

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

$$= \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k$$

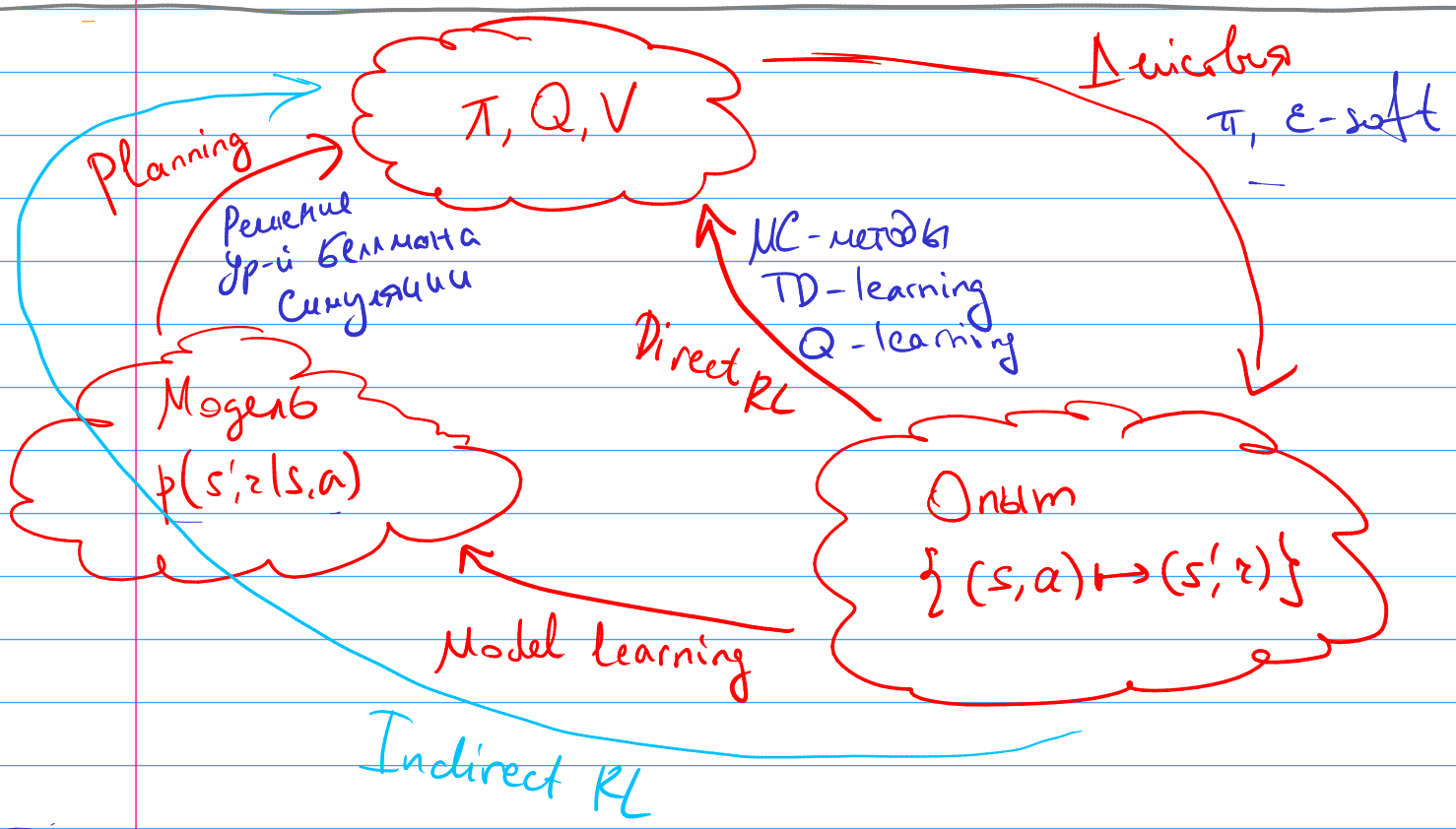
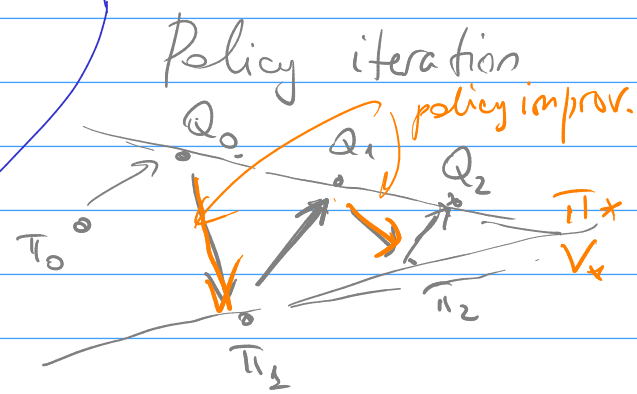
$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

