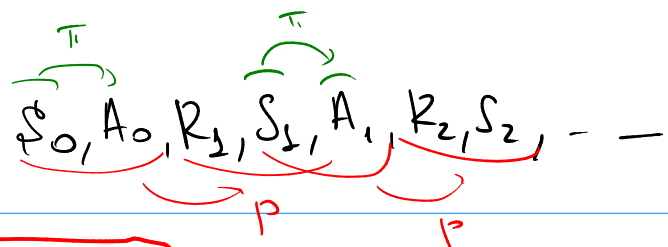
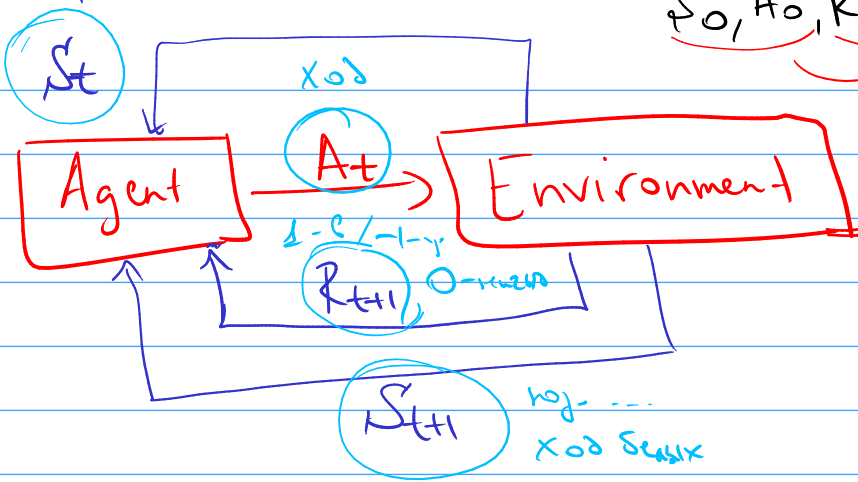


① MDP

$\mathcal{S}, A(s), r \in \mathcal{R}$



вопросы  
покупать  
50 товаров  
3-4 раза  
хотел сказать



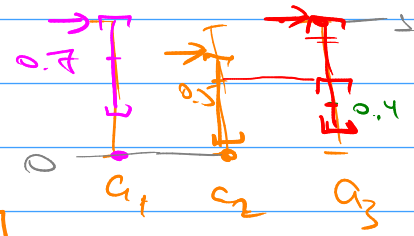
MDP Markov decision process

dynamics: 
$$p(s', z | s, a) = P_z [S_{t+1} = s', R_{t+1} = z | S_t = s, A_t = a]$$

$$\forall s, a \sum_{s', z} p(s', z | s, a) = 1$$

strategy: 
$$\pi(a|s) = P_z [A_t = a | S_t = s], \forall s \sum \pi(a|s) = 1$$

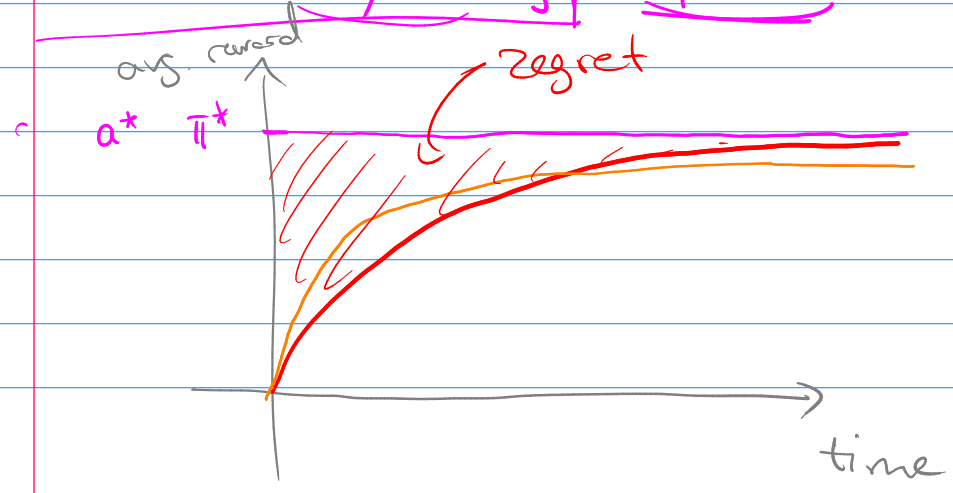
② Multicomed bandits  $|\mathcal{S}| = 1$



