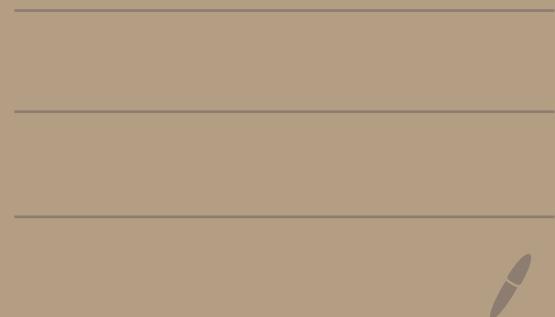


# Lecture 23 RL

---





$r_0$  $r_1 - \dots - r_t$ 

current reward

future

 $\gamma^t$ 

smaller

infinite time



if time is finite

$$\max_{T \text{ steps}} \left[ \sum_{t=0}^T r_t \right]$$

} imitation

policy



efficiency      vs      effectiveness



policy:

transition:

randomness  
 $\pi(a|s)$

$P(s'|s, a)$

$s_0 \rightarrow a_1 \rightarrow s_1 \rightarrow a_2 \rightarrow s_2 \dots$

trajectory

$$\sqrt{\tau}(s)$$


starting in  $s$  execute  $\tau$

$$V^{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t E[RCS_t, \pi(s_t)]$$

deterministic  
Starting from s

$$= R(s, \pi(s)) + \sum_{t=1}^{\infty} \gamma^t E[RCS_t, \pi(s_t)]^{s_{t+1}}$$

$$= R(s, \pi(s)) + V \cdot \sum_{t=0}^{\infty} \gamma^t E[RCS_t, \pi(s_t)].$$

Starting from s

$$P_{CS'}(s, \pi(s)) \cdot V_{CS'}(s')$$

Y.  
 $E_{s'}[V_{CS'}]$   
the next state  
of s, s'

$$V^\pi(s) = \underbrace{\sum_{t=0}^{\infty} \gamma^t E[R(s_t, \pi(s_t)) | s_0=s]}$$

↓                  ↓                  ↗

s                  s

$$V^\pi(s)$$

$S \in \{1, 2, 3\}$

Linear

1.  $\boxed{V^\pi(S=1)} = R(S=1, \pi(S=1)) + \gamma \sum_{S_i \in \{1, 2, 3\}} P_{S_i | S=1} \pi(S=i)$

