

Lecture 3

Review

- $P(s', r | s, a)$
 $= P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$
- $MDP = \langle S, A, R, T, \gamma \rangle$

- Trajectory

$$= \{S_0, A_0, R_1, S_1, A_1, \dots\}$$

- Return

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

- Value

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

- Q-Value

$$Q^\pi(s, a) = E [G_t | S_t = s, A_t = a]$$

- Optimal-Value

$$V^*(s) = \max_\pi V^\pi(s)$$

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

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

$$Q^*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma V^*(s_{t+1}) \mid S_t = s, A_t = a \right]$$

Bellman Eq

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

- Bellman Optimality Eq

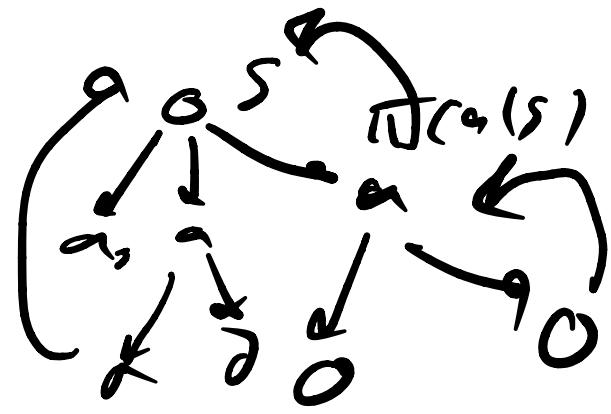
$$V^*(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')]$$

$$Q^*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} Q^*(s', a')]$$

- Dynamic Programming

→ Policy Evaluation

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$$



→ Policy improvement

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

$$\pi' = \arg \max_{\alpha} Q^{\pi}(s, \alpha)$$

+

Policy Iteration

1.) Initialise $V(s)$ and $\pi(s)$

2.) Policy Evaluation ↗
 $V^{\pi}(s)$

3.) Policy Improvement
 $\pi \rightarrow \pi'$

$$\pi' = \arg \max_{\pi'} Q^{\pi}(s, a)$$

4.) Repeat

• Value Iteration

1) Initial $V(s)$

{2) One-step policy evaluation
One-step policy improvement

$$V(S) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) \left[r + \gamma V(s') \right]$$

3) Repeat

Monte Carlo Method

- Previously, we assume knowledge over MDP
i.e. $p(s', r | s, a)$
- Now, let's say we don't know $p(s', r | s, a)$

