

# Reinforcement Learning Pseudocode Collection

September 20, 2025

## Global Notation

States  $s \in \mathcal{S}$ , actions  $a \in \mathcal{A}(s)$ , rewards  $r \in \mathbb{R}$ . Discount  $\gamma \in [0, 1]$ , step-size  $\alpha \in (0, 1]$ , threshold  $\theta > 0$ . Value functions:  $V(s)$ ,  $Q(s, a)$ . Policies: target  $\pi(a|s)$  and behavior  $b(a|s)$ , where  $b$  has coverage of  $\pi$ . Episodes are denoted  $S_0, A_0, R_1, S_1, A_1, \dots, S_T$  with terminal at time  $T$ .

## 1 Bandits

---

**Algorithm 1**  $\varepsilon$ -Greedy  $k$ -Armed Bandit (Incremental Averages)

---

**Input:**  $\varepsilon \in [0, 1]$

- 1: Initialize  $Q(a) \leftarrow 0$ ,  $N(a) \leftarrow 0$  for all arms  $a = 1, \dots, k$
  - 2: **for**  $t = 1, 2, \dots$  **do**
  - 3:   With prob.  $\varepsilon$ : pick  $A \sim \text{Uniform}(\{1, \dots, k\})$ ; else  $A \leftarrow \arg \max_a Q(a)$
  - 4:   Pull arm  $A$  and observe reward  $R$
  - 5:    $N(A) \leftarrow N(A) + 1$
  - 6:    $Q(A) \leftarrow Q(A) + \frac{1}{N(A)}(R - Q(A))$  ▷ incremental mean
- 

---

**Algorithm 2** Incremental Mean Update (Reference)

---

**Input:** Current mean  $Q_n$ , new sample  $R_n$ , count  $n$

- 1:  $Q_{n+1} \leftarrow Q_n + \frac{1}{n}(R_n - Q_n)$
- 

## 2 Markov Decision Processes

We assume a known or unknown transition *kernel*  $p(s', r \mid s, a)$  depending on the algorithm. In model-based Dynamic Programming (DP),  $p$  is known; in model-free methods (MC/TD), we learn from experience.

### 3 Dynamic Programming (Model-Based)

---

**Algorithm 3** Iterative Policy Evaluation (Prediction)

---

**Input:** Policy  $\pi$ , model  $p(s', r \mid s, a)$ , discount  $\gamma$ , threshold  $\theta$

- 1: Initialize  $V(s) \leftarrow 0$  for all  $s$
  - 2: **repeat**
  - 3:      $\Delta \leftarrow 0$
  - 4:     **for** each state  $s$  **do**
  - 5:          $v \leftarrow V(s)$
  - 6:          $V(s) \leftarrow \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$
  - 7:          $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - 8:     **until**  $\Delta < \theta$
  - 9: **return**  $V \approx v_\pi$
- 

---

**Algorithm 4** Policy Iteration

---

**Input:** Model  $p(s', r \mid s, a)$ , discount  $\gamma$ , threshold  $\theta$

- 1: Initialize  $V(s)$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s$
  - 2: **repeat** ▷ Policy Evaluation
  - 3:     **repeat**
  - 4:          $\Delta \leftarrow 0$
  - 5:         **for** each state  $s$  **do**
  - 6:              $v \leftarrow V(s)$
  - 7:              $a \leftarrow \pi(s)$
  - 8:              $V(s) \leftarrow \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$
  - 9:              $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - 10:     **until**  $\Delta < \theta$  ▷ Policy Improvement
  - 11:     policy\_stable  $\leftarrow$  true
  - 12:     **for** each state  $s$  **do**
  - 13:         old  $\leftarrow \pi(s)$
  - 14:          $\pi(s) \leftarrow \arg \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$
  - 15:         **if**  $\pi(s) \neq \text{old}$  **then**
  - 16:             policy\_stable  $\leftarrow$  false
  - 17:     **until** policy\_stable
  - 18: **return**  $(V, \pi) \approx (v^*, \pi^*)$
-

