



Reinforcement Learning (RL) – A Deep Dive

1. What is Reinforcement Learning (RL)?

Reinforcement Learning (RL) is a **learning paradigm** where an **agent** interacts with an **environment** and learns **optimal actions** by maximizing cumulative rewards.

Example Scenario

- Imagine **training a robot to walk**.
- The **agent** (robot) takes **actions** (move left, right, forward).
- The **environment** (ground) responds with **rewards** (+1 for moving forward, -1 for falling).
- The agent **learns a policy** to maximize rewards **over time**.

Key RL Concepts

Concept	Definition
Agent	The decision-maker (e.g., robot, AI playing chess)
Environment	The world the agent interacts with (e.g., chessboard, game, robot's surroundings)
State (s)	A representation of the environment at a given time
Action (a)	A choice the agent makes in a given state
Reward (r)	A signal that tells the agent if an action was good or bad
Policy (π)	The strategy the agent uses to choose actions
Value Function ($V(s)$)	The expected long-term reward from state s
Q-Value ($Q(s, a)$)	The expected reward for taking action a in state s

2. Markov Decision Processes (MDP) – The RL Framework

What is an MDP?

An MDP (Markov Decision Process) provides the **mathematical foundation** for RL. It consists of:

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$$

where:

- $\mathcal{S} \rightarrow$ Set of **states**
- $\mathcal{A} \rightarrow$ Set of **actions**
- $P(s'|s, a) \rightarrow$ **Transition probability** (probability of moving to state s' from s after action a)
- $R(s, a) \rightarrow$ **Reward function** (immediate reward for taking action a in state s)
- $\gamma \rightarrow$ **Discount factor** (how much future rewards matter)

Key MDP Properties

- **Markov Property:** The future state **only depends on the current state and action**, not past history.
 - **Goal:** Find an **optimal policy** $\pi^*(s)$ that **maximizes cumulative rewards** over time.
-

3. Policy & Value Functions

What is a Policy ($\pi(s)$)?

A **policy** is the agent's strategy for choosing actions.

- **Deterministic Policy:** Always picks the best action.

$$\pi(s) = a$$

- **Stochastic Policy:** Picks actions based on probabilities.

$$\pi(a|s) = P(a|s)$$

Value Functions

- **State-Value Function ($V(s)$):** Expected total reward **starting from state s** .

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, \pi \right]$$

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

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

Bellman Equation (Recursion)

The **Bellman equation** expresses $V(s)$ in terms of **future rewards**:

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) [R(s, a) + \gamma V^\pi(s')]$$

This lets us compute **optimal policies** efficiently.

4. Q-Learning: Learning Without a Model

Q-Learning is a **model-free RL algorithm** that learns **Q-values** without knowing the environment dynamics.

Q-Learning Update Rule

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

where:

- α = learning rate.
- r = reward.
- γ = discount factor.

Deep Q-Learning (DQN)

Deep Q-Networks (DQN) use **neural networks** to approximate $Q(s, a)$ when the state space is large.

- **Replace Q-table with a deep neural network.**
 - **Uses Experience Replay** → Stores past experiences and learns from random samples.
 - **Target Network Stabilization** → Uses two networks to improve stability.
-

5. Policy Gradient Methods

Instead of learning $Q(s, a)$, **policy gradient methods** directly optimize the policy.

Policy Gradient Loss

$$J(\theta) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \log \pi_{\theta}(a_t | s_t) \right]$$

This is optimized using **stochastic gradient ascent**.