We want to find

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

- We could estimate
 $P(s', r | s, a)$
from data

This is called

"Model-based" method

- We could learn
 V_{CS} , $\pi(s)$, etc

directly from data

"Model-free" method

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

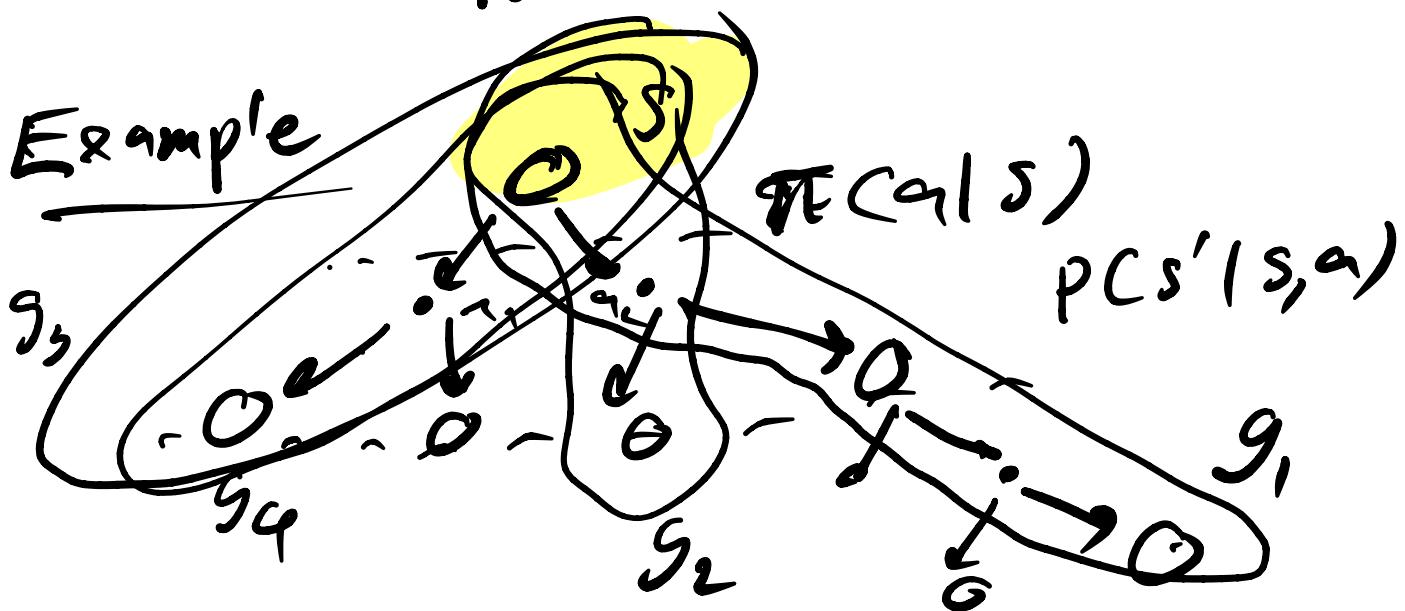
$$V^\pi(s) = E[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots]$$

$E[G] = \sum g_t p(g_t | s)$

- $E[G] \approx \frac{1}{N} \sum_i^n g_i$

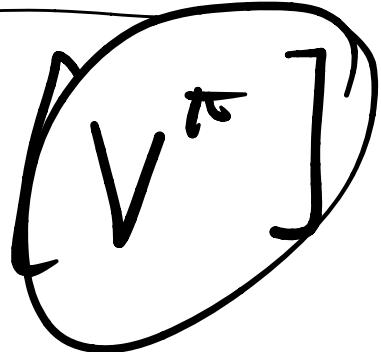
"Monte-Carlo Estimation"

- $V^\pi(s) \approx \frac{1}{N} \sum_i^n (r_t^i + \gamma r_{t+1}^i + \gamma^2 r_{t+2}^i + \dots)$



