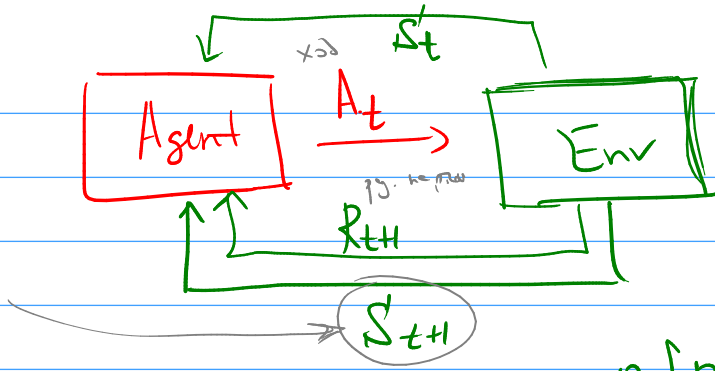


MDP - Markov Decision Process

$$S \rightarrow A$$



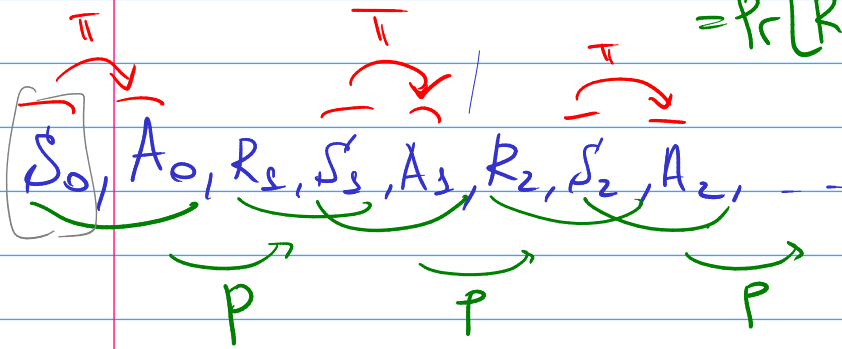
Стратегия:

$$\pi(a|s) = \Pr[A_t = a | S_t = s]$$

Dynamics:

$$p(r, s' | s, a) =$$

$$= \Pr[R_{t+1} = r, S_{t+1} = s' | S_t = s, A_t = a]$$



$$p(r|s,a) = \int p(r,s'|s,a) ds'$$

Episodic tasks

Continuous tasks

$$S_0, A_0, R_1, \dots, (R_T, S_T)$$

$$S_0, A_0, R_1, \dots, S_t, A_t, \dots$$

$$G = R_1 + R_2 + \dots + R_T \rightarrow \max$$

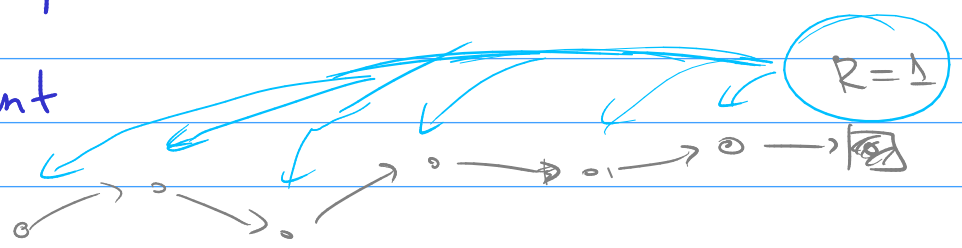
$$G = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots = \sum_{k=1}^{\infty} \gamma^{k-1} R_k, \gamma < 1 \rightarrow \max$$

$$\gamma = 1 \quad G = \sum_{k=1}^{\infty} \gamma^{k-1} R_k$$

$$\mathbb{E}_{\pi, P} \left[\sum_{k=1}^{\infty} \gamma^{k-1} R_k \right] \xrightarrow{\pi} \max$$

