

图 12.1 马尔可夫模型的状态图

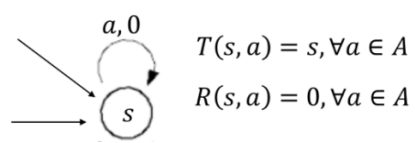


图 12.2 终止状态示意图

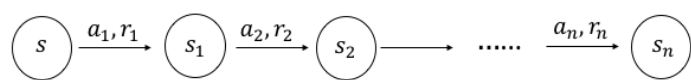


图 12.3 长度为  $n$  的状态路径

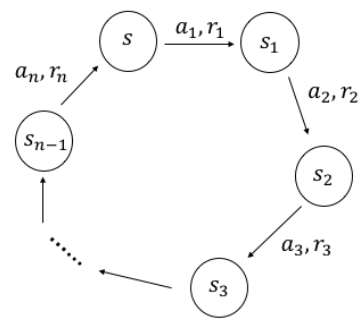


图 12.4 形成一个圈的状态路径

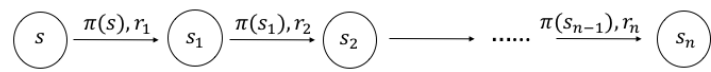


图 12.5 策略  $\pi$  从状态  $s$  发出的  $n$  个行动的状态图

### 值迭代算法

Compute\_environment\_Q\_value(  $S, A, T, R, \gamma$ ):

For all  $s \in S$ :  $V_0(s) = 0$

$n = 0$

While True:

$n \leftarrow n + 1$

For all  $s \in S$ :

For all  $a \in A$ :

$$Q_n(s, a) = R(s, a) + \gamma V_{n-1}(T(s, a))$$

For all  $s \in S$ :

$$V_n(s) = \max_{a \in A} Q_n(s, a)$$

If for all  $s \in S$ :  $V_n(s) = V_{n-1}(s)$ :

Break

Return  $Q_n$

Value\_Iteration (  $S, A, T, R, \gamma$ ):

$Q = \text{Compute\_environment\_Q\_value} (S, A, T, R, \gamma)$

For all  $s \in S$ :

$$\pi(s) = \operatorname{argmax}_{a \in A} Q(s, a)$$

Return  $\pi$

图 12.6 值迭代算法描述

**machine\_learning.reinforcement\_learning.environment**

```
1  import numpy as np
2
3  class Environment:
4      S = []
5      A = []
6      T = []
7      R = []
8      S_end = { }
9
10     def reset(self):
11         s_start = np.random.randint(0, len(self.S))
12         return s_start
```

图 12.7 Envinronment 类

```

machine_learning.reinforcement_learning.value_iteration

1  import numpy as np
2
3  def compute_environment_Q_values(S, A, T, R, gamma):
4      V = np.zeros(len(S))
5      Q = np.zeros((len(S), len(A)))
6      while True:
7          for s in S:
8              for a in A:
9                  Q[s][a] = R[s][a] + gamma * V[int(T[s][a])]
10         converge = True
11         for s in S:
12             if np.max(Q[s]) != V[s]:
13                 converge = False
14                 V[s] = np.max(Q[s])
15         if converge:
16             break
17     return Q
18
19 def value_iteration(env, gamma):
20     S,A,T,R = env.S, env.A, env.T, env.R
21     Q = compute_environment_Q_values(S, A, T, R, gamma)
22     pi = np.argmax(Q, axis = 1)
23     return pi

```

图 12.8 值迭代算法





图 12.9 地雷与宝藏游戏

```

machine_learning.reinforcement_learning.landmine_and_treasure

1  import numpy as np
2  from machine_learning.reinforcement_learning.environment import Environment
3
4  class LandmineAndTreasure(Environment):
5      def __init__(self, m, d):
6          n_states = m
7          n_actions = 2
8          S = range(n_states)
9          S_end = {0, d, n_states - 1}
10         A = range(n_actions)
11         R = np.zeros((n_states, n_actions))
12         R[n_states - 2][1] = 1000
13         R[1][0] = 1000
14         R[d + 1][0] = -1000
15         R[d - 1][1] = -1000
16         T = np.zeros((n_states, n_actions))
17         for s in S:
18             T[s][0] = s - 1
19             T[s][1] = s + 1
20             if s in S_end:
21                 T[s][0] = s
22                 T[s][1] = s
23         self.S, self.A, self.T, self.R, self.S_end = S, A, T, R, S_end

