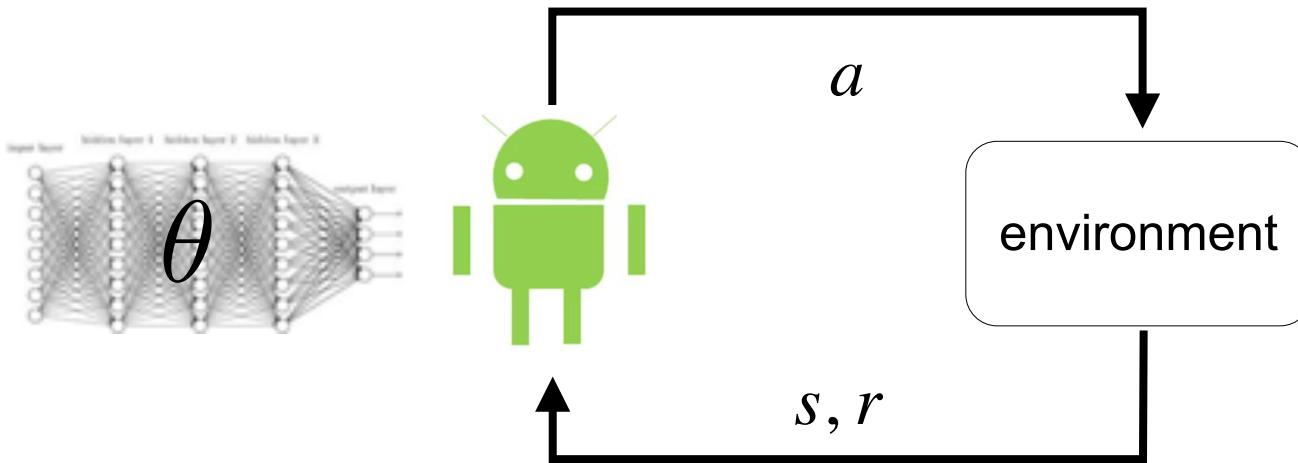


# Lecture 4

# Markov Decision Process

# Reinforcement learning



Agent's goal: learn a policy to maximize the return

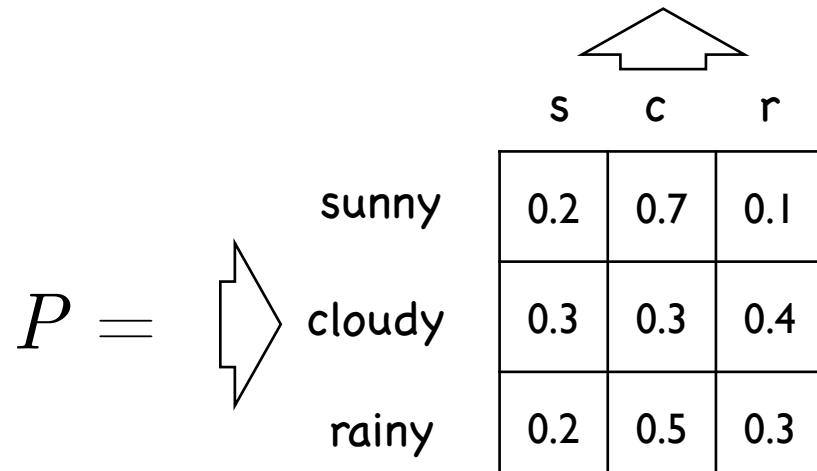
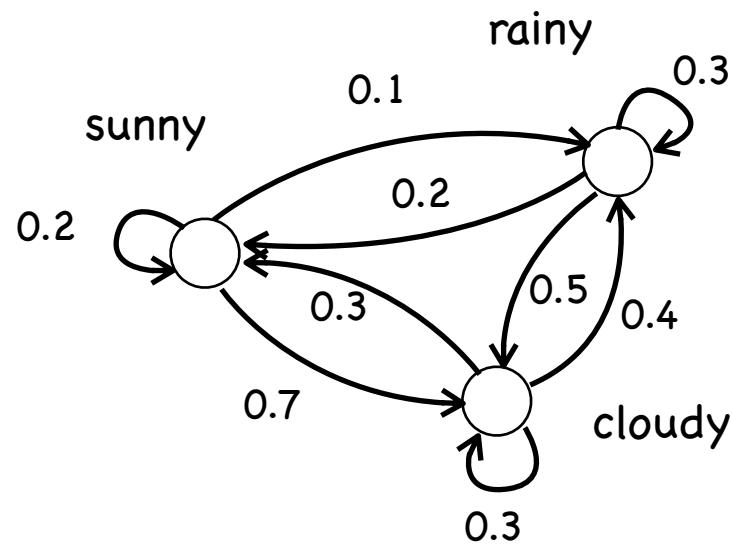
$$\text{T-step: } \sum_{t=1}^T r_t \quad \text{average: } \frac{1}{T} \sum_{t=1}^T r_t \quad \text{discounted: } \sum_{t=1}^{\infty} \gamma^t r_t$$

