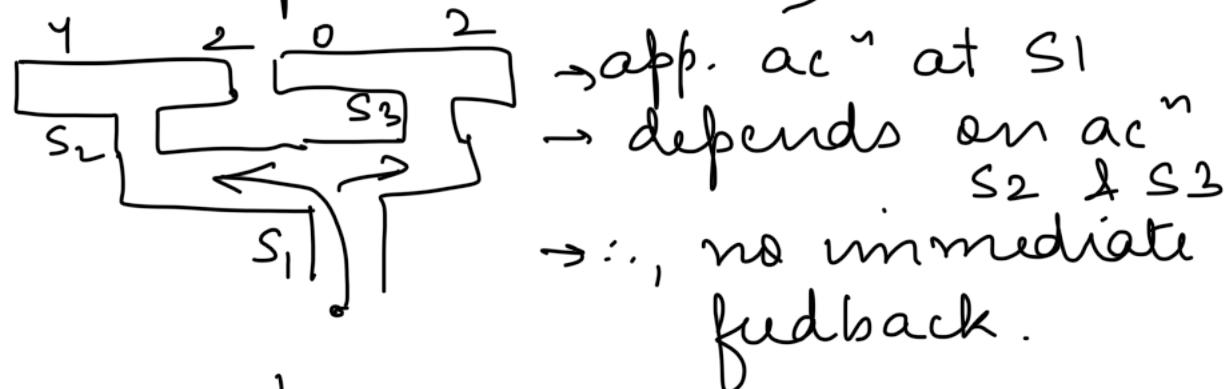


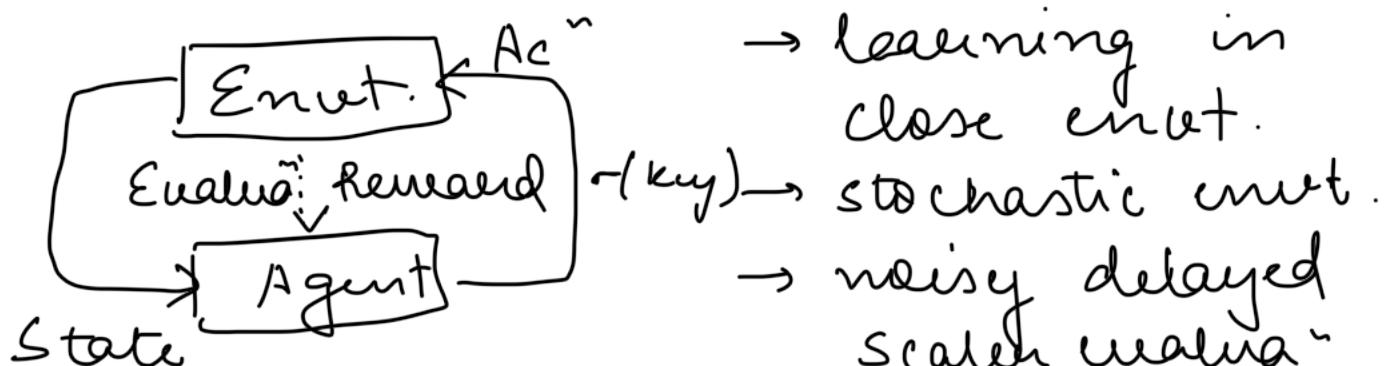
RL Week-2 Notes

T1 Full RL Problem

- Immediate RL, Bandit Problem, contextual bandit problem.
- What abt. Tic-Tac-Toe? Seq. of actions, you get results at-the end.
- your 2nd problem you see, is due to the acⁿ that you took for 1st problem.
- Acⁿ at a (Temporal Distance)



