

# 5. Reinforcement Learning: Other Solution Methods

Based on Sutton and Barto's Reinforcement Learning Book

## Contents

<b>1</b>	<b>Introduction to Other Solution Methods</b>	<b>2</b>
<b>2</b>	<b>Limitations of On-Policy Methods</b>	<b>2</b>
2.1	Definition and Characteristics . . . . .	2
2.2	Limitations . . . . .	2
2.2.1	Exploration-Exploitation Tradeoff . . . . .	2
2.2.2	Sample Inefficiency . . . . .	2
2.2.3	Necessity for Incremental Policy Improvement . . . . .	3
2.3	On-Policy vs. Off-Policy Learning . . . . .	3
<b>3</b>	<b>Monte Carlo Techniques</b>	<b>3</b>
3.1	Properties of Monte Carlo Methods . . . . .	3
3.2	Monte Carlo Prediction (Evaluation) . . . . .	3
3.3	Monte Carlo Control . . . . .	4
3.4	Off-Policy Monte Carlo Methods . . . . .	5
<b>4</b>	<b>Temporal Difference Learning</b>	<b>6</b>
4.1	TD Prediction . . . . .	6
4.2	Advantages of TD Learning . . . . .	7
4.3	SARSA: On-Policy TD Control . . . . .	7
<b>5</b>	<b>Q-Learning</b>	<b>8</b>
5.1	Q-Learning Algorithm . . . . .	8
5.2	Differences Between SARSA and Q-Learning . . . . .	9
5.3	Expected SARSA . . . . .	9
<b>6</b>	<b>Policy Gradient Methods</b>	<b>10</b>
6.1	Policy Parameterization . . . . .	10
6.2	Performance Measure . . . . .	10
6.3	Policy Gradient Theorem . . . . .	10
6.4	REINFORCE Algorithm . . . . .	11
6.5	Advantages of Policy Gradient Methods . . . . .	11
<b>7</b>	<b>Comparison of Methods</b>	<b>12</b>
<b>8</b>	<b>Conclusion</b>	<b>12</b>

# 1 Introduction to Other Solution Methods

In previous sections, we explored policy iteration and value iteration methods for solving reinforcement learning problems, which require a complete model of the environment. However, in many practical scenarios, we either don't have a complete model or it's computationally expensive to use one. In this section, we explore alternative approaches to solving reinforcement learning problems that either relax the requirement for a complete model or approach the problem from different angles.

These alternative methods include:

- Understanding the limitations of on-policy methods
- Monte Carlo techniques for estimating value functions
- Temporal Difference Learning, which combines ideas from Monte Carlo methods and dynamic programming
- Q-Learning as an off-policy temporal difference method
- Policy Gradient methods that directly optimize policy parameters

These alternative approaches are crucial for extending reinforcement learning to a wide variety of practical applications where complete models are unavailable or impractical.

## 2 Limitations of On-Policy Methods

On-policy methods evaluate and improve the same policy that is used to make decisions. While conceptually straightforward, this approach has several limitations:

### 2.1 Definition and Characteristics

On-policy methods learn the value functions and improve the policy based on actions taken according to the current policy. The key characteristic is that the policy being evaluated and the policy generating behavior are the same.

### 2.2 Limitations

#### 2.2.1 Exploration-Exploitation Tradeoff

One fundamental limitation is the exploration-exploitation tradeoff. To find the optimal policy, the agent needs to explore various states and actions. However, to maximize rewards, it should exploit its current knowledge by choosing the actions it believes are best. On-policy methods must balance these competing objectives using the same policy.

For example, in  $\epsilon$ -greedy policies, the agent follows the greedy policy with probability  $1 - \epsilon$  and takes a random action with probability  $\epsilon$ . While this introduces exploration, it also means the policy being learned is not truly optimal but is constrained by the exploration requirements.

#### 2.2.2 Sample Inefficiency

On-policy methods can be sample inefficient because they can only learn from data collected under the current policy. When the policy changes, previously collected experience becomes less relevant or even irrelevant for evaluating the new policy. This limitation is particularly significant in environments where data collection is expensive or time-consuming.

### 2.2.3 Necessity for Incremental Policy Improvement

On-policy methods typically require incremental policy improvements to ensure stability. Large policy changes can lead to drastic changes in the data distribution, making it difficult to accurately estimate value functions.

