This document provides a detailed explanation of on-policy and off-policy methods, focusing on SARSA and Q-learning.

On-Policy Methods: SARSA

Definition

On-policy methods evaluate and improve the policy that the agent is currently following. The agent learns the value of the policy it is executing, including the exploration steps. SARSA (State-Action-Reward-State-Action) is a quintessential example of an on-policy method.

SARSA Algorithm

The SARSA algorithm updates its Q-values based on the action that the agent actually takes in the next state, adhering to the current policy. The update rule for SARSA is given by:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

where:

- $Q(s_t, a_t)$ is the action-value function, representing the expected return for taking action a in state s .
- α is the learning rate.
- r_{t+1} is the reward received after taking action a_t in state s_t .
- γ is the discount factor.
- s_{t+1} and a_{t+1} are the next state and action, respectively, following the current policy.

Example of SARSA

Consider an agent navigating a gridworld where each cell represents a state, and the agent can move in four possible directions: up, down, left, and right. The agent's goal is to reach a target cell with the maximum cumulative reward.

1. The agent starts at an initial state s and selects an action a based on its current policy (e.g., (\cdot) -greedy).

2. The agent executes a , receives a reward r , and transitions to the next state s' .
3. The agent then selects the next action a' based on its current policy.
4. The Q-value for the state-action pair (s, a) is updated using the reward and the Q-value of the next state-action pair (s', a') .

This process repeats, with the agent continuously updating its Q-values based on the actions taken according to its current policy.

Off-Policy Methods: Q-Learning

Definition

Off-policy methods involve learning the value of an optimal policy independently of the agent's actions. The agent can learn the optimal policy while following a different behavior policy for exploration. Q-learning is a prominent example of an off-policy method.

Q-Learning Algorithm

The Q-learning algorithm updates its Q-values based on the maximum reward of the next state, regardless of the action actually taken. The update rule for Q-learning is given by:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]$$

where:

- $Q(s_t, a_t)$ is the action-value function.
- α is the learning rate.
- r_{t+1} is the reward received after taking action a_t in state s_t .
- γ is the discount factor.
- s_{t+1} is the next state.
- $\max_{a'} Q(s_{t+1}, a')$ is the maximum Q-value over all possible actions in the next state.

Example of Q-Learning

Consider the same gridworld example. The Q-learning algorithm proceeds as follows:

1. The agent starts at an initial state s and selects an action a based on its behavior policy (e.g., (\cdot) -greedy).
2. The agent executes a , receives a reward r , and transitions to the next state s' .
3. The Q-value for the state-action pair (s, a) is updated using the reward and the maximum Q-value of the next state s' .

This process repeats, with the agent continuously updating its Q-values based on the maximum possible reward of the next state.

Key Differences Between SARSA and Q-Learning

Policy Evaluation

- **SARSA (On-Policy):** Evaluates and improves the policy that the agent is currently following. The Q-value update is based on the action taken by the agent following its current policy.
- **Q-Learning (Off-Policy):** Learns the value of the optimal policy independently of the agent's behavior. The Q-value update is based on the maximum reward of the next state, regardless of the action taken.

Exploration and Exploitation

- **SARSA:** Balances exploration and exploitation by updating Q-values based on the actions actually taken, which may include exploratory actions.
- **Q-Learning:** Focuses on exploitation by updating Q-values based on the maximum possible reward, which may lead to more aggressive learning.

Convergence Properties

- **SARSA:** More conservative and potentially more stable in stochastic environments. It converges to a policy that balances exploration and exploitation.
- **Q-Learning:** More aggressive and can converge faster in deterministic environments. It aims to learn the optimal policy directly.

Practical Considerations

When to Use SARSA

- **Safety:** SARSA is more conservative and safer, making it suitable for environments where mistakes are costly (e.g., training physical robots).
- **Stochastic Environments:** SARSA performs better in environments with high variability and uncertainty.

When to Use Q-Learning

- **Optimal Policy:** Q-learning is more efficient for learning the optimal policy, making it suitable for simulation environments where exploration is less costly.
- **Deterministic Environments:** Q-learning performs better in deterministic environments where the optimal policy can be learned quickly.

Conclusion

The distinction between on-policy and off-policy methods in reinforcement learning is fundamental and has significant implications for the agent's learning strategy. On-policy methods like SARSA focus on evaluating and improving the current policy, making them more conservative and stable. Off-policy methods like Q-learning aim to learn the optimal policy independently of the agent's behavior, allowing for more aggressive learning and faster convergence. Understanding these differences allows practitioners to select the appropriate algorithm based on the specific requirements and characteristics of the environment.

Applications in Physics

Quantum Physics

RL has been applied to optimize quantum control problems, such as designing better quantum computers and simulators. Classical neural networks can interact with quantum devices to improve their performance.

Trajectory Optimization

In classical mechanics, RL can be used to optimize trajectories in complex systems. For example, controlling the motion of a robotic arm or optimizing the path of a particle in a potential field.

Material Science

RL can assist in discovering new materials by optimizing the synthesis process and predicting material properties.

Advanced Topics

Physics-Informed Reinforcement Learning

Physics-informed reinforcement learning (PIRL) incorporates physical laws and constraints into the RL framework. This approach enhances the learning process by ensuring that the agent's actions adhere to known physical principles, improving both efficiency and accuracy[1][2][3].

Model-Based Reinforcement Learning

Model-based RL involves learning a model of the environment's dynamics and using it to plan actions. This approach can significantly improve sample efficiency by reducing the need for extensive interaction with the environment[2][4].

Deep Reinforcement Learning

Deep RL combines deep learning with RL, allowing the agent to handle high-dimensional state and action spaces. Techniques such as Deep Q-Networks (DQN) and Proximal Policy Optimization (PPO) have shown remarkable success in various applications, including game playing and robotics[2][5][6].

Case Studies

Quantum Circuit Design

RL has been used to optimize quantum circuit designs, addressing challenges such as gate selection and error minimization. This application leverages the RL framework to explore and exploit the vast design space of quantum circuits[7].

Optical Systems

Deep RL has been applied to control and optimize optical systems, such as mode-locked lasers. The RL agent learns to adjust system parameters to achieve optimal performance, demonstrating the potential of RL in complex physical systems[8].

Summary

Reinforcement learning offers powerful tools for solving complex decision-making problems in physics. By understanding the basic concepts and algorithms, physicists can leverage RL to advance their research and explore new frontiers.

Some References

1. Sutton, R. S., & Barto, A. G. (2018). Reinforcement Learning: An Introduction (2nd ed.). MIT Press.
2. Marquardt, F. (2020). Machine Learning for Physicists. Friedrich-Alexander Universität Erlangen-Nürnberg.
3. Martín-Guerrero, J. D., & Lamata, L. (2021). Reinforcement Learning and Physics. Applied Sciences, 11(18), 8589.
4. Banerjee, C., et al. (2023). A Survey on Physics Informed Reinforcement Learning. arXiv:2309.01909.
5. Ramesh, A., & Ravindran, B. (2023). Physics-Informed Model-Based Reinforcement Learning. Proceedings of the 5th Annual Conference on Learning for Dynamics and Control.
6. Sun, C., et al. (2020). Deep reinforcement learning for optical systems: A case study of mode-locked lasers. arXiv:2006.05579.
7. Wiering, M., & van Otterlo, M. (2012). Reinforcement Learning: State-of-the-Art. Springer.
8. Wallscheid, O. (2020). Lecture 02: Markov Decision Processes. University of Paderborn.