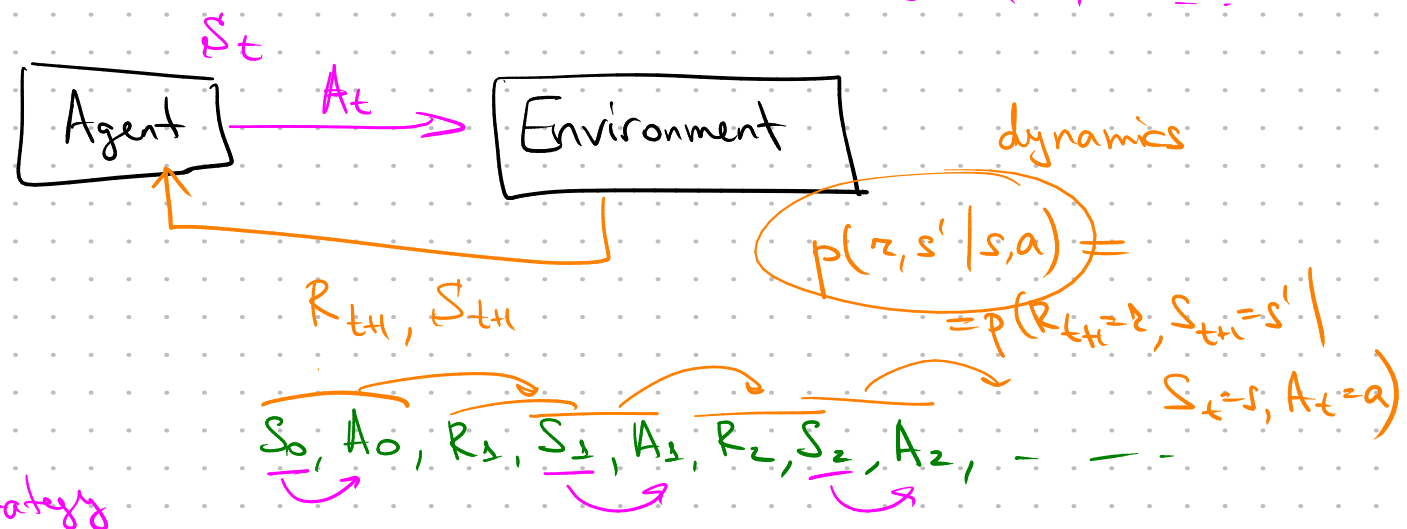


Markov decision process

$t = 0, 1, \dots$



strategy

$\pi: S^t \mapsto \text{Prob}(A(S))$

$\pi(a|s)$

$\pi(s, a) = p(A_t = a | S_t = s)$

Уахмары:

- $S$  - нагунна + счёрники
- $A$  - ходы
- $R$  - прыт + напрыт

Episodic tasks

Continuous tasks

$t = 1, 2, \dots, T$

$t = 1, 2, \dots, T, \dots$

$\sum_{t=1}^T R_t \xrightarrow{\pi} \max$

$\gamma = 1$

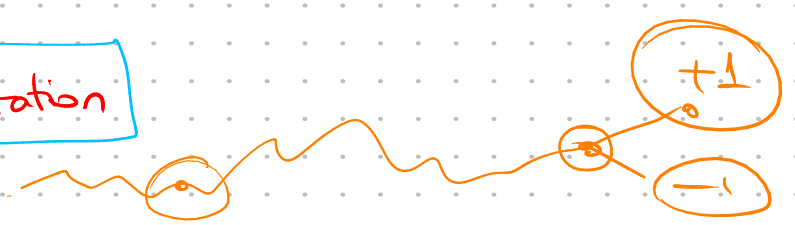
$\sum_{t=1}^{\infty} \gamma^{t-1} R_t \xrightarrow{\pi} \max$