1) Exploration vs. exploitation

2) Credit assignment



Multicarmed bandits

$$|\mathcal{A}| = 2$$

a_1, a_2, \dots, a_n

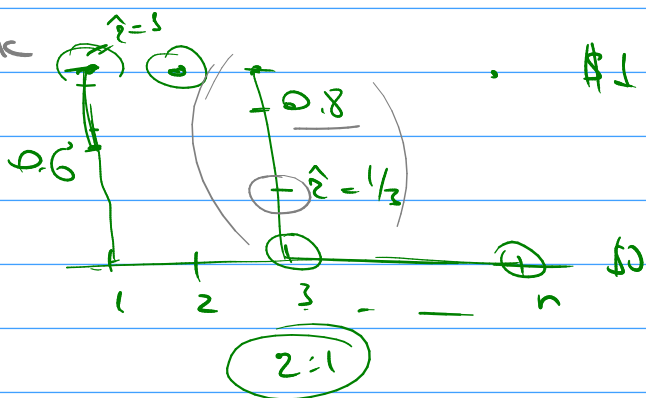
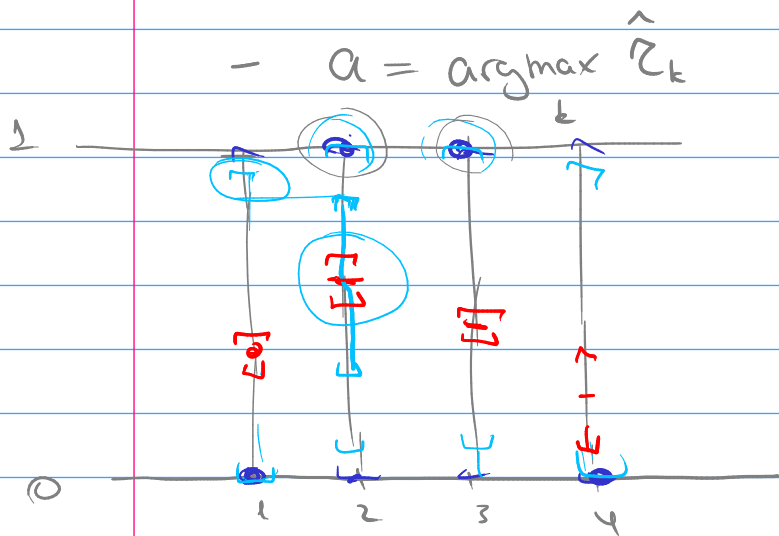
$$\pi(a) = Pr[A_t = a]$$

Greedy

$$\hat{z}_k = \frac{\sum_{t: A_t = a_k} r_t}{\left(\sum_{t: A_t = a_k} 1\right) = n_k}$$

$$z = \frac{r}{r_0}$$

$$z_k = Pr[R=1 | A=a_k]$$



$$\hat{z}_k = \frac{\sum r_t}{n_k} = \frac{r_1 + r_2 + \dots + r_{n_k}}{n_k}$$

$$\hat{z}_k = \frac{(z_{t-1} + z_{n_k}) + r}{n_k + 1} = \frac{1}{n_k + 1} (\hat{z}_k \cdot n_k + r) = \hat{z}_k + \frac{1}{n_k + 1} (r - \hat{z}_k)$$

$$\hat{z}_k = \hat{z}_k + \frac{1}{n_k + 1} (r - \hat{z}_k)$$

Общая оценка = Оценка + Оценка $\left(\text{Число} - \text{Оценка} \right)$

$$\sum_{t=1}^{\infty} h_t = \infty$$

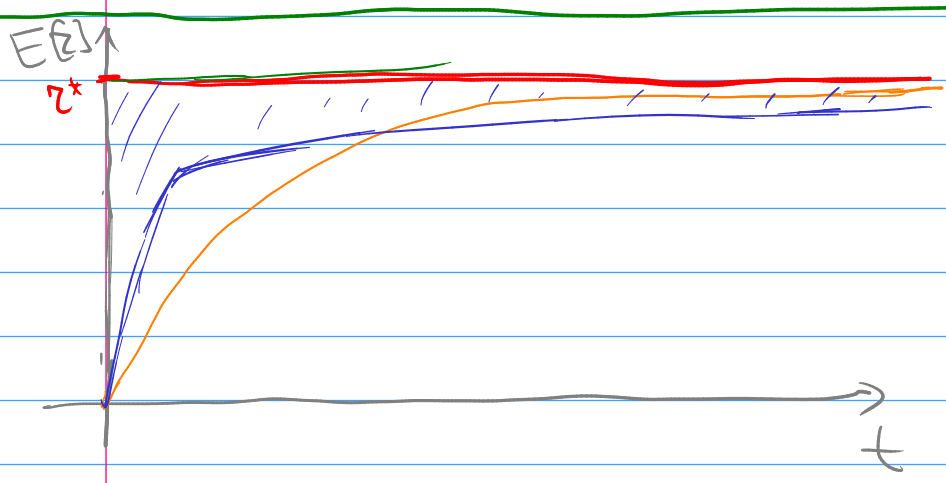
$$\sum_{t=1}^{\infty} h_t^2 < \infty$$

Оценка $\left(\frac{1}{t} \right) \Rightarrow \hat{z}$

$$\hat{a}_{t+1} = \hat{a}_t + h \cdot (r_t - \hat{a}_t) = \dots \rightarrow 0$$

Nonstationary
bandits

$$\begin{aligned} &= h r_t + (1-h) \hat{a}_t = \\ &= h r_t + (1-h) (h r_{t-1} + (1-h) \hat{a}_{t-1}) = \dots \\ &= h r_t + h(1-h) r_{t-1} + h(1-h)^2 r_{t-2} + \dots \\ &= \sum_{k=t}^0 h(1-h)^{k-t} r_k \end{aligned}$$