# Markov Process

(finite) state space  $S$ , transition matrix  $P$

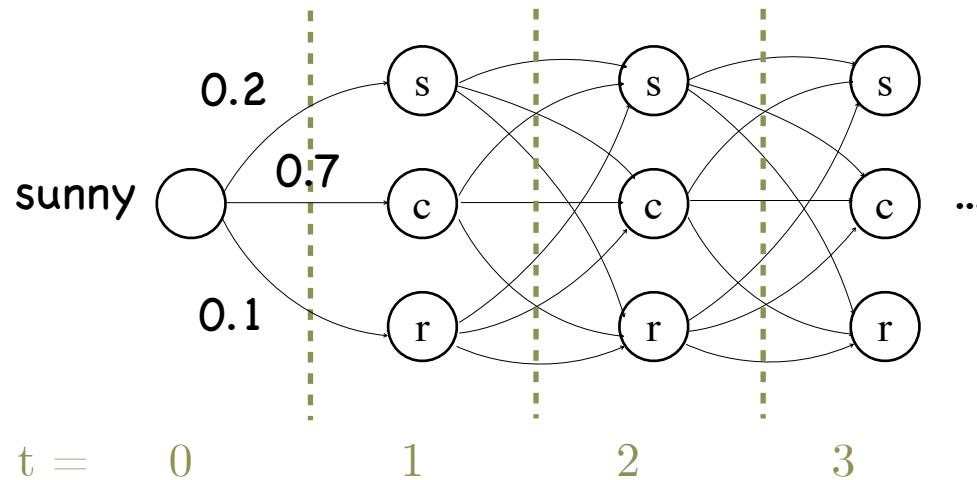
a process  $s_0, s_1, \dots$  is Markov if has no memory

$$P(s_{t+1} \mid s_t, \dots, s_0) = P(s_{t+1} \mid s_t) \quad \text{discrete } S \rightarrow \text{Markov chain}$$



$$s_{t+1} = s_t P = s_0 P^{t+1}$$

horizontal view



stationary distribution:  $s == sP$

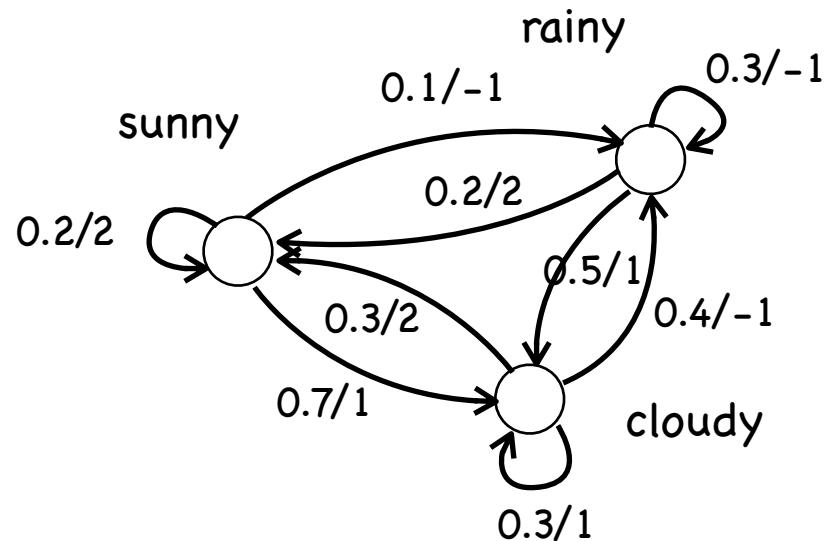
sampling from a Markov process:

s, c, c, r ...

s, c, s, c ...

# Markov Reward Process

introduce reward function  $R$



how to calculate the long-term total reward?

