

Reinforcement Learning

Some definitions

Agent → The agent is the one who takes actions.

Action →  $A = \{a_1, a_2, \dots, a_n\}$  Finite Set of all actions available to agent

Environment → The world through which agent moves. Everything outside agent

$$E(s_t, a_t) = s_{t+1}, r_{t+1}$$

↓      ↓      |      ↓  
current state    current action    next state    reward  
for being in this state

State ( $s$ ) → A concrete and immediate situation which the agent is in. Generally represented as a vector  $s \in \mathbb{R}^d$

Reward ( $r$ ) → Feedback to the agent. Reward can be immediate or delayed.

$s, r \rightarrow$  Random Variables.

Discount Factor ( $\gamma$ ) → How much important future rewards are to the agent.

$\gamma = 1 \rightarrow$  Future reward as important as immediate future. Far-Sighted Agent

$\gamma = 0 \rightarrow$  Don't care about future. Myopic Agent  
Agent looks only one step ahead to maximize reward.  
 $0 \leq \gamma \leq 1$

Example → Hill climbing Algorithm  
(local, greedy, can trap in local sub-optimal solution)

Total Global Reward at time  $t$   
 $\Rightarrow G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} \dots$

$k \rightarrow$  no. of steps in future its reward.

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$G_t \rightarrow$  Sum of infinite terms

→ finite if  $R_{t+k+1}$  is finite and  $\gamma < 1$

Bellman Expectation Equation

Policy ( $\pi$ )  $\rightarrow \pi(s) = a$  strategy that the agent employs to determine the next action based on current state.

Start State ( $s_0$ )  $\rightarrow$  End  $\rightarrow$  horizon  
 Terminal State ( $s_T$ )  $\rightarrow$  End state when only action is "Exit"  
 Episode  $\rightarrow$  One trial run from start to terminal state

$$s_0, a_0, s_1, r_1, a_1, s_2, r_2, a_2, \dots, s_T, r_T, \text{"Exit"}$$

Transition Function  $P(s_{t+1} | s_t, a_t) \rightarrow$  State Transition Probability

$$P(s_{t+1} | s_t, a_t) = \sum_R P(s_{t+1}, R | s_t, a_t)$$

Optimal Policy ( $\pi^*$ )  $\rightarrow$  The goal of RL is to find optimal policy.

$$\pi^* = \arg \max_{\pi} \mathbb{E}[G_T | \pi]$$

where  $G_T \rightarrow$  Total global reward at terminating time  $T$

$$G_T = \sum_{R=0}^{T-1} \gamma^R R(s_K, a_K, s_{K+1})$$

Value Function  $V^\pi(s) : \text{Value (expected return)}$  of being in the state  $s$ .

Expected return when starting from state  $s$  and following policy  $\pi$

$$V^\pi(s) = \mathbb{E}_{\pi} [G_t | s_t = s, \pi] = \mathbb{E} \left[ \sum_{R=0}^{\infty} \gamma^R R(s_t, a_t, s_{t+1}) \mid s_t = s, \pi \right]$$

more specifically if we start at  $t=0$  and terminate at  $t=H \rightarrow \text{Horizon}$ . ②

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{H-1} \gamma^t R(s_t, a_t, s_{t+1}) \mid \pi, s_0 = s \right]$$

Episode  $\rightarrow s_0, a_0, R_1, s_1, a_1, R_2 \dots, R_H, s_H, a_H = \text{"Exit"}$

Optimal state value function

$$V^{\pi^*}(s) = \max_{\pi} V^\pi(s) \quad \forall s \in S$$

Q-value or action-value (Q)

long term return for state  $s$ , taking action  $a$  under policy  $\pi$

$$Q^\pi(s, a) = \mathbb{E}_{\pi} \left[ G_t \mid s_t = s, a_t = a \right]$$

$$= \mathbb{E}_{\pi} \left[ \sum_{t=0}^{H-1} \gamma^t R(s_t, a_t, s_{t+1}) \mid s, a, \pi \right]$$

If  $V^{\pi^*}(s)$  is given starting from  $s$ , optimal policy is

$$\mathbb{E}_{s_{t+1} \sim T(s_{t+1} \mid s_t, a_t)} [V^{\pi^*}(s_{t+1})]$$

choose the ~~next~~ ~~at~~ current action  $a_t$  s.t. the next state  $s_{t+1}$  has max.