Actor-Critic Methods

- **Actor** updates the policy $\pi_{\theta}(s)$.
 - **Critic** estimates $V(s)$ to guide learning.
-

6. Fully Commented Manual Implementation of Q-Learning

python

Copy

Edit

```
import numpy as np
import random

class QLearningAgent:
    def __init__(self, n_states, n_actions, alpha=0.1, gamma=0.9, epsilon=0.1):
        self.q_table = np.zeros((n_states, n_actions)) # Initialize Q-table
        self.alpha = alpha # Learning rate
        self.gamma = gamma # Discount factor
        self.epsilon = epsilon # Exploration rate

    def choose_action(self, state):
        if random.uniform(0, 1) < self.epsilon: # Exploration
            return random.randint(0, self.q_table.shape[1] - 1)
        return np.argmax(self.q_table[state]) # Exploitation
```



```
def update(self, state, action, reward, next_state):  
    best_next_action = np.argmax(self.q_table[next_state]) # Best future action  
    target = reward + self.gamma * self.q_table[next_state, best_next_action]  
    self.q_table[state, action] += self.alpha * (target - self.q_table[state, action])
```

```
# Example: 5 states, 2 actions per state
agent = QLearningAgent(n_states=5, n_actions=2)

# Simulate an experience
state, action, reward, next_state = 0, 1, 10, 3
agent.update(state, action, reward, next_state)

# Print updated Q-table
print(agent.q_table)
```

```

1 """
2 CS 229 Machine Learning
3 Question: Reinforcement Learning – The Inverted Pendulum
4 """
5 from __future__ import division, print_function
6 from env import CartPole, Physics
7 import matplotlib.pyplot as plt
8 import numpy as np
9 from scipy.signal import lfilter
10
11 """
12 Parts of the code (cart and pole dynamics, and the state
13 discretization) are inspired from code available at the RL repository
14 http://www-anw.cs.umass.edu/rlr/domains.html
15
16 Briefly, the cart-pole system is described in `cart_pole.py`. The main
17 simulation loop in this file calls the `simulate()` function for
18 simulating the pole dynamics, `get_state()` for discretizing the
19 otherwise continuous state space in discrete states, and `show_cart()`
20 for display.
21
22 Some useful parameters are listed below:
23
24 `NUM_STATES`: Number of states in the discretized state space
25 You must assume that states are numbered 0 through `NUM_STATES` - 1. The
26 state numbered `NUM_STATES` - 1 (the last one) is a special state that
27 marks the state when the pole has been judged to have fallen (or when
28 the cart is out of bounds). However, you should NOT treat this state
29 any differently in your code. Any distinctions you need to make between
30 states should come automatically from your learning algorithm.
31
32 After each simulation cycle, you are supposed to update the transition
33 counts and rewards observed. However, you should not change either
34 your value function or the transition probability matrix at each
35 cycle.
36
37 Whenever the pole falls, a section of your code below will be
38 executed. At this point, you must use the transition counts and reward
39 observations that you have gathered to generate a new model for the MDP
40 (i.e. transition probabilities and state rewards). After that, you
41 must use value iteration to get the optimal value function for this MDP
42 model.