$$V^\pi(s) \approx \frac{1}{4} (g_1 + g_2 + g_3 + g_4)$$

First-Visit MC estimation



- 0.) ~~Input~~ input π
- 1.) Initialise

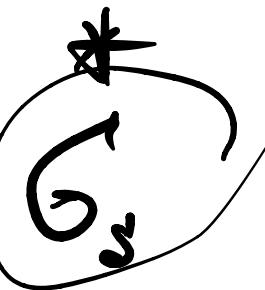
~~$V(s)$~~ , $G_s = []$

2.) Generate trajectory with $\tilde{\pi}$

3.) For each s appear

- $g(s) = \sum_i \gamma^i r$

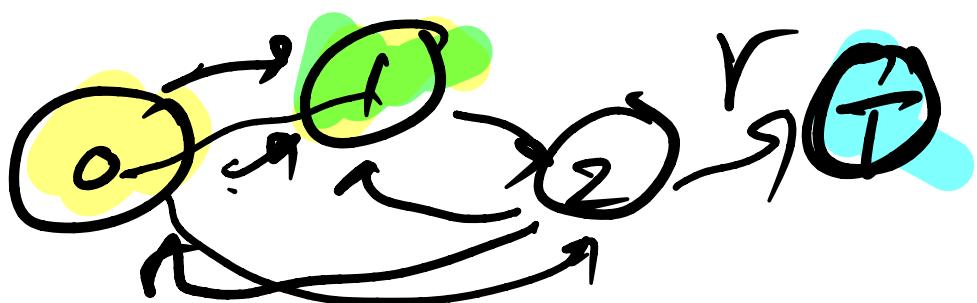
- Append $g(s)$ to G_s



4.) Compute $\underline{V}(s) \leftarrow \text{Average}(G_s)$



Example



$$\mathcal{T} = \{ S_0 = 0, a_0 = 1, r_1 = 1 \}$$

$$S_1 = 1, a_1 = 0, r_2 = 0$$