$\pi(a) = P_r [A_t = a]$   
 $p(z|a) = P_r [R_{t+1} = z | A_t = a]$

$A = a_1, \dots, a_k$

$D = \{(A_t, R_{t+1})\}$  exploration vs. exploitation



### ③ Returns & Value functions

$t: R_{t+1}, R_{t+2}, R_{t+3}, \dots$

episodic:  $G_t = \sum_{\tau=t+1}^T R_\tau$  - return

continuous:  $(?)$   $(T = \infty)$   $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$   
 $= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

discount factor



return

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

$T = \infty \Rightarrow \gamma < 1$   
 $T < \infty \Rightarrow \gamma = 1$

$$G_t = R_{t+1} + \gamma \cdot G_{t+1}$$

### Value functions

$$V_\pi(s) = E_\pi[G_t | S_t = s] = E_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right]$$

$$Q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a] = E_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

$$V_\pi(s) = E_{a \sim \pi} [Q_\pi(s, a)] = \sum_a \pi(a|s) Q_\pi(s, a)$$

$$Q_\pi(s, a) = E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] =$$

$$= \underbrace{E_\pi[R_{t+1} | S_t = s, A_t = a]}_{z(s, a) = E_p[r | s, a]} + \gamma \underbrace{E_\pi[G_{t+1} | S_t = s, A_t = a]}_{E_{s' \sim p(s'|s, a)} [V_\pi(s') | s, a]}$$

$$Q_\pi(s, a) = z(s, a) + \gamma \sum_{s'} p(s' | s, a) V_\pi(s')$$

④ Bellman equations

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s] = E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$V_{\pi}(s) = \sum_a \pi(a|s) \left[ \sum_{s', z} p(s', z | s, a) [z + \gamma V_{\pi}(s')] \right]$$

Ceset. mat. y' - u

$$Q_{\pi}(s, a) = \sum_{s', z} p(s', z | s, a) \left[ z + \gamma \sum_{a'} \pi(a' | s') Q_{\pi}(s', a') \right]$$

Bellman equations

Ceset. mat. y' - u

$\pi_*$  - opt. eq.  $\forall \pi \forall s \quad V_*(s) \geq V_{\pi}(s)$

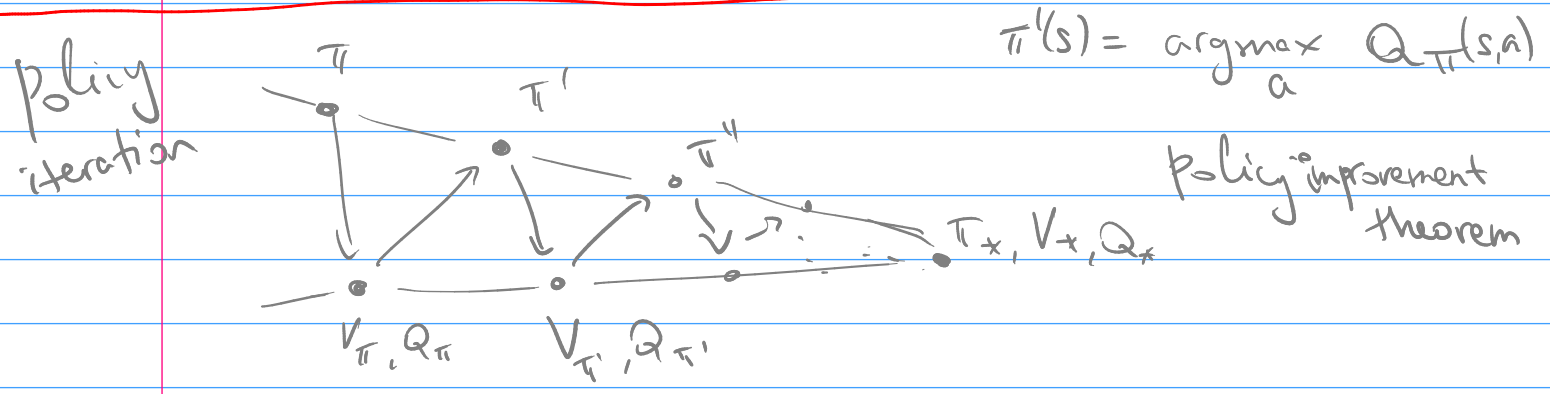
$$V_*(s) = \max_{\pi} V_{\pi}(s), \quad Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