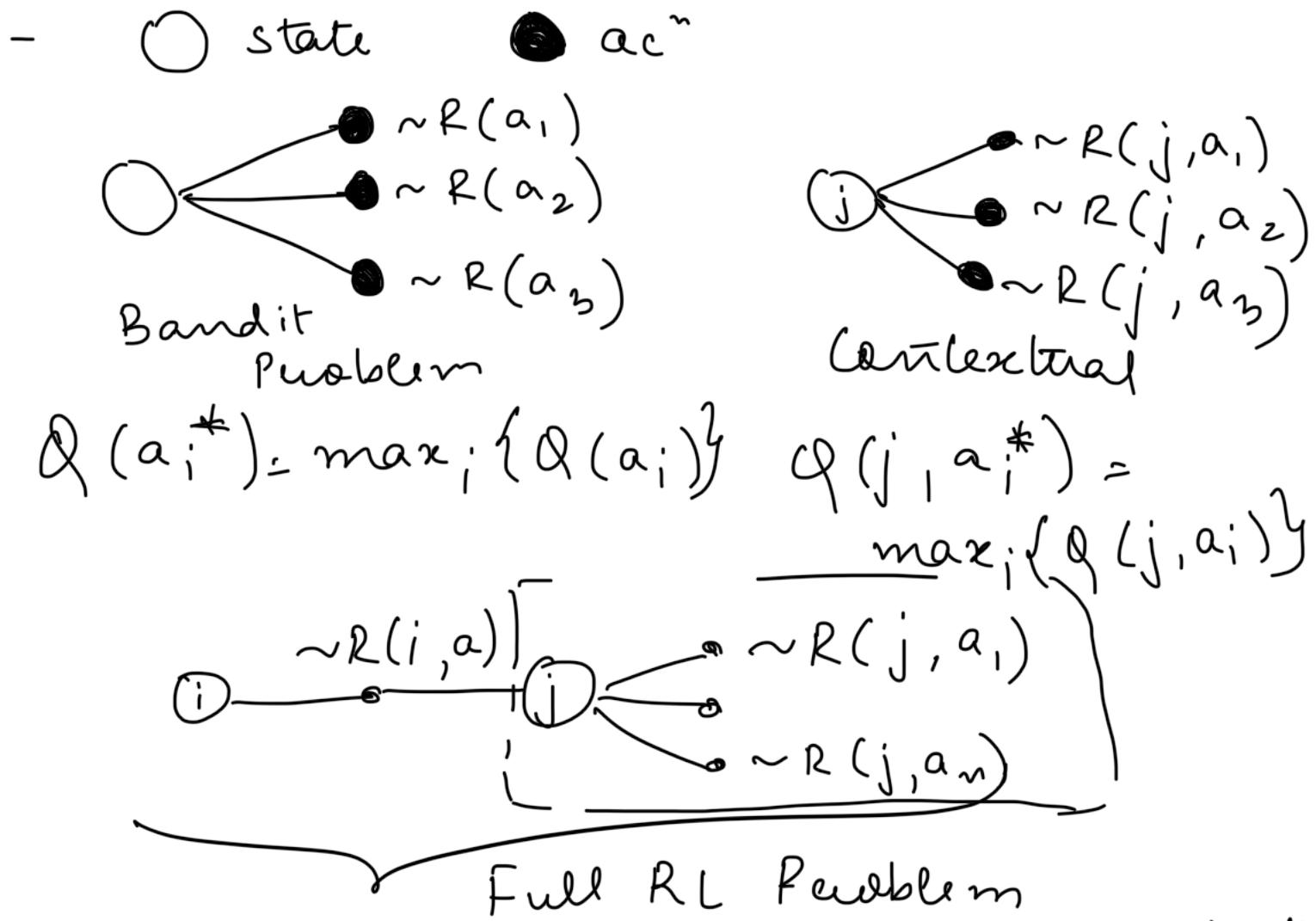
- RL Framework



→ goal is to maximize the long term results.