$$S_2 = 0, a_2 = 1, r_3 = 1$$

$$S_3 = 2, a_3 = 1, r_4 = 1$$

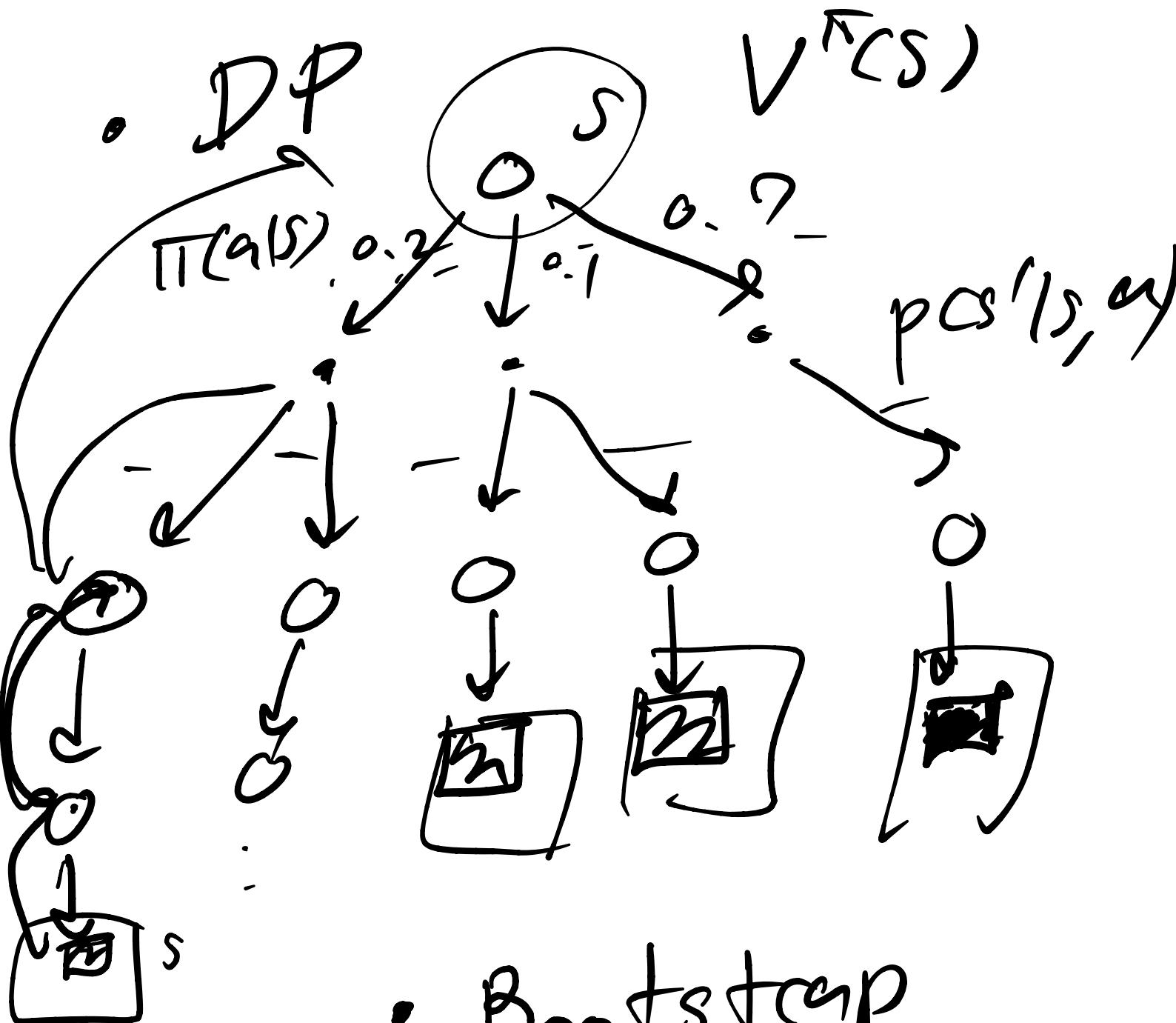
~~S₄~~ ≈ Terminal

$$g(0) = r_1 + \delta r_2 + \delta^2 r_3 + \delta^3 r_4$$

$$= 1 + g(0) + \delta^2 g(1) + \delta^3 g(2)$$

$$G_n = [g(0)]$$

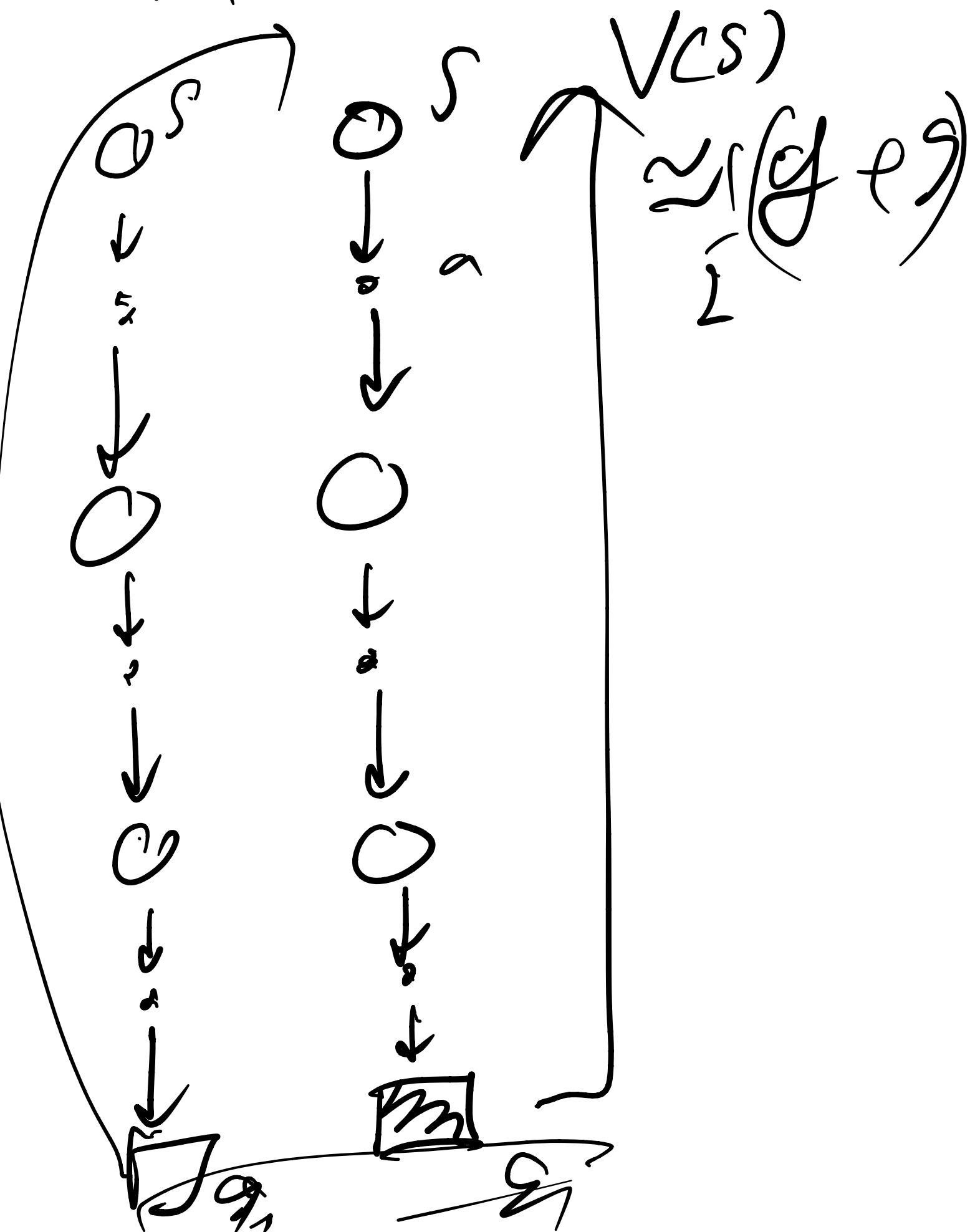
Backup diagram



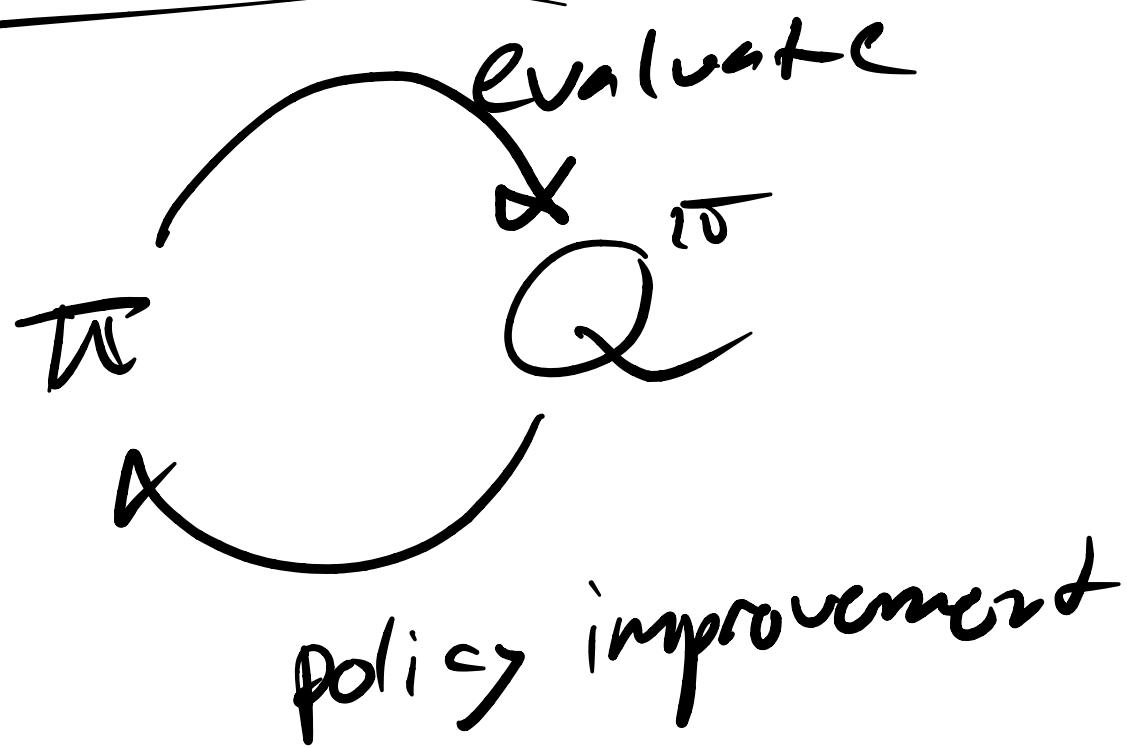
- Bootstrap

- Complete knowledge

• MC back up



MC - Control



1.) Initialise

$$Q(s,a), \pi(s)$$
$$G_{(s,a)} = []$$

2.) Generate

trajectory with π

- For each s, a

- $g(s, a) \leftarrow \sum_i \gamma^i r$

- $\text{Append } g(s, a) \rightarrow G_{s, a}$

- $\text{Compute } Q(s, a) = \text{Avg}[G]$

- For each s in the traj

$$\pi(s) \leftarrow \arg\max_a Q(s, a)$$

Exploration in MC

- We want to estimate

$$V^\pi(s) \text{ or } G^\pi(s, a)$$

$$H_s \text{ or } H(s, a)$$

- For a particular π ,

(s, a) may never be visited.

- One way to cheat

"Exploration Start assumption"

- We could consider
using $\pi(a|s)$

→ e.g. ϵ -greedy policy

$$a^* = \operatorname{argmax}_a Q(s, a)$$

