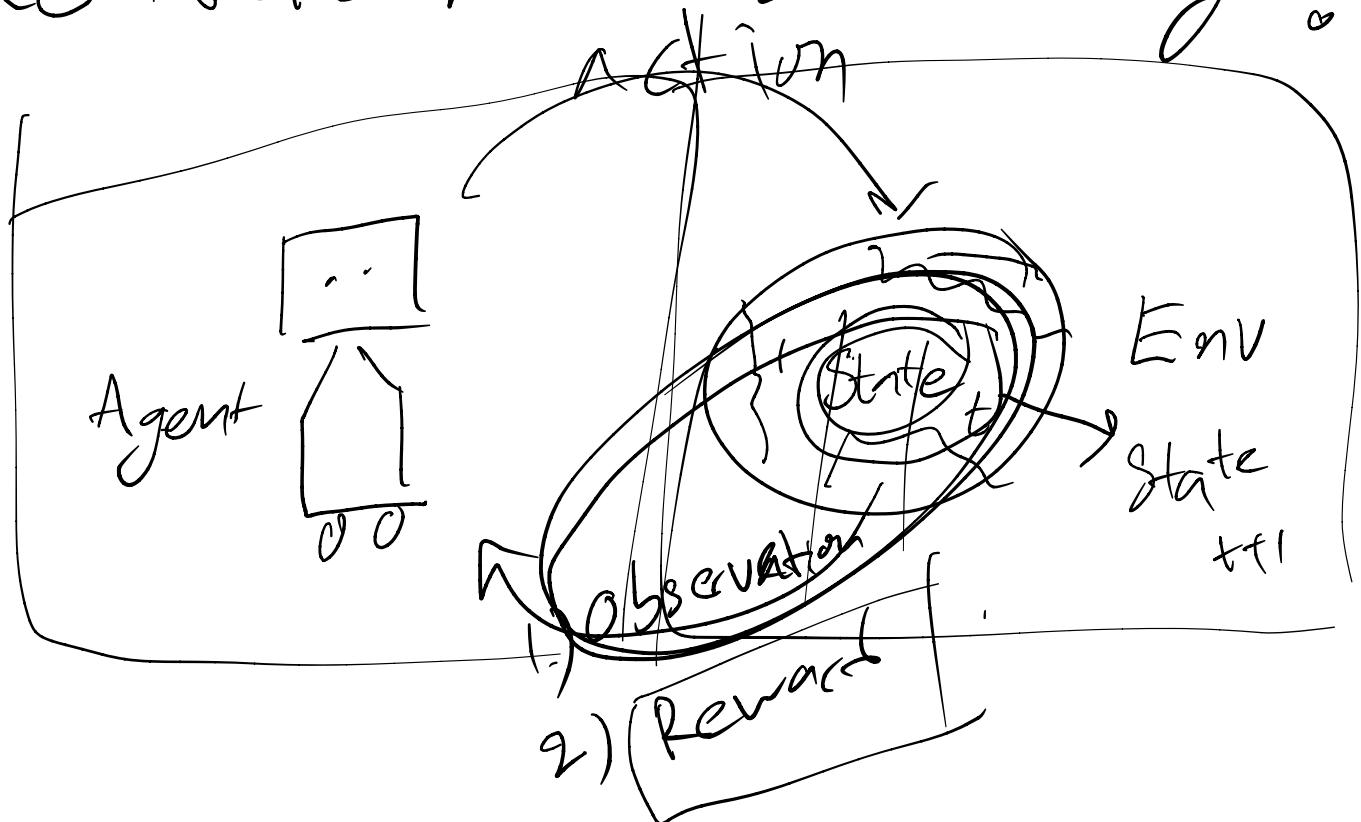


Reinforcement Learning!



What is Environment?

- Environment maintains a State

$$S = \{ s_1, \hat{s}_2, s_3, \dots, s_n \}$$

- State changes

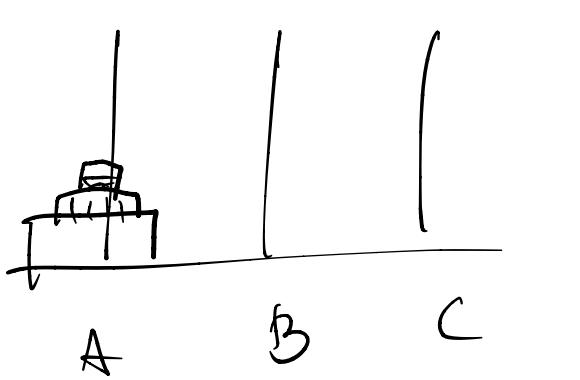
$$s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow \dots$$

- Environment emits reward

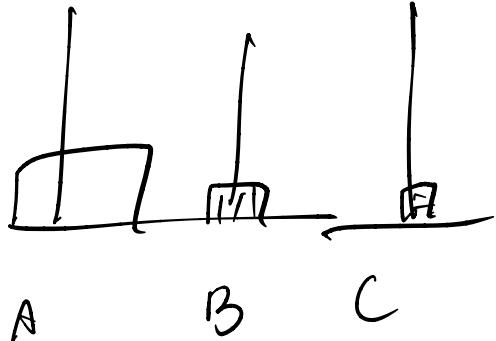
r

Example

Tower of Hanoi



s_1



s_2

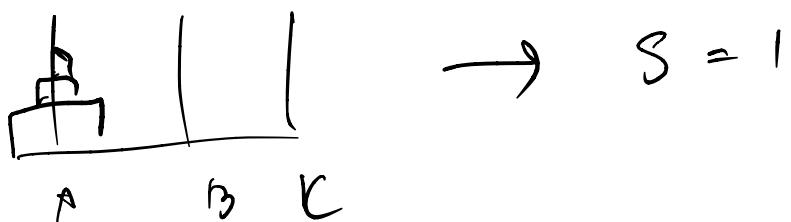
We can

Represent states with

(i) number

(ii) Vector

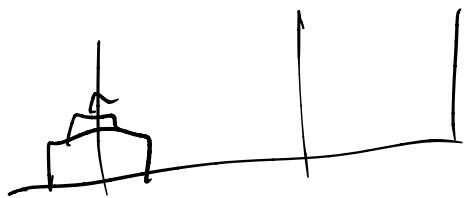
(iii) Matrix etc.