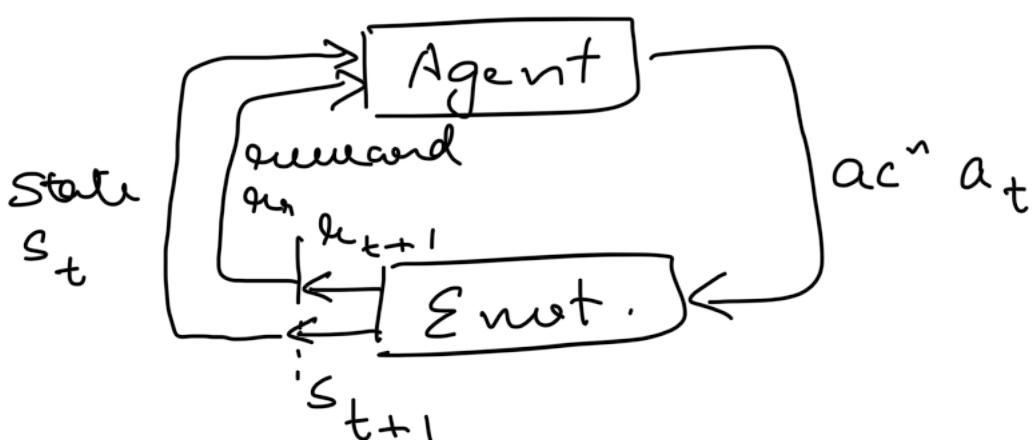
→ Agent learns the mapping from state to acⁿ, to get the best value.

→ inherently though the model is deterministic but currently model does not know, ∴, it is stochastic in nature.



- basically we are trying to use bandit " as surrogate feedback.
- for every state, you have multiple ac["]
→ this basically leads to DP.

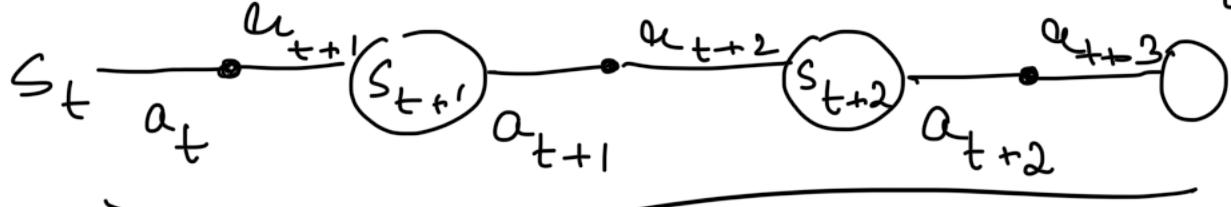
1.2 Markov Decision Process



- agent & env. are working in discrete time steps $t \in 0, 1, 2$ (not time clock)

- state at step t : $s_t \in S$ (discrete steps)
- produces ac["] at step t : $a_t \in A(s_t)$
you get a reward: $r_{t+1} \in R$

then the state becomes $\rightarrow s_{t+1}$



- Markov Assumpⁿ \rightarrow the state \rightarrow all the relevant info. available to the agent about the envt. State can be abstract entity or a seq. of sensa["]. But it should have "essential" info.

$$P_a \{ s_{t+1} = s', a_{t+1} = a' | s_t, a_t \}$$

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

it should \rightarrow not eng. the history.

\hookrightarrow it holds the Markov Property.

- MDP (Markov Decision Process)

$$M = \langle S, A, p, r \rangle \Rightarrow \text{Tuple}$$

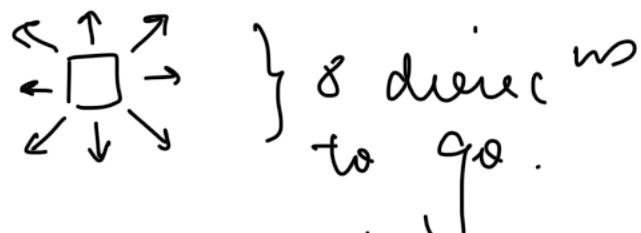
set of \leftarrow set of \uparrow probabilities \downarrow exp. reward
 states ac["] of transition $S \times A \times S \rightarrow [0, 1]$ $S \times A \times S \rightarrow R$

- Policy: $J: S \times A \rightarrow [0, 1]$ (can be state to ac["] map deterministic)

- we want to max. total expected reward, such, policy is cd an optimal policy.