---

**Algorithm 5** Value Iteration

---

**Input:** Model  $p(s', r \mid s, a)$ , discount  $\gamma$ , threshold  $\theta$

```
1: Initialize  $V(s) \leftarrow 0$  for all  $s$ 
2: repeat
3:    $\Delta \leftarrow 0$ 
4:   for each state  $s$  do
5:      $v \leftarrow V(s)$ 
6:      $V(s) \leftarrow \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$ 
7:      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
8: until  $\Delta < \theta$ 
9: Derive  $\pi(s) \leftarrow \arg \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$ 
10: return  $(V, \pi)$ 
```

---

## 4 Monte Carlo Methods (Model-Free)

---

**Algorithm 6** First-Visit Monte Carlo Prediction (State Values)

---

**Input:** Policy  $\pi$ , discount  $\gamma$

```
1: Initialize  $V(s)$  arbitrarily;  $\text{Returns}(s) \leftarrow$  empty list for all  $s$ 
2: for all episodes do
3:   Generate  $S_0, A_0, R_1, \dots, S_T$  using  $\pi$ 
4:    $G \leftarrow 0$ 
5:   for  $t \leftarrow T - 1$  down to 0 do
6:      $G \leftarrow \gamma G + R_{t+1}$ 
7:     if  $S_t$  is first visit to its state in episode then
8:       append  $G$  to  $\text{Returns}(S_t)$ 
9:        $V(S_t) \leftarrow \text{average}(\text{Returns}(S_t))$ 
```

---

---

**Algorithm 7** On-Policy First-Visit MC Control (  $\varepsilon$ -Soft )

---

**Input:** Small  $\varepsilon > 0$ , discount  $\gamma$

```
1: Initialize  $\pi$  as any  $\varepsilon$ -soft policy;  $Q(s, a)$  arbitrarily;  $\text{Returns}(s, a) \leftarrow$  empty lists
2: for all episodes do
3:   Generate  $S_0, A_0, R_1, \dots, S_T$  using  $\pi$ 
4:    $G \leftarrow 0$ 
5:   for  $t \leftarrow T - 1$  down to 0 do
6:      $G \leftarrow \gamma G + R_{t+1}$ 
7:     if  $(S_t, A_t)$  is first visit then
8:       append  $G$  to  $\text{Returns}(S_t, A_t)$ 
9:        $Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$ 
10:       $A^* \leftarrow \arg \max_a Q(S_t, a)$ 
11:      for all  $a \in \mathcal{A}(S_t)$  do
12:        if  $a = A^*$  then
13:           $\pi(a|S_t) \leftarrow 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)|$ 
14:        else
15:           $\pi(a|S_t) \leftarrow \varepsilon/|\mathcal{A}(S_t)|$ 
```

---

---

**Algorithm 8** MC Control with Exploring Starts (MCES)

---

**Input:** Discount  $\gamma$ ; assume exploring starts

- 1: Initialize  $Q(s, a)$  arbitrarily;  $\pi(s)$  arbitrary
  - 2: **for all** episodes **do**
  - 3:   Choose  $(S_0, A_0)$  with nonzero prob. for all state–action pairs
  - 4:   Generate episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_T$
  - 5:    $G \leftarrow 0$
  - 6:   **for**  $t \leftarrow T - 1$  down to 0 **do**
  - 7:      $G \leftarrow \gamma G + R_{t+1}$
  - 8:     **if**  $(S_t, A_t)$  first visit **then**
  - 9:       Update  $Q(S_t, A_t) \leftarrow$  average of returns for  $(S_t, A_t)$
  - 10:     Improve  $\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$
- 

---

**Algorithm 9** Off-Policy MC Prediction (Weighted Importance Sampling)

---

**Input:** Target policy  $\pi$ , behavior policy  $b$  (with coverage), discount  $\gamma$

- 1: Initialize  $Q(s, a)$  arbitrarily;  $C(s, a) \leftarrow 0$  for all  $(s, a)$
  - 2: **for all** episodes **do**
  - 3:   Generate episode using  $b$ :  $S_0, A_0, R_1, \dots, S_T$
  - 4:    $G \leftarrow 0, W \leftarrow 1$
  - 5:   **for**  $t \leftarrow T - 1$  down to 0 **while**  $W \neq 0$  **do**
  - 6:      $G \leftarrow \gamma G + R_{t+1}$
  - 7:      $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
  - 8:      $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} (G - Q(S_t, A_t))$
  - 9:      $W \leftarrow W \cdot \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$
- 