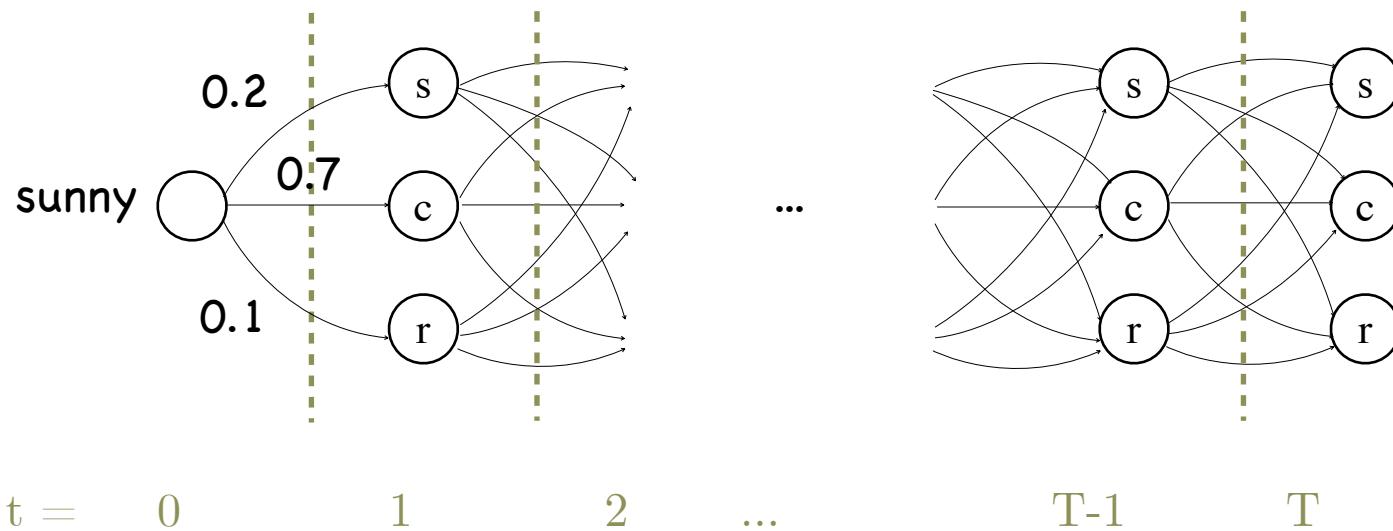
$$V(\text{sunny}) = E\left[\sum_{t=1}^T r_t | s_0 = \text{sunny}\right]$$

$$V(\text{sunny}) = E\left[\sum_{t=1}^{\infty} \gamma^t r_t | s_0 = \text{sunny}\right]$$

value function

# Markov Reward Process

horizontal view: consider T steps



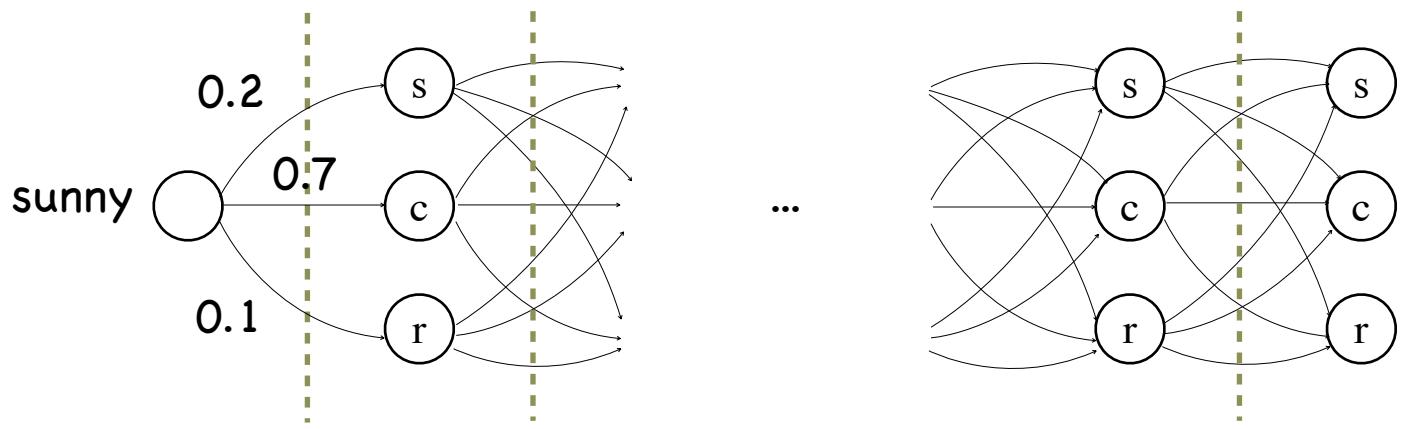
recursive definition:

$$\begin{aligned} V(\text{sunny}) &= P(s|s)[R(s) + V(s)] \\ &\quad + P(c|s)[R(c) + V(c)] \\ &\quad + P(r|s)[R(r) + V(r)] \end{aligned}$$