L3 MDP: Problem to Formula"

- states \rightarrow enough info to take decision, have input (not enough)
- actions \rightarrow control vars., moves in the game \rightarrow can be discrete / continuous
 \hookrightarrow they are not always so easy, learn how to run a car.
- rewards \rightarrow defines the problem game.
- eg. of student defining the reward funcⁿ to learn how to cycle. It ended up doing rounds only. But remember there are numerous success stories as well.
- Eg. 2D Workspace



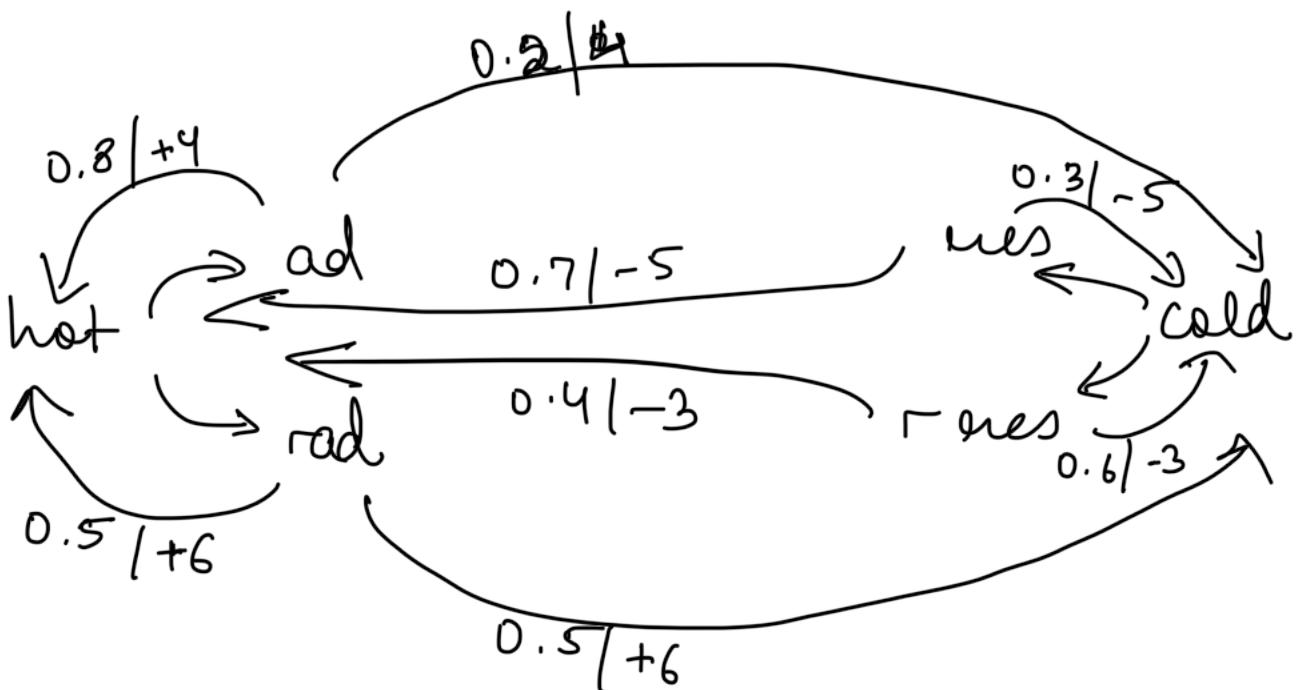
- Eg. Computer Manufacturer Company
 (you should be able to convert words into MDP)

Actions $\rightarrow A = \{ad, \neg ad, aes, \neg aes\}$

States $\rightarrow S = \{\text{hot}, \text{cold}\}$

$$A(\text{hot}) = \{ad, \neg ad\}$$

$$A(\text{cold}) = \{aes, \neg aes\}$$



- Eg. Robot Control \rightarrow input, action, +ve/-ve rewards.