$\gamma < 1$



1) Exploration vs. exploitation

2) ~~Credit assignment~~



# Multiarmed bandits

$$|S| = 1$$

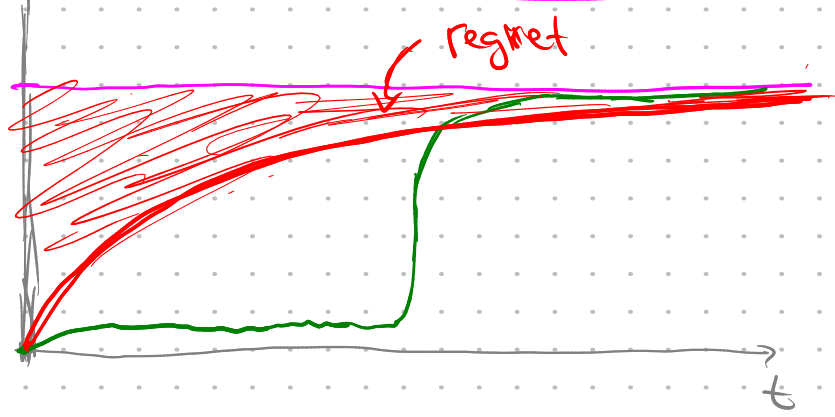
$$A: a_1, a_2, \dots, a_K$$

$$R_{t+1} \sim p(R_{t+1} = z | A_t = a)$$

$$p(a=a^*) \rightarrow 1 \text{ (opt)}$$

$$E[z] =$$

$$a^* = \arg \max_k E_{p(z|k)}[z]$$



Greedy:

- always pick  $a^*$   $\forall a = 1, \dots, K$

- for  $t = K+1, \dots$

$$\hat{\mu}_i(t) = \frac{1}{n_i} \sum_{t: A_t = i} R_t$$

$$A_t = \arg \max_i \hat{\mu}_i(t)$$

binary bandits

$$R_t \in \{0, 1\}$$

$\epsilon$ -Greedy:

$$\hat{\mu}_i(t)$$

$$A_t = \begin{cases} \arg \max_i \hat{\mu}_i(t) & , 1 - \epsilon \\ \text{uniform} & , \epsilon \end{cases}$$

$$\pi(a) = (1 - \epsilon) \cdot [a = a^*] + \frac{\epsilon}{K}$$

$$R_t \quad \hat{\mu}_i(n+1) = \frac{1}{n+1} \sum_{t=1}^{n+1} R_t = \frac{1}{n+1} \left( R_{n+1} + \sum_{t=1}^n R_t \right) =$$

$$= \frac{1}{n+1} R_{n+1} + \frac{n}{n+1} \hat{\mu}_i(n) \rightarrow \underline{1} - \frac{1}{n+1}$$

$$\hat{\mu}(n+1) = \hat{\mu}(n) + \frac{1}{n+1} (R_{n+1} - \hat{\mu}(n))$$

$$\sum_{n=1}^{\infty} \alpha_n = \infty$$

$$\sum_{n=1}^{\infty} \alpha_n^2 < \infty$$

Новая оценка = Старая оценка + Уар. (Мноб. знак. - Старая оценка)

$\alpha_n = \alpha$  Nonstationary bandits

$$\begin{aligned} \hat{\mu}(n+1) &= \hat{\mu}(n) + \alpha (R_{n+1} - \hat{\mu}(n)) = \\ &= (1-\alpha)\hat{\mu}(n) + \alpha R_{n+1} = \\ &= \alpha R_{n+1} + (1-\alpha)(\alpha R_n + \hat{\mu}(n-1) - (1-\alpha)) = \\ &= \dots = \alpha R_{n+1} + \alpha(1-\alpha)R_n + \alpha(1-\alpha)^2 R_{n-1} + \dots \end{aligned}$$

Binary  $R_t \in \{0, 1\}$   $(n_k, w_k)$   $k=1, \dots, K$   $0 \leq w_k \leq n_k$   
 $t=1, \dots, T$   $(n_1, w_1, n_2, w_2, \dots, n_K, w_K)$ ,  $\sum n_k \leq T$

$$V_*(n_1, \dots, w_K) = E \left[ \sum_{t=Z_{n_i+1}}^T R_t \right]$$

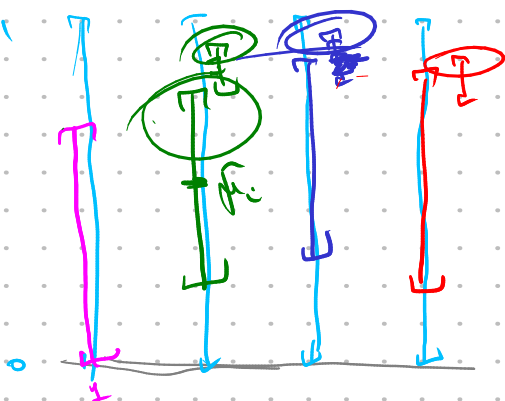
Gittins index

Базис:  $\sum n_k = T \Rightarrow V_* = 0$

Рекурсия:  $V_*(n_1, \dots, n_i, w_i, \dots, w_K) =$

$$\operatorname{argmax}_i \gamma(n_i, w_i)$$

$$= \max_{i=1}^K \left[ \begin{aligned} &= \frac{w_i+1}{n_i+2} \\ & \mu_i \cdot (1 + V_*(n_1, \dots, \underline{n_i+1}, w_i+1, \dots, w_K)) \\ & + (1-\mu_i) \cdot V_*(n_1, \dots, \underline{n_i+1}, w_i, \dots, w_K) \end{aligned} \right]$$



