

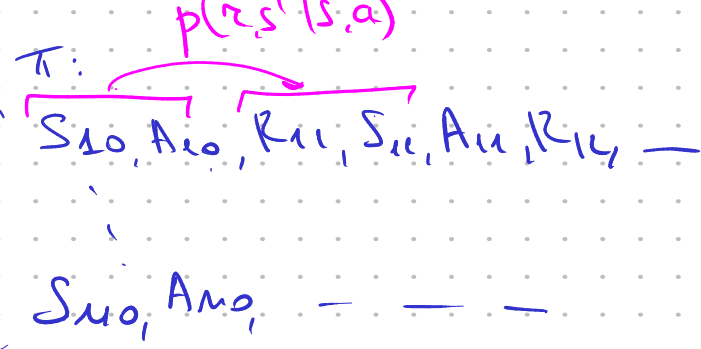
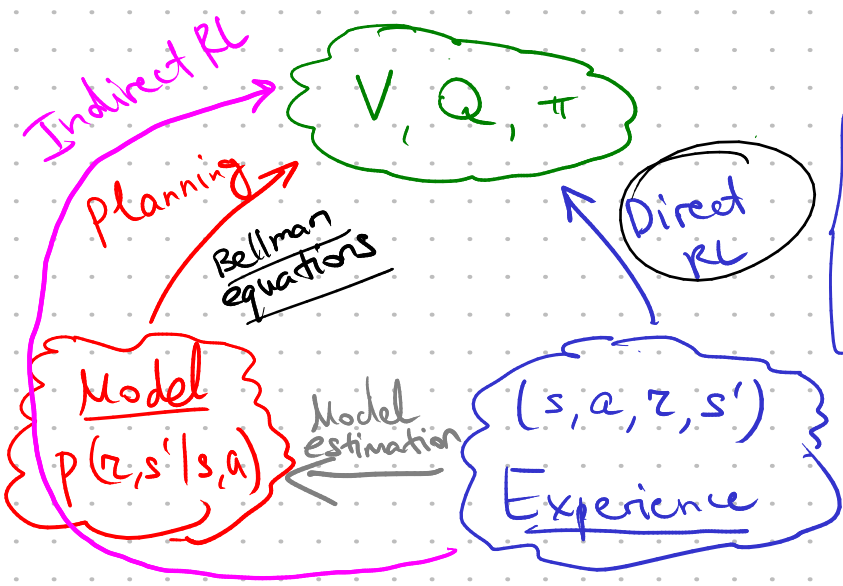
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$V_\pi(s) = E_\pi [G_t | S_t = s]$$

$$Q_\pi(s, a) = E_\pi [G_t | S_t = s, A_t = a]$$

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

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



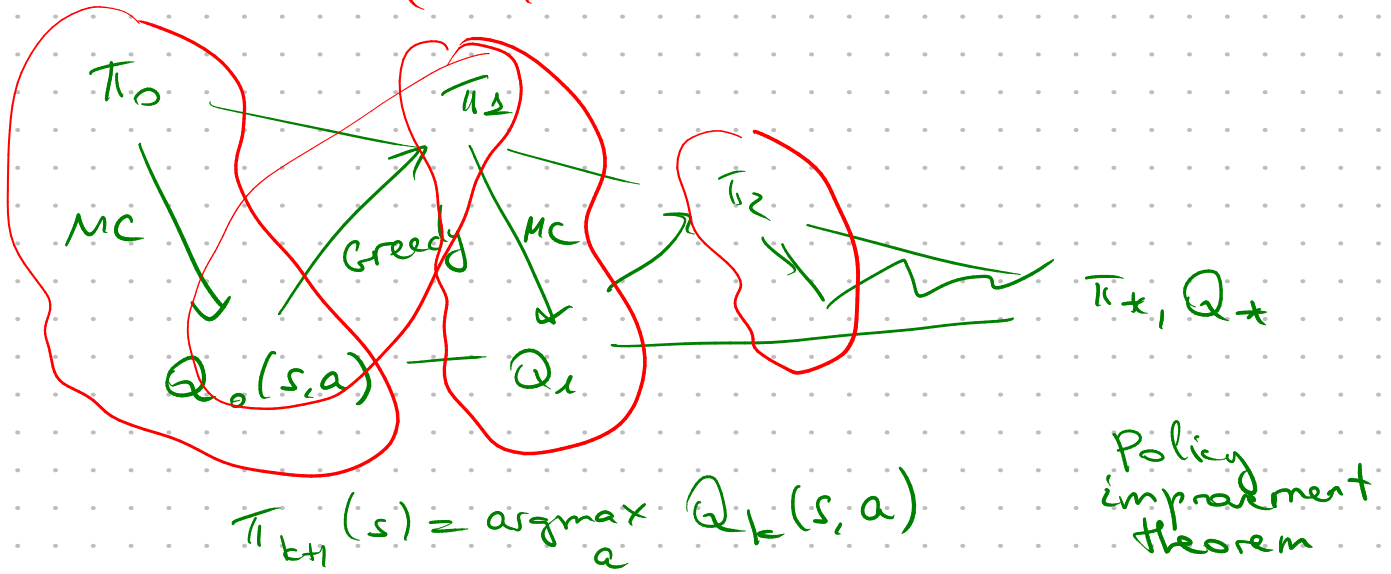
$$V_\pi(s) = E_\pi [G_t | S_t = s] \approx \text{Avg} (G_t | S_t = s, \pi)$$