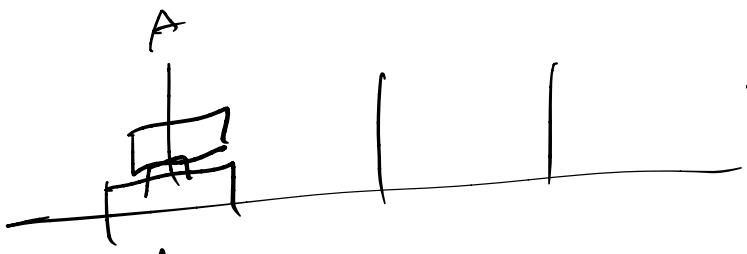
$$\rightarrow s = 1$$



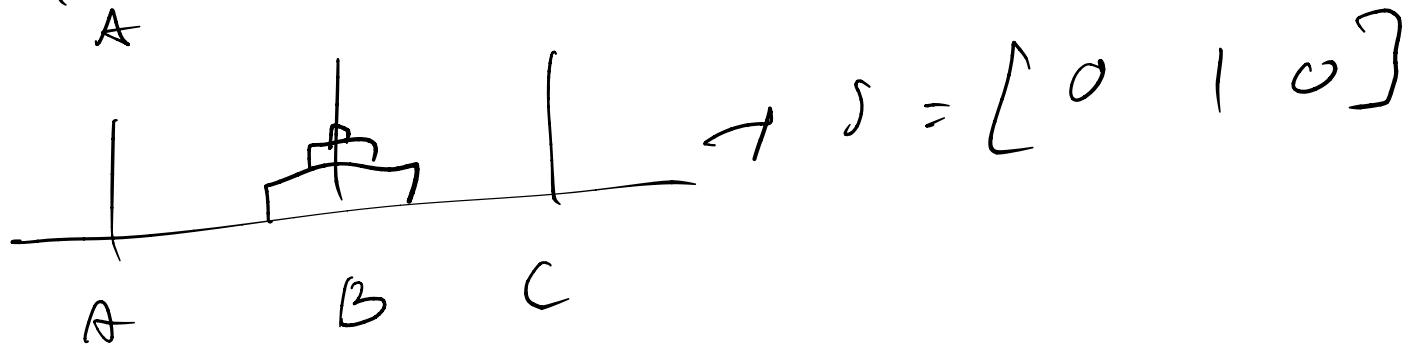
;



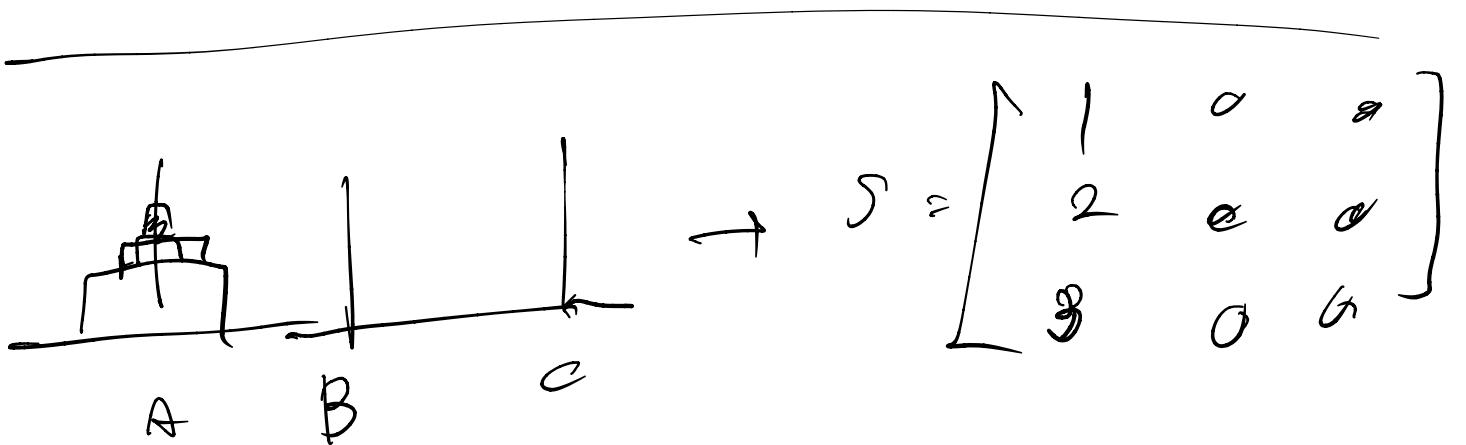
$$\rightarrow S = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$