$$= \sum_s P(s|\text{sunny})(R(s) + V(s))$$

# Markov Reward Process

horizontal view: consider T steps



$t =$  0 1 2 ...  $T-1$   $T$

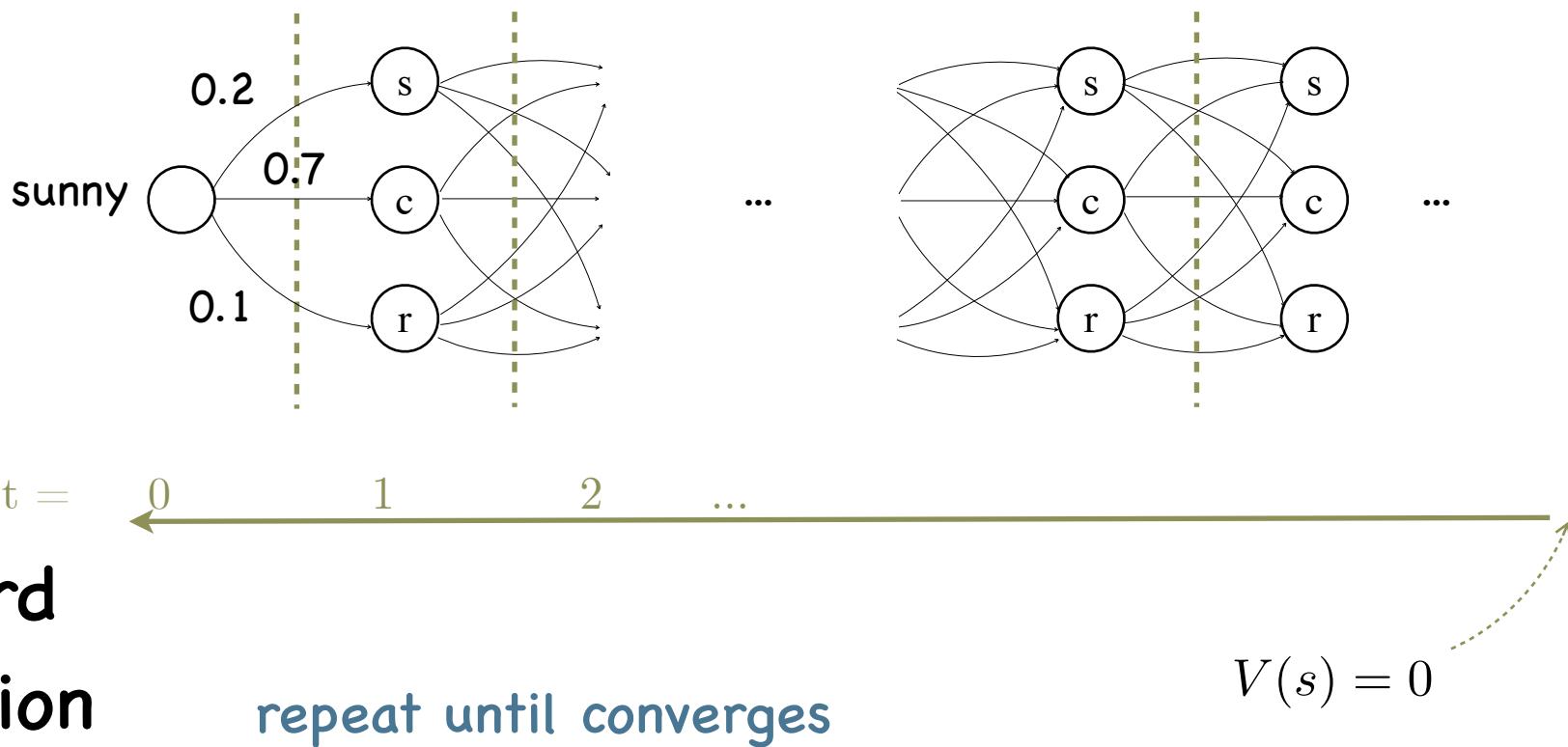
backward  
