

ML: $D \leadsto ?$

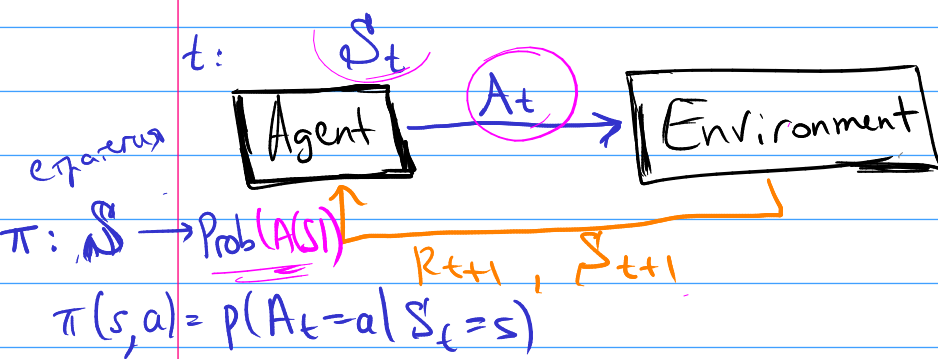
Supervised Learning
 $\bar{x} \mapsto y?$
 $D = \{(\bar{x}, y)\}$
 $p(y|\bar{x})$

Unsupervised learning
 \bar{x} , $D = \{\bar{x}\}$
 $p(\bar{x})$

Reinforcement learning

Markov decision process (MDP)

S_t & maximize:
 - $\sum_{t=0}^{\infty} \gamma^t R_t$
 - ...



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

Episodic task

Episodic task diagram and equation:

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

Diagram shows a sequence of states $0 \rightarrow 0 \rightarrow 0 \rightarrow 0 \rightarrow \dots \rightarrow R_t \rightarrow R_{t+1} = 0$ with a terminal state indicated by a double circle.

Continuous task

Continuous task diagram and equation:

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

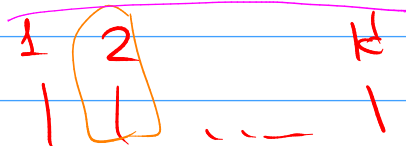
Diagram shows a sequence of states $0 \rightarrow 0 \rightarrow 0 \rightarrow 0 \rightarrow \dots$ with a discount factor $\gamma < 1$ indicated.

- 1) Exploration vs. exploitation
- 2) Credit assignment

Sutton, Barto 1998 RL: An Introduction 2018

Multiarmed bandits

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



stationary bandits

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

At $A_t \in \{1, 2, \dots, K\}$
 $R_{t+1} \sim p(R_{t+1} = r | A_t = a)$

$$i_t = \arg \max_{i=1..K} \mathbb{E}_{p(z|i)} [z]$$



regret

Greedy:

- выбираю no 1 погу же наибольшего погугу
- for $t = k+1, \dots$:

$$A_t = \operatorname{argmax}_i \hat{\mu}_i(t)$$

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

binary bandits
 $\mu_i = p(x=1|i)$

1 1 0 1 0 1 0 0

Exploitation
~~Exploration~~

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

" $n_i \cdot \hat{\mu}_i(t)$

$$= \frac{1}{n_{i+1}} \left(R_{n_{i+1}} + (n_{i+1} - 1) \hat{\mu}_i(t) \right)$$

$$\hat{\mu}_i(t+1) = \hat{\mu}_i(t) + \frac{1}{n_{i+1}} \cdot (R_{n_{i+1}} - \hat{\mu}_i(t))$$

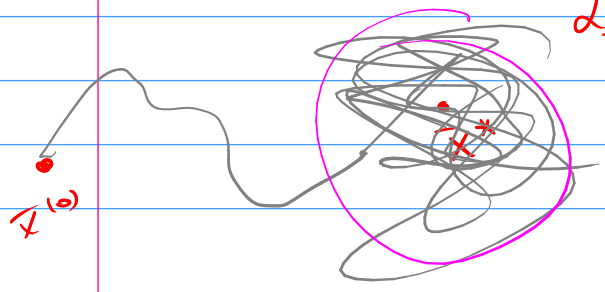
2/ Новая оценка = Старая оценка + Шаг \cdot (Множ. знамен. - Старая оценка)

$$\|F(\bar{x}^*) - F(\bar{x}^{(k)})\| \leq \frac{R^2 + G^2 \cdot \sum_{t=1}^k d_t^2}{2 \sum_{t=1}^k d_t} < \infty$$