$$V_*(s) = \max_{\pi} V_{\pi}(s) = \max_{\pi} \left[ \sum_a \pi(a|s) \sum_{s', z} p(s', z | s, a) [z + \gamma \max_{\pi} V_{\pi}(s')] \right]$$

Bellman eq

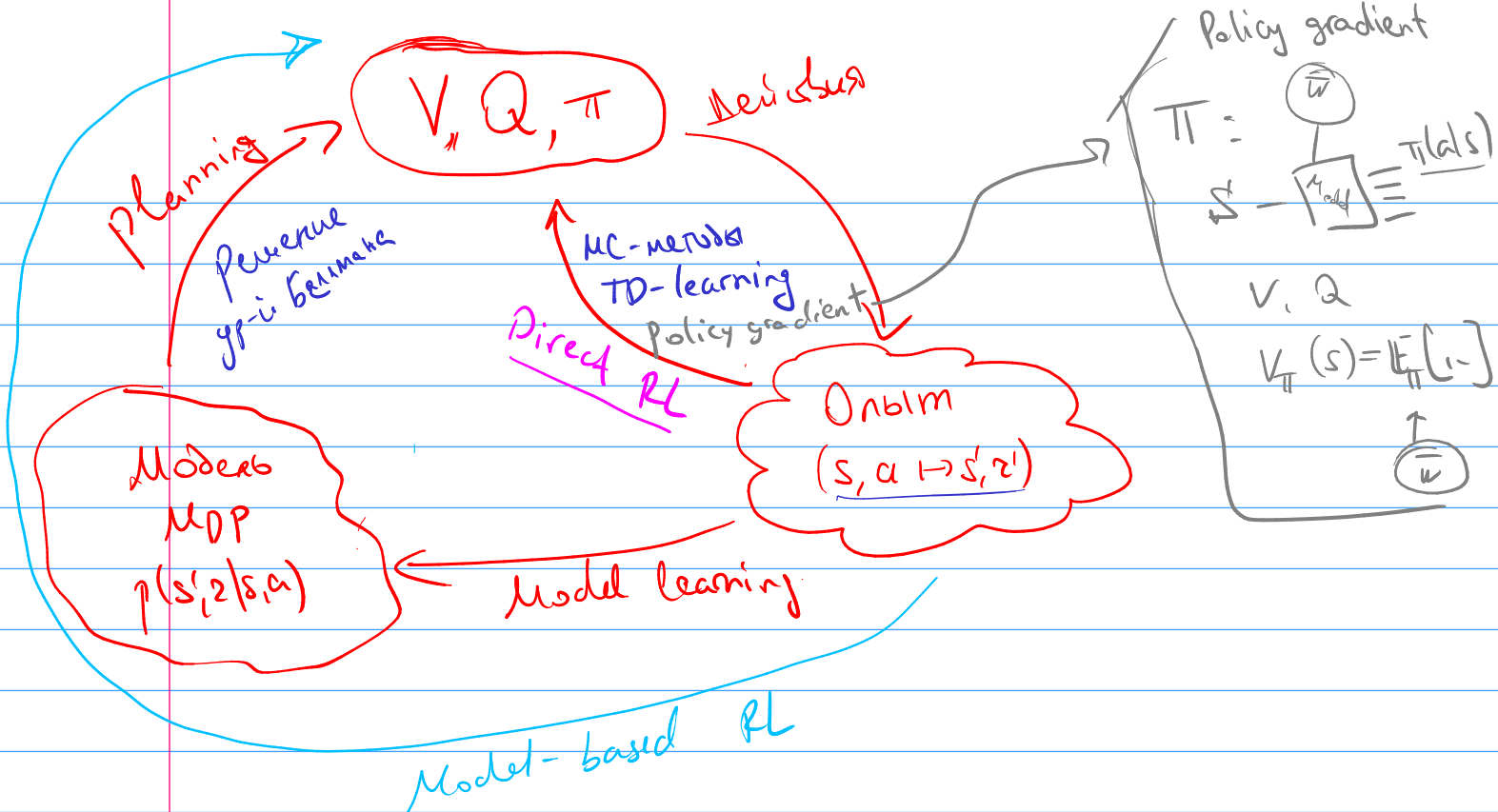
$$V_*(s) = \max_a \left( \sum_{s', z} p(s', z | s, a) [z + \gamma V_*(s')] \right)$$