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

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

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

Bellman equations



Monte Carlo

① MC estimation

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s] \approx \frac{1}{N} \sum_{n=1}^N G_t^{(n)}$$

$\pi: s_0, A_0, R_1, S_1, A_1, \dots$

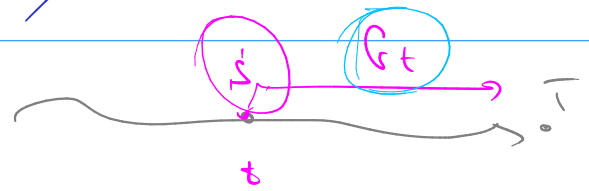
- init $\pi, V_{\pi}(s)$

- loop:

- попарно. эмульг по сфер π :

$s_0, A_0, R_1, S_1, A_1, \dots, S_T$

MC estimation



- central $G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad \forall t: \quad G_T = 0$
- $\forall t$ го состояния G_t $G_{t-1} = R_t$
- $\forall t$ го состояния G_t $G_{t-2} = \gamma R_t + R_{t-1}$
- $V_\pi(S_t) = \text{Avg}(\text{Returns}[S_t])$ $G_t = \gamma G_{t+1} + R_{t+1}$

② On-policy MC control

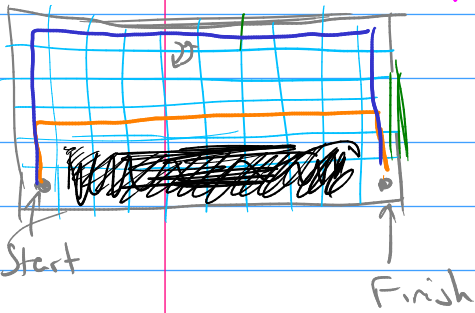
- init π -

- loop

- " " " " $G_t, \text{Returns}[S_t, A_t]$

- $Q(S_t, A_t) = \text{Avg}(\text{Returns}[S_t, A_t])$

- $\forall t \quad \pi(a|S_t) = \begin{cases} 1-\epsilon, & a = \text{argmax}_a Q(S_t, a) \\ \epsilon, & a \sim \text{Unif} \end{cases}$



$\begin{cases} \pi(a|S_t) = \text{argmax}_a Q(S_t, a) \\ b(a|S_t) = \begin{cases} 1-\epsilon \\ \epsilon \end{cases} \end{cases}$ behaviour

③ Off-policy MC control : - no policy, no b

- обьект не Q_θ, a, Q_π

Importance sampling

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

$$Q_\theta(S_t, a) = E_\theta[G_t | \dots]$$

$$E_{p(\bar{x})}[f(\bar{x})]$$

$$\bar{x}_n \sim q(\bar{x})$$

$$\int q(\bar{x}) \dots d\bar{x}$$

$$\int p(\bar{x}) f(\bar{x}) d\bar{x} = \int \frac{p(\bar{x})}{q(\bar{x})} f(\bar{x}) q(\bar{x}) d\bar{x} =$$

$$= E_{q(\bar{x})} \left[f(\bar{x}) \frac{p(\bar{x})}{q(\bar{x})} \right] \approx \frac{1}{N} \sum_n \frac{p(\bar{x}_n)}{q(\bar{x}_n)} f(\bar{x}_n)$$