```

图 12.10 生成 LandmineAndTreasure 子环境的算法

```
1 from machine_learning.reinforcement_learning.value_iteration import value_iteration
2 import machine_learning.reinforcement_learning.landmine_and_treasure as game
3
4 env = game.LandmineAndTreasure(m = 100, d = 30)
5 pi = value_iteration(env, 0.95)
6 print(pi)
```

图 12.11 小机器人的最优行走策略的值迭代算法

### 策略迭代算法

Compute\_policy\_Q\_value (  $S, A, T, R, \gamma, \pi$ ):

For all  $s \in S : V_0^\pi(s) = 0$

$n = 1$

While True:

$n \leftarrow n + 1$

For all  $s \in S :$

For all  $a \in A :$

$$Q_n^\pi(s, a) = R(s, a) + \gamma V_{n-1}^\pi(T(s, a))$$

For all  $s \in S :$

$$V_n^\pi(s) = Q_n^\pi(s, \pi(s))$$

If  $V_n^\pi = V_{n-1}^\pi :$

Break

Return  $Q_n^\pi$

Policy\_Iteration(  $S, A, T, R, \gamma$ ):

For all  $s \in S :$

$$\pi_0(s) = \operatorname{argmax}_{a \in A} R(s, a)$$

$n = 0$

While True:

$$Q^{\pi_n} = \text{Compute\_policy\_Q\_value}(S, A, T, R, \gamma, \pi_n)$$

For all  $s \in S :$

$$\pi_{n+1}(s) = \operatorname{argmax}_{a \in A} Q^{\pi_n}(s, a)$$

If  $\pi_{n+1} = \pi_n :$

Break

$n \leftarrow n + 1$

Return  $\pi_n$

图 12.12 策略迭代算法描述

### **machine\_learning.reinforcement\_learning.policy\_iteration**

```
1  import numpy as np
2
3  def compute_policy_Q_values(S, A, T, R, gamma, pi):
4      V = np.zeros(len(S))
5      Q = np.zeros((len(S), len(A)))
6      while True:
7          for s in S:
8              for a in A:
9                  Q[s][a] = R[s][a] + gamma * V[int(T[s][a])]
10         converge = True
11         for s in S:
12             if Q[s][pi[s]] != V[s]:
13                 converge = False
14             V[s] = Q[s][pi[s]]
15         if converge:
16             break
17     return Q
18
19 def policy_iteration(env, gamma):
20     S,A,T,R = env.S, env.A, env.T, env.R
21     pi = np.argmax(R, axis = 1)
22     while True:
23         Q = compute_policy_Q_values(S, A, T, R, gamma, pi)
24         converge = True
25         for s in S:
26             if np.argmax(Q[s]) != pi[s]:
27                 converge = False
28             pi[s] = np.argmax(Q[s])
29         if converge:
30             break
31     return pi
```

图 12.13 策略迭代算法

### $\epsilon$ 贪心探索策略

Epsilon\_greedy ( $\tilde{Q}, s, A, \epsilon$ ):

$r = \text{random number} \in [0,1)$

If  $r < \epsilon$ :

Return random  $a \in A$

Else:

Return  $\operatorname{argmax}_{a \in A} \tilde{Q}(s, a)$

图 12.14  $\epsilon$  贪心探索策略描述

### Sarsa 算法

Sarsa( $S, A, T, R, N, \gamma, \epsilon, \eta$ ):

For all  $s \in S, a \in A$ :  $\tilde{Q}(s, a) = 0$

For  $i = 1, 2, \dots, N$ :

$s_{\text{cur}} = \text{initial state in } S$

While  $s_{\text{cur}} \notin \{\text{terminal states of } S\}$ :

$a_{\text{cur}} = \text{Epsilon\_greedy}(\tilde{Q}, s_{\text{cur}}, A, \epsilon)$

$s_{\text{next}} = T(s_{\text{cur}}, a_{\text{cur}})$

$a_{\text{next}} = \text{Epsilon\_greedy}(\tilde{Q}, s_{\text{next}}, A, \epsilon)$