Regret

$$z^* - \mathbb{E}_{\pi} [z] = \sum_{i=1}^K n_i$$

UCB1
upper confidence bound

$$\text{Priority}_k = \hat{z}_k + \frac{c \sqrt{\log n}}{\sqrt{n_k}}$$

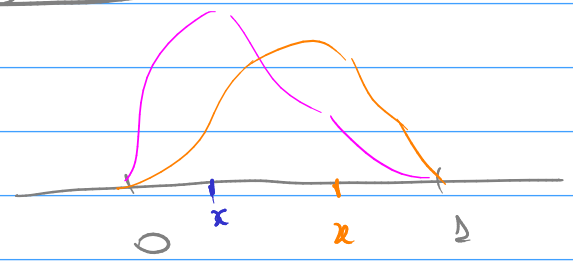
= seznam je k

$T = O(\log T)$ po definiciji je udoben prijem

$$\text{regret} \leq O(\sqrt{K \cdot T \cdot \log T})$$

Thompson sampling

$$\pi(a) = p(a = a^* | D)$$



$$\begin{aligned}
 \underbrace{G_t}_{\text{return}} &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
 &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = \\
 &= R_{t+1} + \gamma G_{t+1}
 \end{aligned}$$

Value function

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

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

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

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

$$\begin{aligned}
 Q_{\pi}(s, a) &= E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] = \\
 &= \sum_{z, s'} p(z, s' | s, a) \cdot (z + \gamma E_{\pi} [G_{t+1} | S_{t+1} = s']) \\
 &= r(s, a) + \gamma \cdot \sum_{s'} p(s' | s, a) \cdot V_{\pi}(s')
 \end{aligned}$$

Bellman equations

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

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

$$\pi(a|s)$$

$$\bar{x} = f(\bar{x})$$

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

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

Blackjack

$S = ?$

($\sum_{i=1, \dots, n}$ ^{wable} Ace, ^{kapia} Queen)

$$\pi: S \rightarrow a$$

Hit / Stand

(Hand, Ace, Dealer card,, Score)

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

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

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

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

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

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

$$\bar{x} = f(\bar{x})$$

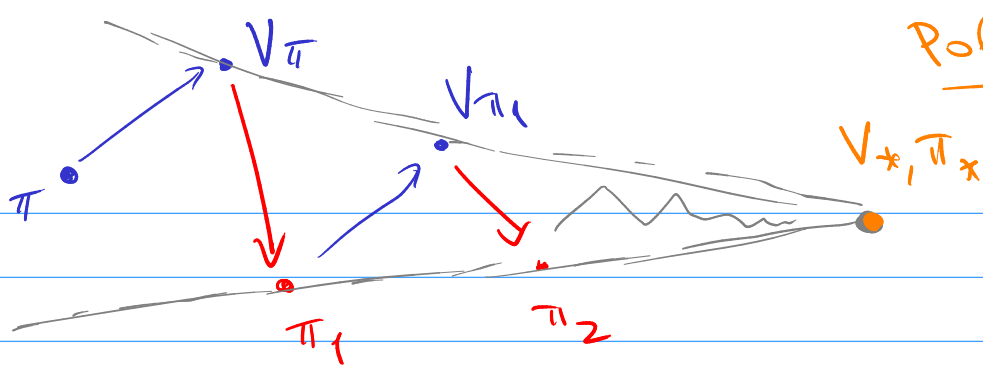
$$\pi^*(a|s) = \begin{cases} 1, & a = \arg \max_a Q_*(s, a) \\ 0, & \text{---} \end{cases}$$

Policy improvement

$$\pi \rightarrow V_{\pi}, Q_{\pi}$$

~~V_*~~

Policy iteration



$$\pi' \geq \pi \iff \forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

Thm. (Policy Improvement thm). Policy π, π' - gen. eq., \forall

$$\forall s \quad Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s)$$

$$\text{Then } \forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

$$V_{\pi'}(s) \leq Q_{\pi}(s, \pi'(s)) =$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | S_t = s, A_t = \pi'(s)] =$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | S_t = s] \leq$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | S_t = s] =$$

$$= \dots \leq \dots \quad \mathbb{E}_{\pi'} [G_t | S_t = s] = V_{\pi'}(s)$$

$$\pi^{(k+1)}(s) = \operatorname{argmax}_a Q_{\pi^{(k)}}(s, a)$$

$$\text{Exam } \exists \pi'; \forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

$$\exists s_0 \quad V_{\pi'}(s_0) > V_{\pi}(s_0)$$

$$\Rightarrow Q_{\pi}(s_0, \pi'(s_0)) > V_{\pi}(s_0)$$