Екау $q(\bar{x}) = 0, \pi$
и $p(\bar{x}) = 0$

$$s_t, A_t : R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots, s_T \leftarrow \text{Traj}$$

$$\begin{aligned} \Pr[\text{Traj} | \pi, s_t, A_t] &= p(s_{t+1}, R_{t+1} | s_t, A_t) \pi(A_{t+1} | s_{t+1}) \\ &\quad p(s_{t+2}, R_{t+2} | s_{t+1}, A_{t+1}) \pi(A_{t+2} | s_{t+2}) \dots \\ &= \prod_{k=t+1}^T \underbrace{p(s_k, R_k | s_{k-1}, A_{k-1})}_{\text{orange underline}} \pi(A_k | s_k) \end{aligned}$$

$$\Pr[\dots | b \dots] = \prod_{k=t+1}^T \underbrace{p(s_k, R_k | s_{k-1}, A_{k-1})}_{\text{orange underline}} b(A_k | s_k)$$

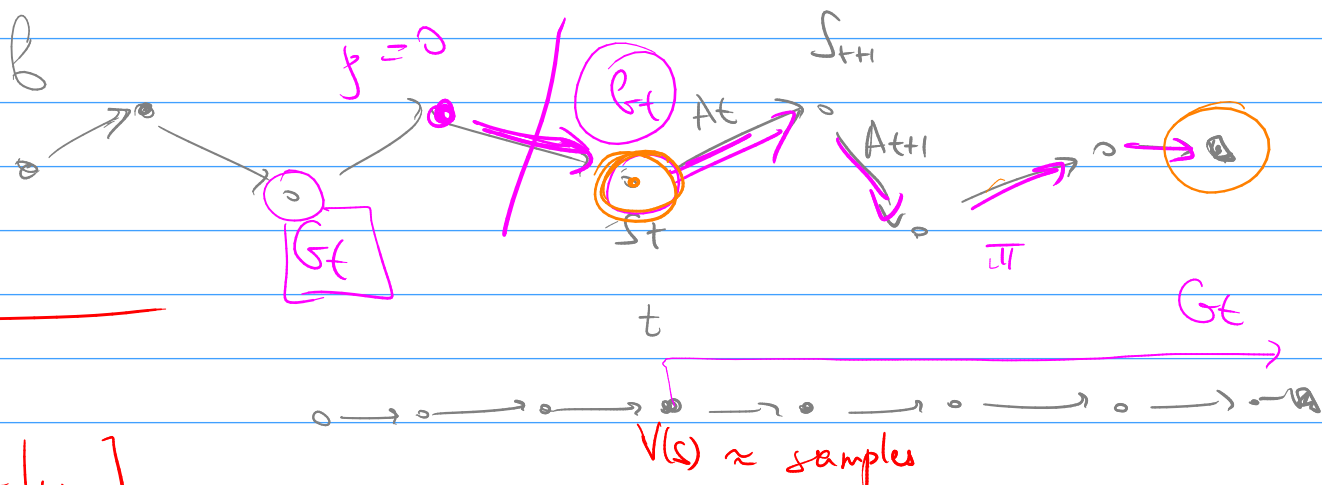
$$J_{t:T-1} = \frac{\Pr[\text{Traj} | \pi]}{\Pr[\text{Traj} | b]} = \prod_k \frac{\pi(A_k | s_k)}{b(A_k | s_k)}$$

$$G_t : G_T = 0, \quad G_t = \gamma G_{t+1} + R_{t+1}$$

$$J_T = 1, \quad J_t = J_{t+1} \cdot \frac{\pi(A_{t+1} | s_{t+1})}{b(A_{t+1} | s_{t+1})}$$

