

Machine Learning

Supervised learning

$$D = \{d = (\bar{x}, y)\}$$

$$f: \bar{x} \mapsto y$$

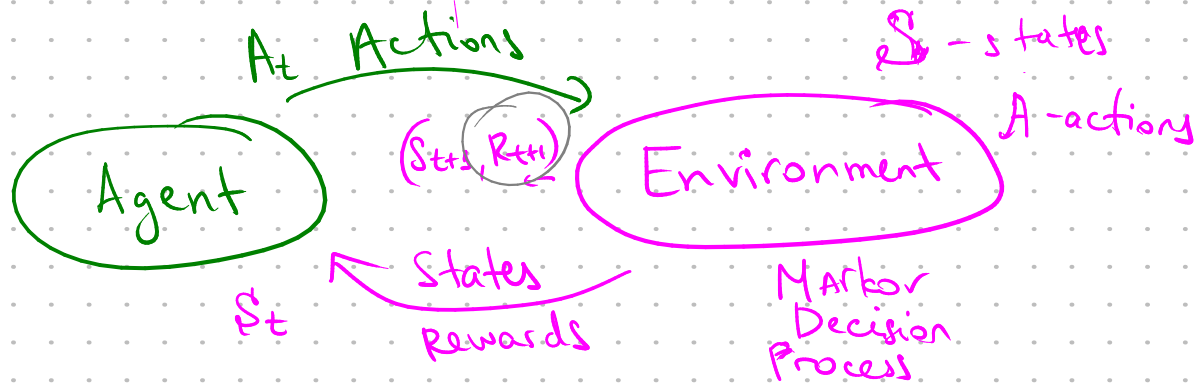
$$p(y|\bar{x})$$

Unsupervised learning

$$D = \{\bar{x}\}$$

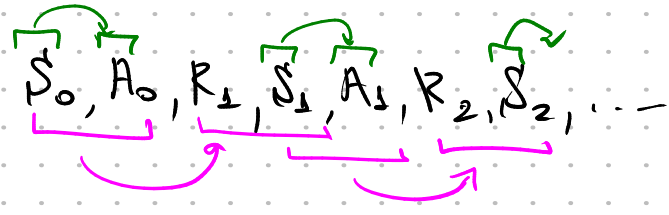
$$p(\bar{x})$$

Reinforcement learning



Agent strategy $\pi: S \rightarrow \text{Prob}[A]$

Environment dynamics $p(s', r | s, a)$



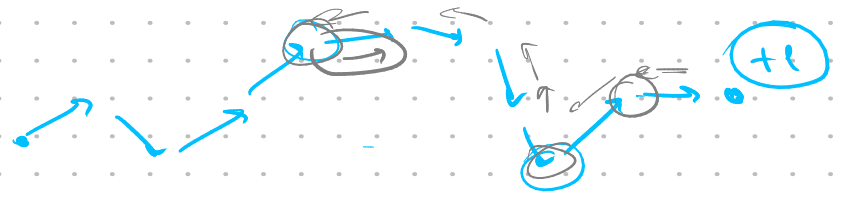
Указания

state - положение в какой области, время, покрытие
 action - ход
 reward - $+1$ бонусом, -1 штрафом, 0 нулем
 en passant 50 ходов, покрыт. обн.

Episodic tasks

$$R_t = 0, R_T \neq 0?$$

Credit assignment



Exploration vs. exploitation

Multiarmed bandits

$$|S| = 1 \quad A = (1, \dots, M) \quad (p_1, \dots, p_M)$$

binary
 $R_t \in \{0, 1\}$
 $R_t \in \mathbb{R}$

$$A_t \rightarrow \text{Env} \rightarrow R_t$$

$a_1, a_2, a_3, \dots, a_t$
 $r_1, r_2, r_3, \dots, r_t$

$$a^* = \operatorname{argmax} E[R_t]$$

Greedy:

$$\hat{p}_i = \frac{1}{n_i} \sum_{t: a_t=i} r_t$$

$$a_t = \operatorname{argmax}_i \hat{p}_i$$

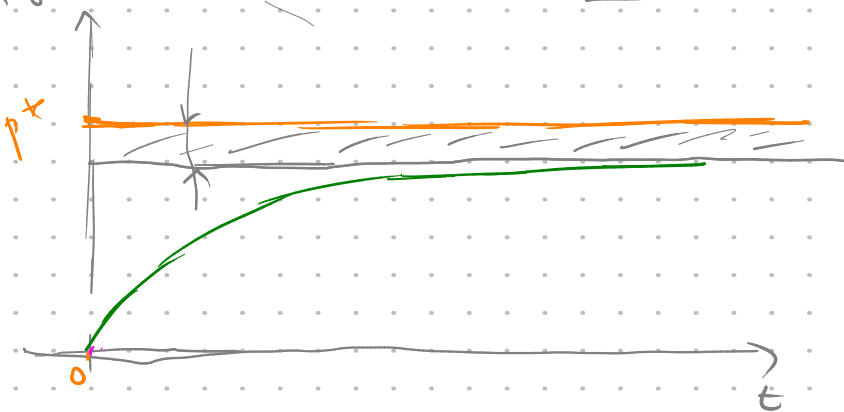
ϵ -Greedy:

$$a_t = \begin{cases} \operatorname{argmax}_i \hat{p}_i, & 1-\epsilon \\ \text{uniform}, & \epsilon \end{cases}$$

