# Lecture 11

Reinforcement learning :
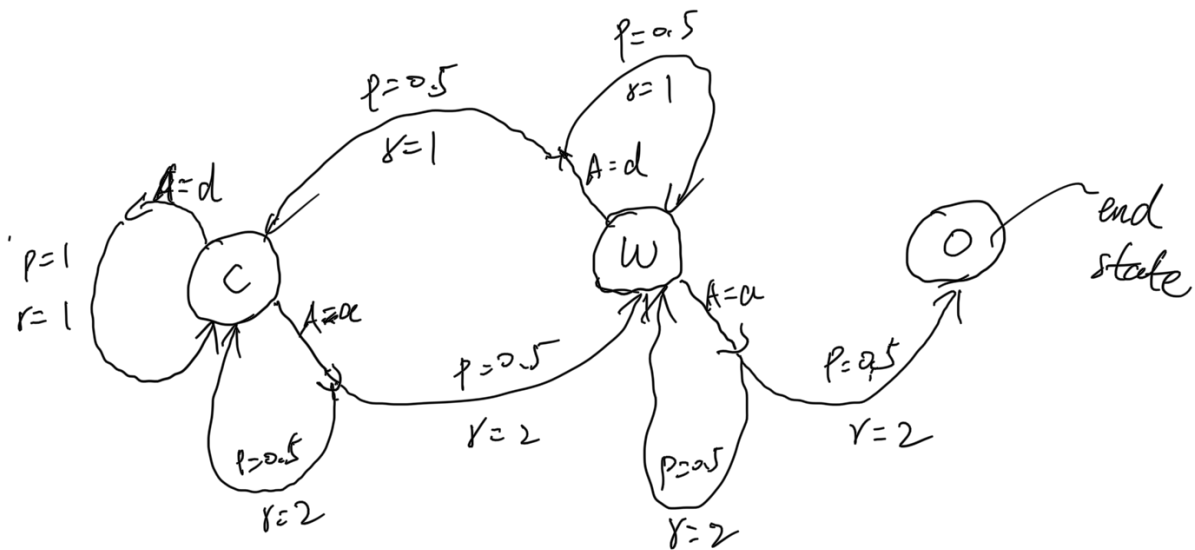
A car (robot car)

States = { cool, warm, overheated } = { c, w, o }.

Actions = { Acc , De-acc } = { $\underline{a}$ , $\underline{d}$ }.

standard time step

reward (A) : $r(ce) = 2$ , $r(d) = 1$ , $r(o) = -10$

$r(A=a) = 2$ , $r(A=d) = 1$ , $r(s=o) = -10$



Markov Decision Process :

⎰ set of state S

⎱ Start state $s_0$

Set of actions A

Transition Prob $P(s'|s,a)$ (or $T(s,a,s')$ )

reward $R(s,a,s')$  ⎰ receive after an action

( reward ( ... ) , ... )                  └ receive after reach ce state
reward discount $\gamma$.

receive a sequence of reward over 3 steps:

   case1 : $r_1 = [1, 2, 3]$

   case 2 : $r_2 = [3, 2, 1]$

without discount : tot. reward $R_1 = 6$
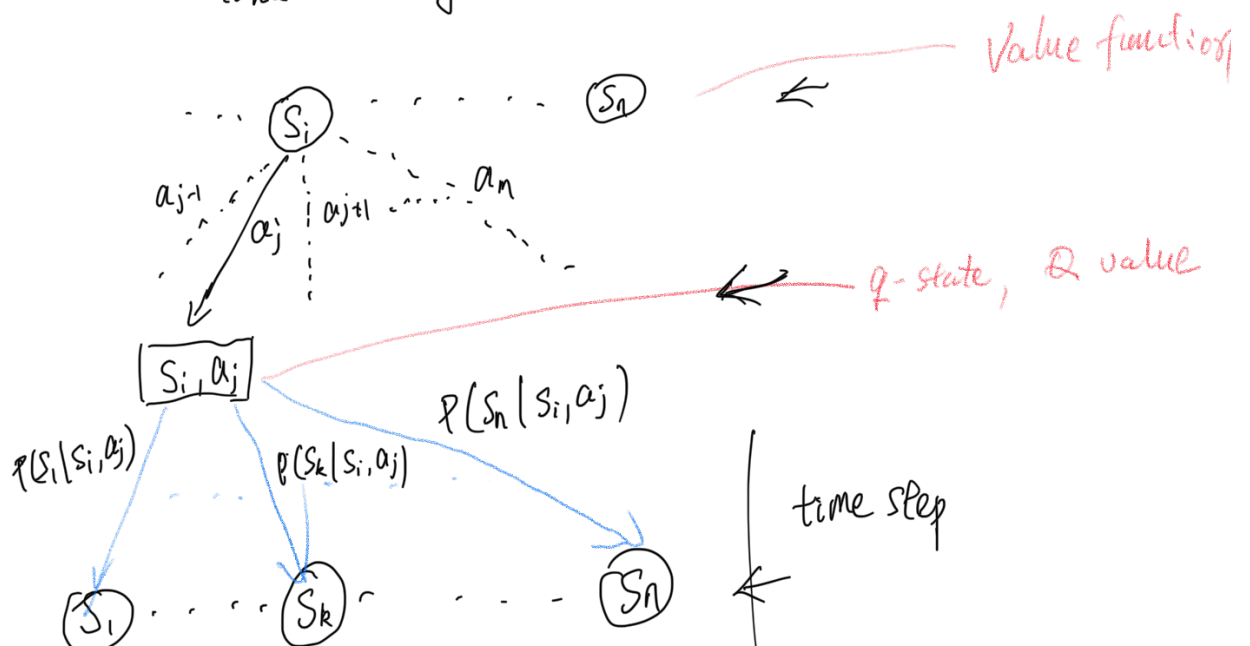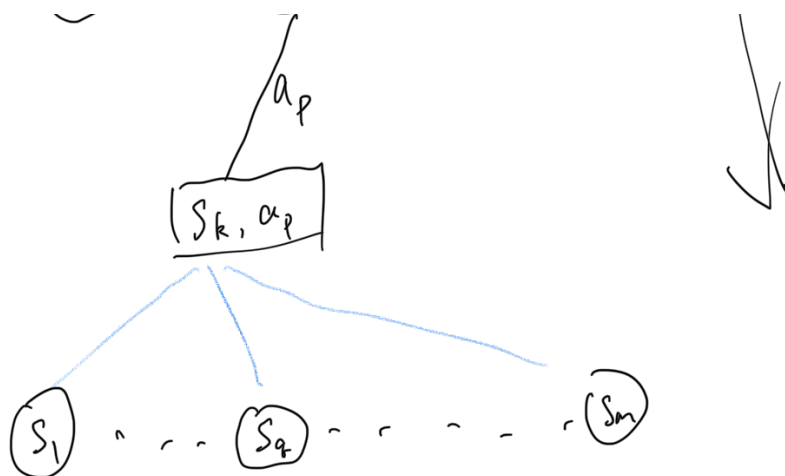$$R_2 = 6$$

discount : 0.5
$$R_1 = 1 + 2 \times 0.5 + 3 \times 0.5 \times 0.5 = 2.75$$
$$R_2 = 3 + 2 \times 0.5 + 1 \times 0.5 \times 0.5 = 4.25$$

Goal : find or generate a policy that can get the max possible sum of discounted rewards.

Policy is a handbook that tells the robot which action to take at a given state.



Value function

$S_i$      $S_n$

$a_{j-1}$    $a_j$   $a_{j+1}$   $a_n$

q-state, Q value

$S_i, a_j$

$P(S_n | S_i, a_j)$

$P(S_i | S_i, a_j)$    $P(S_k | S_i, a_j)$

time step

$S_i$     $S_k$     $S_n$

$a_p$

$[S_k, a_p]$

$(S_1) \quad \cdots \quad (S_q) \quad - \quad - \quad - \quad - \quad (S_m)$

$V(s)$ —— expected all (sum of) future rewards in state $s$.

$Q(s, a)$ —— expected all (sum of) future rewards, after taking action $a$, at state $s$.