---

**Algorithm 10** Off-Policy MC Control (Weighted IS)

---

**Input:** Behavior  $b$  (soft, coverage), discount  $\gamma$

- 1: Initialize  $Q(s, a)$  arbitrarily;  $C(s, a) \leftarrow 0$ ;  $\pi(s) \leftarrow \arg \max_a Q(s, a)$
  - 2: **for all** episodes **do**
  - 3:   Generate episode using  $b$ :  $S_0, A_0, R_1, \dots, S_T$
  - 4:    $G \leftarrow 0, W \leftarrow 1$
  - 5:   **for**  $t \leftarrow T - 1$  down to 0 **while**  $W \neq 0$  **do**
  - 6:      $G \leftarrow \gamma G + R_{t+1}$
  - 7:      $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
  - 8:      $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} (G - Q(S_t, A_t))$
  - 9:      $\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$
  - 10:    **if**  $A_t \neq \pi(S_t)$  **then**
  - 11:      **break**
  - 12:     $W \leftarrow W \cdot \frac{1}{b(A_t|S_t)}$
-

---

**Algorithm 11** Off-Policy Incremental Update (Scalar Form)

---

**Input:** Stream of returns  $\{G_n\}$  with weights  $\{W_n\}$

- 1: Initialize  $v \leftarrow 0$ ,  $c \leftarrow 0$
  - 2: **for**  $n = 1, 2, \dots$  **do**
  - 3:      $c \leftarrow c + W_n$
  - 4:      $v \leftarrow v + \frac{W_n}{c}(G_n - v)$
- 

## 5 Temporal-Difference Methods (Model-Free)

---

**Algorithm 12** TD(0) Prediction

---

**Input:** Policy  $\pi$ , step-size  $\alpha$ , discount  $\gamma$

- 1: Initialize  $V(s)$  arbitrarily with  $V(\text{terminal}) = 0$
  - 2: **for all** episodes **do**
  - 3:     Initialize  $S$
  - 4:     **while**  $S$  not terminal **do**
  - 5:         Choose  $A \sim \pi(\cdot|S)$ ; take  $A$ , observe  $R, S'$
  - 6:          $V(S) \leftarrow V(S) + \alpha(R + \gamma V(S') - V(S))$
  - 7:          $S \leftarrow S'$
- 

---

**Algorithm 13** SARSA (On-Policy TD Control)

---

**Input:** Step-size  $\alpha$ , discount  $\gamma$ ; behavior=target policy (e.g.,  $\varepsilon$ -greedy w.r.t.  $Q$ )

- 1: Initialize  $Q(s, a)$  arbitrarily with  $Q(\text{terminal}, \cdot) = 0$
  - 2: **for all** episodes **do**
  - 3:     Initialize  $S$ ; choose  $A$  from  $S$  using policy derived from  $Q$
  - 4:     **while**  $S$  not terminal **do**
  - 5:         Take  $A$ , observe  $R, S'$
  - 6:         Choose  $A'$  from  $S'$  using policy derived from  $Q$
  - 7:          $Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A))$
  - 8:          $S \leftarrow S'$ ;  $A \leftarrow A'$
- 

---

**Algorithm 14** Q-Learning (Off-Policy TD Control)

---

**Input:** Step-size  $\alpha$ , discount  $\gamma$ ; behavior policy (e.g.,  $\varepsilon$ -greedy w.r.t.  $Q$ )

- 1: Initialize  $Q(s, a)$  arbitrarily with  $Q(\text{terminal}, \cdot) = 0$
  - 2: **for all** episodes **do**
  - 3:     Initialize  $S$
  - 4:     **while**  $S$  not terminal **do**
  - 5:         Choose  $A$  using behavior policy; take  $A$ , observe  $R, S'$
  - 6:          $Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma \max_a Q(S', a) - Q(S, A))$
  - 7:          $S \leftarrow S'$
-