$$b : s_0, \dots, R_T, s_T$$

$$Q_t = E_{\pi}[\dots] = E_b[J] \approx \frac{1}{N} \sum_{i=1}^N G_t \cdot J_t$$



V, Q

$E_{\pi}[G_t | \dots]$

$V(s) \approx \text{samples}$

$V = E_{\pi}[\dots] \approx \text{samples}$

TD Learning



$$V(s_t) = \mathbb{E}_{\pi} [G_{t+1} | s_t] = \mathbb{E}_{\pi} [\gamma G_{t+1} + R_{t+1} | s_t] = \mathbb{E}_{\pi} [R_{t+1} | s_t] + \gamma \mathbb{E}_{\pi} [G_{t+1} | s_t]$$

Sample step (t+1)



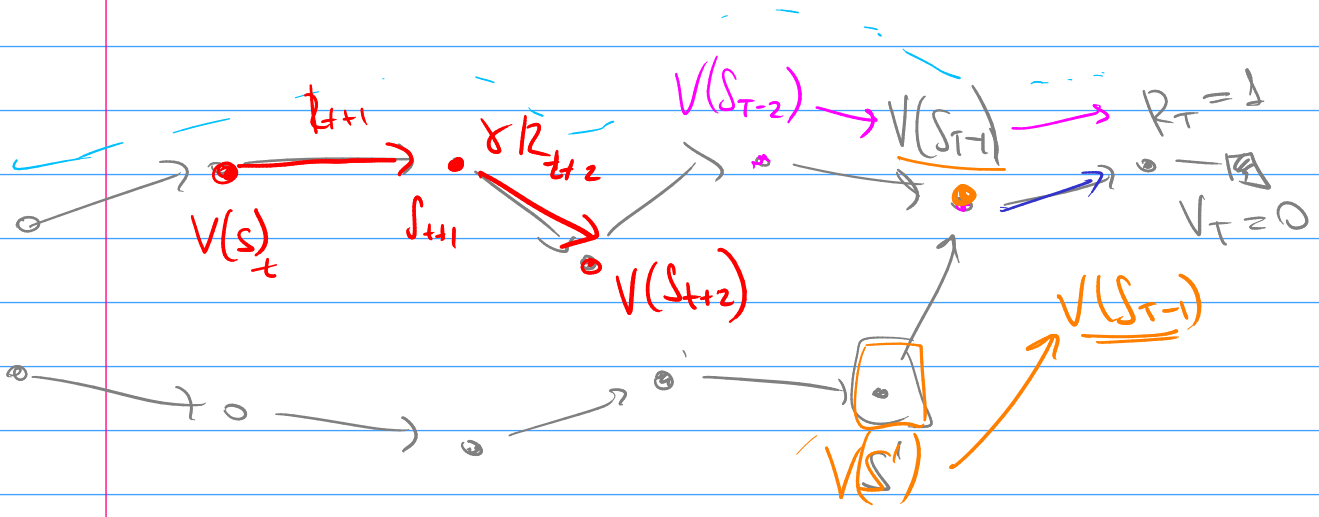
$$V(s_{t+1}) \approx R_{t+1} + \gamma \mathbb{E}_{\pi} [G_{t+1} | s_{t+1}] \approx R_{t+1} + \gamma V(s_{t+1})$$

TD estimation

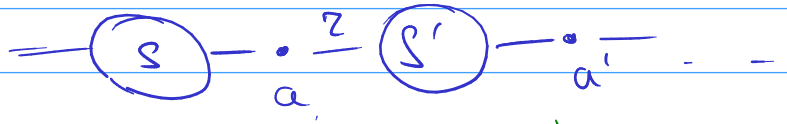
$$V(s_t) := V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

target

$$V(s_t) \approx G_t = R_{t+1} + \gamma (R_{t+2} + \dots) \approx V(s_{t+1})$$



On-policy TD control



Source

- init
- loop

$$Q(s, a, r, s', a') - Q(s, a) := Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$$

