

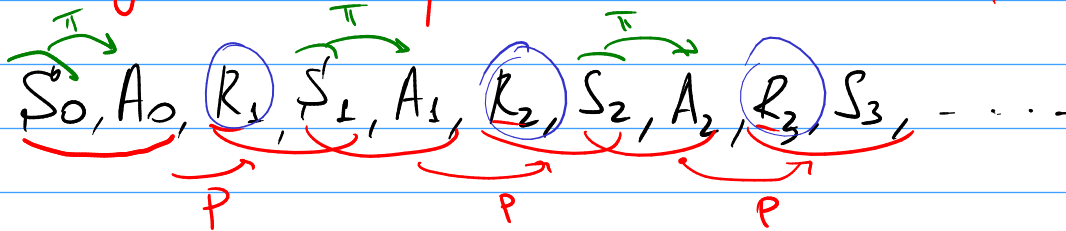
MDP ← Markov Decision Process:

$S, A$

Strategy:  $\pi: S \rightarrow \text{Prob}(A)$

$$\pi(s, a) = \text{Pr}[A_t = a | S_t = s]$$

Dynamics:  $p(r, s' | s, a) = \text{Pr}[R_{t+1} = r, S_{t+1} = s' | S_t = s, A_t = a]$

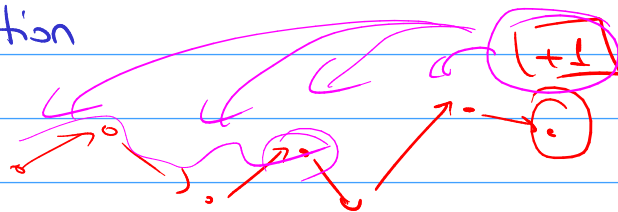


$$\sum_{t=1}^{\infty} R_t \cdot \gamma^{t-1} \rightarrow \max$$

$|\gamma| < 1$

1) Exploration vs. exploitation

2) Credit assignment



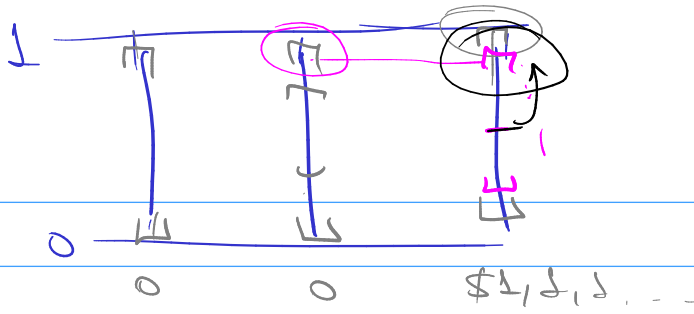
Multiarmed bandits  $|\mathcal{S}| = 1$

$A$   $a_1, \dots, a_n$

~~exploration~~



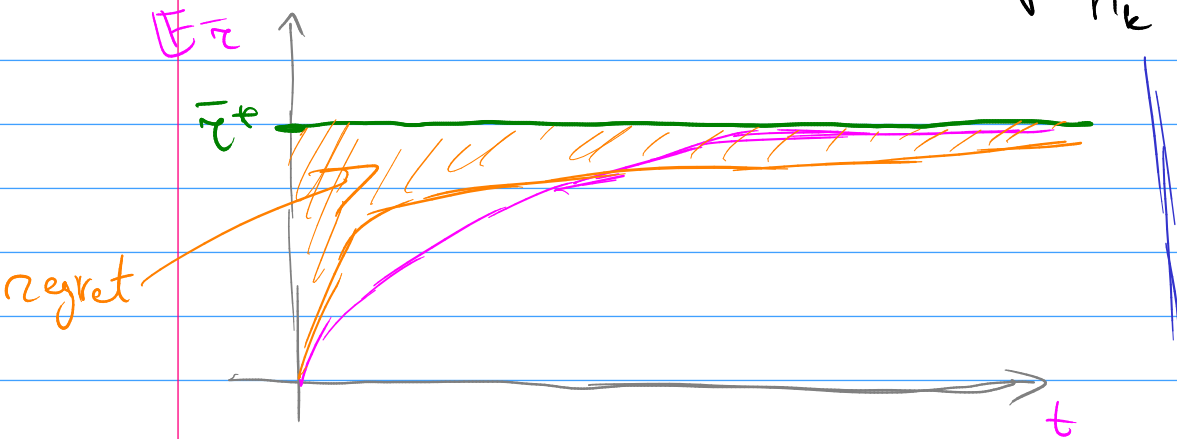
~~exploration~~



$$Pr_i(a_k) = \bar{r}_k + \dots$$

$$N, n_k, \sum n_k = N,$$

$$\text{UCB1: } Pr_i(a_k) = \bar{r}_k + c \sqrt{\frac{\log N}{n_k}}$$



Thm  $c = \sqrt{2}$   
 $O(\log T)$

$$G_t = \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k \quad \text{return}$$