## 2.3 On-Policy vs. Off-Policy Learning

To address these limitations, we can use off-policy methods, which allow the agent to learn from data that was collected using a different policy (behavior policy) than the one being evaluated and improved (target policy). This separation offers more flexibility but introduces its own challenges.

## 3 Monte Carlo Techniques

Monte Carlo (MC) methods learn directly from complete episodes of experience without requiring a model of the environment's dynamics. They are based on averaging complete returns from episodes of experience.

### 3.1 Properties of Monte Carlo Methods

- Only applicable to episodic tasks (tasks with a definite ending point)
- Learn from complete episodes: no bootstrapping
- Can be off-policy, learning about a target policy while following a behavior policy
- Do not require knowledge about environment dynamics
- Simple conceptually and can handle environments with unknown dynamics

### 3.2 Monte Carlo Prediction (Evaluation)

Monte Carlo prediction estimates the state-value function  $v_\pi$  for a given policy  $\pi$  by averaging returns observed after visits to a state.

---

**Algorithm 1** First-visit Monte Carlo Prediction

---

```
1: Input: policy  $\pi$  to be evaluated
2: Initialize  $V(s)$  for all  $s \in \mathcal{S}$ 
3: Initialize  $Returns(s)$  as an empty list for all  $s \in \mathcal{S}$ 
4: for each episode do
5:   Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
6:    $G \leftarrow 0$ 
7:   for  $t = T - 1, T - 2, \dots, 0$  do
8:      $G \leftarrow \gamma G + R_{t+1}$ 
9:     if  $S_t$  appears for the first time in the episode then
10:      Append  $G$  to  $Returns(S_t)$ 
11:       $V(S_t) \leftarrow average(Returns(S_t))$ 
12:     end if
13:   end for
14: end for
```

---

### Notation Overview

- $\pi$  - policy to be evaluated
- $V(s)$  - estimated value function for state  $s$
- $Returns(s)$  - list of returns observed after visiting state  $s$
- $S_t$  - state at time  $t$
- $A_t$  - action at time  $t$
- $R_t$  - reward at time  $t$
- $G$  - return (cumulative discounted reward)
- $\gamma$  - discount factor
- $T$  - final time step of the episode

The first-visit MC method estimates  $v_\pi(s)$  as the average of returns following the first visits to state  $s$ , while the every-visit MC method averages returns following all visits to  $s$ .

### 3.3 Monte Carlo Control

Monte Carlo control methods find the optimal policy by alternating between policy evaluation and policy improvement, similar to policy iteration but using Monte Carlo methods for the evaluation step.

---

**Algorithm 2** Monte Carlo Control with Exploring Starts

---

```
1: Initialize  $Q(s, a)$  arbitrarily for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
2: Initialize  $\pi(s)$  arbitrarily for all  $s \in \mathcal{S}$ 
3: Initialize  $Returns(s, a)$  as empty lists for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
4: for each episode do
5:   Choose  $S_0 \in \mathcal{S}$  and  $A_0 \in \mathcal{A}(S_0)$  such that all pairs have probability  $> 0$ 
6:   Generate an episode from  $S_0, A_0$  following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
7:    $G \leftarrow 0$ 
8:   for  $t = T - 1, T - 2, \dots, 0$  do
9:      $G \leftarrow \gamma G + R_{t+1}$ 
10:    if  $(S_t, A_t)$  pair appears for the first time in the episode then
11:      Append  $G$  to  $Returns(S_t, A_t)$ 
12:       $Q(S_t, A_t) \leftarrow average(Returns(S_t, A_t))$ 
13:       $\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$ 
14:    end if
15:  end for
16: end for
```

---

### Notation Overview

- $Q(s, a)$  - action-value function for state  $s$  and action  $a$
- $\pi(s)$  - policy that maps states to actions
- $Returns(s, a)$  - list of returns observed after taking action  $a$  in state  $s$
- $S_t$  - state at time  $t$
- $A_t$  - action at time  $t$
- $R_t$  - reward at time  $t$
- $G$  - return (cumulative discounted reward)
- $\gamma$  - discount factor
- $T$  - final time step of the episode
- $\mathcal{S}$  - set of all states
- $\mathcal{A}(s)$  - set of all actions available in state  $s$

### 3.4 Off-Policy Monte Carlo Methods

Off-policy Monte Carlo methods allow learning about one policy (the target policy) while following another policy (the behavior policy). This approach addresses the exploration-exploitation tradeoff by using a separate exploratory policy for data collection.