$\tilde{Q}(s_{\text{cur}}, a_{\text{cur}}) \leftarrow (1 - \eta)\tilde{Q}(s_{\text{cur}}, a_{\text{cur}}) + \eta \left( R(s_{\text{cur}}, a_{\text{cur}}) + \gamma \tilde{Q}(s_{\text{next}}, a_{\text{next}}) \right)$

$s_{\text{cur}} = s_{\text{next}}$

For all  $s \in S$ :

$\pi(s) = \operatorname{argmax}_{a \in A} \tilde{Q}(s, a)$

Return  $\pi$

图 12.15 Sarsa 算法描述

**machine\_learning.reinforcement\_learning.epsilon\_greedy**

```
1 import numpy as np
2
3 def epsilon_greedy(Q_s, n_actions, epsilon):
4     if np.random.rand() < epsilon:
5         return np.random.randint(n_actions)
6     else:
7         return np.argmax(Q_s)
```

图 12.16  $\epsilon$  贪心探索策略



#### **machine\_learning.reinforcement\_learning.sarsa**

```
1 import numpy as np
2 from machine_learning.reinforcement_learning.epsilon_greedy import epsilon_greedy
3
4 def sarsa(env, N, gamma, epsilon, eta):
5     S, A, T, R = env.S, env.A, env.T, env.R
6     n_states, n_actions = len(S), len(A)
7     Q = np.zeros((n_states, n_actions))
8     for iter in range(N):
9         env.reset()
10        s_cur = env.s_start
11        while s_cur not in env.S_end:
12            a_cur = epsilon_greedy(Q[s_cur], n_actions, epsilon)
13            s_next = int(T[s_cur][a_cur])
14            a_next = epsilon_greedy(Q[s_next], n_actions, epsilon)
15            q = (1 - eta) * Q[s_cur][a_cur] + eta * (R[s_cur][a_cur] + gamma * Q[s_next][a_next])
16            Q[s_cur][a_cur] = q
17            s_cur = s_next
18    return np.argmax(Q, axis = 1)
```

图 12.17 Sarsa 算法

## Q 学习算法

Q\_learning ( $S, A, T, R, N, \gamma, \epsilon, \eta$ ):

For all  $s \in S, a \in A$ :  $\tilde{Q}(s, a) = 0$

For  $i = 1, 2, \dots, N$ :

$s_{\text{cur}} = \text{initial state in } S$

While  $s_{\text{cur}} \notin \{\text{terminal states of } S\}$ :

$a_{\text{cur}} = \text{Epsilon\_greedy}(\tilde{Q}, s_{\text{cur}}, A, \epsilon)$

$s_{\text{next}} = T(s_{\text{cur}}, a_{\text{cur}})$

$a_{\text{next}} = \operatorname{argmax}_{a \in A} \tilde{Q}(s_{\text{next}}, a)$

$\tilde{Q}(s_{\text{cur}}, a_{\text{cur}}) \leftarrow (1 - \eta)\tilde{Q}(s_{\text{cur}}, a_{\text{cur}}) + \eta \left( R(s_{\text{cur}}, a_{\text{cur}}) + \gamma \tilde{Q}(s_{\text{next}}, a_{\text{next}}) \right)$

$s_{\text{cur}} = s_{\text{next}}$

For all  $s \in S$ :

$\pi(s) = \operatorname{argmax}_{a \in A} \tilde{Q}(s, a)$

Return  $\pi$

图 12.18 Q 学习算法描述

```

machine_learning.reinforcement_learning.q_learning

1  import numpy as np
2  from machine_learning.reinforcement_learning.epsilon_greedy import epsilon_greedy
3
4  def q_learning(env, N, gamma, epsilon, eta):
5      S, A, T, R = env.S, env.A, env.T, env.R
6      n_states, n_actions = len(S), len(A)
7      Q = np.zeros((n_states, n_actions))
8      for iter in range(N):
9          env.reset()
10         s_cur = env.s_start
11         while s_cur not in env.S_end:
12             a_cur = epsilon_greedy(Q[s_cur], n_actions, epsilon)
13             s_next = int(T[s_cur][a_cur])
14             a_next = np.argmax(Q[s_next])
15             
$$q = (1 - \text{eta}) * Q[s\_cur][a\_cur] + \text{eta} * (R[s\_cur][a\_cur] + \text{gamma} * Q[s\_next][a\_next])$$

16             Q[s_cur][a_cur] = q
17             s_cur = s_next
18     return np.argmax(Q, axis = 1)