Value function

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} \left[ \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k | S_t = s \right]$$

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[ \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k | S_t = s, A_t = a \right]$$

$$V_{\pi}(s) = \sum_a \pi(a|s) Q(s, a)$$

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t | s, a] = \mathbb{E}_{\pi} \left[ R_{t+1} + \left( \sum_{k=t+2}^{\infty} \gamma^{k-t-1} R_k \right) | s, a \right]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | s, a] = \underbrace{\mathbb{E}_{\pi} [R_{t+1} | s, a]}_{r(s, a)} + \gamma \underbrace{\mathbb{E}_{\pi} [G_{t+1} | s, a]}_{V_{\pi}(s')}$$

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

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

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

Bellman equations

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

$$V_{*}(s) = \max_{\pi} \mathbb{E}_{\pi} [G_t | S_t = s]$$

$$\pi_{*}(s) = \operatorname{argmax}_a Q_{*}(s, a)$$

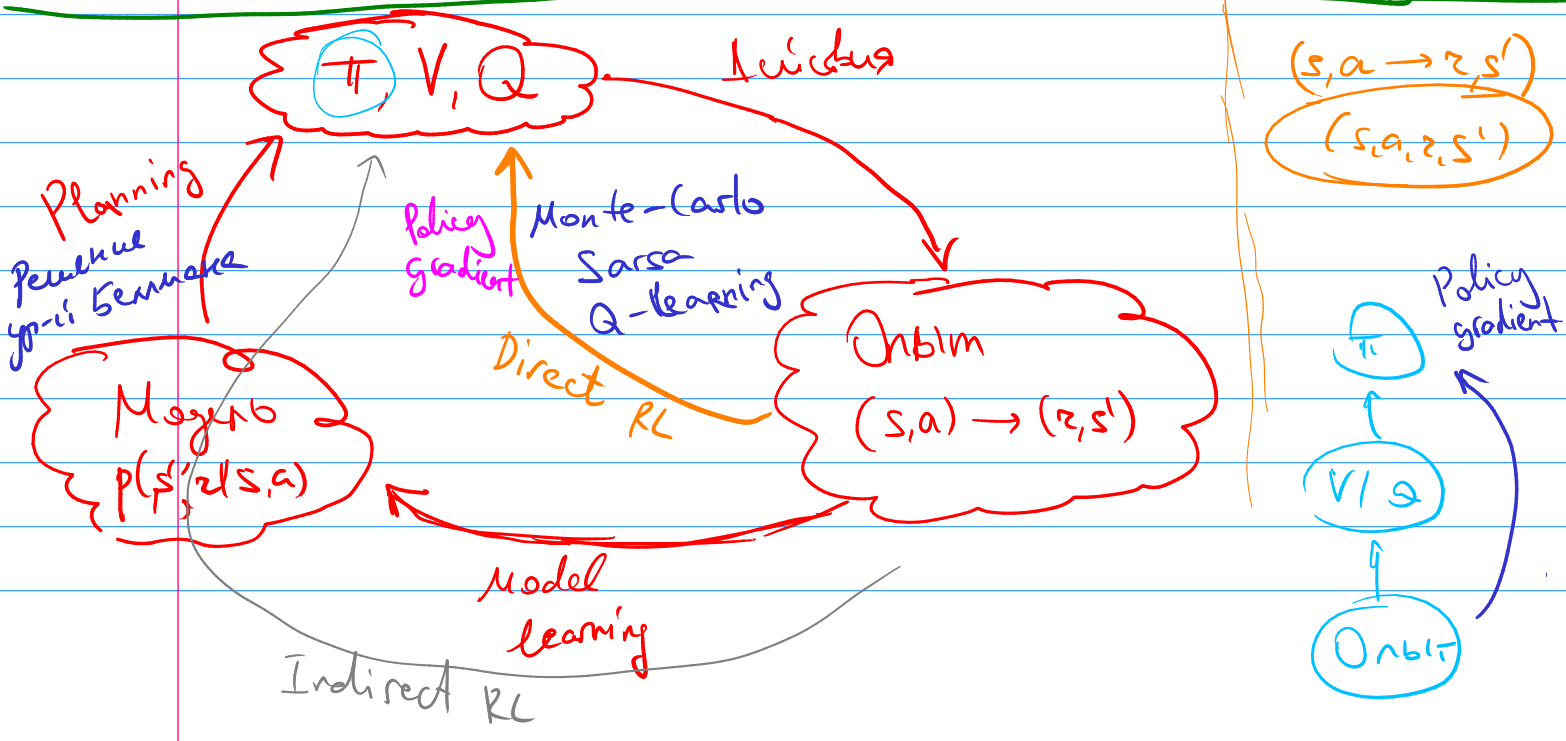
$$Q_{*}(s, a) = \max_{\pi} \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a]$$