---

**Algorithm 3** Off-policy Monte Carlo Prediction

---

```
1: Input: target policy  $\pi$  to be evaluated
2: Initialize  $V(s)$  for all  $s \in \mathcal{S}$ 
3: Initialize  $C(s)$  to 0 for all  $s \in \mathcal{S}$ 
4: for each episode do
5:   Generate an episode using behavior policy  $b$ :  $S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
6:    $G \leftarrow 0$ 
7:    $W \leftarrow 1$ 
8:   for  $t = T - 1, T - 2, \dots, 0$  do
9:      $G \leftarrow \gamma G + R_{t+1}$ 
10:     $C(S_t) \leftarrow C(S_t) + W$ 
11:     $V(S_t) \leftarrow V(S_t) + \frac{W}{C(S_t)}[G - V(S_t)]$ 
12:     $W \leftarrow W \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$ 
13:    if  $W = 0$  then
14:      break
15:    end if
16:  end for
17: end for
```

---

## Notation Overview

- $\pi$  - target policy to be evaluated
- $b$  - behavior policy used to generate episodes
- $V(s)$  - estimated value function for state  $s$
- $C(s)$  - cumulative weight for state  $s$
- $W$  - importance sampling weight
- $S_t$  - state at time  $t$
- $A_t$  - action at time  $t$
- $R_t$  - reward at time  $t$
- $G$  - return (cumulative discounted reward)
- $\gamma$  - discount factor
- $T$  - final time step of the episode
- $\pi(A_t|S_t)$  - probability of taking action  $A_t$  in state  $S_t$  under policy  $\pi$
- $b(A_t|S_t)$  - probability of taking action  $A_t$  in state  $S_t$  under policy  $b$

The term  $\frac{\pi(A_t|S_t)}{b(A_t|S_t)}$  is the importance sampling ratio, which corrects for the difference in action selection probabilities between the target and behavior policies.

## 4 Temporal Difference Learning

Temporal Difference (TD) learning combines ideas from Monte Carlo methods and dynamic programming. Like Monte Carlo methods, TD methods can learn directly from experience without a model of the environment. Like DP, TD methods update estimates based in part on other learned estimates, without waiting for a final outcome (they bootstrap).

### 4.1 TD Prediction

The simplest TD method, known as TD(0), updates the value estimate for a state based on the observed reward and the estimated value of the next state:

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \quad (1)$$

### Notation Overview

- $V(S_t)$  - estimated value of state  $S_t$
- $R_{t+1}$  - reward received after taking action in state  $S_t$
- $S_{t+1}$  - next state after taking action in state  $S_t$
- $\alpha$  - step-size parameter (learning rate)
- $\gamma$  - discount factor
- $R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  - TD error, often denoted as  $\delta_t$

The term  $R_{t+1} + \gamma V(S_{t+1})$  is called the TD target, and  $R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  is called the TD error.

---

**Algorithm 4** Tabular TD(0) for estimating  $v_\pi$ 

---

```
1: Input: policy  $\pi$  to be evaluated
2: Parameters: step size  $\alpha \in (0, 1]$ 
3: Initialize  $V(s)$  arbitrarily for all  $s \in \mathcal{S}$ 
4: for each episode do
5:   Initialize  $S$ 
6:   for each step of episode do
7:     Take action  $A$  according to  $\pi(\Delta|S)$ 
8:     Observe reward  $R$  and next state  $S'$ 
9:      $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ 
10:     $S \leftarrow S'$ 
11:   end for
12: end for
```

---

## 4.2 Advantages of TD Learning

TD learning has several advantages over both Monte Carlo methods and dynamic programming:

- TD methods don't require a model of the environment (unlike DP)
- TD methods can learn online, after each step, without waiting for the end of an episode (unlike MC)
- TD methods can learn from incomplete episodes, making them applicable to continuing (non-terminating) tasks
- TD methods are generally more sample efficient than Monte Carlo methods

## 4.3 SARSA: On-Policy TD Control

SARSA (State-Action-Reward-State-Action) is an on-policy TD control method that learns action-value functions. The update rule for SARSA is:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)] \quad (2)$$

#### Notation Overview

- $Q(S_t, A_t)$  - estimated value of taking action  $A_t$  in state  $S_t$
- $R_{t+1}$  - reward received after taking action  $A_t$  in state  $S_t$
- $S_{t+1}$  - next state after taking action  $A_t$  in state  $S_t$
- $A_{t+1}$  - next action taken in state  $S_{t+1}$
- $\alpha$  - step-size parameter (learning rate)
- $\gamma$  - discount factor

