

Thompson Sampling

Александра Харь, 774

Многорукие бандиты, постановка задачи:

Пусть у нас есть автомат с N ручками, в каждый момент времени $t = 1, 2, 3, \dots$ игрок выбирает одну из этих ручек. Для каждой ручки есть какое-то свое распределение выигрыша (игроку неизвестное). Сразу после того, как игрок сыграл конкретную ручку, он получает выигрыш.

Пусть μ_i - матожидание выигрыша на ручке i ; $\mu^* := \max_i \mu_i$, $\Delta_i := \mu^* - \mu_i$. Тогда определим суммарный regret:

$$\mathbb{E}[R(T)] = \mathbb{E}\left[\sum_{t=1}^T (\mu^* - \mu_{i(t)})\right] = \sum_i \Delta_i \cdot \mathbb{E}[k_i(T+1)] \quad (1)$$

Мы будем рассматривать Бернуллиевских бандитов: распределения выигрышей на ручек являются бернуллиевскими. Без ограничения общности, будем считать, что первая ручка имеет максимальный средний выигрыш ($\mu^* = \mu_1$)

Algorithm 1 (Thompson Sampling using Beta prior)

For each arm $i = 1, \dots, N$ set $S_i = 0, F_i = 0$
For each $t = 1, 2, \dots$ **do**
 For each arm $i = 1, \dots, N$ sample $\theta_i(t)$ from the $Beta(S_i + 1, F_i + 1)$ distribution
 Play arm $i(t) := \operatorname{argmax}_i \theta_i(t)$ and observe reward r_t
 If $r_t = 1$, then $S_{i(t)} = S_{i(t)} + 1$,
 else $F_{i(t)} = F_{i(t)} + 1$

Теорема. Для N -рукого стохастического бандита, используя Thompson Sampling using Beta priors, имеет место следующая оценка:

$$\mathbb{E}[R(T)] \leq O(\sqrt{NT \ln T}) \quad (2)$$

Для начала введем некоторые обозначения:

- $i(t)$ - ручка, которую играл игрок в момент времени t
- $k_i(t)$ - количество игр ручки i до момента времени $t - 1$
- $S_i(t)$ - количество успешных игр ручки i до момента времени $t - 1$
- эмпирическое среднее $\hat{\mu}_i(t) = \frac{\sum_{\tau=1}^{t-1} \mathbb{1}_{i(\tau)=i} r_i(\tau)}{k_i(t) + 1}$
- величины x_i, y_i впоследствии будем выбирать таким образом, чтобы $\forall i \hookrightarrow \mu_i < x_i < y_i < \mu_1$
- обозначим $E_i^\mu(t)$ событие: $\hat{\mu}_i(t) \leq x_i$, $E_i^\theta(t)$ событие: $\theta_i(t) \leq y_i$
- $\mathcal{F}_{t-1} = \{i(w), r_{i(w)}(w), w = 1, \dots, t-1\}$ - история до момента времени $t - 1$
- $p_{i,t} = P(\theta_1(t) > y_i | \mathcal{F}_{t-1})$

Лемма 1. Для любого $t \in [1, T]$ и $t \neq 1$:

$$P(i(t) = i, E_i^\mu(t), E_i^\theta(t) | \mathcal{F}_{t-1}) \leq \frac{1 - p_{i,t}}{p_{i,t}} P(i(t) = 1, E_i^\mu(t), E_i^\theta(t) | \mathcal{F}_{t-1}), \quad p_{i,t} = P(\theta_1(t) > y_i | \mathcal{F}_{t-1}) \quad (3)$$

Доказательство: Положим, что $E_i^\mu(t)$ - верно (иначе левая часть равно 0 и неравенство выполнено). Тогда достаточно доказать, что

$$P(i(t) = i | E_i^\theta(t), \mathcal{F}_{t-1}) \leq \frac{1 - p_{i,t}}{p_{i,t}} P(i(t) = 1 | E_i^\theta(t), \mathcal{F}_{t-1}) \quad (4)$$

$E_i^\theta(t), i(t) = i$ только если $\theta_j(t) \leq y_i, \forall j$.