```

图 12.19 Q 学习算法

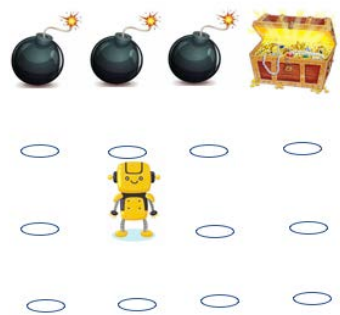


图 12.20 埋藏着宝藏的沼泽地

```

machine_learning.reinforcement_learning.landmine_and_treasure_2d

1  import numpy as np
1  from machine_learning.reinforcement_learning.environment import Environment
2
4  class LandmineAndTreasure2d(Environment):
5      def __init__(self, m, n):
6          n_states = m * n
7          n_actions = 4
8          S = range(n_states)
9          S_end = range(n)
10         A = range(n_actions)
11         R = np.zeros((n_states, n_actions))
12         for j in range(n-1):
13             R[n + j][0] = -1000
14             R[2 * n - 1][0] = 100
15         T = np.zeros((n_states, n_actions))
16         for i in range(m):
17             for j in range(n):
18                 s = n * i + j
19                 T[s][0] = s - n if i > 0 else s
20                 T[s][1] = s + 1 if j < n - 1 else s
21                 T[s][2] = s + n if i < m - 1 else s
22                 T[s][3] = s - 1 if j > 0 else s
23         for s in S_end:
24             for a in A:
25                 T[s][a] = s
26         self.S, self.A, self.T, self.R, self.S_end = S, A, T, R, S_end

```

图 12.21 地雷与宝藏游戏的环境模型

```
1 from machine_learning.reinforcement_learning.td.sarsa import sarsa
2 from machine_learning.reinforcement_learning.td.q_learning import q_learning
3 import machine_learning.reinforcement_learning.landmine_and_treasure_2d as game
4
5 m = 4
6 n = 4
7 env = game.LandmineAndTreasure2d(m, n)
8 pi_sarsa = sarsa(env, 1000, 0.95, 0.2, 0.1)
9 pi_q = q_learning(env, 1000, 0.95, 0.2, 0.1)
10
11 map = {0: "U", 1: "R", 2: "D", 3: "L"}
12 for i in range(m):
13     print([map[pi_sarsa[s]] for s in range(i * n, (i + 1) * n)])
14 for i in range(m):
15     print([map[pi_q[s]] for s in range(i * n, (i + 1) * n)])
```

图 12.22 Sarsa 算法和 Q 学习算法求解寻宝最优策略

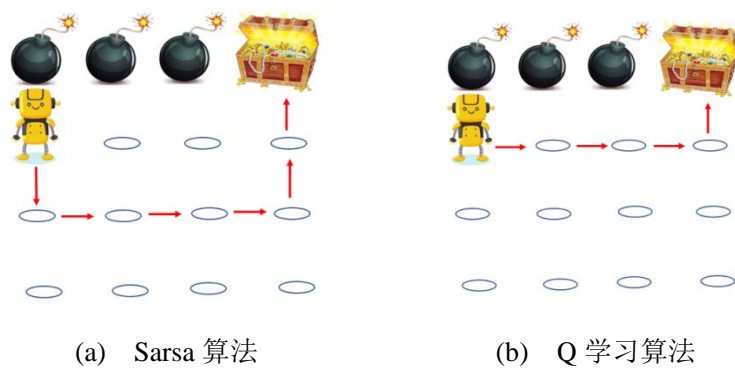


图 12.23 两个算法的寻宝路线图

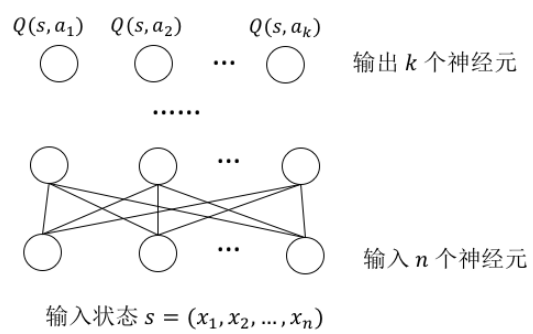


图 12.24 DQN 模型



### DQN 算法

Deep\_Q\_Network ( $S, A, T, R, N, \gamma, \epsilon$ ):