$$\rightarrow S = \begin{bmatrix} 2 & 0 & 0 \end{bmatrix}$$



$$\rightarrow S = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$



$$\rightarrow S = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

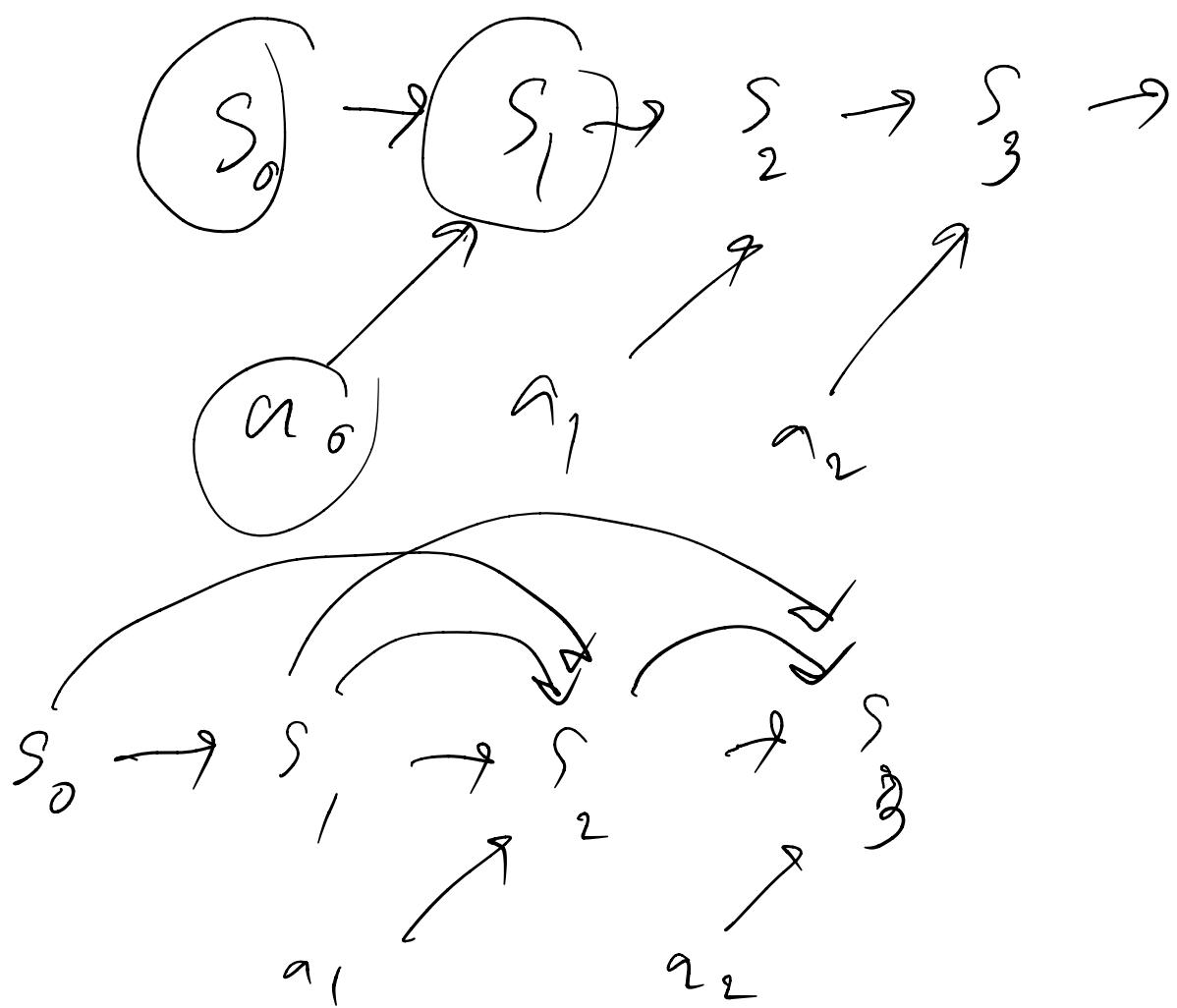
What is Reward ?

- Reward is a scalar number
 - At each time-step, a reward is emitted.
- * Reward encapsulates the objective of the task .

What make the state change ?

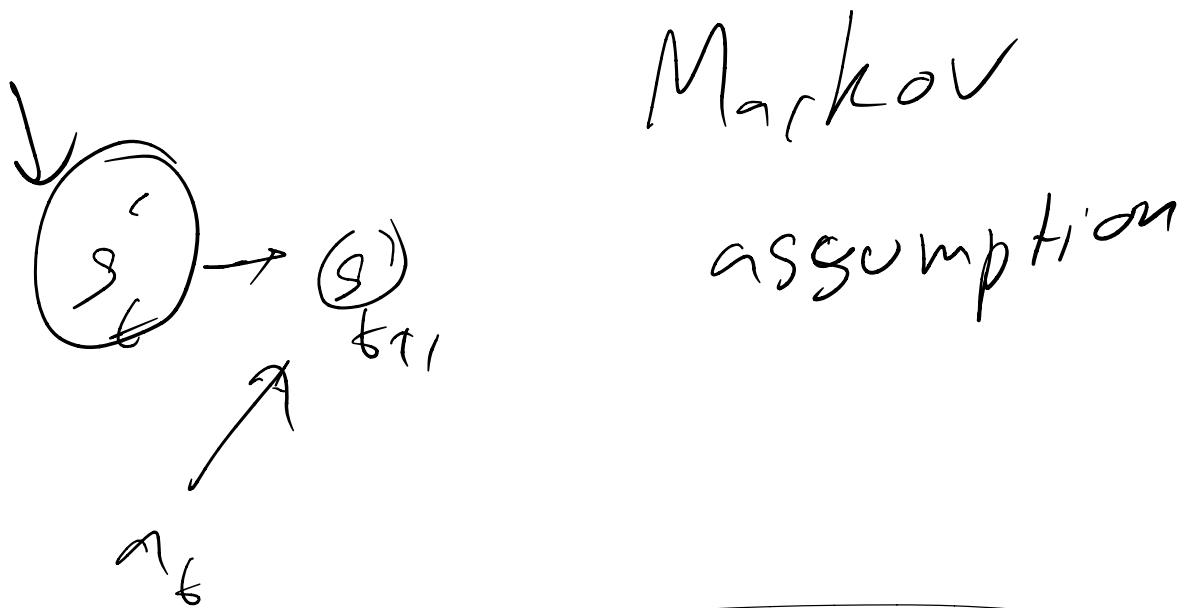
1.) Action

2.) Time



Let's define a transition function

$$\mathcal{T} = P(S_{t+1} = s' \mid S_t = s, A_t = a)$$



Markov Decision Process (MDP)