```

```
43
44 `TOLERANCE`: Controls the convergence criteria for each value iteration
45 run. In value iteration, you can assume convergence when the maximum
46 absolute change in the value function at any state in an iteration
47 becomes lower than `TOLERANCE`.
48
49 You need to write code that chooses the best action according
50 to your current value function, and the current model of the MDP. The
51 action must be either 0 or 1 (corresponding to possible directions of
52 pushing the cart)
53
54 Finally, we assume that the simulation has converged when
55 `NO_LEARNING_THRESHOLD` consecutive value function computations all
56 converged within one value function iteration. Intuitively, it seems
57 like there will be little learning after this, so we end the simulation
58 here, and say the overall algorithm has converged.
59
60
61 Learning curves can be generated by calling a code snippet at the end
62 (it assumes that the learning was just executed, and the array
63 `time_steps_to_failure` that records the time for which the pole was
64 balanced before each failure is in memory). `num_failures` is a variable
65 that stores the number of failures (pole drops / cart out of bounds)
66 till now.
67
68 Other parameters in the code are described below:
69
70 `GAMMA`: Discount factor to be used
71
72 The following parameters control the simulation display; you dont
73 really need to know about them:
74
75 `pause_time`: Controls the pause between successive frames of the
76 display. Higher values make your simulation slower.
77 `min_trial_length_to_start_display`: Allows you to start the display only
78 after the pole has been successfully balanced for at least this many
79 trials. Setting this to zero starts the display immediately. Choosing a
80 reasonably high value (around 100) can allow you to rush through the
81 initial learning quickly, and start the display only after the
82 performance is reasonable.
83 """"
```

```

114 def initialize_mdp_data(num_states):
115     """
116     Return a variable that contains all the parameters/state you need for your MDP.
117     Feel free to use whatever data type is most convenient for you (custom classes, tuples, dicts, etc)
118
119     Assume that no transitions or rewards have been observed.
120     Initialize the value function array to small random values (0 to 0.10, say).
121     Initialize the transition probabilities uniformly (ie, probability of
122     |   transitioning for state x to state y using action a is exactly
123     |   1/num_states).
124     Initialize all state rewards to zero.
125
126     Args:
127     |   num_states: The number of states
128
129     Returns: The initial MDP parameters
130     """
131     transition_counts = np.zeros((num_states, num_states, 2))
132     transition_probs = np.ones((num_states, num_states, 2)) / num_states
133     #Index zero is count of rewards being -1 , index 1 is count of total num state is reached
134     reward_counts = np.zeros((num_states, 2))
135     reward = np.zeros(num_states)
136     value = np.random.rand(num_states) * 0.1
137
138     return {
139         'transition_counts': transition_counts,
140         'transition_probs': transition_probs,
141         'reward_counts': reward_counts,
142         'reward': reward,
143         'value': value,
144         'num_states': num_states,
145     }
146
147 def sample_random_action():
148     return 0 if np.random.uniform() < 0.5 else 1
149
150

```

```

121 def choose_action(state, mdp_data):
122     """
123     Choose the next action (0 or 1) that is optimal according to your current
124     mdp_data. When there is no optimal action, return a random action using
125     sample_random_action.
126
127     Args:
128         state: The current state in the MDP
129         mdp_data: The parameters for your MDP. See initialize_mdp_data.
130
131     Returns:
132         0 or 1 that is optimal according to your current MDP
133     """
134
135     # *** START CODE HERE ***
136     pi0 = np.sum(mdp_data['transition_probs'][state, :, 0] * mdp_data['value'][:])
137     pi1 = np.sum(mdp_data['transition_probs'][state, :, 1] * mdp_data['value'][:])
138     if pi0 == pi1:
139         return sample_random_action()
140     else:
141         if pi0 > pi1:
142             return 0
143         else:
144             return 1
145
146     # *** END CODE HERE ***
147
148 def update_mdp_transition_counts_reward_counts(mdp_data, state, action, new_state, reward):
149     """
150     Update the transition count and reward count information in your mdp_data.
151     Do not change the other MDP parameters (those get changed later).
152
153     Record the number of times `state, action, new_state` occurs.
154     Record the rewards for every `new_state`
155     (since rewards are -1 or 0, you just need to record number of times reward -1 is seen in 'reward_counts' index new_state,0)
156     Record the number of time `new_state` was reached (in 'reward_counts' index new_state,1)
157
158     Args:
159         mdp_data: The parameters of your MDP. See initialize_mdp_data.
160         state: The state that was observed at the start.
161         action: The action you performed.