Значит, для любого $i \neq 1$:

$$\begin{aligned} P(i(t) = i | E_i^\theta(t), \mathcal{F}_{t-1}) &\leq P(\theta_j(t) \leq y_i, \forall j | E_i^\theta(t), \mathcal{F}_{t-1}) = P(\theta_1(t) \leq y_i | \mathcal{F}_{t-1}) P(\theta_j(t) \leq y_i, \forall j \neq 1 | E_i^\theta(t), \mathcal{F}_{t-1}) = \\ &= (1 - p_{i,t}) P(\theta_j(t) \leq y_i, \forall j \neq 1 | E_i^\theta(t), \mathcal{F}_{t-1}) \end{aligned} \quad (5)$$

$$\begin{aligned} P(i(t) = 1 | E_i^\theta(t), \mathcal{F}_{t-1}) &\geq P(\theta_1(t) > y_i \geq \theta_j(t), \forall j \neq 1 | E_i^\theta(t), \mathcal{F}_{t-1}) = \\ &= P(\theta_1(t) > y_i | \mathcal{F}_{t-1}) P(\theta_j(t) \leq y_i, \forall j \neq 1 | E_i^\theta(t), \mathcal{F}_{t-1}) = p_{i,t} P(\theta_j(t) \leq y_i, \forall j \neq 1 | E_i^\theta(t), \mathcal{F}_{t-1}) \end{aligned} \quad (6)$$

$$\Rightarrow P(i(t) = i | E_i^\theta(t), \mathcal{F}_{t-1}) \leq \frac{1 - p_{i,t}}{p_{i,t}} P(i(t) = 1 | E_i^\theta(t), \mathcal{F}_{t-1}). \quad \square$$

Доказательство Теоремы:

$$\mathbb{E}[k_i(T)] = \sum_{t=1}^T P(i(t) = i) = \sum_{t=1}^T P(i(t) = i, E_i^\mu(t), E_i^\theta(t)) + \sum_{t=1}^T P(i(t) = i, E_i^\mu(t), \overline{E_i^\theta(t)}) + \sum_{t=1}^T P(i(t) = i, \overline{E_i^\mu(t)}) \quad (7)$$

Пусть τ_k - это момент времени, когда k -ый раз была сыграна первая ручка, $\tau_0 = 0$. Заметим, что для любого i , для любого $k > k_i(T)$, верно $\tau_k > T$ (также $\tau_T \geq T$).

Используем Лемму 1:

$$\begin{aligned} \sum_{t=1}^T P(i(t) = i, E_i^\mu(t), E_i^\theta(t)) &= \sum_{t=1}^T \mathbb{E}[P(i(t) = i, E_i^\mu(t), E_i^\theta(t) | \mathcal{F}_{t-1})] \leq \sum_{t=1}^T \mathbb{E}\left[\frac{1 - p_{i,t}}{p_{i,t}} P(i(t) = 1, E_i^\mu(t), E_i^\theta(t) | \mathcal{F}_{t-1})\right] = \\ &= \sum_{t=1}^T \mathbb{E}\left[\mathbb{E}\left[\frac{1 - p_{i,t}}{p_{i,t}} I(i(t) = 1, E_i^\mu(t), E_i^\theta(t) | \mathcal{F}_{t-1})\right]\right] = \sum_{t=1}^T \mathbb{E}\left[\frac{1 - p_{i,t}}{p_{i,t}} I(i(t) = 1, E_i^\mu(t), E_i^\theta(t))\right] \leq \\ &\leq \sum_{k=0}^{T-1} \mathbb{E}\left[\frac{1 - p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \sum_{t=\tau_k+1}^{\tau_{k+1}} I(i(t) = 1)\right] = \sum_{k=0}^{T-1} \mathbb{E}\left[\frac{1}{p_{i,\tau_k+1}} - 1\right] \end{aligned} \quad (8)$$

Лемма 2.