---

**Algorithm 15**  $n$ -Step TD Prediction

---

**Input:** Policy  $\pi$ , step-size  $\alpha$ , discount  $\gamma$ , integer  $n \geq 1$

```
1: Initialize  $V(s)$  arbitrarily;  $T \leftarrow \infty$ 
2: for all episodes do
3:   Initialize and store  $S_0 \neq \text{terminal}$ 
4:   for  $t = 0, 1, 2, \dots$  do
5:     if  $t < T$  then
6:       Take  $A_t \sim \pi(\cdot|S_t)$ ; observe  $R_{t+1}, S_{t+1}$ 
7:       if  $S_{t+1}$  terminal then
8:          $T \leftarrow t + 1$ 
9:        $\tau \leftarrow t - n + 1$  ▷ time whose estimate to update
10:      if  $\tau \geq 0$  then
11:         $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
12:        if  $\tau + n < T$  then
13:           $G \leftarrow G + \gamma^n V(S_{\tau+n})$ 
14:           $V(S_\tau) \leftarrow V(S_\tau) + \alpha(G - V(S_\tau))$ 
15:      if  $\tau = T - 1$  then
16:        break
```

---

---

**Algorithm 16**  $n$ -Step SARSA (On-Policy)

---

**Input:** Step-size  $\alpha$ , discount  $\gamma$ , integer  $n \geq 1$ ,  $\varepsilon$  for  $\varepsilon$ -greedy

```
1: Initialize  $Q(s, a)$  arbitrarily; define  $\pi$  as  $\varepsilon$ -greedy w.r.t.  $Q$ ;  $T \leftarrow \infty$ 
2: for all episodes do
3:   Initialize  $S_0$ ; choose  $A_0 \sim \pi(\cdot|S_0)$ 
4:   for  $t = 0, 1, 2, \dots$  do
5:     if  $t < T$  then
6:       Take  $A_t$ , observe  $R_{t+1}, S_{t+1}$ 
7:       if  $S_{t+1}$  terminal then
8:          $T \leftarrow t + 1$ 
9:       else
10:        Choose  $A_{t+1} \sim \pi(\cdot|S_{t+1})$ 
11:       $\tau \leftarrow t - n + 1$ 
12:      if  $\tau \geq 0$  then
13:         $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
14:        if  $\tau + n < T$  then
15:           $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
16:           $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha(G - Q(S_\tau, A_\tau))$ 
17:          Update  $\pi(\cdot|S_\tau)$  to be  $\varepsilon$ -greedy w.r.t.  $Q(S_\tau, \cdot)$ 
18:      if  $\tau = T - 1$  then
19:        break
```

---

---

**Algorithm 17**  $n$ -Step Off-Policy SARSA (Importance Sampling)

---

**Input:** Step-size  $\alpha$ , discount  $\gamma$ , integer  $n \geq 1$ ; target policy  $\pi$ ; behavior policy  $b$

```
1: Initialize  $Q(s, a)$  arbitrarily;  $T \leftarrow \infty$ 
2: for all episodes do
3:   Initialize  $S_0$ ; choose  $A_0 \sim b(\cdot|S_0)$ 
4:   for  $t = 0, 1, 2, \dots$  do
5:     if  $t < T$  then
6:       Take  $A_t$ , observe  $R_{t+1}, S_{t+1}$ 
7:       if  $S_{t+1}$  terminal then
8:          $T \leftarrow t + 1$ 
9:       else
10:        Choose  $A_{t+1} \sim b(\cdot|S_{t+1})$ 
11:       $\tau \leftarrow t - n + 1$ 
12:      if  $\tau \geq 0$  then
13:         $\rho \leftarrow \prod_{k=\tau+1}^{\min(\tau+n-1, T-1)} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}$ 
14:         $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
15:        if  $\tau + n < T$  then
16:           $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
17:           $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha \rho (G - Q(S_\tau, A_\tau))$ 
18:        if  $\tau = T - 1$  then
19:          break
```

---

## 6 Greedy Policy Improvement (Helper)

Given action-values  $Q(s, \cdot)$ , a deterministic improvement is  $\pi(s) \leftarrow \arg \max_a Q(s, a)$ . An  $\varepsilon$ -soft version at state  $s$  sets the greedy action prob. to  $1 - \varepsilon + \varepsilon/|\mathcal{A}(s)|$  and all others to  $\varepsilon/|\mathcal{A}(s)|$ .