MDP is one way to model environment mathematically

MDP

$$\langle S, A, T, R, \gamma \rangle$$

$$S = \{s_0, s_1, s_2, \dots, s_n\}$$

$$A = \{a_0, a_1, \dots, a_m\}$$

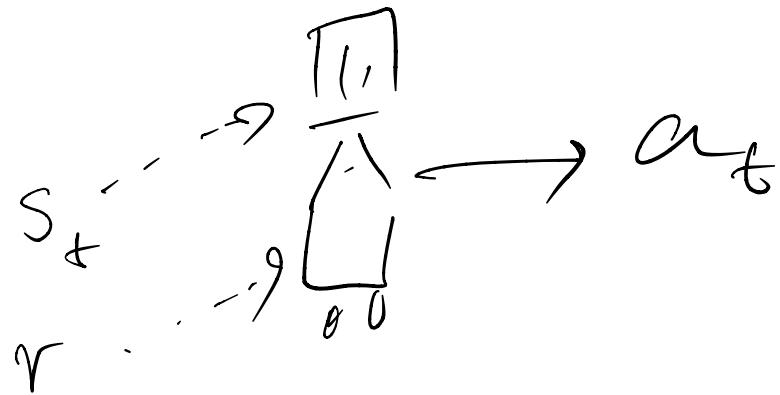
$$T = P(s_{t+1} | s_t, a_t)$$

$$R = R(s_t, a_t) \rightarrow r$$

γ = discount factor ; $\gamma \in [0, 1]$

$$P(R | s_t, a_t)$$

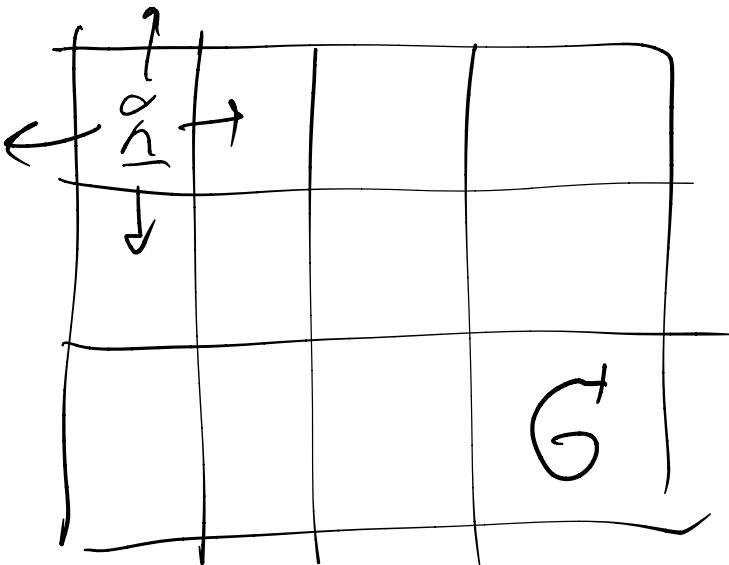
Agent



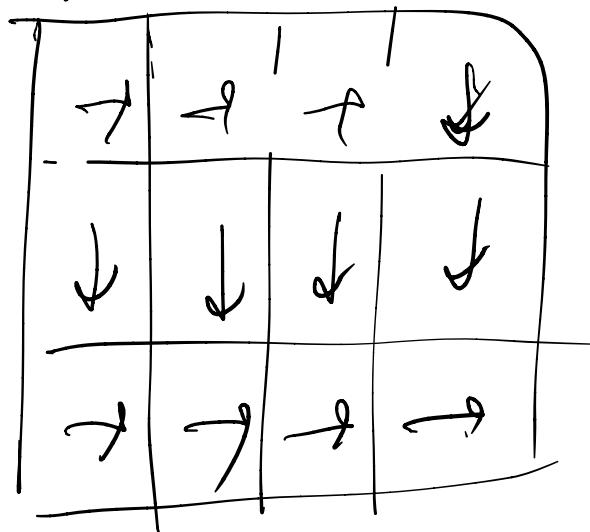
It follows a policy, π

$$\pi = p(A|S)$$

Example