$x = V^\pi(S=1)$   
 $y = V^\pi(S=2)$   
 $z = V^\pi(S=3)$

$V^\pi(S_1)$

2.  $V^\pi(S=2) = R(S=2) - - - - -$

3.  $V^\pi(S=3) - - - - -$

$$\sqrt{\bar{\pi}}(s)$$

$$\max_{a \in A}$$

$$\overbrace{\pi(a|s)}$$

$$\overbrace{\pi^*(a|s)}$$

$$\overbrace{\pi^*}$$

$$V^{(est)}(s) \rightarrow V^*$$

$\boxed{\pi(s)}$  deterministic;  $Q(s, a) = V(s)$

$$\pi^*$$

$$Q(s, a) = R(s, a) + \gamma \dots$$

$V(s)$

$V^{\bar{a}}(s)$  given  $\bar{e}$

$R(s, a), P(s' | s, a)$  known

0.1 0.2 0.7

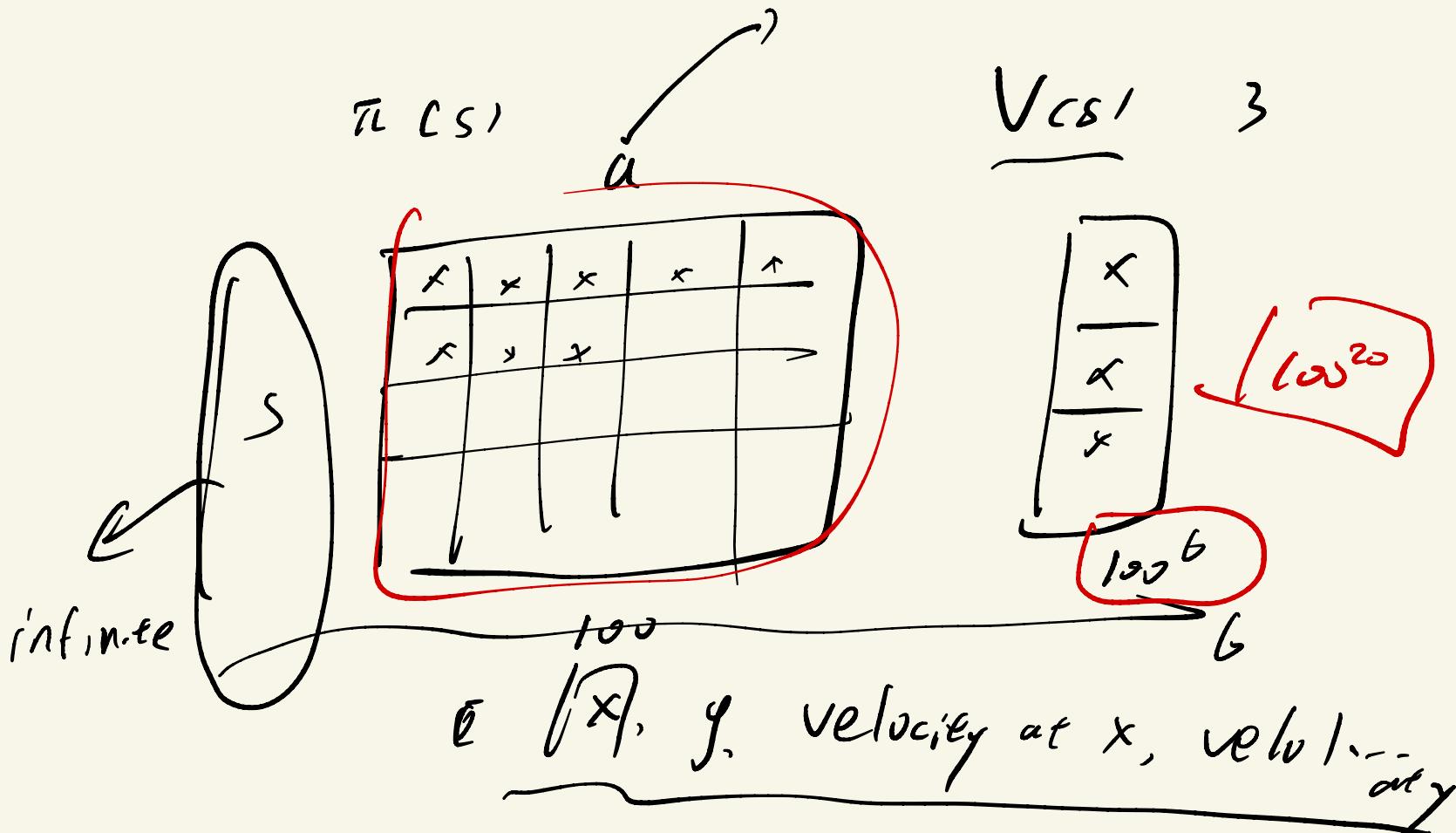
{0, 1}

Strength  $\{1, 2, 3, 4, 5\}$

1000 times. 10 times chair fall down.

$$R(\text{Strength}=1) = \frac{10}{1000}$$

$$= 0.01$$



$$\left[ \begin{array}{c} c_s, \\ \sqrt{c_s} \end{array} \right]$$

$$\pi(a|s)$$

continuous

Policy model

$$\gamma_\theta = \bar{E}_{z \sim p_\theta} \left[ \sum_{t=0}^{T-1} V^t R(s_t, o_t) \right]$$

unrelated to  $\phi$

$$\cdot \bar{E}_{z \sim q_\phi(c_t|x_t)} [B_\phi(\theta, \phi)]$$

$z \sim \text{Gaussian}$

$z \quad \underline{s_0, a_1, s_1, a_2, s_2, a_3}$

$$\mathbb{E}_{\tau \sim P_\theta} [f(\tau)]$$

$$= \mathbb{E}_\theta \int P_\theta(\tau) f(\tau) d\tau$$

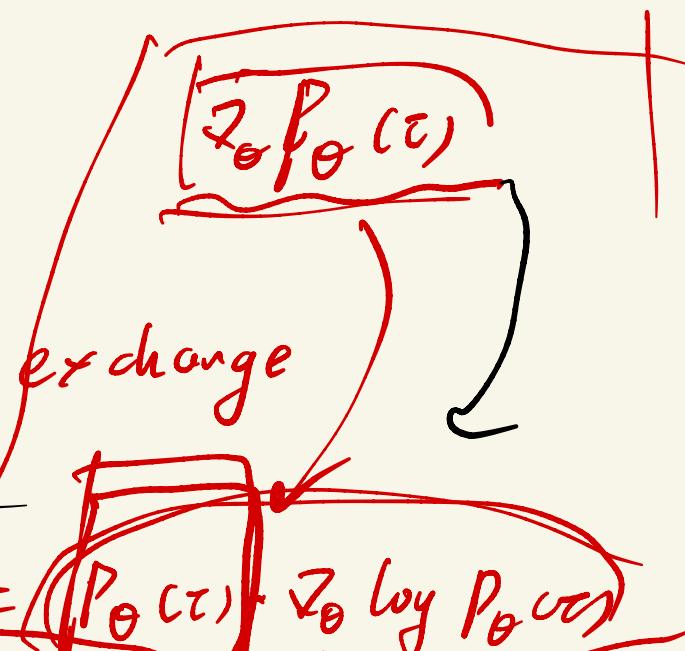
$$= \int_{\mathcal{T}} \mathbb{E}_\theta P_\theta(\tau) f(\tau) d\tau$$

$$= \int_{\mathcal{T}} P_\theta(\tau) \mathbb{E}_\theta \log P_\theta(\tau) f(\tau) d\tau$$

$$= \mathbb{E}_{\tau \sim P_\theta} [\mathbb{E}_\theta \log P_\theta(\tau) / f(\tau)]$$

$$= P_\theta(\tau) \cdot \frac{1}{P_\theta(\tau)} \cdot \mathbb{E}_\theta P_\theta(\tau)$$

$$= \mathbb{E}_\theta P_\theta(\tau)$$



$z$  can be discrete

1. policy gradient
2. Gumbel softmax

$\mathbb{E}_{z \sim q(z)} \mathcal{L}(\theta)$

discrete