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



⑤ Monte-Carlo

$$p(z, s' | s, a) = ?$$



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

- loop:

ME evaluation - no policy. no  $\pi$   $S_0, A_0, R_1, S_1, A_1, \dots$

- current  $G_t = \gamma G_{t+1} + R_{t+1}$

-  $G_t$  is cumulative Returns ( $S_t$ )

-  $V_\pi(S_t) := \text{Avg}(\text{Returns}(S_t))$

MC control

$\pi \rightarrow \pi' \rightarrow \dots$

Q( $S_t, A_t$ ) := Avg(Returns( $S_t, A_t$ ))

$\pi(a|S_t)$  =  $\epsilon$ -Greedy(Q)

=  $\begin{cases} 1-\epsilon, & \text{argmax}_a Q \\ \epsilon, & \text{arbitrary} \end{cases}$

On-policy control

updates w.r.t  $\pi$   
и оценок  $Q_\pi$

importance sampling

off-policy control

updates w.r.t  $\pi$   
оценок  $Q_\pi$  и  $Q_{\pi'}$

# ⑥ TD-learning (temporal difference)

- envy  $S_0, A_0, R_1$  -

$$V(S_t) \approx G_t = R_{t+1} + \gamma(R_{t+2} + \gamma^2 R_{t+3} + \dots)$$

TD( $\lambda$ )

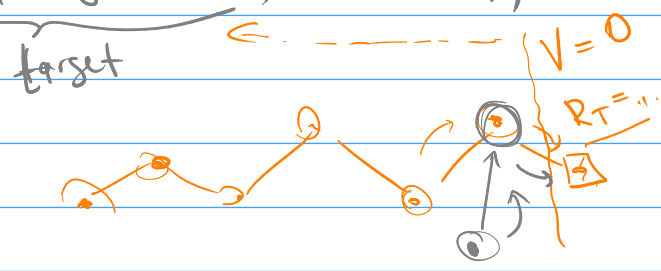
$$V(S'_t) \rightarrow R_{t+1} + \gamma \cdot V(S_{t+1})$$

$$G_{t+1} \approx V(S_{t+1})$$

TD(0)

