

# MDP Markov decision process

t:

Agent

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

$s_t$

$$\pi(s_t, a_t)$$

$$a_t \in A_t$$

Environment

$r_t$  - reward

t+1:

$s_{t+1}$

$$r_t \sim p(r_t | s_t, a_t)$$

$$s_{t+1} \sim p(s_{t+1} | s_t, a_t)$$

Episodic tasks

$1, \dots, T$

$$\sum_{t=1}^T r_t \rightarrow \max$$

Continuous tasks

$$\sum_{t=1}^{\infty} r_t \rightarrow \max$$

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

$$D = \frac{1}{N} \sum_{n=1}^N (\bar{x}_n, y_n)$$

$$p(y|\bar{x})$$

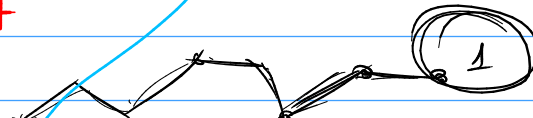
$$D = \frac{1}{N} \sum_{n=1}^N \bar{x}_n$$

$$p(\bar{x})$$

$$p(\bar{x} | \bar{\alpha}, \bar{\mu}_1, \dots, \bar{\mu}_K, \bar{\Sigma}_1, \dots, \bar{\Sigma}_K)$$

① Exploration vs. exploitation

② Credit assignment



Reward -  $r_t \in \{-1, 0, 1\}$

$$V(s) \approx E[R | \text{начало и позиция } s]$$

$$Q(s, a) \approx E[R | \text{начало и } s \text{ и chosen действие } a]$$

$$a^* = \arg \max_a Q(s, a)$$

Multiarmed bandits

$a_1 \quad a_2 \quad \dots \quad a_K$

$$t, \underline{a_i} \sim \underline{r_t(i)} \sim p_i(z), \quad E_{p_i} z = \bar{R}_i$$

$$i^* = \operatorname{argmax}_i R_i$$

$$\hat{R}_i = \frac{1}{n_i} \sum_{t=1}^{n_i} r_t$$

Greedy: — Алгоритм 1 шаг за шагом прыты

0/1

$$\forall t = 1, \dots, T:$$

$$I_t = \operatorname{argmax}_i \hat{R}_i$$

$$\hat{R}_{I_t} = \frac{1}{n_{I_t}} \sum r_t$$

$$r_t \in [0, 1]$$

— „бупі-жэн.“ —  $n_i = 1$ ,  $\hat{R}_i = 5$

— Greedy

$$t: \underbrace{Q_t(a)}_{r_{t+1}} = \frac{\sum_{i=1}^{n_a} r_{ti}}{n_a}$$

$t_i$  — бэне, кале сэрону а

$$Q_t(a) = \frac{1}{n_a} (r_{t1} + \dots + r_{tna})$$

$$Q_{t+1}(a) = \frac{1}{n_a+1} ((r_{t1} + \dots + r_{tna}) + r_{t+1}) =$$

$$= \frac{r_{t1} + \dots + r_{tna}}{n_a+1} + \frac{r_{t+1}}{n_a+1} =$$

$$1 - \frac{1}{n_a+1}$$

$$= \frac{n_a}{n_a+1} Q_t(a) + \frac{1}{n_a+1} r_{t+1}$$

$$Q_{t+1}(a) = Q_t(a) + \frac{1}{n_a+1} [r_{t+1} - Q_t(a)]$$

Новаа оцэнка = Стара оцэнка + Вар. [r - Стара оцэнка]