calculation

$$V(s) = 0$$

$$V(s) = \sum_{s'} P(s'|s) (R(s') + V(s'))$$

# Markov Reward Process

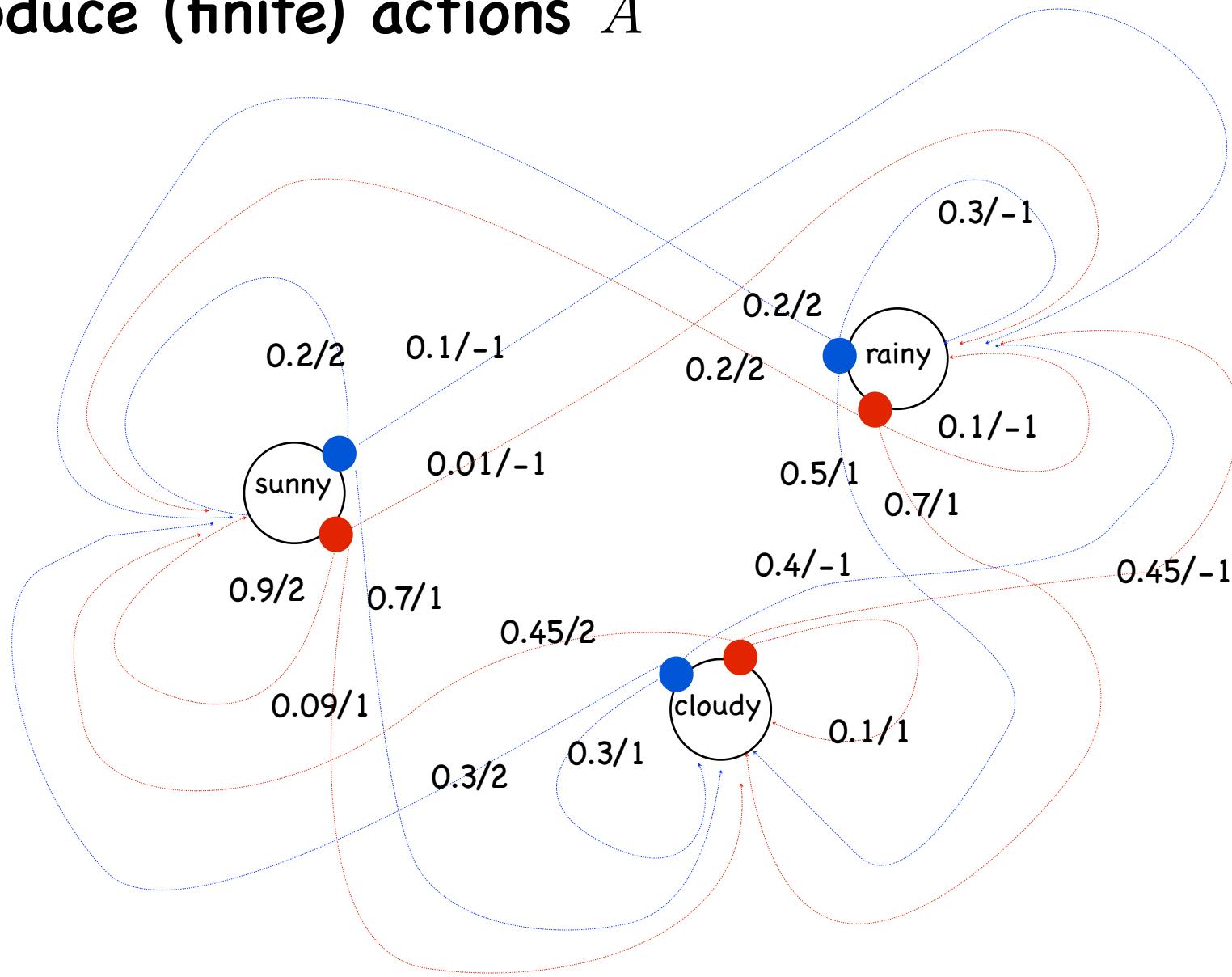
horizontal view: consider discounted infinite steps