$DQN = \text{initial DQN model}$

For  $i = 1, 2, \dots, N$ :

$s_{\text{cur}} = \text{initial state of } S$

While  $s_{\text{cur}} \notin \{\text{terminal states of } S\}$ :

$a_{\text{cur}} = \text{Epsilon\_greedy}(DQN, s_{\text{cur}}, A, \epsilon)$

$s_{\text{next}} = T(s_{\text{cur}}, a_{\text{cur}})$

$a_{\text{next}} = \operatorname{argmax}_{a \in A} DQN(s_{\text{next}}, a)$

$\text{loss} = (R(s_{\text{cur}}, a_{\text{cur}}) + \gamma DQN(s_{\text{next}}, a_{\text{next}}) - DQN(s_{\text{cur}}, a_{\text{cur}}))^2$

$\text{BackProp}(DQN, \text{loss})$

$s_{\text{cur}} = s_{\text{next}}$

For all  $s \in S$ :

$\pi(s) = \operatorname{argmax}_{a \in A} DQN(s, a)$

Return  $\pi$

图 12.25 DQN 算法描述

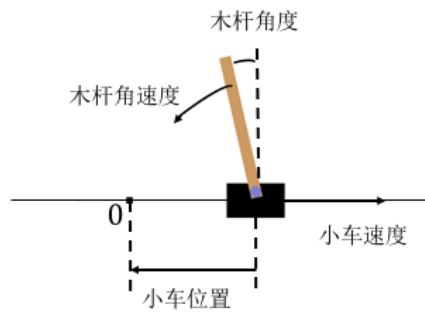


图 12.26 木杆平衡的环境

```
1 import gym
2
3 env = gym.make("CartPole-v0")
4 state = env.reset()
5 action = 1
6 while True:
7     state, reward, done, info = env.step(action)
8     env.render(mode = "rgb_array")
9     if done:
10         break
```

图 12.27 OpenAI 中的木杆平衡虚拟环境

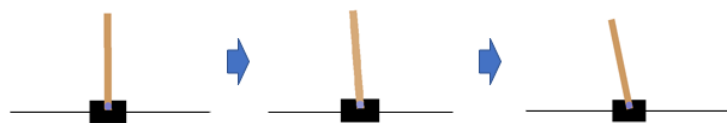


图 12.28 木杆状态变化过程

```

1 import numpy as np
2 import gym
3 import tensorflow as tf
4 import machine_learning.reinforcement_learning.epsilon_greedy as eg
5
6 state_size = 4
7 n_hidden1 = 24
8 n_hidden2 = 24
9 n_actions = 2
10 State = tf.placeholder(tf.float32, shape=[None, state_size])
11 hidden1 = tf.layers.dense(State, n_hidden1, activation = tf.nn.relu)
12 hidden2 = tf.layers.dense(hidden1, n_hidden2, activation = tf.nn.relu)
13 Q_values = tf.layers.dense(hidden2, n_actions)
14
15 Target = tf.placeholder(tf.float32)
16 Action = tf.placeholder(tf.int32)
17 Q_value = tf.reduce_sum(Q_values * tf.one_hot(Action, n_actions))
18 loss = tf.reduce_mean(tf.square(Target - Q_value))
19 optimizer = tf.train.AdamOptimizer(learning_rate = 0.001)
20 training_op = optimizer.minimize(loss)
21
22 env = gym.make("CartPole-v0")
23 gamma = 0.95
24 n_iterations = 1000
25 epsilon_max = 1.0
26 epsilon_min = 0.01
27 epsilon_decay = 0.99
28 stop_penalty = -100
29 with tf.Session() as sess:
30     tf.global_variables_initializer().run()
31     epsilon = epsilon_max
32     for iteration in range(n_iterations):
33         state = env.reset()
34         done = False
35         steps = 0
36         while not done:
37             steps += 1
38             s_cur = state.reshape(1, state_size)