$$V(S_t) := V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

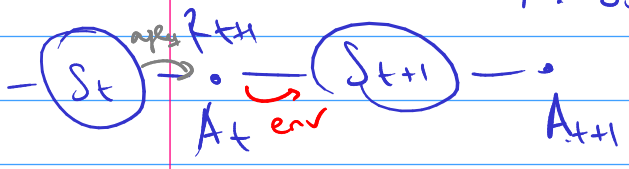
TD evaluation



## TD control

ON-policy - SARSA

$$\pi: S_0, A_0, R_1, S_1, A_1, \dots, G$$



$$Q(S_t, A_t) := Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

$$\pi := \epsilon\text{-Greedy}(Q)$$

## off-policy Q-learning

$$(s, a) \mapsto (s', r)$$

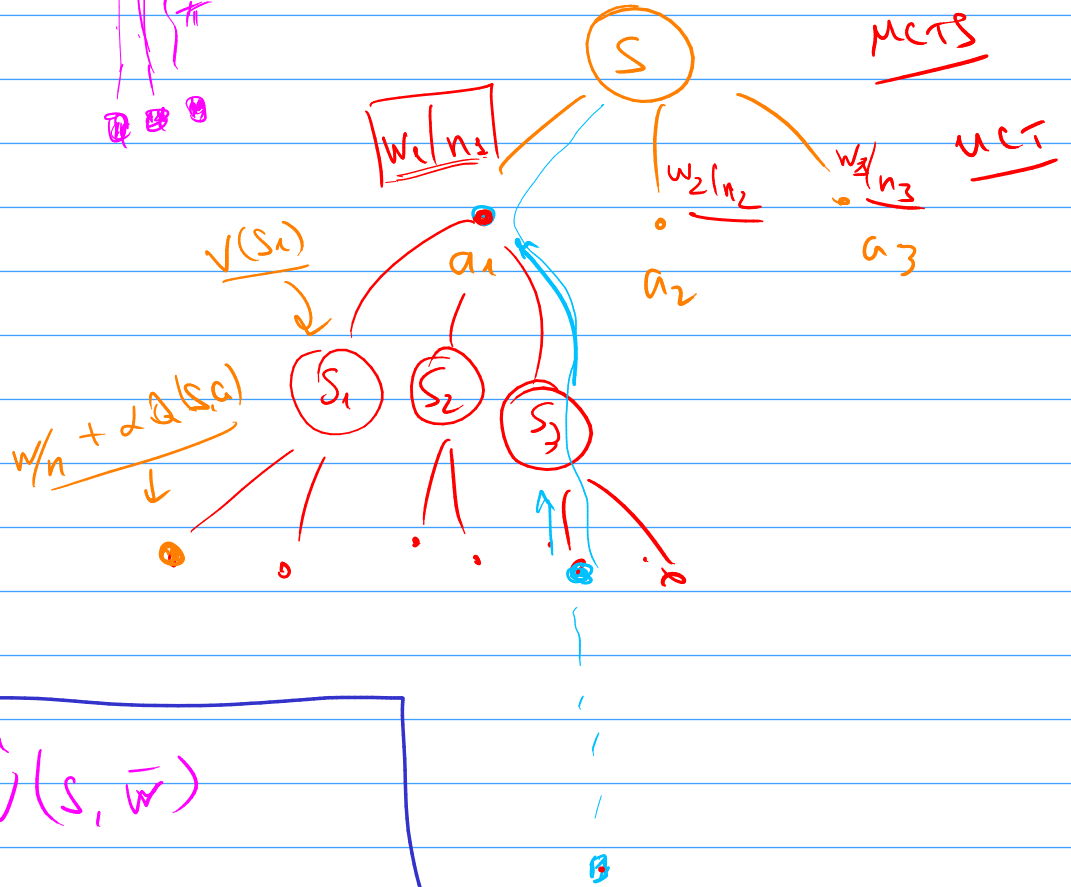
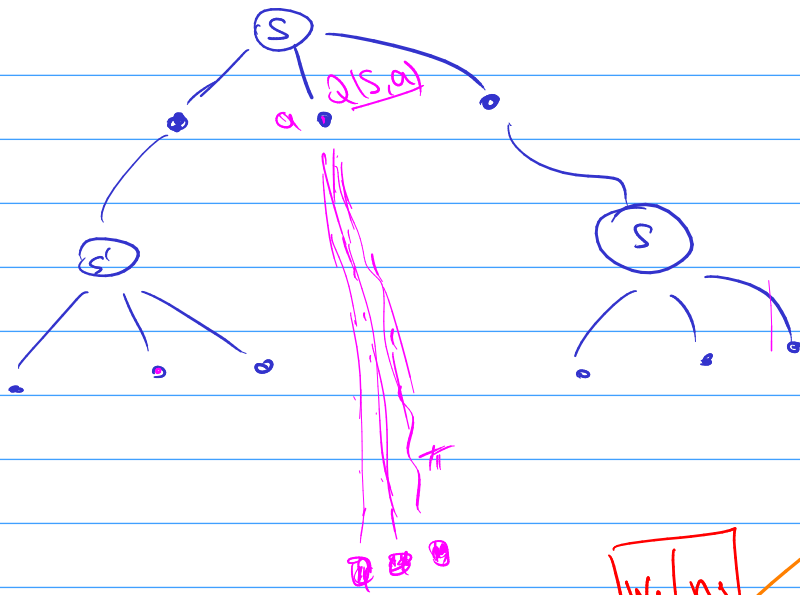
$$Q(s, a) = Q(s, a) + \alpha (r + \gamma \cdot \max_{a'} Q(s', a') - Q(s, a))$$