$V^{\pi^*}(s_{t+1})$  more computationally extensive

$$V^\pi(s) = \sum_a \pi(a \mid s) Q^\pi(s, a)$$

If  $Q^\pi(s, a)$  is given optimal policy is  $\arg \max_a Q^\pi(s, a)$

$$V^\pi(s) = \max_a Q^\pi(s, a)$$

more memory extensive,

$$Q^{\pi^*}(s, a) = \max_{\pi} Q^{\pi}(s, a)$$

$$Q^{\pi^*}(s, a) =$$

$$\mathbb{E}_{\pi} \left[ R_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t, a_t, \pi \right]$$

# Bellman Equations

## value iteration

$$V^\pi(s) = \mathbb{E}_\pi [G_t | s_t = s]$$

we know

$$G_t = R_{t+1} + \gamma G_{t+1}$$

current reward      future reward

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

$$\mathbb{E}_\pi [R_{t+1} | s_t = s] =$$

$$\sum_a \pi(a | s_t = s) \mathbb{E}_{s'} [\dots]$$

$$\sum_{s'} p(s_{t+1} = s' | s_t = s, a = a) = \sum_{s'} p(s' | s, a) \gamma \sum_{a'} \mathbb{E}_{\pi} [\dots]$$

$$\mathbb{E}_\pi [\gamma G_{t+1} | s_t = s] =$$

$$\sum_a \pi(a | s_t = s)$$

$$\sum_{s'} p(s_{t+1} = s' | s, a) \mathbb{E}_\pi [\dots]$$

$$V^\pi(s) = \sum_a \pi(a | s)$$

$$\sum_{s'} p(s' | s, a)$$

$$[R(s, a, s') +$$

$$\gamma V^\pi(s')]$$

$$V_t^\pi(s) = \max_a \sum_{s'} p(s' | s, a)$$

$$[R(s, a, s') + \gamma V_{t-1}^\pi(s')]$$

## Q-value Iteration

$$Q^\pi(s, a) = \mathbb{E}_\pi [G_t | s_t = s, a_t = a]$$

$$Q^\pi(s, a) = \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | s_t = s, a_t = a]$$

$$\mathbb{E}_\pi [R_{t+1} | s_t, a_t] = \sum_{s'} p(s' | s, a) R(s, a, s')$$

$$\mathbb{E}_\pi [\gamma G_{t+1} | s, a] = \sum_{s'} p(s' | s, a) \gamma \mathbb{E}_\pi [\dots]$$

$$= \sum_{s'} p(s' | s, a) \gamma \sum_{a'} \mathbb{E}_\pi [\dots]$$

$$\gamma \mathbb{E}_\pi [R_{t+2} | s', a']$$

$$Q^\pi(s, a) = \sum_{s'} p(s' | s, a) [R(s, a, s') + \gamma \mathbb{E}_\pi [R_{t+2} | s', a']]$$

$$\sum_{a'} \mathbb{E}_\pi [\dots]$$

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

$$\max_{a'} Q_t^\pi(s', a')$$

$$Q^\pi(s, a) = \sum_{s'} p(s' | s, a) [R(s, a, s') + \gamma \max_{a'} Q_t^\pi(s', a')]$$

## Bellman Equations

(3)

$$V^\pi(s) = \mathbb{E}_\pi [R_{t+1}(s, a, s') + \gamma G_{t+1} | s, a]$$

$$V_k^\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|a,s)$$

$$+ \sum_{s'} [R(s, a, s') + \gamma V_{k-1}^\pi(s')]$$

$$\hookrightarrow \mathbb{E}_\pi [R_{t+1}(s, a, s')] = \sum_a \pi(a|s) \sum_{s'} p(s'|a,s) R(s, a, s')$$

value iteration

we can also have:

$$V^\pi(s) = \mathbb{E}_\pi [R_{t+1}(s, a, s') + \gamma G_{t+1} | s, a]$$

$$V^\pi(s) = R(s, a, s') + \gamma \sum_{s'} p(s'|a,s) V^\pi(s')$$

$$V^\pi(s) = (I - \gamma T^\pi)^{-1} R_{t+1}(s, a, s')$$

$$V^\pi(s) = (I - \gamma T^\pi)^{-1} R^\pi$$

dxd  
Transition Matrix

d → no. of states

Linear system  
of equations

## Value Iteration Algorithm

1. Initialize  $V_0(s) = 0 \forall s \in S$

2. For  $t = 1 \dots H$  (one episode)

For all  $s \in S$

For all  $a$  possible at state  $s$

$$\text{update value } V_t^\pi(s) \leftarrow \max_a \sum_{s'} p(s'|s,a) [R(s,a,s') + \gamma V_{t-1}^\pi(s')]$$

update policy  $\pi_t^\pi(s) \leftarrow \text{that max action } a$

Can be solved using Dynamic Programming

Note → for each time step, update values of all states, update policy. Repeat time steps till convergence.