```

```
39     Q_s_cur = Q_values.eval(feed_dict = {State: s_cur})
40     a_cur = eg.epsilon_greedy(Q_s_cur, n_actions, epsilon)
41     if epsilon > epsilon_min:
42         epsilon *= epsilon_decay
43     state, reward, done, info = env.step(action)
44     s_next = state.reshape(1, state_size)
45     if done:
46         target = stop_penalty
47         feed_dict = {State: s_cur, Action: a_cur, Target: target}
48         sess.run(training_op, feed_dict = feed_dict)
49         print(steps)
50         break
51     else:
52         Q_s_next = Q_values.eval(feed_dict = {State: s_next})
53         target = reward + gamma * np.max(Q_s_next, axis=1)
54         feed_dict = {State: s_cur, Action: a_cur, Target: target}
55         sess.run(training_op, feed_dict = feed_dict)
56 env.close()
```

图 12.29 木杆平衡环境的 DQN 算法

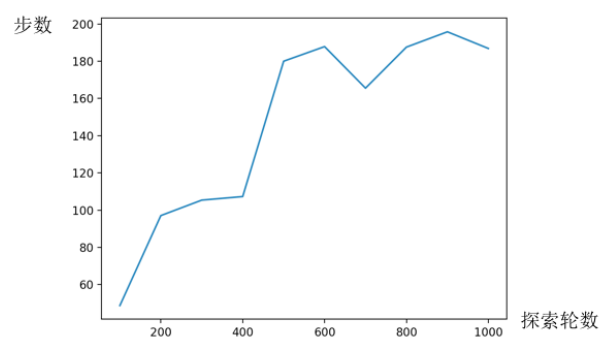


图 12.30 DQN 算法环境探索轮数与木杆平衡步数的关系

### REINFORCE 算法

REINFORCE ( $S, A, T, R, N, \gamma, \eta$ ):

$\mathbf{W}$  = random initial model parameters

For  $i = 1, 2, \dots, N$ :

$\mathbf{s}_{\text{cur}}$  = initial state in  $S$

$t = 0$

While  $\mathbf{s}_{\text{cur}} \notin \{\text{terminal states of } S\}$ :

$a_{\text{cur}} \sim h_{\mathbf{W}}(\mathbf{s}_{\text{cur}})$

$\mathbf{s}_{\text{next}} = T(\mathbf{s}_{\text{cur}}, a_{\text{cur}})$

$\mathbf{W} \leftarrow \mathbf{W} + \eta \cdot \gamma^t \cdot R(\mathbf{s}_{\text{cur}}, a_{\text{cur}}) \nabla \log h_{\mathbf{W}}(\mathbf{s}_{\text{cur}}, a_{\text{cur}})$

$\mathbf{s}_{\text{cur}} = \mathbf{s}_{\text{next}}$

$t \leftarrow t + 1$

Return  $h_{\mathbf{W}}$

图 12.31 REINFORCE 算法描述



```

1 import numpy as np
2 import gym
3
4 def softmax(scores):
5     e = np.exp(scores)
6     s = e.sum()
7     return e / s
8
9 def REINFORCE(env, state_size, n_actions, n_iter, gamma, eta):
10     W = np.random.rand(state_size, n_actions)
11     for iter in range(n_iter):
12         state = env.reset()
13         done = False
14         discount = 1
15         steps = 0
16         while not done:
17             steps += 1
18             s = state.reshape(1, state_size)
19             probs = softmax(s.dot(W))
20             a = np.random.choice(n_actions, p = probs.reshape(-1))
21             state, reward, done, info = env.step(a)
22             y = np.zeros(n_actions)
23             y[a] = 1
24             y = y.reshape(1, n_actions)
25             gradient = -s.T.dot(probs - y)
26             if done:
27                 reward = -100
28                 print("iteration { } lasts for { } steps".format(iter, steps))
29                 W = W + eta * discount * reward * gradient
30                 discount *= gamma
31
32 env = gym.make("CartPole-v0")
33 REINFORCE(env, 4, 2, 3000, 0.95, 0.1)

```

图 12.32 木杆平衡问题的 REINFORCE 算法

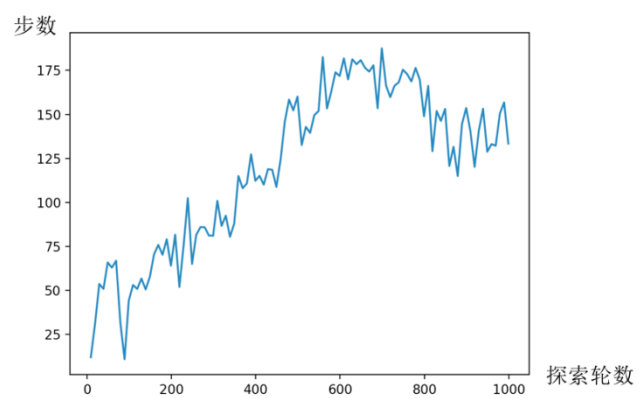


图 12.33 REINFORCE 算法探索轮数与木杆平衡步数的关系

### Actor-Critic 算法

Actor\_Critic ( $S, A, T, R, N, \gamma, \eta_u, \eta_w$ ):