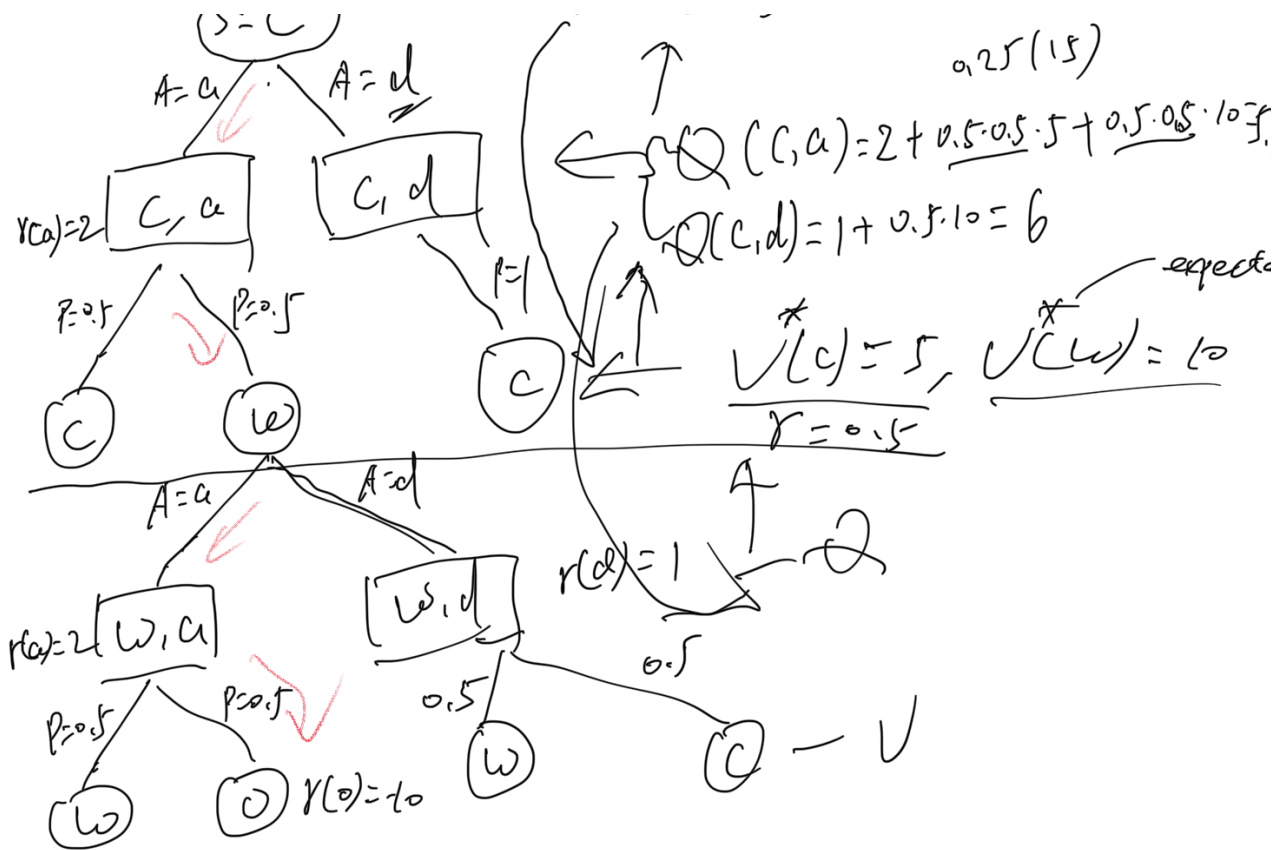$$\begin{cases} V(s) = \max_{a_i} Q(s, a_i) \qquad i = 1 \rightarrow m \\ Q(s, a) = r_a + \sum_{s'} P(s' | s, a) \left( r(s') + \overset{\lambda}{\gamma} V(s') \right) \end{cases}$$

$$V(s) = \max_{a_i} Q(s, a_i) \qquad i = 1 \rightarrow m$$

$\vdots$

Bellman equations.

$S \triangleq \{c, \omega, o\}, \quad A \triangleq \{a, d\}, \quad r(a)=2, \; r(d)=1, \; r(o)=-10.$

$P(c|c, a) = 0.5, \quad P(c|c, d) = 1, \quad P(o|\omega, a) = 0.5, \; P(c|\omega, d) = 0.5$

$P(\omega|c, a) = 0.5, \quad P(\omega|c, d) = 0, \quad P(\omega|\omega, a) = 0.5, \; P(\omega|\omega, d) = 0.5$

$\longleftarrow V(c) = 6$

$(S = C)$

$(S=C)$

$A=a$    $A=d$

$r(a)=2$ [C, a]    [C, d]

$P=0.5$    $P=0.5$    $P=1$

(C)    (W)    (C)

$Q(C,a)=2+0.5\cdot0.5\cdot5+0.5\cdot0.5\cdot10\}$

$Q(C,d)=1+0.5\cdot10=6$

$a.25(15)$

$V^*(c)=5, \quad V^*(W)=10$ — expecta

$\gamma=0.5$

$A=a$    $A=d$

$r(a)=2$ [W, a]    [W, d]    $r(d)=1$

$P=0.5$    $P=0.5$    $0.5$    $0.5$

(W)    (O) $r(o)=10$    (W)    (C) — V

Q — learning algorithm,    step size $\alpha \in [0,1]$

initialize $Q(S,a)$, for $S$, $a$,    $\lambda$ — discount

Loop for each episode:

    choose A at S based on Q-value ($\varepsilon$-greedy)

Take action A, get reward, reach a state $S'$

compute Q.

$\varepsilon$-greedy:

pick a $\varepsilon$ value

$\underset{\text{percentage}}{(1-\varepsilon)}$ of steps: act greed.

$\varepsilon$ of steps: random

$$Q(S,A) \leftarrow Q(S,A) + \alpha\left[ r(a) + \lambda \max_a Q^*(S',a) - Q(S,A) \right]$$

$\underbrace{\quad\quad}_{\text{new Q}}$    $\underbrace{\quad}_{\text{old Q}}$

new Q    $Q = r(a) + \sum_{S'} P\cdot\lambda V(S') + r(S')$

$V(S') = \max Q(S',a)$

$$\text{new } Q \quad Q = r(a) + \sum_{s'} P \lambda \max_a \overset{*}{Q}(s', a) + \underline{r(s')}$$

from $s$, take action $a$, reach $s'$

$$s \Leftarrow s'$$

until $s$ reach end state, or for certain # of steps.

learns $Q(s, a)$

|       | $S_1$ - $S_i$ - $S_n$ |   |   |
|-------|------------------------|---|---|
| $a_1$ | $Q(s_1, a_1)$ |   | $Q(s_n, a_1)$ |
| $\vdots$ |   |   |   |
| $a_m$ | $Q(s_1, a_m)$ |   | $Q(s_n, a_m)$ |