$$V(s) = \sum_{s'} P(s'|s) (R(s') + \gamma V(s'))$$

# Markov Decision Process

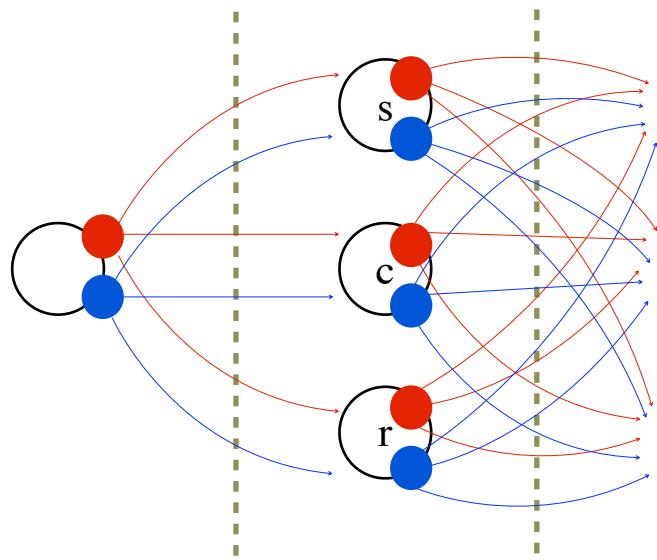
introduce (finite) actions  $A$



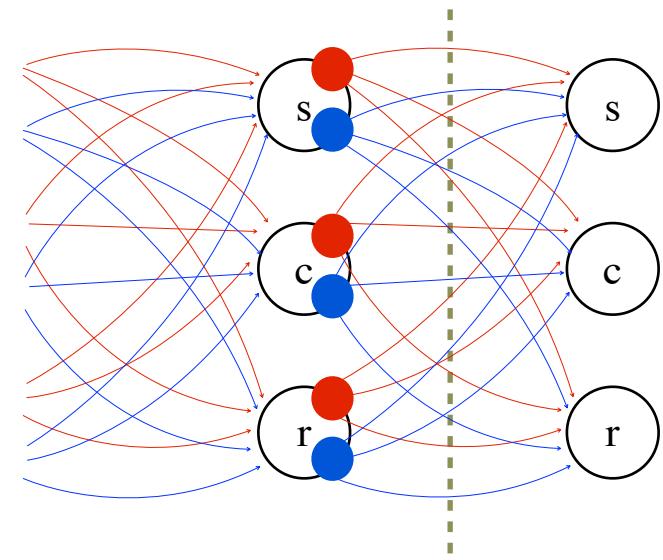
# Markov Decision Process

horizontal view

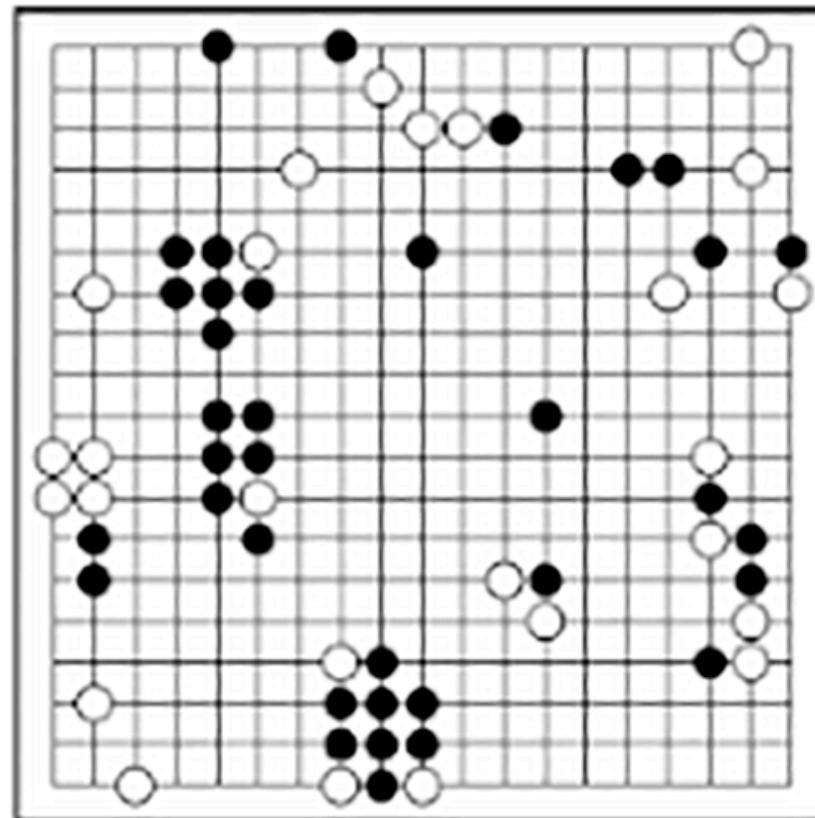
sunny



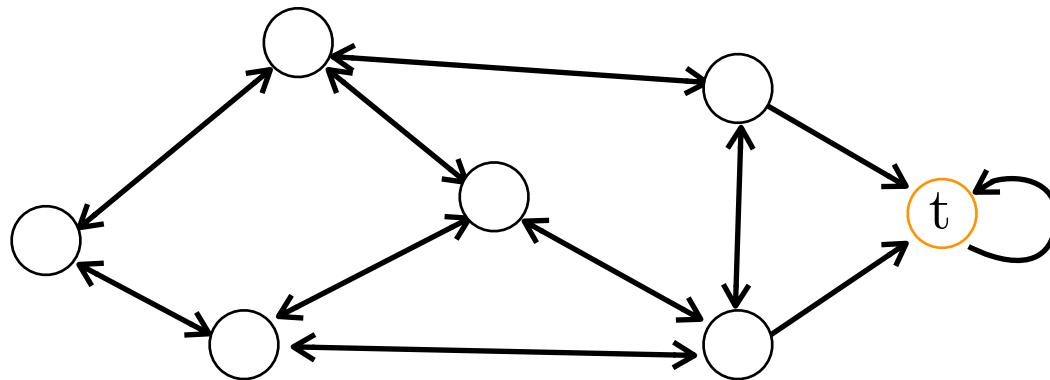
...