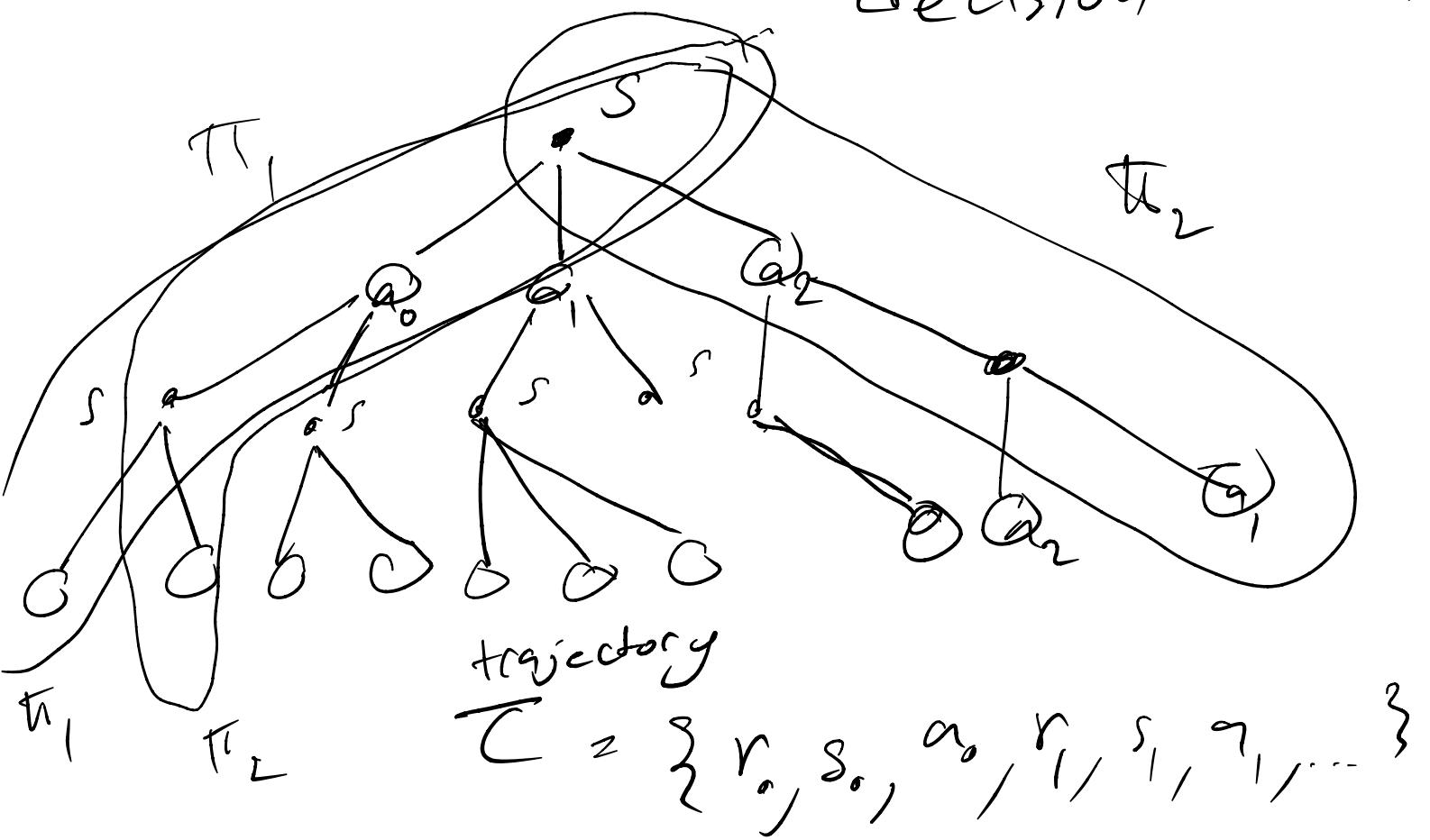
$$\pi_0 =$$



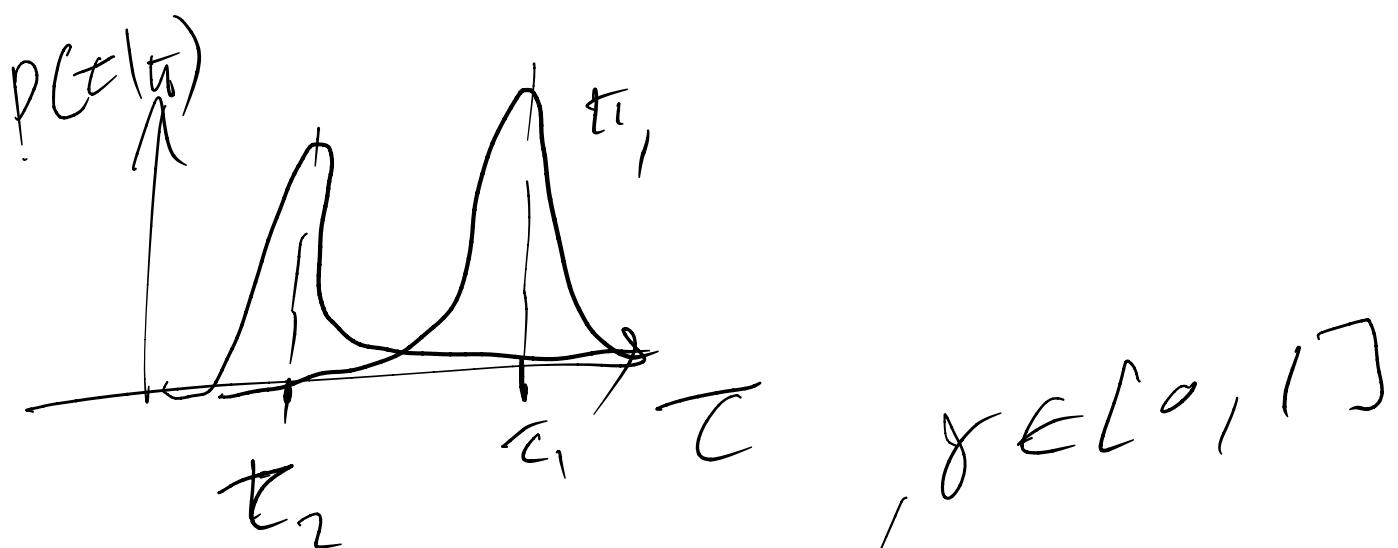
Objective of RL

We want to find
"the best" policy.

Different policies take you
to different path in a
"decision tree".



$P(\tau | \pi)$



Let's define a

Return G^+

$$g = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^T r_T$$

$$\tau_1 = \{r_0, s_0, a_0, r_1, s_1, a_1, \dots, r_T\}$$

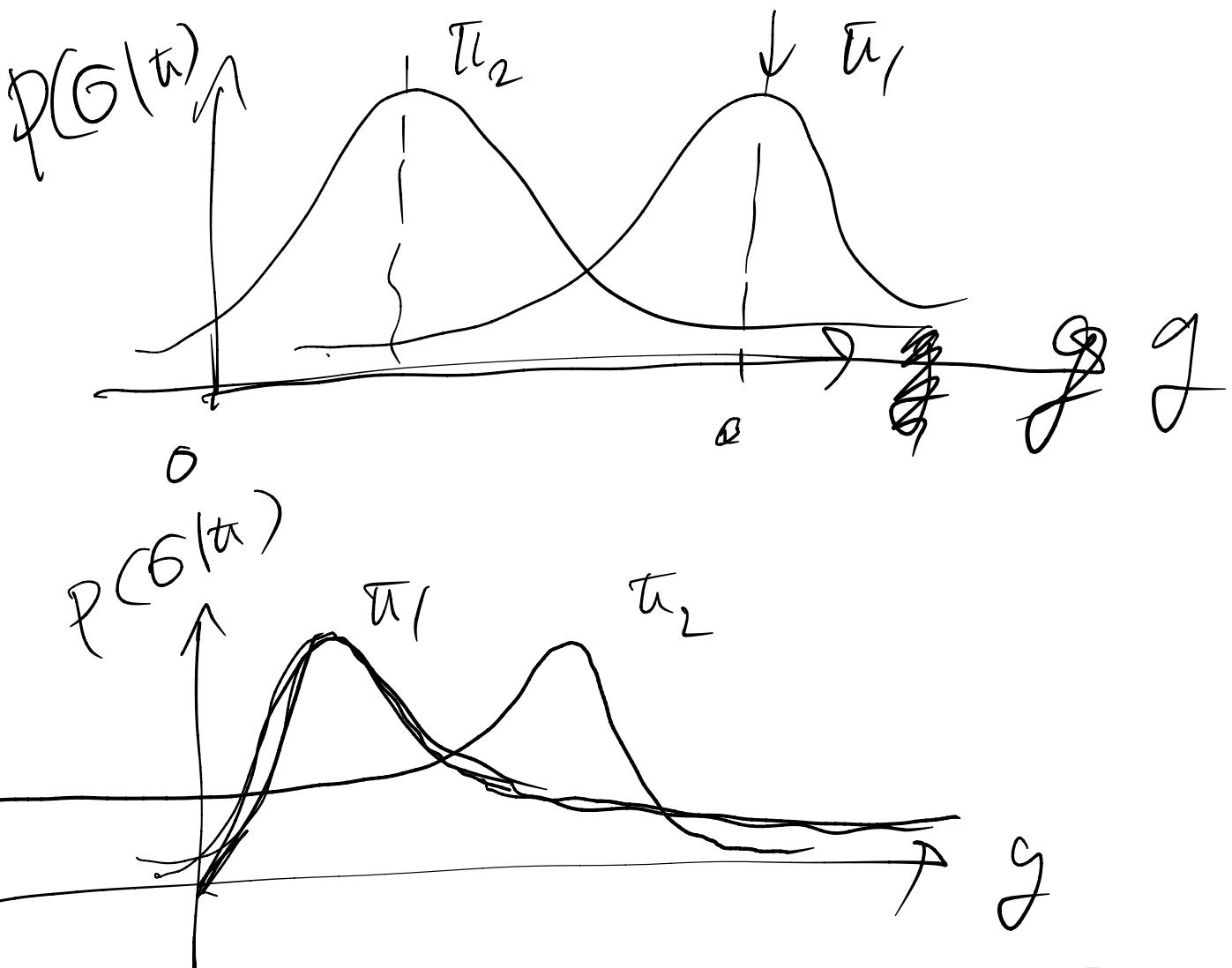
$$g_1 = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^T r_T$$