$$\rightarrow Q_*(s, a)$$

$$\pi := \epsilon\text{-Greedy}(Q)$$

# MCTS - Monte-Carlo Tree Search

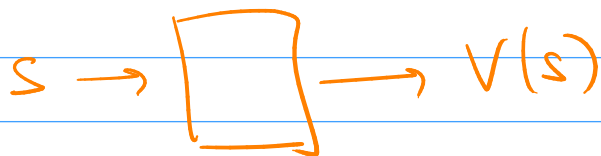
Rollouts



$$V(s) \approx \hat{V}(s, \bar{w})$$

$$\bar{w} := \bar{w} + \alpha \left( R_{t+1} + \gamma \hat{V}(S_{t+1}, \bar{w}) - \hat{V}(S_t, \bar{w}) \right) \cdot \nabla_{\bar{w}} \hat{V} |_{S_t}$$

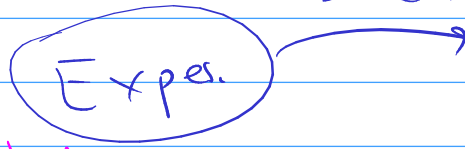
$$\bar{w} := \bar{w} + \alpha \left( R_{t+1} + \gamma \max_{a'} \hat{Q}(S_{t+1}, a', \bar{w}) - \hat{Q}(S_t, a_t, \bar{w}) \right) \nabla_{\bar{w}} \hat{Q}$$



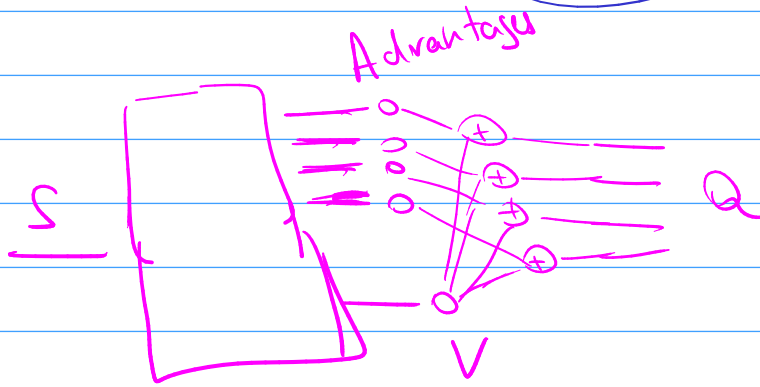
# DQN

deep Q-networks

1) Experience replay



2)



3) Double DQN