TD-target

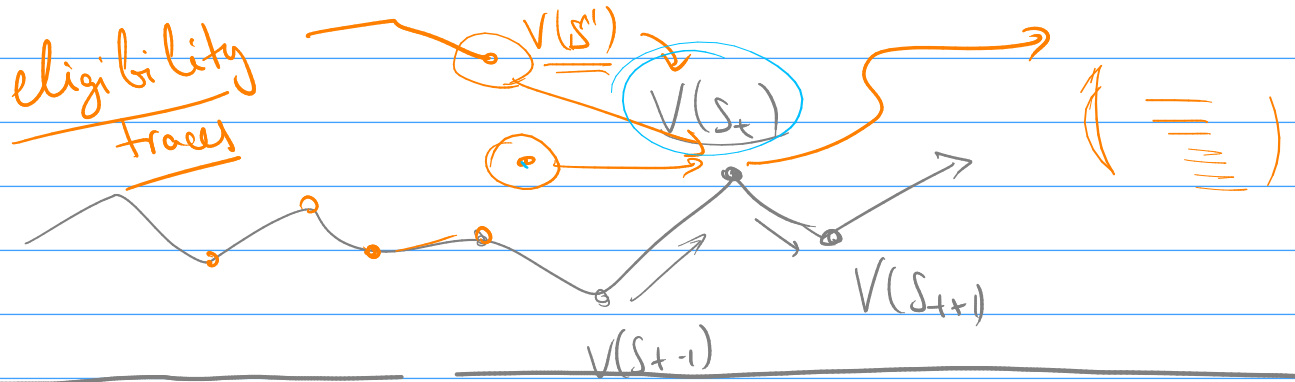
- π - ϵ -HCGH - $q \approx Q$

Off-policy TD-control — Q-learning

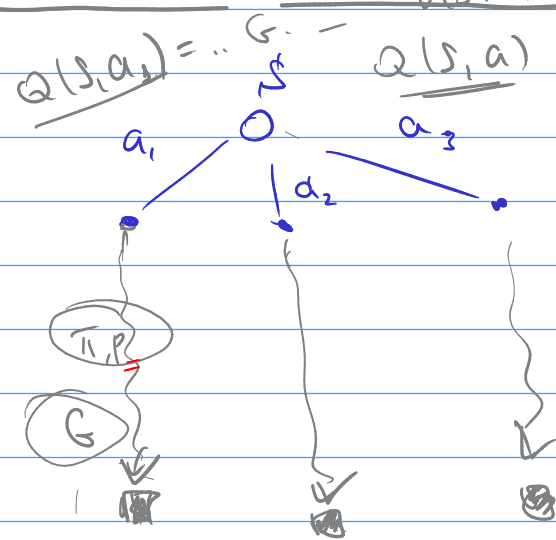
1989

Q_x

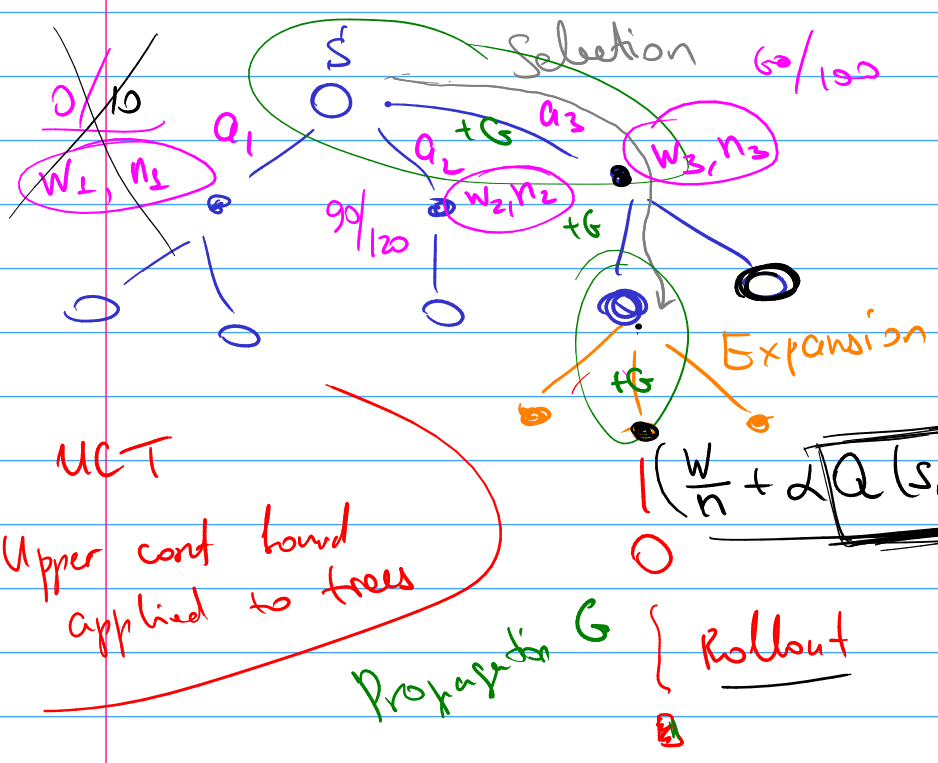
$$Q(s,a) := Q(s,a) + \alpha [r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$