$$V_{*}(s) = \max_{\pi} \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] =$$

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

Bellman equation

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



# 1) Monte-Carlo

Estimation

$$V_{\pi} = \mathbb{E}_{\pi} [G_t | S_t = s] \approx \frac{1}{R} \sum_{z=1}^R G_t^{(z)}$$

- no env. anyd no  $\pi$

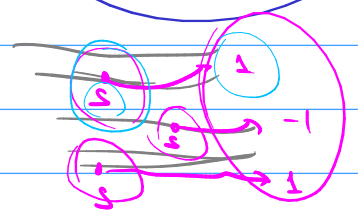
-  $\forall t \quad G_T = 0$

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

-  $\forall t \quad \text{Returns}[S_t]. \text{appear}(G_t)$

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

$Q_{\pi}(S_t, A_t)$



Control

-  $\pi$

-  $\pi$

on-policy

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

MC control

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

$\epsilon$ -soft,  $\epsilon$ -greedy

off-policy

- no env. no  $b$  (behaviours)

- env.  $V_{\pi}(s), Q_{\pi}(s, a)$

Importance sampling

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] \approx \frac{1}{R} \sum_{z=1}^R G_t^{(z)}$$

control no  $\pi, a$  no  $b$

$q(\bar{x})$

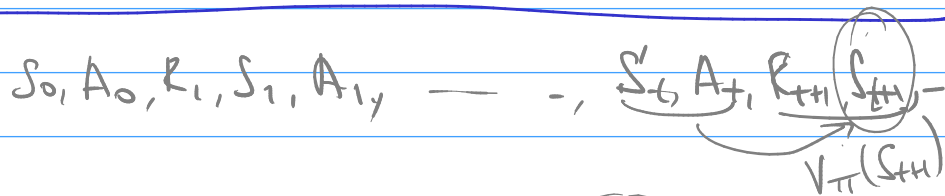
$$\mathbb{E}_{p(\bar{x})} [f(\bar{x})] = \int f(\bar{x}) p(\bar{x}) d\bar{x} =$$

$$= \int f(\bar{x}) \frac{p(\bar{x})}{q(\bar{x})} q(\bar{x}) d\bar{x} = \mathbb{E}_q \left[ f(\bar{x}) \cdot \frac{p(\bar{x})}{q(\bar{x})} \right]$$

importance weights

# 2) TD-learning

temporal difference

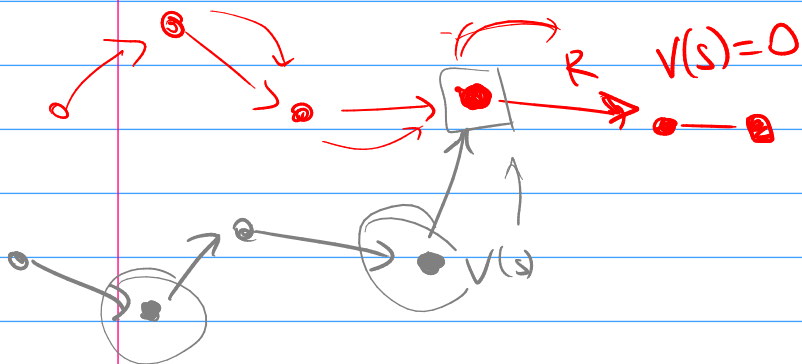


$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

TD(x)

$$V_{\pi}(S_t) \approx R_{t+1} + \gamma V_{\pi}(S_{t+1})$$

$$TD(0): V_{\pi}(S_t) := V_{\pi}(S_t) + \alpha [ \underbrace{R_{t+1} + \gamma V_{\pi}(S_{t+1})}_{\text{target}} - \underbrace{V_{\pi}(S_t)} ]$$