UCB - upper confidence bound

UCB1

$$\text{Priority}_i(t) = \hat{\mu}_i(t) + c \cdot \sqrt{\frac{\log t}{n_i}}$$

Thm  $c = \sqrt{2} \Rightarrow \text{Regret} = O(\sqrt{KT \log T})$

$$\Delta_i = |\mu_{*} - \mu_i| \sqrt{T \log T}$$

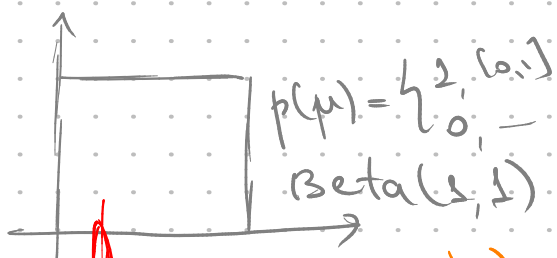
UCB-Tuned

$$\hat{\mu}_i + \sqrt{\frac{\log t}{n_i} \cdot \min\left(\frac{1}{4}, V_i(n_i)\right)}$$

$$V_i(n_i) = \frac{1}{n_i} \sum z^2 - \left(\frac{1}{n_i} \sum z\right)^2 + \sqrt{\frac{2 \log t}{n_i}}$$

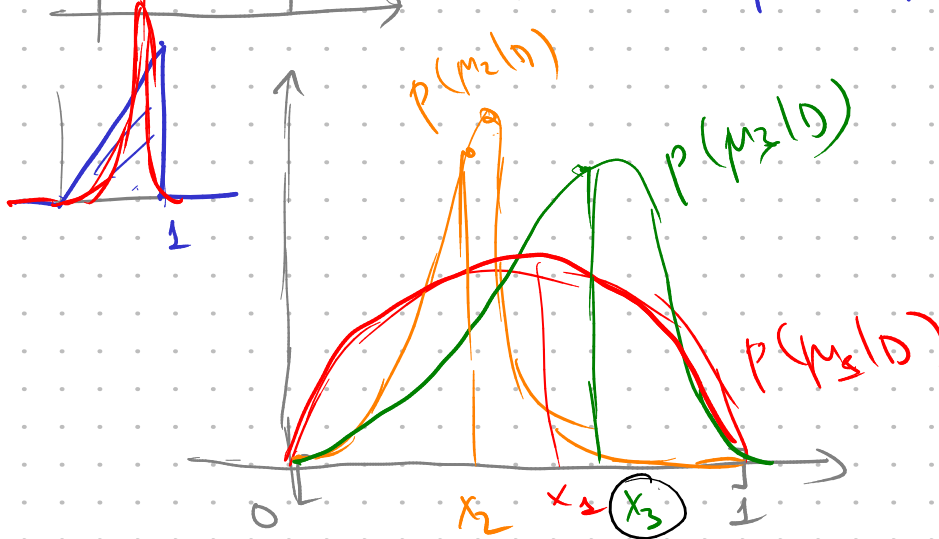
Thompson sampling

binary  $\mu_1, \mu_2, \dots, \mu_K \in [0, 1]$



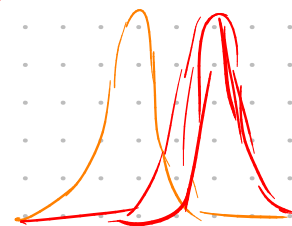
$$x \mu^w (1-\mu)^{n-w}$$

$$p(\mu | D) = \text{Beta}(w+1, n-w+1)$$



$$x_i \sim p(\mu_i | D)$$

$$A_t = \arg \max x_i$$



A/B testing

$a_1, a_2, a_3$

$t = 1, \dots, T$

$\{x_1, \dots, x_M\}$

$x_t$  - correct

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

$$R_t \sim p(R_t | A_t, x_t)$$

$\pi_1, \pi_2, \dots, \pi_A$

~~$p(x_t | x_t)$~~