L4 Returns

- Agent is learning a policy.
Policy at t, π :
- $\pi_t(s, a) = P(a_t = a, s_t = s)$
mapping states to action probabilities
- RL methods define how the policy is for the agent based on experience! Max. the reward goal remains same.
- Episodic tasks \rightarrow interactⁿ breaks into naturally broken tasks.

$$G_t = \mu_{t+1} + \mu_{t+2} + \dots + \mu_T$$

T is the final time step, when you reach the terminal state.

- Continuing Task - w/o any natural eps.

↳ Discounted return

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

$$\text{discounted value} = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

short sighted $\leq \gamma \leq$ far sighted

- we are maximizing the expected return, $E\{G_t\}$, for each step t . If you are certain, the summa^n won't go to ∞ , then you can use $\gamma = 1$

LS Value Func["]

- expected future rewards starts from state t following policy π .

State - value func["] for π .

$$v_{\pi}(s) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s \right]$$

Ac["] value func["]

$$q_{\pi}(s, a) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right]$$

$$v_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$

$$= \mathbb{E}_{\pi} [q_{\pi}(s, a)]$$

L6 MDP: Recycling Robot

- robot \rightarrow collect empty cans, sensors to detect cans, arms to collect it
- room \rightarrow empty cans, humans, recharge stat " (for robot)
- Interface
 1. Navigation \rightarrow go to cans, ΔT to motors, +ve if bot goes to can, -ve if bot topples.
 2. Pick & Place \rightarrow pick the can, put to bin, ΔT to motors, +ve/-ve respectively.
 3. Search \rightarrow search for cans, search, wait for recharge, +ve if bot collects the can, -ve for bot goes 0.
- charge level. $\rightarrow \{ \text{high}, \text{low} \}$