On-policy TD control - Sarsa

$$(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$$

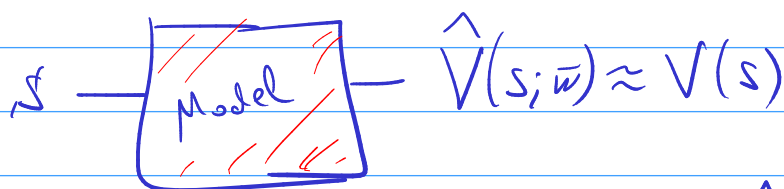
$$Q(S_t, A_t) := Q(S_t, A_t) + \alpha [ R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) ]$$

$$\pi(S_t) = \epsilon\text{-greedy exp. no } Q$$

Off-policy TD-control Q-Learning

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

$$\rightarrow Q_*(S_t, A_t)$$



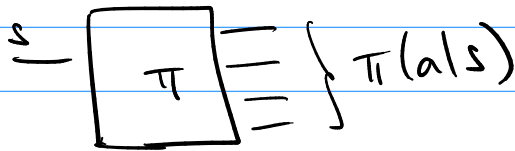
Gradient<sub>MC</sub> 
$$\bar{w} := \bar{w} + \alpha (G_t - \hat{V}(S_t, \bar{w})) \cdot \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$$

Semi-gradient TD

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

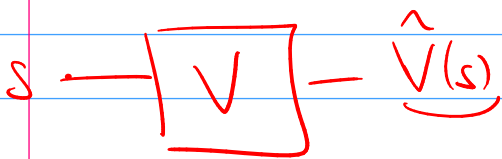
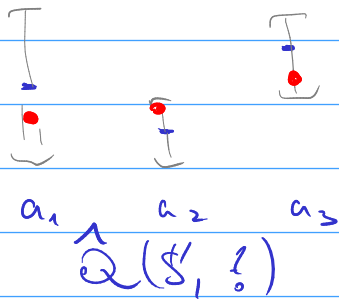
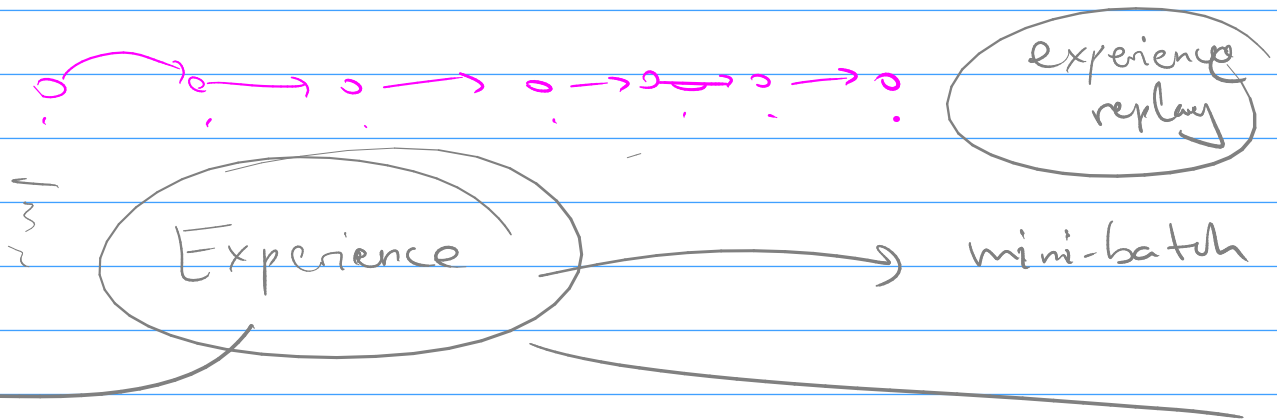
Policy gradient

$$\pi(a|s, \theta)$$



$$J(\theta) = V_{\pi_{\theta}}(s_0)$$

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \alpha \cdot \nabla_{\bar{\theta}} J(\bar{\theta}_t)$$



Dueling networks

advantage

