

Finite armed bandits $\underbrace{\quad}_K$ \leq UCB successive elimination \sqrt{KT}

Continuous armed bandits / Structured bandits
 $f \in \mathbb{R}^d$ - Lipschitz \leq reduction to finite arm $T \frac{d+1}{d+2}$
 zooming (adaptive binning) $T \frac{d+1}{d+2}$

- linear θ - UCB $d\sqrt{T}$
 $x_t \leftarrow \arg \max_x UCB_t(x) \rightarrow \max_{\theta \in \Theta} \langle x, \theta \rangle$

$$C_t = d\theta: \|\theta - \hat{\theta}_{t-1}\| \leq \sqrt{\frac{d}{N_{t-1}}}$$

How to leverage prior information?

- Finite
- Continuous

$\begin{pmatrix} \mu(1) \\ \vdots \\ \mu(K) \end{pmatrix} = \mu$ - mean reward vector $\mu \sim P$ (prior distⁿ) $\mu \in [0,1]^K$

Bayesian Regret $E_{\mu \sim P} \left[\mu^{*T} - \sum_{t=1}^T \mu(a_t) \right]$

For simplicity,

- single parameter reward family
- finite support for reward & $P(\mu)$
- a^* is unique for each μ .

$\mu \sim$ Bernoulli(θ)
 unit-var Gaussian
 zero-mean Gaussian

Fact: History $H_t = \{a_1, \dots, a_t\}$

Posterior distribution $P_H(\mu) = P(\mu | H_t = H)$ doesn't depend on algorithm given history.

Thompson Sampling

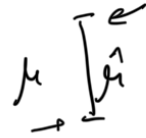
Observe $H_{t-1} = H$

Sample mean vector from posterior $\mu_t \sim P_H(\mu)$ ←

Choose best arm according to μ_t , $a_t = \arg \max_a \mu_t(a)$ ←

Update $H_t \leftarrow (a_t, r_t) \cup H_{t-1}$

If prior is correct, $O(\sqrt{KT \log T})$



Thm: Bayesian regret ↓

Proof: At round t , $BR_t = E_{\mu \sim P} [\mu(a^*) - \mu(a_t)]$

$$= E_{\mu \sim P} [E_{H_t} [\mu(a^*) - \mu(a_t) | H_t]]$$

$$= E_{H_t} [E_{\mu \sim P} [\mu(a^*) - \mu(a_t) | H_t]] \quad *$$

Claim: a_t and a^* have same distribution given H_t .

$$a_t = \arg \max_a \mu_t(a), \quad \mu_t \sim P_r(\mu | H_{t-1}) \propto P(\mu) P(H_{t-1} | \mu)$$

a_t & a^* have same distr given H_{t-1}

$$\begin{aligned} \Pr(a_t = a | H_{t-1}) &= \Pr(\mu_t(a) \geq \mu_t(a^*) | H_{t-1}) \quad \mu_t \sim P_r(\mu | H_{t-1}) \\ &= \Pr(\mu(a) \geq \mu(a^*) | H_{t-1}) \\ &= \Pr(a^* = a | H_{t-1}) \end{aligned}$$

$$\begin{aligned} * \quad BR_t &= E_{H_{t-1}} [E_{\mu \sim P} [\mu(a^*) - \text{UCB}(a^*, H_{t-1}) + \text{UCB}(a_t, H_{t-1}) - \mu(a_t) | H_{t-1}]] \\ &\because E[\text{UCB}(a^*, H_{t-1}) | H_{t-1}] = E[\text{UCB}(a_t, H_{t-1}) | H_{t-1}] \\ &\quad \text{as } a^*, a_t \text{ have same distr given } H_{t-1}. \end{aligned}$$

$$= E[\mu(a^*) - \text{UCB}(a^*, H_{t-1})] + E[\text{UCB}(a_t, H_{t-1}) - \mu(a_t)]$$

$$\leq 0$$

a^*, a are random

whp $\mu(a^*) - \text{UCB}(a^*, H_{t_i}) \leq 0$

$$\text{UCB}(a_t, H_{t_i}) - \mu(a_t) = \sum_a \mathbb{1}_{a_t=a} (\text{UCB}(a, H_{t_i}) - \mu(a))$$

$$\leq \sum_a \mathbb{1}_{a_t=a} \sqrt{\frac{\log T}{n_{t_i}(a)}}$$

$$\text{BR} \leq \mathbb{E} \left[\sum_t \sum_a \mathbb{1}_{a_t=a} \sqrt{\frac{\log T}{n_{t_i}(a)}} \right]$$

$$= \mathbb{E} \left[\sum_a \sum_t \mathbb{1}_{a_t=a} \sqrt{\frac{\log T}{n_{t_i}(a)}} \right]$$

$$= \mathbb{E} \left[\sum_a \underbrace{\sum_{j=1}^{n_T(a)} \frac{1}{\sqrt{j}}}_{\sqrt{n_T(a)}} \sqrt{\log T} \right]$$

$$= O \left(\mathbb{E} \left[\sum_a \sqrt{n_T(a)} \right] \sqrt{\log T} \right)$$

$$\leq O \left(\sqrt{\log T} \sqrt{\underbrace{\sum_a 1}_K \cdot \underbrace{\sum_a n_T(a)}_T} \right) \quad \text{whp}$$

misspecified ε Standard
Carm. regret $+ \varepsilon T$
sublinear - linear

Note: Can be improved to $O(\sqrt{KT \text{entropy}(a^*)})$

$\sqrt{KT \log K}$ $a^* = \arg \max_a P(\mu)$

K - information gain

Continuous action space

$$\mu(a) \sim \text{GP}(m, K)$$

\uparrow \uparrow \nwarrow
 mean covariance
 function function
=
 $K(a, a')$

$$\begin{pmatrix} \mu(1) \\ \vdots \\ \mu(K) \end{pmatrix} = \mu$$

Posterior :

$$m_t(x) = m_{t-1}(x) + K(x, H_{t-1}) K(H_{t-1}, H_{t-1})^{-1} \begin{pmatrix} y_t \\ x_t \end{pmatrix} - m(H_{t-1})$$

$$K_t(x, x) = K(H_{t-1}, H_{t-1}) - K(x, H_{t-1}) K(H_{t-1}, H_{t-1})^{-1} K(H_{t-1}, x)$$

μ, Σ

$$\mu(x) \rightarrow \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$\mu(x|H) \quad \Sigma_{x|H}$$

$$\mu_{1|2} \quad \Sigma_{1|2}$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (y_2 - \mu_2)$$

↓
observed value