$$\mathbb{E}\left[\frac{1}{p_{i,\tau_k+1}}\right] \leq \begin{cases} 1 + \frac{3}{\Delta'_i}, & \text{for } k < \frac{8}{\Delta'_i} \\ 1 + \Theta(e^{-\Delta_i'^2 k/2} + \frac{1}{(k+1)\Delta_i'^2} e^{-D_i k} + \frac{1}{e^{\Delta_i'^2 k/4} - 1}), & \text{for } k \geq \frac{8}{\Delta'_i}, \end{cases} \quad (9)$$

где $\Delta'_i = \mu_1 - y_i$, $D_i = y_i \ln \frac{y_i}{\mu_1} + (1 - y_i) \ln \frac{1 - y_i}{1 - \mu_1}$

Лемма 3.

$$\sum_{t=1}^T P(i(t) = i, \overline{E_i^\mu(t)}) \leq \frac{1}{d(x_i, \mu_i)} + 1 \quad (10)$$

, где $d(x_i, \mu_i) = x_i \log \frac{x_i}{\mu_i} + (1 - x_i) \log \frac{1 - x_i}{1 - \mu_i}$

Лемма 4.

$$\sum_{t=1}^T P(i(t) = i, \overline{E_i^\theta(t)}, E_i^\mu(t)) \leq L_i(T) + 1, \quad (11)$$

где $L_i(T) = \frac{\ln T}{d(x_i, y_i)}$

Пусть

$$x_i = \mu_i + \frac{\Delta_i}{3}, \quad y_i = \mu_1 - \frac{\Delta_i}{3} \quad (12)$$

$$\Delta_i'^2 = (\mu_1 - y_i)^2 = \frac{\Delta_i^2}{9}$$

$$d(x_i, \mu_i) \geq 2(x_i - \mu_i)^2 = \frac{2\Delta_i^2}{9}, \quad d(x_i, y_i) \geq 2(y_i - x_i)^2 \geq \frac{2\Delta_i^2}{9}. \text{ Тогда:}$$

$$L_i(T) = \frac{\ln T}{d(x_i, y_i)} \leq \frac{9\ln T}{2\Delta_i^2} \text{ и } \frac{1}{d(x_i, \mu_i)} \leq \frac{9}{2\Delta_i^2}$$

$$\begin{aligned} \mathbb{E}[k_i(T)] &\leq \frac{24}{\Delta_i'^2} + \sum_{j=0}^{T-1} \Theta\left(e^{-\Delta_i'^2 j/2} + \frac{1}{(j+1)\Delta_i'^2} e^{-D_i j} + \frac{1}{e^{\Delta_i'^2 j/4} - 1}\right) + L_i(T) + 1 + \frac{1}{d(x_i, \mu_i)} + 1 \leq \\ &\leq \sum_{j=0}^{T-1} \Theta\left(e^{\Delta_i'^2 j/2} + \frac{1}{(j+1)\Delta_i'^2} + \frac{4}{\Delta_i'^2 j}\right) + O\left(\frac{\ln T}{\Delta_i^2}\right) = \Theta\left(\frac{1}{\Delta_i'^2} + \frac{\ln T}{\Delta_i'^2}\right) + O\left(\frac{\ln T}{\Delta_i^2}\right) = O\left(\frac{\ln T}{\Delta_i^2}\right) \end{aligned} \quad (13)$$

Для каждой ручки i с $\Delta_i \geq \sqrt{\frac{N\ln T}{T}}$ выполнено $\Delta_i \mathbb{E}[k_i(T)] = O\left(\sqrt{\frac{T\ln T}{N}}\right)$. Для ручек с $\Delta_i \leq \sqrt{\frac{N\ln T}{T}}$ выполнено $\Delta_i \mathbb{E}[k_i(T)] = O(\sqrt{NT\ln T}) \Rightarrow \mathbb{E}[R(T)] = O(\sqrt{NT\ln T})$

Further Optimal Regret Bounds for Thompson Sampling, Shipra Agrawal, Navin Goyal