high $\xrightarrow{\leftarrow}$ search, wait

low $\xrightarrow{\rightarrow}$ search, wait, recharge

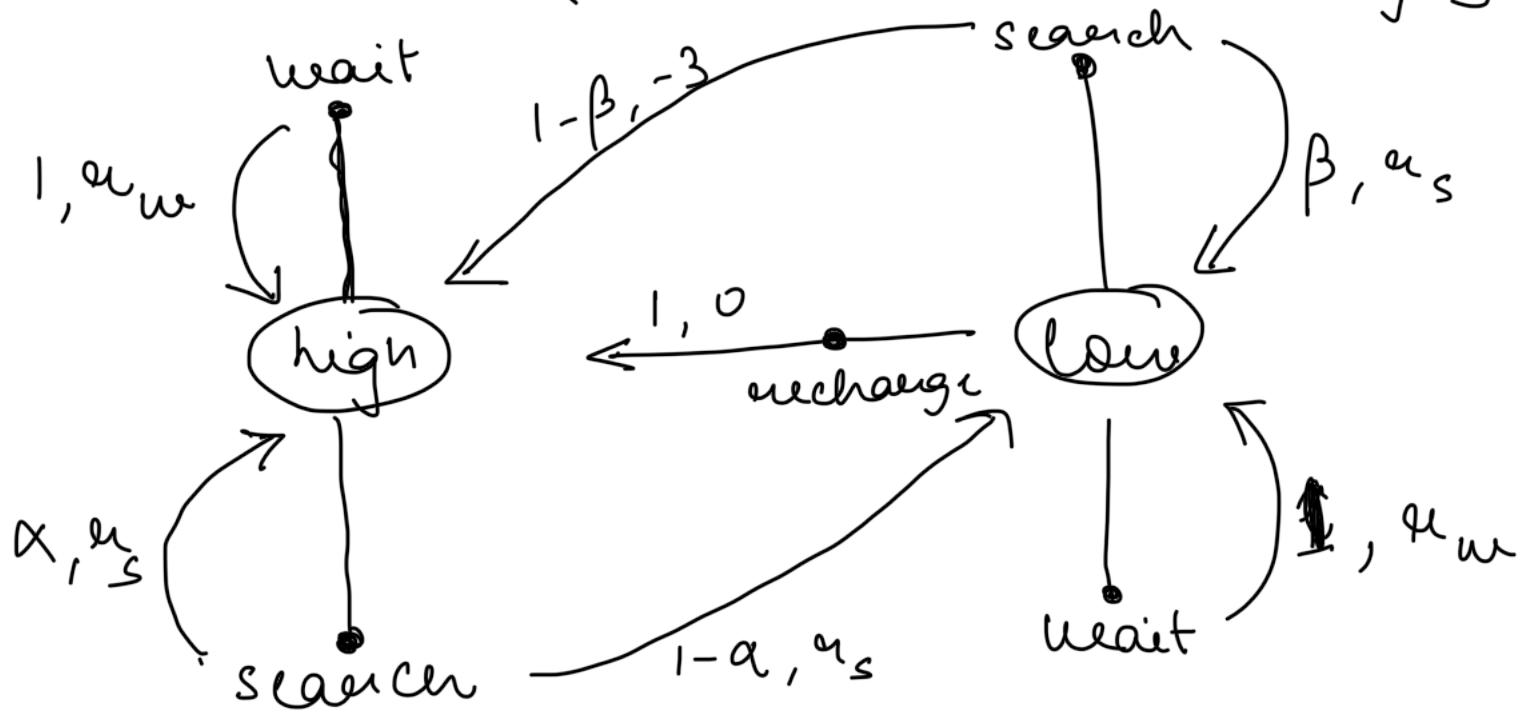
outcome

- + \rightarrow collects a can
- \rightarrow battery runs down.

$S = \{ \text{high}, \text{low} \}$

$A(\text{high}) = \{ \text{search}, \text{wait} \}$

$A(\text{low}) = \{ \text{search}, \text{wait}, \text{recharge} \}$



Transition Graph

○ state
● Acⁿ

$u_w \rightarrow u_{\text{wait}}$

$u_s \rightarrow u_{\text{search}}$

s	a	s'	$p(s' s, a)$	$u(s, a, s')$
high	search	high	α	u_s
high	wait			

Finite MDP $\xrightarrow{\quad}$

Deterministic Policy

$\pi(a|s)$ $\xrightarrow{\quad}$ Stochastic

s	a	$\pi(a s)$	0.8
high	search	1	
high	wait	0	

0.3

trajectory ↓

$s_0 \ a_0 \ s_1 \ a_1 \ a_1 \ s_2 \ a_2 \ a_2 \ s_3 \dots$
high search high & search low, recharge high

$$G_0 = 2 + \gamma | + \gamma^2 0 + \dots$$

$$v_{\pi}(s) = E_{\pi} [G_t \mid s_t = s]$$

$$v_{\pi}(\cdot) = \begin{bmatrix} v_{\pi}(\text{low}) \\ v_{\pi}(\text{high}) \end{bmatrix}$$