- ① Monte Carlo estimation for  $V$
- init
  - repeat:
    - gen  $S_0, A_0, R_1, S_1, A_1, \dots, S_t, A_t, R_{t+1}, S_{t+1}, \dots, S_T$
    - no episode  $\tau$
    - $G := 0$
    - for  $t = T-1, T-2, \dots, 0$ 
      - $G := \gamma G + R_{t+1}$
      - each time, so go through  $G$  to returns( $S_t$ )
      - $V(S_t) := \text{Avg}(\text{Returns}(S_t))$

— " — for Q:

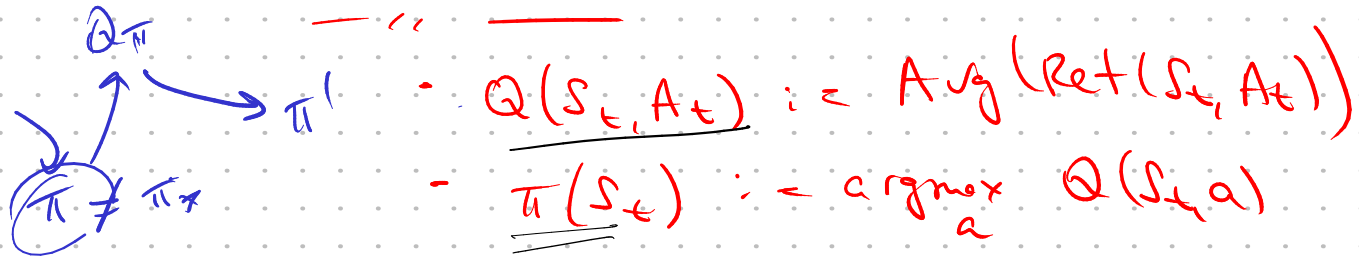
Exploring starts  
 $\forall (s, a) \quad p(s, a) > 0$

— " —  $Q(s_t, A_t)$  —



② MC control with exploring starts

— " —



③ On-policy MC control

— " — without Expl. starts

— " —  $Q(s_t, A_t) := \text{Avg}(\text{Ret}(\dots))$

$\pi(s_t) = \begin{cases} \operatorname{argmax}_a Q(s_t, a) & \text{if } \text{Rep.} > 1 - \epsilon \\ \text{uniform} & \text{if } \text{Rep.} \leq 1 - \epsilon \end{cases}$

$\epsilon$ -soft

$$\pi(a|s) = \frac{\epsilon}{|A|} + (1 - \epsilon) \cdot [a = \operatorname{argmax}_{a'} Q(s, a')] ]$$

Imp. guarantee Policy Improvement. Then also  $\epsilon$ -soft cooperation

$\pi_{k+1}$ -on-poli.  
cyesu  $\epsilon$ -soft

# ④ Off-policy MC control

policy:  $\pi$  →  $\pi'$   
 обучаем  $Q_\pi$   $\pi \neq \pi'$

$$Q_\pi(s, a) = E_{\pi'} [G_t | S_t = s, A_t = a]$$

Итерация / сгорание  
 мы хотим сгорать

## Importance sampling

нужно: если  $p(x) > 0$ , то  $q(x) > 0$

$$E_{p(x)} [f(x)] = \int f(x) p(x) dx = \int f \cdot \frac{p}{q} \cdot q dx = E_{q(x)} [f \cdot \frac{p}{q}]$$

