Two potential approaches.

① Find the value of state

② Find the value of action

Two RL schemes:

1. Policy iteration [monte carlo, Temporal differencing]

evaluate policy to
improve it

emphasis on Sampling
and sampling as much
as we can

What is improving TD(λ)
(what is now and next)

2. Value iteration [Q-learning]

maximize accumulate
reward for (s,a)

Deal with rewards?

MDP - markov decision process

MDP —┬— keep everything        $Q_K = \dfrac{r_1 + r_2 + \dots + r_k}{k}$
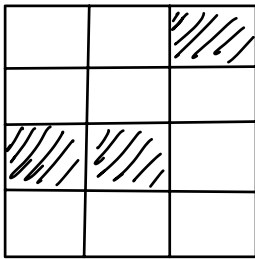      │   (order n)
      └— Do it incrementally

(Avg) $Q_{K+1} = $ (Avg) $Q_K + \dfrac{1}{k+1} \left[ r_{K+1} - Q_K \right]$

Common update rule:

New estimate = Old estimate + step size $*$ [Target - Old estimate]

Caution: State-space problem

elevator system



$$4 \times 4 \times 4 = 64 \longleftarrow \text{total states}$$

2  2  4  $\longrightarrow$ eg: state 62

## Design RL agents :

① Discretize states

② Define actions

③ Determine the reward / punishment

④ Stablish the action policy (taking action)

## How to take actions :

① Random

② Greedy (action with maximum reward)

③ Epsilon $\varepsilon$-greedy (max action with $P = 1 - \varepsilon$
and random with $\varepsilon$, $\varepsilon > 0$)

At the begining $\varepsilon$ is large

④ Softmax action selection

Softmax : We use Gibbs or Boltzmann Distribution

accumulated reward from the table

$$P(a) = \frac{\exp\left(\dfrac{Q(s,a)}{\tau}\right)}{\overline{\sum_{b=1}^{n} \left(\exp Q(s,b) / \tau\right)}} \quad \} \text{ for normalization}$$

(In reality)

The RL agent learns a policy $\pi$.

eg: Grid world

target

$\pi$

robot

need ai when problem is stochastic, not deterministic

exploiting policy, explore

At step $t$, $\pi_t(s,a)$ is the probability that $a_t = a$ when $S_t = S$.

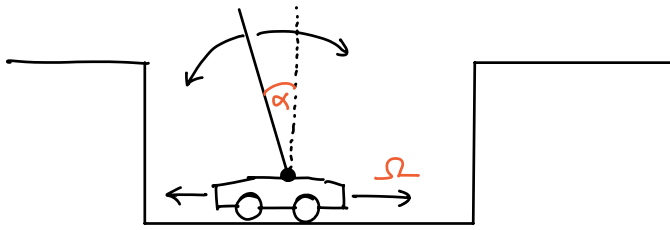RL methods enable a change of policy based on experience.

Maximise reward :

1) Episodic task : Return at $t$, $R_t = r_{t+1} + r_{t+2} + \ldots r_T$

Terminal

2) Continuous task : $R_t = r_{t+1} + \gamma \cdot r_{t+2} + \gamma^2 \cdot r_{t+3} + \dots$

$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \qquad \gamma - \text{discount factor}$$

0 —— $\gamma$ —— 1
short sighted          far sighted

## Inverted pendulum



Avoid failure

1) Poll falls
2) Car hit the wall

Understand as an episodic task.

$r = +1$ (for each step prior to fall)

$R$ = number of steps before fall

Understand as a continuous task

$r = -1$ (for falling), 0 otherwise

$R = -\gamma^k$ for k steps before the fall


Temporal Differencing methods - TD

TD Prediction = Policy evaluation

$\longrightarrow$ Compute the state value function $V^\pi$ for a given policy $\pi$

# Policy evaluation

① Simple every visit Monte-Carlo

$$V(s_t) \longleftarrow V(s_t) + \alpha (\underset{\text{Return}}{\boxed{R_t}} - V(s_t))$$

↖ value of the state

② The simplest TD $\longrightarrow$ TD(0)   estimate of the return

$$V(s_t) \longleftarrow V(s_t) + \alpha [r_{t+1} + \boxed{\gamma V(s_{t+1})} - V(s_t)]$$

## Learning an action value function   estimate

$$Q(s_t, a_t) \longleftarrow Q(s_t, a_t) + \alpha [\underset{\text{reward}}{\boxed{r_{t+1}}} + \boxed{\gamma Q(s_{t+1}, a_{t+1})} - Q(s_t, a_t)]$$

↗ accumulated reward

target
(reward together)



## Q-Learning   design of the agent is difficult, learning algorithm is simple.

(matrix)

① Initialize $Q(s,a)$ randomly.

② for each episode,
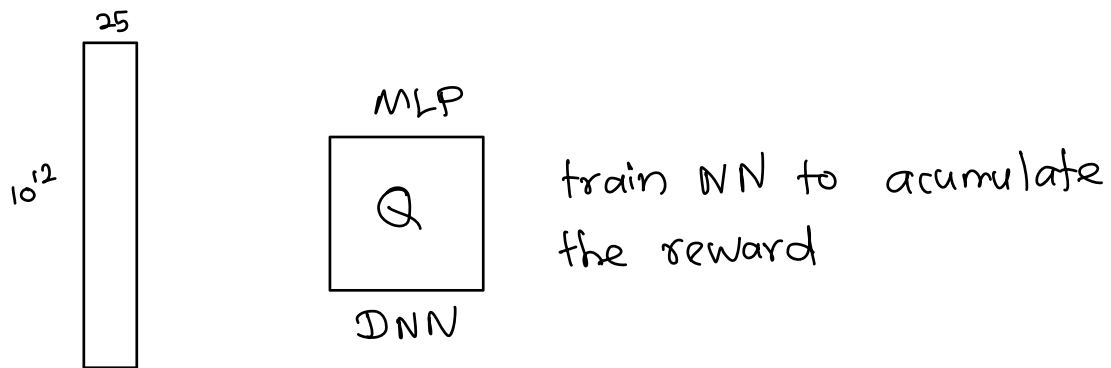
    - initialize state $s$

    - For each step

        - choose $a$ from $s$ using action policy
        - take action $a$

— observe  $r,\ s' = s_{t+1}$

— update  $Q(s,a) \longleftarrow Q(s,a) + \alpha \ldots$

— update  $s \longleftarrow s'$

— until s terminal

## Convergence

Every state has to be visited multiple times.

25

$10^{12}$

MLP

Q

DNN

train NN to acumulate the reward

Deep RL $\longrightarrow$ Use DNN to hold/model Q-matrix