Trajectory τ_1 is better

than τ_2 , if g_1 is better than g_2

A policy π gives us a distribution of G ; $P(G|\pi)$



In Classical RL,

We compare the mean

If $E[G|\pi_1] > E[G|\pi_2]$

then π_1 is better than π_2

How to find
the best π ?

π^* with the highest $E[\text{lot}_t]$
is the best π . π^*

How to find π^*

Let's define "Value"

$$V^\pi(s_t) = E[G_t | \pi, s_t = s]$$

state-value function

$$g_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t} r_T$$
$$g_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^{T-(t+1)} r_T$$

- Given a policy π

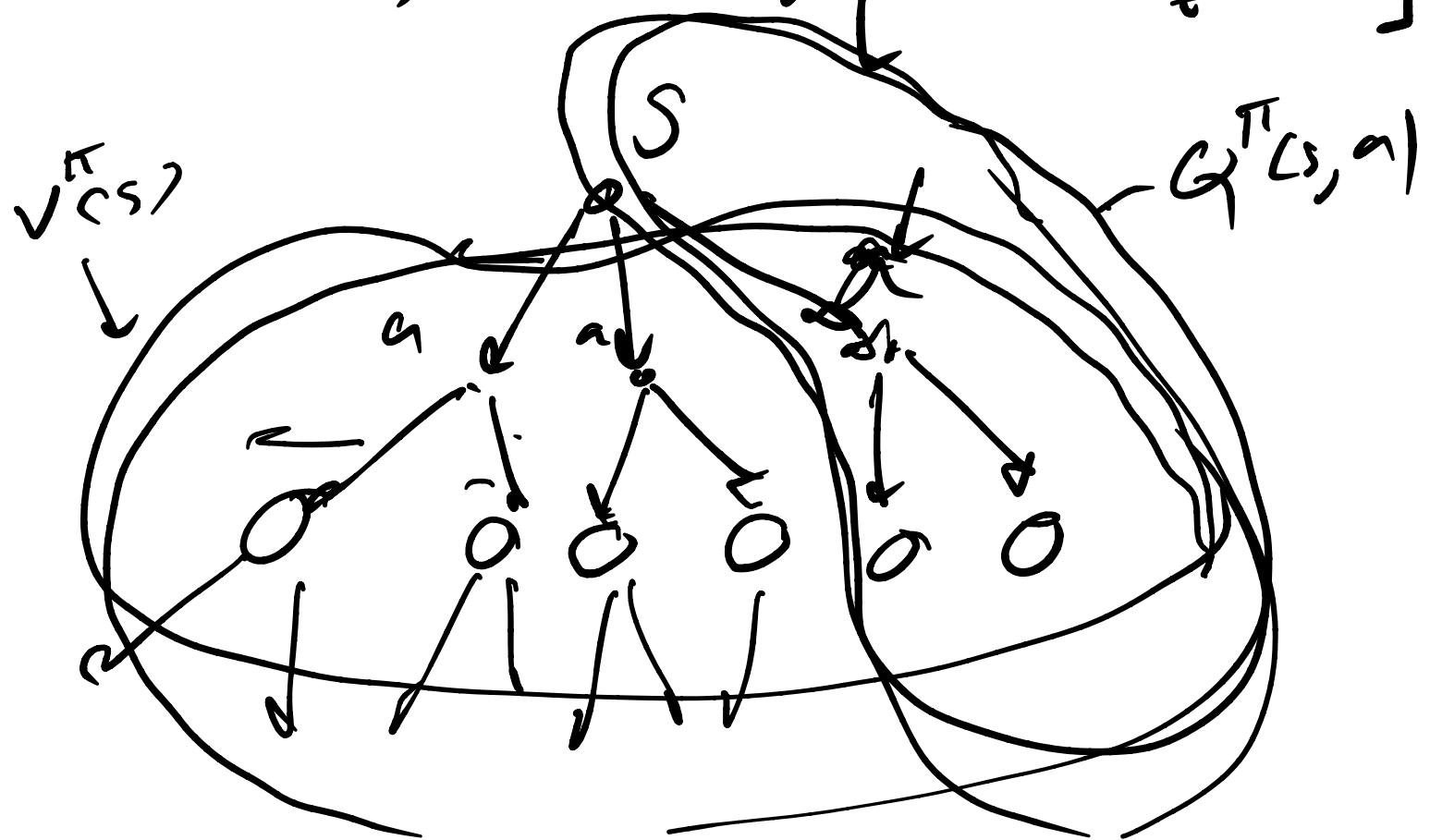
$$\text{if } V^\pi(s_1) > V^\pi(s_2)$$

We prefer s_1 to s_2

Similarly, we define

Q-value

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



Optimal Policies & Optimal Value Function

$$\pi \text{ is better than } \pi'$$

if and only if

$$V^\pi(s) \geq V^{\pi'}(s)$$

$\forall s \in S$

There is always at least one policy

that is better than or equal to

all other policies. π^* (optimal policy)

π^* (optimal policy)

v^* (optimal value function)

$$v^* = v^*(s) = \max_{\pi} V^\pi(s)$$

$\forall s \in S$

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

$\forall s, a \in S, A$

$$v^* = E[G_t | S_t]$$

$$v^* = E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots]$$

$$v^* = E[r_{t+1} + \gamma E[r_{t+2} + \gamma r_{t+3} + \dots]]$$

$$V^\pi(s_t) = E[R_{t+1} + \gamma V^\pi(s_{t+1})]$$

$$V^*(s_t) = E^* [R_{t+1} + \gamma V^*(s_{t+1})] \\ | S_t = s]$$

$$Q^*(s_t, a_t) = E^* [R_{t+1} + \gamma V^*(s_{t+1})] \\ | S_t = s, A_t = a]$$

Bellman Equation for V^π

$$V^\pi(s) = E_{\pi, \tau} [G_t \mid S_t = s]$$

$$= E_{\pi, \tau} \left[R_{t+1} + \gamma G_{t+1} \mid S_t = s \right]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r P(S_{t+1} = s') R_{t+1} = r \left[S_t = s \right]$$