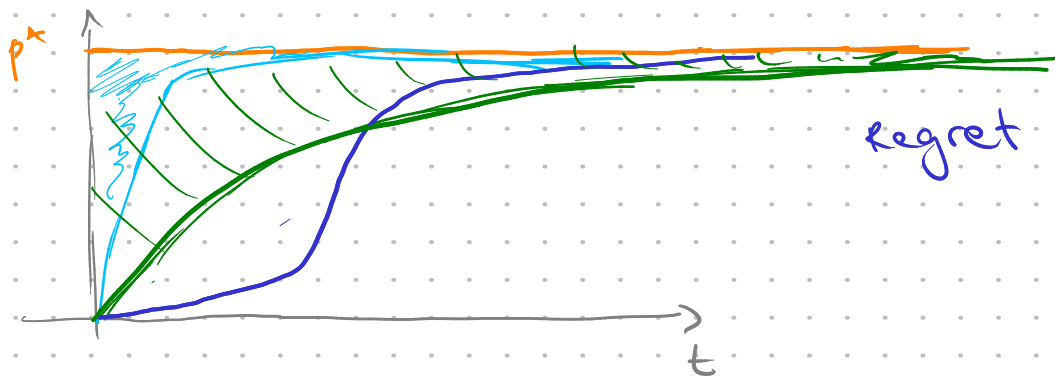
$$\epsilon_t \xrightarrow{t \rightarrow \infty} 0$$

$$\pi(a_t=i) = \frac{\epsilon}{M} + (1-\epsilon) \cdot [i = \operatorname{argmax}_i \hat{p}_i]$$

Avg reward



$$p^* - \epsilon \cdot E[\dots]$$



regret - цена ошибки

$$\hat{p}_i = \frac{r_1 + r_2 + \dots + r_n}{n} = \frac{1}{n} (r_n + \sum_{i=1}^{n-1} r_i) = \frac{1}{n} (r_n + (n-1) \hat{p}_i^{\text{old}})$$

$$\hat{p}_i = \hat{p}_i^{\text{old}} + \frac{1}{n} (r_n - \hat{p}_i^{\text{old}})$$

$$\boxed{\text{New} = \text{Old} + \text{Stepsize} (\text{Target} - \text{Old})}$$

$$\hat{p}_{t+1} = \hat{p}_t + \alpha_t (r_t - \hat{p}_t)$$

$$\alpha_t = \frac{1}{t} \Rightarrow \text{Avg}$$

$\alpha_t = \alpha = \text{Const}$ ← nonstationary bandits

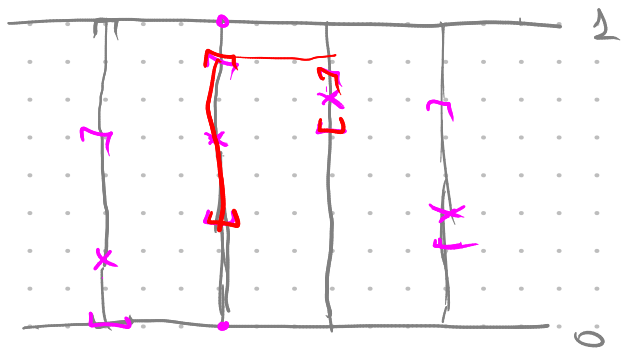
$$\begin{aligned} \hat{p}_{t+1} &= \hat{p}_t + \alpha (r_t - \hat{p}_t) = \alpha r_t + (1-\alpha) \hat{p}_t = \\ &= \alpha r_t + (1-\alpha) (\alpha r_{t-1} + (1-\alpha) \hat{p}_{t-1}) = \dots \\ &= \alpha (r_t + (1-\alpha) r_{t-1} + (1-\alpha)^2 r_{t-2} + \dots + (1-\alpha)^{k-1} r_{t-k+1}) \end{aligned}$$

Thm Nowey. $\hat{p}_{t+1} = \hat{p}_t + \alpha_t (r_t - \hat{p}_t)$ $\alpha_t < \log^{-1}$,

ecce $\sum_{t=1}^{\infty} \alpha_t = \infty$ u $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$

Optimism under uncertainty

- гобекривеннае унепбана - бепхнае ганауна
 - optimistic initial values
- $\tau_0 = 10$



- ucb upper confidence band

Priority $[i] = \hat{p}_i + \dots$

ucb1: $a_t = \text{argmax}_i \left[\hat{p}_i + c \sqrt{\frac{\ln t}{n_i}} \right]$

$[a_t = i]$

Thm. $c = \sqrt{2}$ \Rightarrow ucb1 бадыпер цыдон. гикривна $O(\log T)$ пог жс T wand, u regret = $O(\sqrt{MT \log T})$ $\Delta_i = p^* - p_i$

ucb

Thompson sampling

- sample $p'_i \sim p(p_i | D)$
- $a_t = \operatorname{argmax}_i p'_i$