$S_t, A_t$ :  $R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots, S_T$

Траектория  $Traj$   $\pi(A_{t+1} | S_{t+1})$

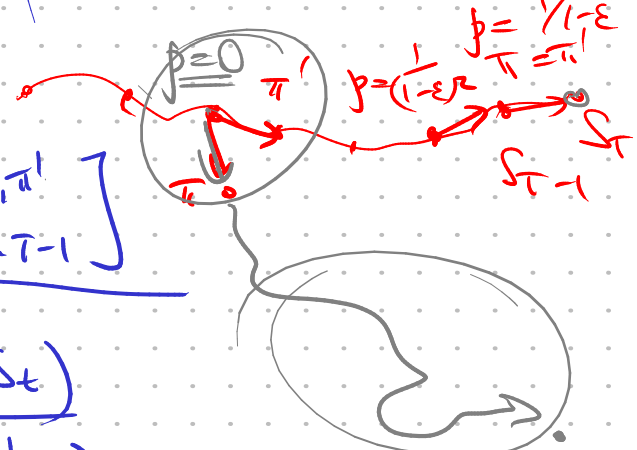
$$p(Traj | \pi) = \underbrace{p(R_{t+1}, S_{t+1} | S_t, A_t)}_{\text{динамика MDP}} \cdot \underbrace{p(A_{t+1} | S_{t+1})}_{\pi(A_{t+1} | S_{t+1})} \cdot p(R_{t+2}, S_{t+2} | S_{t+1}, A_{t+1}) \cdot \dots \cdot p(R_T, S_T | S_{T-1}, A_{T-1})$$

$$p(Traj | \pi') = \underbrace{p(R_{t+1}, S_{t+1} | S_t, A_t)}_{\text{динамика MDP}} \cdot \underbrace{p(A_{t+1} | S_{t+1})}_{\pi'(A_{t+1} | S_{t+1})} \cdot \dots$$

$$\frac{p(Traj | \pi)}{p(Traj | \pi')} = \frac{\pi(A_{t+1} | S_{t+1}) \cdot \pi(A_{t+2} | S_{t+2}) \cdot \dots \cdot \pi(A_{T-1} | S_{T-1})}{\pi'(A_{t+1} | S_{t+1}) \cdot \pi'(A_{t+2} | S_{t+2}) \cdot \dots \cdot \pi'(A_{T-1} | S_{T-1})}$$

## Off-policy estimation:

$$V_\pi(s) = E_{\pi'} [G_t \cdot J_{t:T-1}^{\pi, \pi'}]$$



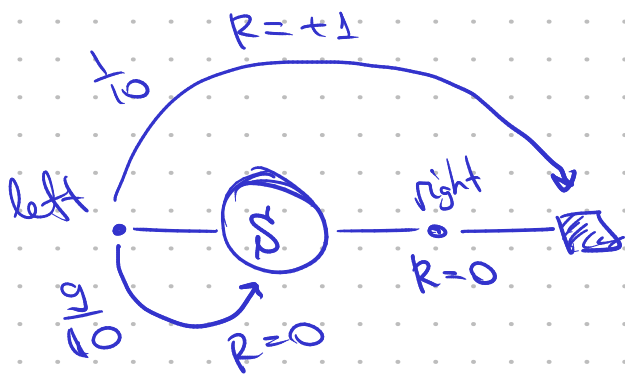
$$J_{t:T-1}^{\pi, \pi'} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{\pi'(A_k | S_k)}$$

$$\text{Returns}(S_t) \leftarrow G_t \cdot J_{t:T-1}^{\pi, \pi'}$$

$$Q_{\pi}(s, a) = E_{\pi'} \left[ G_t \cdot \gamma_{t+1:T-1}^{\pi, \pi'} \right]$$

$$V_{\pi}(s) = \frac{1}{N} \sum_{t=1}^N \gamma_{t:T-1}^{\pi, \pi'} G_t$$

$$V_{\pi}(s) = \frac{\sum_{t=1}^N \gamma_{t:T-1}^{\pi, \pi'} G_t}{\sum_{t=1}^N \gamma_{t:T-1}^{\pi, \pi'}}$$



$$\pi(\text{left}|s) = 1$$

$$\pi'(\text{left}|s) = \pi'(\text{right}|s) = \frac{1}{2}$$

$$V_{\pi}(s) = E_{\pi'} \left[ G \cdot \frac{\prod_t \pi(A_t|S_t)}{\prod_t \pi'(A_t|S_t)} \right]$$

$$\text{Var}[X] = E[X^2] - (EX)^2$$

$$E_{\pi'} \left[ \left( G \cdot \prod_{t=0}^{T-1} \frac{\pi(A_t|S_t)}{\pi'(A_t|S_t)} \right)^2 \right] = \sum_{\text{Traj}} \pi'(\text{Traj}) \cdot \left( \dots \right)^2 =$$

$$= \frac{1}{2} \cdot \frac{1}{10} \cdot 1^2 \cdot \left( \frac{1}{1/2} \right)^2 + \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{1}{10} \cdot 1^2 \cdot \left( \frac{1 \cdot 1}{\frac{1}{2} \cdot \frac{1}{2}} \right)^2 +$$

$$+ \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{1}{10} \cdot 1 \cdot \left( \frac{1}{\left(\frac{1}{2}\right)^3} \right)^2 =$$

$$= \frac{1}{10} \cdot \sum_{k=0}^{\infty} \left( \frac{9}{10} \right)^k \cdot \left( \frac{1}{2} \right)^{k+1} \cdot 2^{2(k+1)} = \frac{1}{5} \sum_{k=0}^{\infty} \left( \frac{9}{5} \right)^k = \infty$$