```

```
161         action: The action you performed.
162         new_state: The state after your action.
163         reward: The reward after your action (i.e. reward corresponding to new_state).
164
165     Returns:
166     |     Nothing
167     """
168
169     # *** START CODE HERE ***
170     mdp_data['transition_counts'][state, new_state, action] += 1
171     mdp_data['reward_counts'][new_state, 1] += 1
172     if reward == -1:
173         mdp_data['reward_counts'][new_state, 0] += 1
174     # *** END CODE HERE ***
175
176     # This function does not return anything
177     return
178
```

```

179 def update_mdp_transition_probs_reward(mdp_data):
180     """
181     Update the estimated transition probabilities and reward values in your MDP.
182
183     Make sure you account for the case when a state-action pair has never
184     been tried before, or the state has never been visited before. In that
185     case, you must not change that component (and thus keep it at the
186     initialized uniform distribution).
187
188     Args:
189     |     mdp_data: The data for your MDP. See initialize_mdp_data.
190
191     Returns:
192     |     Nothing
193
194     """
195
196     # *** START CODE HERE ***
197     for i in range(mdp_data['transition_counts'].shape[0]):
198         if mdp_data['reward_counts'][i, 1] == 0:
199             continue
200         else:
201             for j in range(2):
202                 denom = 0
203                 denom = np.sum(mdp_data['transition_counts'][i, :, j])
204                 if denom == 0:
205                     continue
206                 else:
207                     for k in range(mdp_data['transition_counts'].shape[0]):
208                         mdp_data['transition_probs'][i, k, j] = mdp_data['transition_counts'][i, k, j] / denom
209             mdp_data['reward'][i] = mdp_data['reward_counts'][i, 0] / mdp_data['reward_counts'][i, 1]
210
211     # *** END CODE HERE ***
212
213     # This function does not return anything
214     return

```



```
374     if plot:
375         # plot the learning curve (time balanced vs. trial)
376         log_tstf = np.log(np.array(time_steps_to_failure))
377         plt.plot(np.arange(len(time_steps_to_failure)), log_tstf, 'k')
378         window = 30
379         w = np.array([1/window for _ in range(window)])
380         weights = lfilter(w, 1, log_tstf)
381         x = np.arange(window//2, len(log_tstf) - window//2)
382         plt.plot(x, weights[window:len(log_tstf)], 'r--')
383         plt.xlabel('Num failures')
384         plt.ylabel('Log of num steps to failure')
385         plt.savefig('./control.pdf')
386
387     return np.array(time_steps_to_failure)
388
389 if __name__ == '__main__':
390     main()
391
```

```

338 # Recompute MDP model whenever pole falls
339 # Compute the value function V for the new model
340 if new_state == NUM_STATES - 1:
341     update_mdp_transition_probs_reward(mdp_data)
342
343     converged_in_one_iteration = update_mdp_value(mdp_data, TOLERANCE, GAMMA)
344
345     if converged_in_one_iteration:
346         consecutive_no_learning_trials = consecutive_no_learning_trials + 1
347     else:
348         consecutive_no_learning_trials = 0
349
350
351 # Do NOT change this code: Controls the simulation, and handles the case
352 # when the pole fell and the state must be reinitialized.
353 if new_state == NUM_STATES - 1:
354     num_failures += 1
355     if num_failures >= max_failures:
356         break
357     print(['INFO'] Failure number {}).format(num_failures))
358     time_steps_to_failure.append(time - time_at_start_of_current_trial)
359     # time_steps_to_failure[num_failures] = time - time_at_start_of_current_trial
360     time_at_start_of_current_trial = time
361
362     if time_steps_to_failure[num_failures - 1] > min_trial_length_to_start_display:
363         display_started = 1
364
365     # Reinitialize state
366     # x = 0.0
367     x = -1.1 + np.random.uniform() * 2.2
368     x_dot, theta, theta_dot = 0.0, 0.0, 0.0
369     state_tuple = (x, x_dot, theta, theta_dot)
370     state = cart_pole.get_state(state_tuple)
371 else:
372     state = new_state
373

```

```

304 mdp_data = initialize_mdp_data(NUM_STATES)
305
306 # This is the criterion to end the simulation.
307 # You should change it to terminate when the previous
308 # 'NO_LEARNING_THRESHOLD' consecutive value function computations all
309 # converged within one value function iteration. Intuitively, it seems
310 # like there will be little learning after this, so end the simulation
311 # here, and say the overall algorithm has converged.
312
313 consecutive_no_learning_trials = 0
314 while consecutive_no_learning_trials < NO_LEARNING_THRESHOLD:
315
316     action = choose_action(state, mdp_data)
317
318     # Get the next state by simulating the dynamics
319     state_tuple = cart_pole.simulate(action, state_tuple)
320     # x, x_dot, theta, theta_dot = state_tuple
321
322     # Increment simulation time
323     time = time + 1
324
325     # Get the state number corresponding to new state vector
326     new_state = cart_pole.get_state(state_tuple)
327     # if display_started == 1:
328     #     cart_pole.show_cart(state_tuple, pause_time)
329
330     # reward function to use - do not change this!
331     if new_state == NUM_STATES - 1:
332         R = -1
333     else:
334         R = 0
335
336     update_mdp_transition_counts_reward_counts(mdp_data, state, action, new_state, R)
337