$\mathbf{W}, \mathbf{u}$  = random initial model parameters

For  $i = 1, 2, \dots, N$ :

$\mathbf{s}_{\text{cur}}$  = initial state in  $S$

$t = 0$

While  $\mathbf{s}_{\text{cur}} \notin \{\text{terminal states of } S\}$ :

$a_{\text{cur}} \sim h_{\mathbf{W}}(\mathbf{s}_{\text{cur}})$

$\mathbf{s}_{\text{next}} = T(\mathbf{s}_{\text{cur}}, a_{\text{cur}})$

$\delta = R(\mathbf{s}_{\text{cur}}, a_{\text{cur}}) + \gamma V_{\mathbf{u}}(\mathbf{s}_{\text{next}}) - V_{\mathbf{u}}(\mathbf{s}_{\text{cur}})$

$\mathbf{W} \leftarrow \mathbf{W} + \eta_w \cdot \gamma^t \cdot \delta \cdot \nabla \log h_{\mathbf{W}}(\mathbf{s}_{\text{cur}}, a_{\text{cur}})$

$\mathbf{u} \leftarrow \mathbf{u} + \eta_u \cdot \gamma^t \cdot \delta \cdot \nabla V_{\mathbf{u}}(\mathbf{s}_{\text{cur}})$

$\mathbf{s}_{\text{cur}} = \mathbf{s}_{\text{next}}$

$t \leftarrow t + 1$

图 12.34 Actor-Critic 算法描述

```

1  import numpy as np
2  import gym
3
4  def softmax(scores):
5      e = np.exp(scores)
6      s = e.sum()
7      return e / s
8
9  def actor_critic(env, state_size, n_actions, n_iter, gamma, eta_u, eta_W):
10     W = np.random.rand(state_size, n_actions)
11     u = np.random.rand(state_size, 1)
12     for iter in range(n_iter):
13         state = env.reset()
14         done = False
15         discount = 1
16         steps = 0
17         while not done:
18             steps += 1
19             s_cur = state.reshape(1, state_size)
20             probs = softmax(s_cur.dot(W))
21             a_cur = np.random.choice(n_actions, p = probs.reshape(-1))
22             state, reward, done, info = env.step(a_cur)
23             if done:
24                 print("iteration { } lasts for { } steps".format(iter, steps))
25                 reward = -100
26                 delta = reward - s_cur.dot(u)
27             else:
28                 s_next = state.reshape(1, state_size)
29                 delta = reward + gamma * s_next.dot(u) - s_cur.dot(u)
30             y = np.zeros(n_actions)
31             y[a_cur] = 1
32             y = y.reshape(1, n_actions)
33             gradient_W = s_cur.T.dot(probs - y)
34             W = W - eta_W * discount * delta * gradient_W
35             gradient_u = s_cur.T
36             u = u + eta_u * discount * delta * gradient_u
37             discount *= gamma
38
39     env = gym.make("CartPole-v0")
40     actor_critic(env, 4, 2, 1000, 0.95, 0.1, 0.1)

```

图 12.35 木杆平衡问题的 Actor-Critic 算法

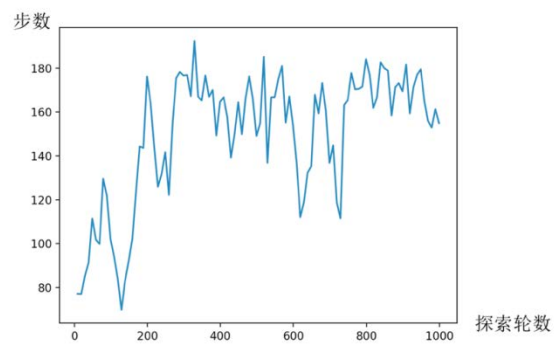


图 12.36 Actor-Critic 算法探索轮数与木杆平衡步数的关系