$R^2 \geq \|x^{(k)} - x^*\|^2$
 $G^2 \geq \|g_k\|^2$

$d_t = h$

$\frac{R^2}{2h \cdot k} + \frac{G^2 \cdot h}{2}$



Nonstationary bandits

$$\begin{aligned}\hat{\mu}_i(t+1) &= \hat{\mu}_i(t) + \alpha (R_t - \hat{\mu}_i(t)) = \\ &= \alpha R_t + (1-\alpha) \hat{\mu}_i(t) = \\ &= \alpha R_t + (1-\alpha) (\alpha R_{t-1} + (1-\alpha) \hat{\mu}_i(t-1)) = \\ \dots &= \alpha R_t + \alpha(1-\alpha) R_{t-1} + \alpha(1-\alpha)^2 R_{t-2} + \alpha(1-\alpha)^3 R_{t-3} + \dots\end{aligned}$$

ϵ -Greedy:

$$A_t = \begin{cases} \operatorname{argmax}_i \hat{\mu}_i(t) & \text{c. prob. } 1-\epsilon, \\ \text{uniform} & \text{c. prob. } \epsilon. \end{cases}$$

Binary $R_t \in \{0, 1\}$ $(n_1, w_1, n_2, w_2, \dots, n_K, w_K)$, $0 \leq w_i \leq n_i$, $\sum n_i \leq T$

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

Soln: $\sum n_i = T \Rightarrow V_* = 0$

recursion: $V_*(n_1, w_1, \dots, n_K, w_K, \dots) =$

$$= \max_{i=1}^K \left(\mu_i (1 + V_*(\dots, n_i+1, w_i+1, \dots)) + (1-\mu_i) \cdot V_*(\dots, n_i+1, w_i, \dots) \right)$$

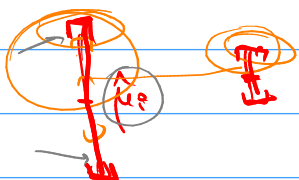
$$\frac{w_i+1}{n_i+2}$$

Gittins index

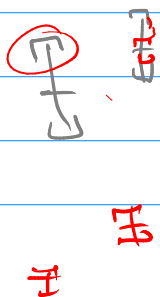
$$\operatorname{argmax}_i g(n_i, w_i)$$

$$\hat{\mu}_i, i=1, \dots, K$$

$$\hat{\mu}_i \rightarrow \operatorname{softmax} \rightarrow _$$



$$\operatorname{argmax}_i [\hat{\mu}_i + \dots]$$



Upper confidence bound

UCB 1

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

ϵ -greedy $O(T)$

Thm 1

Let $c = \sqrt{2}$

$$\hat{\mu}_i + \sqrt{\frac{2 \log t}{n_i}}$$

$$\text{regret} = O(\sqrt{KT \log T})$$

$O(\log t)$ bound upon generic

UCB 2

UCB-tuned

$$\sqrt{\frac{\log t}{n_i}} \cdot \min\left(\frac{1}{4}, V_j(n_j)\right), \text{ or}$$

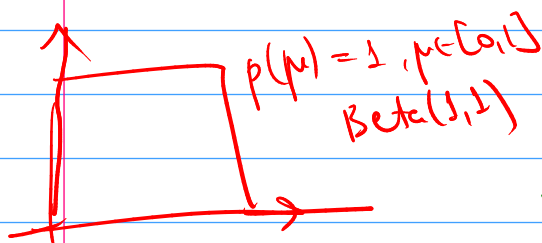
$$V_j(n_j) = \left(\frac{1}{n_j} \sum z^2\right) + \left(\frac{1}{n_j} \sum z^2\right)^2 + \sqrt{\frac{2 \log t}{n_j}}$$

a. $\sum_{i \neq i^*} \frac{\log t}{\Delta_i}$

$$\Delta_i = \mu_* - \mu_i$$

Thompson sampling:

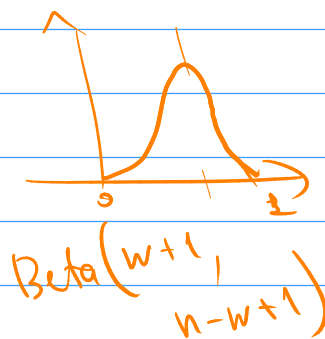
$$\mu_1, \mu_2, \dots, \mu_K \in [0, 1]$$



$p(p)$

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

$p(D|\mu)$



$p(\mu|D)$



Probability matching:

$$i \sim p(\mu_i = \max_j \mu_j | D)$$