```

```

264 def main(plot=True):
265     # Seed the randomness of the simulation so this outputs the same thing each time
266     np.random.seed(0)
267
268     # Simulation parameters
269     pause_time = 0.0001
270     min_trial_length_to_start_display = 100
271     display_started = min_trial_length_to_start_display == 0
272
273     NUM_STATES = 163
274     GAMMA = 0.995
275     TOLERANCE = 0.01
276     NO_LEARNING_THRESHOLD = 20
277
278     # Time cycle of the simulation
279     time = 0
280
281     # These variables perform bookkeeping (how many cycles was the pole
282     # balanced for before it fell). Useful for plotting learning curves.
283     time_steps_to_failure = []
284     num_failures = 0
285     time_at_start_of_current_trial = 0
286
287     # You should reach convergence well before this
288     max_failures = 500
289
290     # Initialize a cart pole
291     cart_pole = CartPole(Physics())
292
293     # Starting `state_tuple` is (0, 0, 0, 0)
294     # x, x_dot, theta, theta_dot represents the actual continuous state vector
295     x, x_dot, theta, theta_dot = 0.0, 0.0, 0.0, 0.0
296     state_tuple = (x, x_dot, theta, theta_dot)
297
298     # `state` is the number given to this state, you only need to consider
299     # this representation of the state
300     state = cart_pole.get_state(state_tuple)
301     # if min_trial_length_to_start_display == 0 or display_started == 1:
302     #     cart_pole.show_cart(state_tuple, pause_time)
303

```

```

237 # *** START CODE HERE ***
238 flag = False
239 flag_convergence = False
240 cont_convergence = 0
241 deltas = np.zeros(mdp_data['transition_counts'].shape[0])
242
243 while flag_convergence == False:
244     cont_convergence += 1
245     for i in range(mdp_data['transition_counts'].shape[0]):
246         value = mdp_data['value'][i]
247         sum0 = np.sum(mdp_data['transition_probs'][i, :, 0] * mdp_data['value'][:])
248         sum1 = np.sum(mdp_data['transition_probs'][i, :, 1] * mdp_data['value'][:])
249         V0 = mdp_data['reward'][i] + gamma * sum0
250         V1 = mdp_data['reward'][i] + gamma * sum1
251         mdp_data['value'][i] = max(V0, V1)
252         delta = abs(mdp_data['value'][i] - value)
253         deltas[i] = delta
254
255     delta_max = np.max(deltas)
256     if delta_max < tolerance:
257         if cont_convergence == 1:
258             flag = True
259             flag_convergence = True
260
261 return flag
262 # *** END CODE HERE ***

```

```
216 def update_mdp_value(mdp_data, tolerance, gamma):
217     """
218     Update the estimated values in your MDP.
219
220
221     Perform value iteration using the new estimated model for the MDP.
222     The convergence criterion should be based on `TOLERANCE` as described
223     at the top of the file.
224
225     Return true if it converges within one iteration.
226
227     Args:
228         mdp_data: The data for your MDP. See initialize_mdp_data.
229         tolerance: The tolerance to use for the convergence criterion.
230         gamma: Your discount factor.
231
232     Returns:
233         True if the value iteration converged in one iteration
234
235     """
```

Final Takeaways

- ✓ RL is about maximizing cumulative rewards by learning an optimal policy.
 - ✓ Q-Learning is a fundamental technique for model-free RL.
 - ✓ Deep Q-Networks (DQN) use neural networks to approximate Q-values.
 - ✓ Policy Gradient methods learn policies directly.
 - ✓ Actor-Critic methods combine value-based and policy-based learning.
-