L7 Haunted House Example

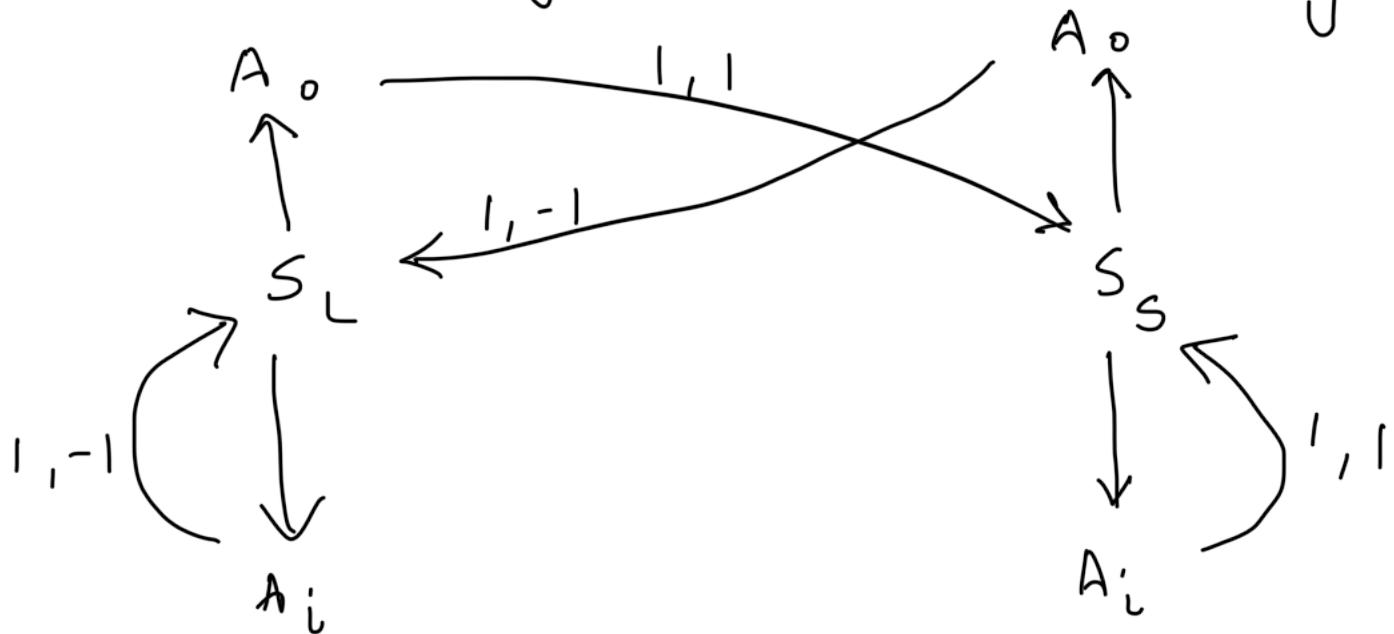
- Read the para queen in lecture.

states — { laughter s_L , silence s_S }
acts — { organ A_0 , incense A_i }

$$\gamma = 0.9$$

$$R \rightarrow \{+1, -1\}$$

$s_0 = s_L$, goal = silent always



s	a	s'	$p(s' s, a)$	$r(s, a, s')$
s_L	A_0	s_L	0	-
s_L	A_i	s_S	1	1
s_L	A_i	s_L	1	-1
s_L	A_i	s_S	0	-
:			$\pi(a s)$	
s_L	A_0		1	
s_L	A_i		0	
s_S	A_0		0	
s_S	A_i		1	

Policy

$s_0 \ a_0 \ s_1 \ a_1 \ s_2 \ a_2 \ s_3 \ a_3 \dots$
 $s_L \ A_0 \ s_S \ | \ A_i \ s_S \ | \ A_i \ s_S \ | \dots$

trajectory

$$G_0 = 1 + \gamma \cdot 1 + \gamma^2 \cdot 1 + \dots$$

$$= \frac{1}{1 - \gamma} = \frac{1}{1 - 0.9} = \frac{1}{0.1} = 10$$

$$V_{\pi}(s_L) = E_{\pi}[G_T \mid S_t = s_L] = \frac{1}{1 - \gamma} = 10$$

$$V_{\pi}(s_S) = E_{\pi}[G_t \mid S_t = s_S] = \frac{1}{1 - \gamma} = 10$$