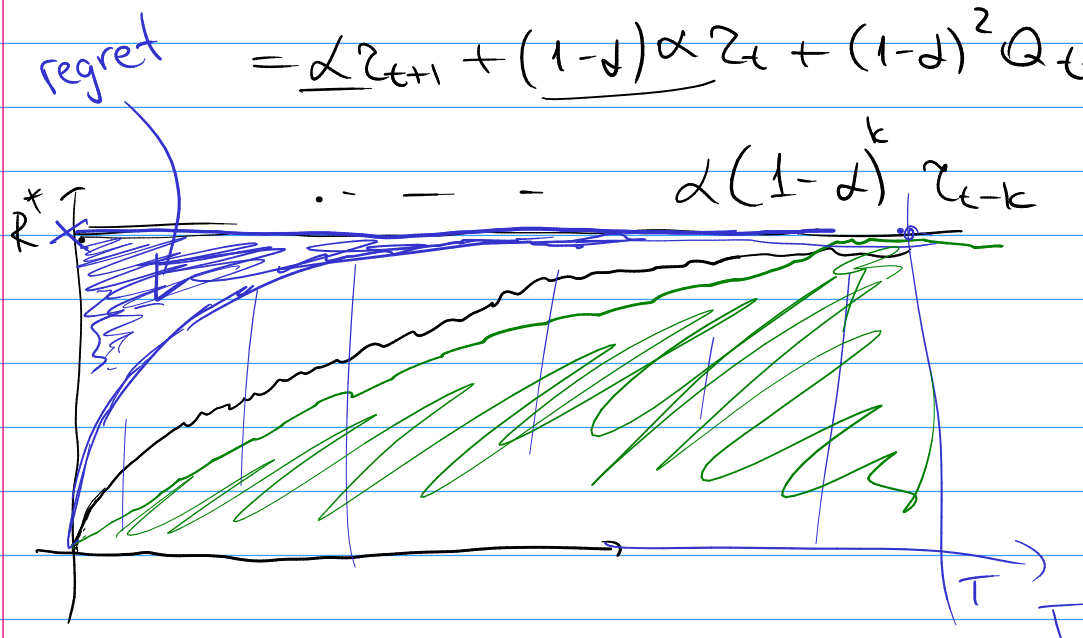
$$\|F(x^*) - F(x_k)\| \leq \frac{R^2 + G^2 \cdot \sum d_t^2}{2 \cdot \sum d_t} \rightarrow \infty$$

$$\sum_{t=1}^{\infty} \alpha_t = \infty \quad \sum \alpha_t^2 < \infty$$

$$Q_{t+1}(a) = Q_t(a) + \alpha \cdot [z_{t+1} - Q_t(a)] =$$

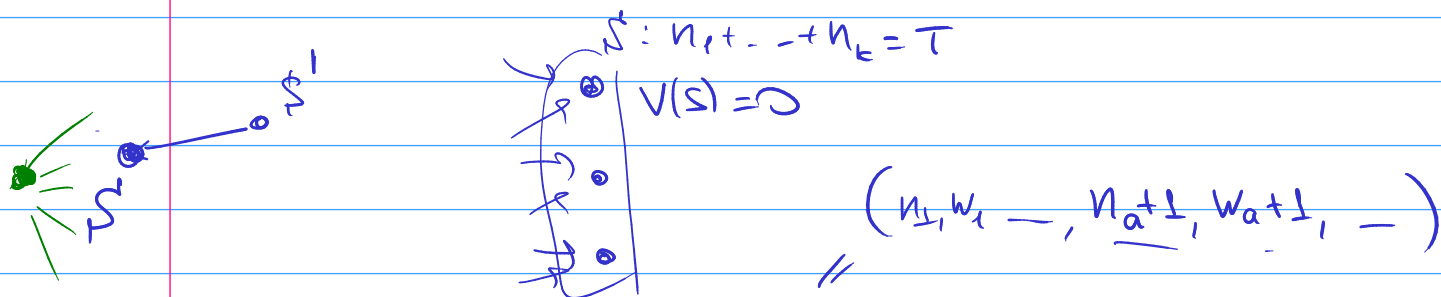
$$= \alpha z_{t+1} + (1-\alpha) Q_t(a) =$$

$$= \alpha z_{t+1} + (1-\alpha) \alpha z_t + (1-\alpha)^2 Q_{t-1}(a)$$



$$S = \{n_1, w_1, (n_2, w_2), \dots, n_k, w_k\} \quad n = n_1 + \dots + n_k$$

$$V(S) = \mathbb{E}[\text{harpa to kupa } T]$$



$$Q(S, a) = p_a \cdot (1 + V(S')) + (1 - p_a) V(S'')$$

$$V(S) = \max_a Q(S, a)$$

$$p_a = \frac{w_{a+1}}{n_{a+2}}$$

Gittins indices

# Thompson sampling

$$z_t \in \{0, 1\}, \quad \theta_i = p(z_t = 1 | I_t = i)$$

$$t: \quad i \rightarrow (D_{i,t}) \quad \{n_{i,t}, w_{i,t}\}$$

$$p(\theta_i) = B(\theta_i | 1, 1)$$

$$p(\theta_i | D_{i,t}) =$$

$$p(\underline{\theta_i} | w_{i,t}, n_{i,t}) = B(\theta_i | w_{i,t} + 1, n_{i,t} - w_{i,t} + 1)$$

$$- t = 1, \dots, T:$$

$$- \theta_{i,t} \sim p(\theta_i | w_{i,t}, n_{i,t})$$

$$- \underline{I_t} = \arg \max_i \theta_{i,t}$$

↑  
zu  $w_i$  separieren  
ke wane  $t$