---

#### Algorithm 5 SARSA (on-policy TD control) for estimating $Q \approx q_*$

---

```

1: Parameters: step size  $\alpha \in (0, 1]$ , small  $\epsilon > 0$ 
2: Initialize  $Q(s, a)$  arbitrarily for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 
3: for each episode do
4:   Initialize  $S$ 
5:   Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
6:   for each step of episode do
7:     Take action  $A$ , observe  $R$ ,  $S'$ 
8:     Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
9:      $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
10:     $S \leftarrow S'$ 
11:     $A \leftarrow A'$ 
12:   end for
13: end for

```

---

## 5 Q-Learning

Q-Learning is an off-policy TD control method. It directly approximates the optimal action-value function,  $q_*$ , independent of the policy being followed.

### 5.1 Q-Learning Algorithm

The update rule for Q-Learning is:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)] \quad (3)$$

#### Notation Overview

- $Q(S_t, A_t)$  - estimated value of taking action  $A_t$  in state  $S_t$
- $R_{t+1}$  - reward received after taking action  $A_t$  in state  $S_t$
- $S_{t+1}$  - next state after taking action  $A_t$  in state  $S_t$
- $\max_a Q(S_{t+1}, a)$  - maximum estimated value obtainable from state  $S_{t+1}$
- $\alpha$  - step-size parameter (learning rate)
- $\gamma$  - discount factor



---

**Algorithm 6** Q-Learning (off-policy TD control) for estimating  $\pi \approx \pi_*$ 

---

```
1: Parameters: step size  $\alpha \in (0, 1]$ , small  $\epsilon > 0$ 
2: Initialize  $Q(s, a)$  arbitrarily for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ 
3: for each episode do
4:   Initialize  $S$ 
5:   for each step of episode do
6:     Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
7:     Take action  $A$ , observe  $R$ ,  $S'$ 
8:      $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
9:      $S \leftarrow S'$ 
10:   end for
11: end for
```

---

## 5.2 Differences Between SARSA and Q-Learning

The key difference between SARSA and Q-Learning lies in their update targets:

- SARSA uses the Q-value of the next state-action pair actually taken:  $Q(S_{t+1}, A_{t+1})$ , making it on-policy.
- Q-Learning uses the maximum Q-value for the next state:  $\max_a Q(S_{t+1}, a)$ , regardless of what action is actually taken next, making it off-policy.

This means Q-Learning can learn the optimal policy even while following an exploratory policy, while SARSA learns the optimal policy subject to the constraint of the exploration strategy being used.

## 5.3 Expected SARSA

Expected SARSA is a variant of SARSA that uses the expected value of the next state-action pair instead of the sample:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a) - Q(S_t, A_t)] \quad (4)$$

### Notation Overview

- $Q(S_t, A_t)$  - estimated value of taking action  $A_t$  in state  $S_t$
- $R_{t+1}$  - reward received after taking action  $A_t$  in state  $S_t$
- $S_{t+1}$  - next state after taking action  $A_t$  in state  $S_t$
- $\pi(a|S_{t+1})$  - probability of taking action  $a$  in state  $S_{t+1}$  under policy  $\pi$
- $\sum_a \pi(a|S_{t+1})Q(S_{t+1}, a)$  - expected value of the next state-action pair
- $\alpha$  - step-size parameter (learning rate)
- $\gamma$  - discount factor

Expected SARSA can be more computationally expensive than SARSA but can reduce variance in the updates and lead to better performance in some environments.

## 6 Policy Gradient Methods

Policy gradient methods take a different approach by directly parameterizing the policy and updating the parameters to maximize the expected return.

### 6.1 Policy Parameterization

The policy is represented as a parameterized function:

$$\pi(a|s, \theta) = P(A_t = a | S_t = s, \theta_t = \theta) \quad (5)$$

where  $\theta$  is a vector of policy parameters that we adjust to optimize performance.

#### Notation Overview

- $\pi(a|s, \theta)$  - probability of taking action  $a$  in state  $s$  under the policy parameterized by  $\theta$
- $\theta$  - vector of policy parameters
- $A_t$  - action at time  $t$
- $S_t$  - state at time  $t$
- $\theta_t$  - policy parameters at time  $t$

### 6.2 Performance Measure

The performance measure to be maximized is the expected return:

$$J(\theta) = \mathbb{E}_{\pi_\theta}[G_0] \quad (6)$$

where  $G_0$  is the return starting from the initial state.