## V-value iteration

$$V^\pi_t(s) = \max_a \sum_{s'} p(s'|a,s) [R(s,a,s') + \gamma V_{t+1}^\pi(s')]$$

## Q-value iteration

$$Q_t^\pi(s,a) = \sum_{s'} p(s'|a,s) [R(s,a,s') + \gamma \max_{a'} Q_{t+1}^\pi(s',a')]$$

## Policy Iteration

① We start with an initial policy  $\pi_t$

② we find values until convergence.

$\xrightarrow{\text{value update}} V_i^{\pi_t}(s) = \max_a \sum_{s'} p(s'|a,s) [R(s,a,s') + \gamma V_{i-1}^{\pi_t}(s')]$

(loop over  $s \in S$ )

③ Update policy based on  $V_i^{\pi_t}(s)$  using one step look ahead.

$\xrightarrow{\text{one policy update}} \pi_{t+1}(s) \leftarrow \arg \max_a \sum_{s'} p(s'|s,a) [R(s,a,s') + \gamma V_{i-1}^{\pi_t}(s')]$

④ Repeat until policy converges

\* Dynamic Programming  
 If we know  $p(s'|s,a)$ , the transition function and  $R(s,a,s')$ , the reward system, we can use dynamic programming to compute  $\pi^*(s)$  (optimal policy)

Policy and value iterations are iterative  
 hence, intermediate outputs can be saved to reduce computational complexity.

# Reinforcement Learning

## Model Free Approach

Goal → Find optimal policy

By finding optimal value functions

$$V^{\pi}(s), Q^{\pi}(s, a)$$

condition →  $p(s' | s, a)$  Transition Function  
 $R(s, a, s')$  Reward → UNKNOWN

Learn by interacting with Environment.

2 approaches.

- ① Monte Carlo Approach
- ② Temporal Difference Approach.

### Monte Carlo Approach

Few concepts

1. Incremental mean:

$$\mu_k = \frac{1}{k} \sum_{j=1}^k x_j \quad \text{mean of } j \text{ terms}$$

$$= \frac{1}{k} (x_k + \sum_{j=1}^{k-1} x_j)$$

$$= \frac{1}{k} (x_k + (k-1) \mu_{k-1})$$

$$= \mu_{k-1} + \frac{x_k}{k} + \mu_{k-1} - \frac{1}{k} \mu_{k-1}$$

$$\mu_k = \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

$$\text{new Mean} = \text{old mean} + \frac{1}{\text{Total no. of terms now}} (\text{new value} - \text{old mean})$$

$$\mu_k = f(\mu_{k-1}, k, x_k)$$

$$V_t^{\pi}(s) = V_t^{\pi}(s) + \frac{1}{N(s)} (G_t - V_t^{\pi}(s))$$

total reward at step  $t$

## Monte Carlo Approach

$$* V^\pi(s) = \underset{\pi}{\mathbb{E}} [G_t | S_t = s]$$

we learn  $V^\pi(s)$  by taking average over  $G_t$

\* Evaluated over one each episode.

### Algorithm

Random initialize  $V(s) \forall s \in \mathcal{S}$

1. Loop forever (for each episode)

2. Generate an episode ( $E$ ) from  $S_0 \rightarrow S_T$

$$G_t = 0$$

3. For each  $t$  from  ~~$T-1 \rightarrow 0$~~

$$G_t = \gamma G_t + R_{t+1}$$

we can  
get reward  
from first occurrence  
or every occurrence  
of state  $s$  in  $E$

for  $s$  in  $E$  (state visited in  
current episode)  
update  $N(s_t) \leftarrow N(s_t) + 1$

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t))$$

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t))$$

For non-stationary  
problems.