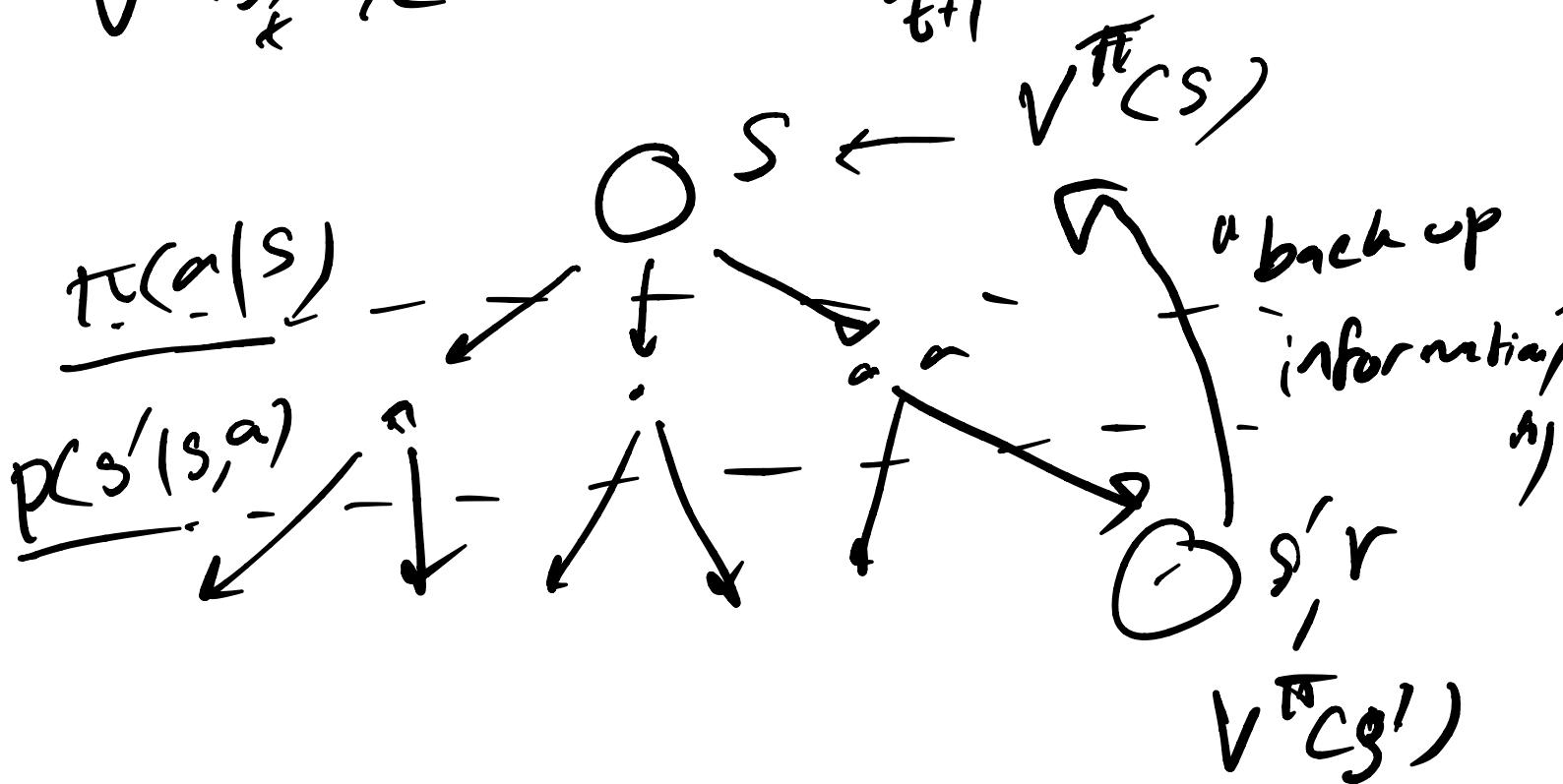
$$\cdot [r + \gamma E_{\pi} [G_{t+1} \mid S_{t+1} = s']]$$