#### Notation Overview

- $J(\theta)$  - expected return under the policy parameterized by  $\theta$
- $\mathbb{E}_{\pi_\theta}$  - expectation over trajectories generated by following policy  $\pi_\theta$
- $G_0$  - return (cumulative discounted reward) starting from the initial state
- $\pi_\theta$  - policy parameterized by  $\theta$

### 6.3 Policy Gradient Theorem

The policy gradient theorem provides a way to compute the gradient of  $J(\theta)$  without requiring differentiation of the state distribution:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi(A_t | S_t, \theta) \cdot G_t] \quad (7)$$

#### Notation Overview

- $\nabla_{\theta} J(\theta)$  - gradient of the expected return with respect to policy parameters
- $\mathbb{E}_{\pi_{\theta}}$  - expectation over trajectories generated by following policy  $\pi_{\theta}$
- $\nabla_{\theta} \log \pi(A_t|S_t, \theta)$  - gradient of the log probability of taking action  $A_t$  in state  $S_t$
- $G_t$  - return (cumulative discounted reward) from time step  $t$
- $\pi_{\theta}$  - policy parameterized by  $\theta$

## 6.4 REINFORCE Algorithm

REINFORCE is a Monte Carlo policy gradient method that implements the policy gradient theorem:

---

**Algorithm 7** REINFORCE (Monte Carlo Policy Gradient)

---

- 1: **Parameters:** step size  $\alpha > 0$
  - 2: Initialize policy parameters  $\theta$  arbitrarily
  - 3: **for** each episode **do**
  - 4:   Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$  following  $\pi(\Delta|\Delta, \theta)$
  - 5:   **for**  $t = 0, 1, \dots, T-1$  **do**
  - 6:      $G \leftarrow$  return from step  $t$
  - 7:      $\theta \leftarrow \theta + \alpha \gamma^t G \nabla_{\theta} \log \pi(A_t|S_t, \theta)$
  - 8:   **end for**
  - 9: **end for**
- 

#### Notation Overview

- $\theta$  - policy parameters
- $\alpha$  - step-size parameter (learning rate)
- $\gamma$  - discount factor
- $G$  - return (cumulative discounted reward)
- $\nabla_{\theta} \log \pi(A_t|S_t, \theta)$  - gradient of the log probability of taking action  $A_t$  in state  $S_t$
- $S_t$  - state at time  $t$
- $A_t$  - action at time  $t$
- $R_t$  - reward at time  $t$
- $T$  - final time step of the episode

## 6.5 Advantages of Policy Gradient Methods

Policy gradient methods have several advantages:

- They can learn stochastic policies, which is important in partially observable environments
- They naturally handle continuous action spaces

- They can incorporate prior knowledge by initializing the policy to specific behaviors
- Convergence properties are often better than value-based methods
- Function approximation is more stable with policy gradients than with value-based methods

## 7 Comparison of Methods

Method	Model Required	On/Off-Policy	Bootstrapping	Updates
Dynamic Programming	Yes	On-policy	Yes	Sweep
Monte Carlo	No	Both	No	Episode
TD Learning	No	Both	Yes	Step
Q-Learning	No	Off-policy	Yes	Step
SARSA	No	On-policy	Yes	Step
Policy Gradient	No	Typically on-policy	Varies	Episode/Step

Table 1: Comparison of Reinforcement Learning Methods

## 8 Conclusion

We have explored various alternative solution methods for reinforcement learning problems beyond the classical dynamic programming approaches. Each method has its strengths and weaknesses, making them suitable for different types of problems and environments.

Monte Carlo methods are conceptually simple and learn directly from complete episodes, but they can only be applied to episodic tasks and may require more samples. Temporal Difference methods combine ideas from Monte Carlo and dynamic programming, allowing learning from incomplete sequences and bootstrapping from existing estimates. Q-Learning provides an off-policy approach that can learn the optimal policy while following an exploratory policy. Finally, policy gradient methods offer a direct approach to optimizing policies, which is particularly useful for continuous action spaces and stochastic policies.

The choice of method depends on the specific problem characteristics, such as whether the environment is episodic or continuing, whether a model is available, and whether the state and action spaces are discrete or continuous.