$\bar{\alpha}$   
remember  
rate  
Learning Rate

### Key Points

\* Value is estimated as mean over  $G_t$

\* we loop over those  $s$  which are in  $E$ , and not

over all possible  $s$  as in dynamic programming

\* Above is value iteration, we can also do  
policy iteration on control tasks.

$$* \text{For Q-value} \\ N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G_t - Q(s_t, a_t))$$

\* First generate  $E$  and then evaluate.

# (5)

## Temporal Difference Approach (T-D)

MC + DP — Dynamic Programming

Monte Carlo

- \* Don't wait till end of episode ( $E$ ) to update the expected future reward estimation ( $V$ )
- \* Wait only till next step.

In MC

$$V^\pi(s) = \mathbb{E}_\pi [G_t | s_t = s]$$

$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$

Actual Return — computed at end of episode.

In TD & DP

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi [G_t | s_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | s_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s] \end{aligned}$$

So,

$$V^\pi(s_t) \leftarrow V^\pi(s_t) + \alpha (R_{t+1} + \gamma V^\pi(s_{t+1}) - V^\pi(s_t))$$

learning rate

An estimate of return

(in MC it was  $G_t$  → Actual Return)

$\downarrow$

Simplest TD(0)  
look ahead one-step

$$Q^\pi(s_t, a_t) \leftarrow Q^\pi(s_t, a_t) + \alpha (R_{t+1} + \gamma \max_{a \in A} Q^\pi(s_{t+1}, a) - Q^\pi(s_t, a_t))$$

We can do more steps look ahead

example ↴

$$\alpha [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 \max_{a \in A} Q^\pi(s_{t+3}, a) - Q^\pi(s_t, a_t)]$$

Algorithm (On-Policy method) SARSA

1. Random initialize  $Q(s, a)$  & say,  $\pi \in A$
  2. loop forever (for each episode)
  3. for each step in episode
    4. choose  $a$  from  $s$  in  $E$  using policy derived from  $Q$
    - Take action ' $a$ ' & observe  $R(s, a, s')$
    - $Q(s, a) \leftarrow Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$
- $s \leftarrow s'$       choose action based on epsilon same greedy method policy current  $Q(s, a)$

Q-learning (off-Policy method)

Agent may not follow the greedy policy.

Agent may not choose  $\max_{a'} Q(s', a')$  as next action (Exploration)

Agent can choose random action (Exploration)

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



choose action based on epsilon greedy method policy

## Reinforcement Learning

(16)

Don't know  $p(s'|a, s)$  (Transition matrix)  
and/or  $R(s, a, s')$  Rewards.

### Q-Learning

Here, we learn Q-Table

Rows  $\rightarrow$  state ( $s$ )

Columns  $\rightarrow$  action ( $a$ )

$$Q(s, a) = \mathbb{E} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t, a_t]$$

expected discounted cumulative award given current state and action

Q-value for a given state & action pair

\* we don't know policy ( $\pi$ )

In Q-learning, we explore the environment to learn Q-Table (update)

Update Q-Table

$$Q(s, a) = Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

new Q-value

learning rate  
 reward for taking current action  
 current Q-value  
 discount factor

Highest Q-value we can achieve in state  $s'$

### Deep Q-learning

$$Q(s, a) = (1-\alpha) Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a')]$$

$$\text{Target} = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

$$Q(s, a) = (1-\alpha) Q(s, a) + \alpha \text{Target}$$

## Q-Learning Algo

1. Initialize Q-Table  $Q(s, a) = 0 \forall s \in \text{states}, a \in \text{actions}$
2. Sample an action  $a \rightarrow$  How to sample action??
3. Perform action  $a$
4. update  $Q(s, a)$
5. reached state  $s'$ , repeat  $\Rightarrow$

How to Sample action - ??

1. choose  $\arg \max_a Q(s, a) \rightarrow$  Greedy Exploitation
2. choose randomly  $\rightarrow$  Exploration

## Approximate Q-Learning

what if  $Q(s, a)$  is a weighted sum of various features computed over  $s, a$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

$$Q(s, a) = \sum_{i=1}^N w_i f_i(s, a)$$

*approximation of  $Q(s, a)$*

MSE

$$\text{error}(w) = (Q(s, a) - \sum_i w_i f_i(s, a))^2$$

$$\frac{\partial E(w)}{\partial w_m} = - (Q(s, a) - \sum_i w_i f_i(s, a)) f_m(s, a)$$

update weight

$$w_m = w_m + \alpha (Q(s, a) - \sum_i w_i f_i(s, a)) f_m(s, a)$$

$\theta$