$$E[X] = \sum_x \underbrace{P(X=x) \cdot x}_{P(x)}$$

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

Bellman Eqn $\rightarrow [V^\pi(s) = E_{\pi, \tau} [R_t + \gamma V^\pi(s')]]$

$$V^\pi(s_t) \approx r + \gamma V^\pi(s_{t+1})$$



Bellman Optimality Equation

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

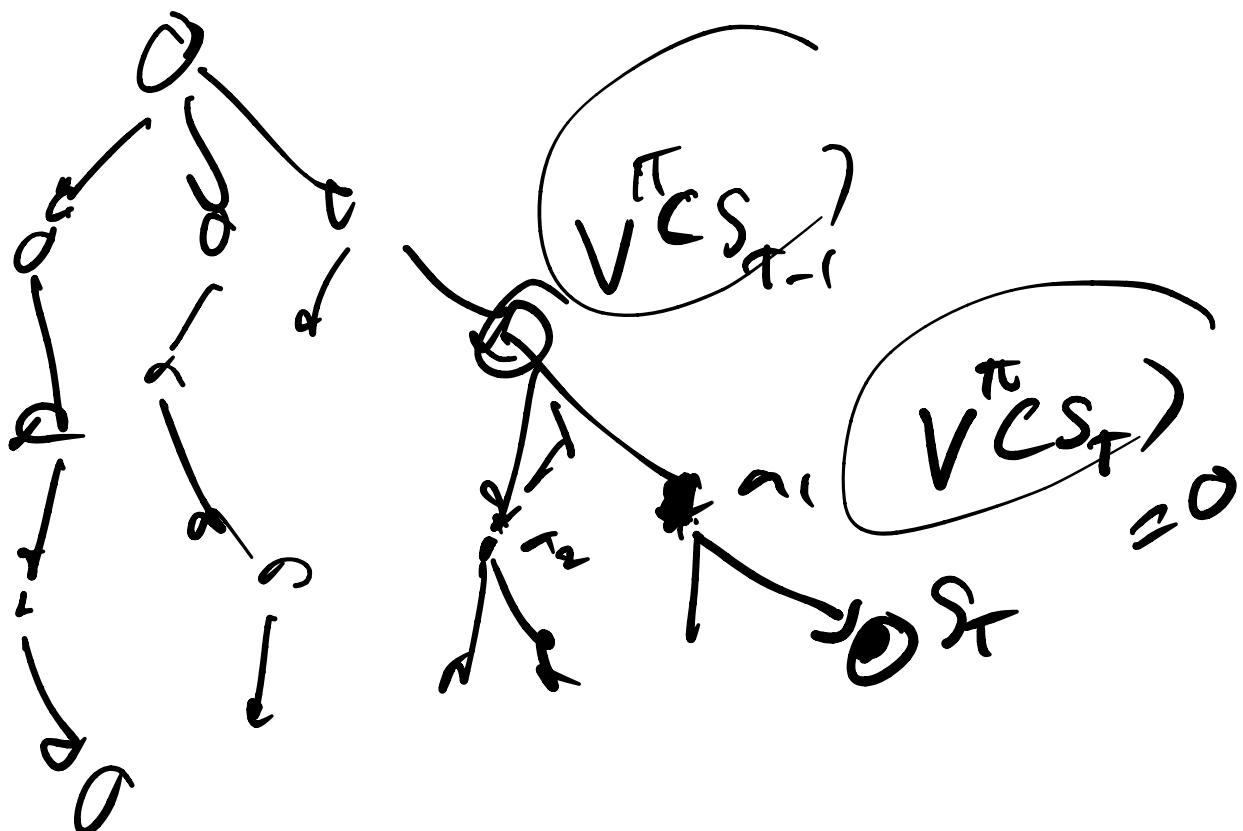
$$V^*(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma \underline{V^*(s')}]$$

$$Q^*(s, a) = \max_{\pi} E_{\pi^*} [R + \gamma Q^*(s', a')]$$

How to find $V^\pi(s)$, $V^*(s)$,
 $Q^*(s, a)$

for all s and a

DP = breaking problem into sub-problem
in a recursive manner.



Policy Evaluation

← Algorithm name

Find $\pi(s)$

$$Q^\pi(s, a)$$

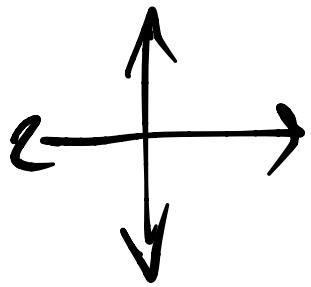
$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V^\pi(s')]$$

1.) Start with an initial $\hat{V}_0(s) \forall s$

2.) Update $\hat{V}(s)$ with

$$\hat{V}_{K+1}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma \hat{V}_K(s')]$$

3.) Repeat until Converge



action

| | | |
|----|----|---|
| 1 | 2 | 3 |
| 4 | 5 | 6 |
| 7 | 8 | 9 |
| 10 | 11 | |

$$R_t = -1$$

π_{random}

$\sqrt{\pi_{\text{random}}}$

$K=0$

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$\frac{1}{k^0} \pi(s)$

$$V = \sum \pi \sum P(S, r | s, a) [r + \gamma V_{K+1}]$$

$K=1$

| | | | |
|------|------|------|------|
| 0.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

$\sqrt{\pi_{\text{random}}}$

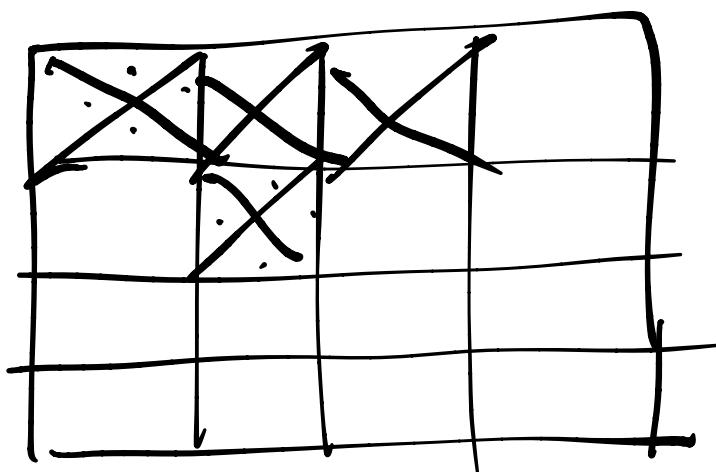
$K=2$

| | | | | |
|---------------|---------------|---------------|---------------|--------|
| $0;0$ | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{3}{4}$ | -2.0 |
| $\frac{1}{4}$ | -1.7 | -2.0 | -2.0 | |
| -1.9 | -2 | -2 | -2 | |
| -2 | -2 | -2 | -1.7 | |
| -2 | -2 | -1.9 | 0.0 | |

$$\begin{aligned}
 \hat{V}_{k=2}(S=1) &= \frac{1}{4} \cdot \left[-1.0 + \gamma \hat{V}_{k=1}(S=2) \right] \\
 &+ \frac{1}{4} \left[-1.0 + \gamma \hat{V}_{k=1}(S=S) \right] \\
 &+ \frac{1}{4} \left[-1.0 + \gamma \hat{V}_{k=1}(S=T) \right] \\
 &+ \frac{1}{4} \left[-1.0 + \gamma \hat{V}_{k=1}(S=1) \right] \\
 &= -\frac{1}{2} - \frac{1}{2} - \frac{1}{4} - \frac{1}{2} \\
 &= -1.75
 \end{aligned}$$

$K = 2$

| | | | |
|-----|-----|-----|-----|
| 0.0 | -14 | -20 | -22 |
| -14 | -18 | -20 | -20 |
| -20 | -20 | -18 | -14 |
| -22 | -20 | -14 | 0.0 |



$\hat{Q}(s, a)$

Policy Improvement

theorem

$$\text{if } Q^{\pi}(s, \pi'(s)) \geq V^{\pi}(s)$$

then π' must be better than π

Policy Iteration

= Policy Evaluation

+ Policy Improvement

Value Iteration

= Policy Evaluation Only (-step)
+ Policy Improvement