Rollout



Decision-time planning



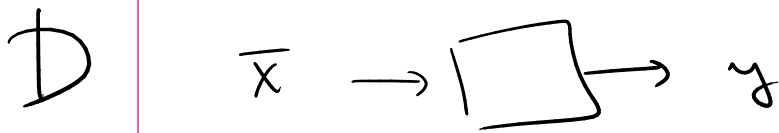
MCTS
Monte-Carlo
Tree Search

800 rollouts

UCT

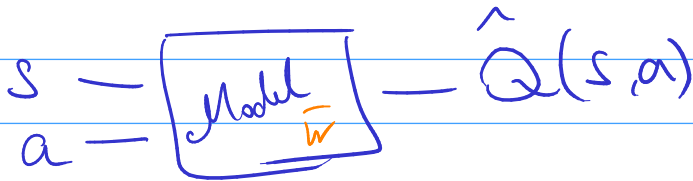
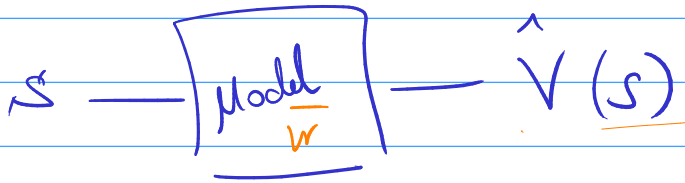
Upper conf bound applied to trees

Propagation G } Rollout



$$s = \bar{x}$$

$$\hat{V}(s) = \sigma(\bar{w}^T \bar{x})$$



$$s: V(s) \quad E = \frac{1}{2} (V(s) - \hat{V}(s))^2 \xrightarrow{\bar{w}} \min$$

$$\nabla_{\bar{w}} E = - (V(s) - \hat{V}(s)) \cdot \nabla_{\bar{w}} \hat{V}$$

$$\bar{w} := \bar{w} - \alpha (V - \hat{V}) \nabla_{\bar{w}} \hat{V}$$

Gradient MC:

$$\bar{w} := \bar{w} - \alpha (G_t - \hat{V}(s'_t, \bar{w})) \nabla_{\bar{w}} \hat{V}(s_t, \bar{w})$$

Semi-gradient TD(0):

$$\bar{w} := \bar{w} - \alpha (R_{t+1} + \gamma \hat{V}(s_{t+1}, \bar{w}) - \hat{V}(s_t, \bar{w})) \nabla_{\bar{w}} \hat{V}(s_t, \bar{w})$$

Semi-gradient Q-learning

$$\bar{w} := \bar{w} - \alpha (R_{t+1} + \gamma \max_{a'} \hat{Q}(s_{t+1}, a') - \hat{Q}(s_t, A_t)) \nabla_{\bar{w}} \hat{Q}(s_t, A_t, \bar{w})$$