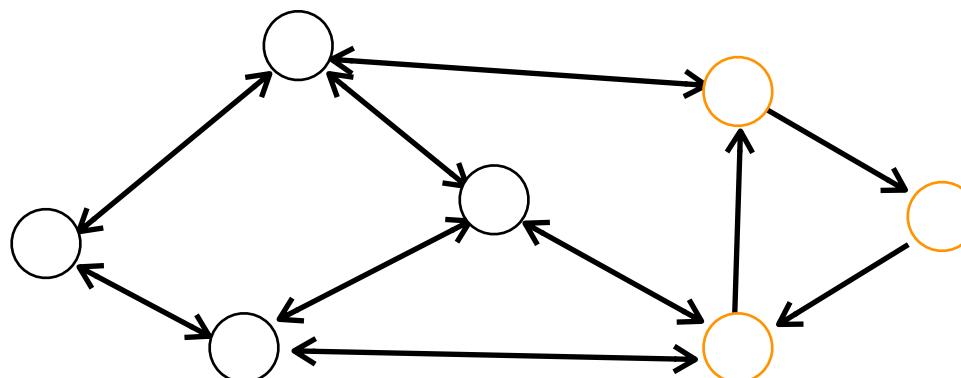
horizontal view of the game of Go



goal-directed



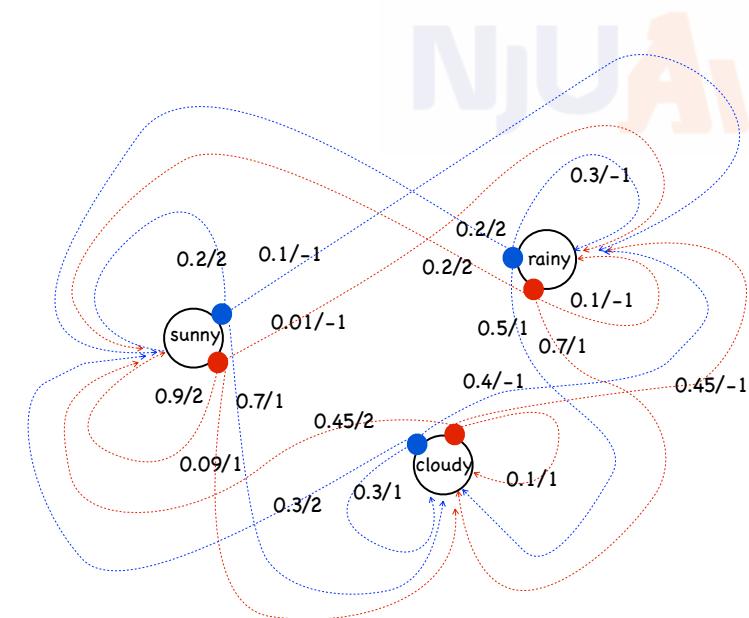
stationary distribution



# Markov Decision Process

**MDP**  $\langle S, A, R, P \rangle$  (often with  $\gamma$ )

essential model for RL  
but not all of RL



**policy**

**stochastic**

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

**deterministic**

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

$|A|^{|S|}$  deterministic policies

**tabular representation**

$\pi =$

	0	0.3
s	1	0.7
	0	0.6
c	1	0.4
	0	0.1
r	1	0.9

# Expected return

how to calculate the expected total reward of a policy?

similar with the Markov Reward Process

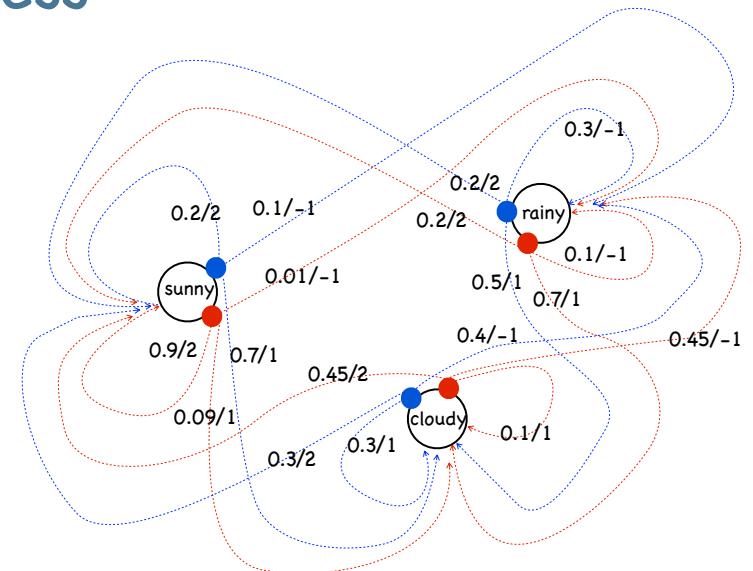
MRP:

$$V(s) = \sum_{s'} P(s'|s)(R(s') + V(s'))$$

MDP:

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

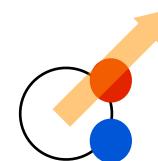
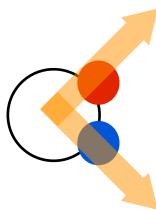
expectation over actions  
with respect to the policy



# Q-function

state value function

$$V^\pi(s) = E\left[\sum_{t=1}^T r_t | s\right]$$



state-action value function

$$Q^\pi(s, a) = E\left[\sum_{t=1}^T r_t | s, a\right] = \sum_{s'} P(s'|s, a) (R(s, a, s') + V^\pi(s'))$$

consequently,

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

Q-function => policy

# Optimality

	0	0.3
s	1	0.7
c	0	0.6
	1	0.4
r	0	0.1
	1	0.9

there exists an optimal policy  $\pi^*$

$$\forall \pi, \forall s, V^{\pi^*}(s) \geq V^\pi(s)$$

optimal value function

$$\forall s, V^*(s) = V^{\pi^*}(s)$$

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

# Bellman optimality equations

	0	0.3
s	1	0.7
c	0	0.6
	1	0.4
r	0	0.1
	1	0.9

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

from the relation between V and Q

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$$

we have

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_a Q^*(s', a))$$

$$V^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$$

the unique fixed point is the optimal value function

# Solve optimal policy in MDP

**idea:**

how is the current policy      **policy evaluation**  
improve the current policy      **policy improvement**

**policy evaluation:**      backward calculation

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

**policy improvement:**      from the Bellman optimality equation

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

# Solve optimal policy in MDP

**policy improvement:** from the Bellman optimality equation

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

let  $\pi'$  be derived from this update

$$\begin{aligned} V^\pi(s) &\leq Q^\pi(s, \pi'(s)) \\ &= \sum_{s'} P(s'|s, \pi'(s))(R(s, \pi'(s), s') + \gamma V^\pi(s')) \\ &\leq \sum_{s'} P(s'|s, \pi'(s))(R(s, \pi'(s), s') + \gamma Q^\pi(s', \pi'(s))) \\ &= \dots \\ &= V^{\pi'} \end{aligned}$$

so the policy is improved

## Policy iteration algorithm:

loop until converges

policy evaluation: calculate V

policy improvement: choose the action greedily

$$\pi_{t+1}(s) = \arg \max_a Q^{\pi_t}(s, a)$$

converges:  $V^{\pi_{t+1}}(s) = V^{\pi_t}(s)$

$$Q^{\pi_{t+1}}(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_a Q^{\pi_t}(s', a))$$

recall the optimal value function about Q

# Solve optimal policy in MDP

embed the policy improvement in evaluation

Value iteration algorithm:

$$V_0 = 0$$

for  $t=0, 1, \dots$

    for all  $s$   $\leftarrow$  synchronous v.s. asynchronous

$$V_{t+1}(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_t(s))$$

    end for

    break if  $\|V_{t+1} - V_t\|_\infty$  is small enough

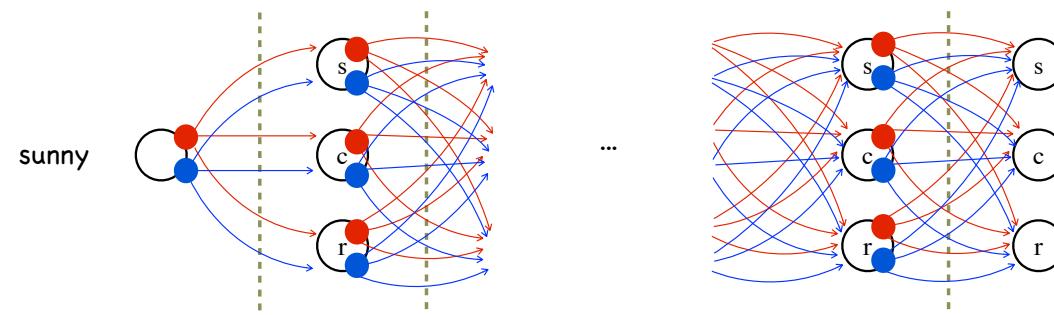
end for

recall the optimal value function about  $V$

# Solve optimal policy in MDP

$$Q^{\pi_{t+1}}(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_a Q^{\pi_t}(s', a))$$

$$V_{t+1}(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_t(s'))$$



Dynamic programming



R. E. Bellman  
1920-1984

Complexity

needs  $\Theta(|S| \cdot |A|)$  iterations to converge on deterministic MDP

[O. Madani. Polynomial Value Iteration Algorithms for Deterministic MDPs. UAI'02]

curse of dimensionality: Go board 19x19,  $|S|=2.08 \times 10^{170}$

[<https://github.com/tromp/golegal>]