

Stochastic Bandits

→ K arms/actions, T rounds

For $t=1, \dots, T$

Algorithm choose a_t action

Receive reward $r_t \in [0, 1]$

$$E[r_t(a)] = \mu(a)$$

Corresponding to action $a_t \leftarrow$

Goal: $\sum_{t=1}^T [\mu(a^*) - \mu(a_t)] =: R(T) \leftarrow$ explore-exploit tradeoff

Non-adaptive exploration

① Uniform exploration

Try each of the K actions N times \leftarrow explore

Play highest empirical reward action rest $T - NK$ times \leftarrow exploit

$$E[R(T)] = O(T^{2/3} (K \log T)^{1/3}) \text{ for } N \leq \left(\frac{T}{K}\right)^{2/3} (\log T)^{1/3} \checkmark$$

② eps-greedy

Toss coin w.p. $\epsilon_t = t^{-1/3} (K \log t)^{1/3} \leftarrow$

If success choose action uniformly @ random

else choose empirically best action

$$E[R(t)] = O(t^{2/3} (K \log t)^{1/3}) \text{ for any } t \leq T.$$

Adaptive Exploration + no. of times its explored

Action taken \wedge depend on previous rounds.

whp $|\hat{\mu}_t(a) - \mu(a)| \leq \sigma_t(a) = O\left(\sqrt{\frac{\log T}{n_t(a)}}\right)$ \leftarrow no. of times action a is taken before t.

$$\mu(a) \in \hat{\mu}_t(a) \pm \sigma_t(a)$$

confidence bound.

$$r(a, t) \quad \begin{bmatrix} r(1,1) & \dots & r(1,T) \\ \vdots & & \vdots \\ r(K,1) & \dots & r(K,T) \end{bmatrix}$$

Fix n . $w_p \geq 1-\delta$ $|\hat{\mu}(a) - \mu(a)| \leq \sqrt{\frac{\log 1/\delta}{n}}$ ✓

$$P(| \quad | \geq \epsilon) \leq \delta T^2 \leftarrow$$

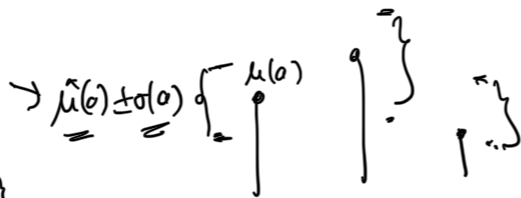
Union bound over all possible actions & all possible values of n .

$$\underbrace{K}_{=} \leq T \quad \underbrace{T}_{=}$$

$$\forall n, \forall a, \quad |\hat{\mu}(a) - \mu(a)| = O\left(\sqrt{\frac{\log 1/\delta}{n_t(a)}}\right) \quad w_p \geq 1 - T^2 \delta$$

$$\rightarrow = O\left(\sqrt{\frac{\log T}{n_t(a)}}\right) =: \sigma_t(a) \quad \delta = \frac{1}{T^4}$$

② Successive elimination



$$\mathcal{A} = \{1, \dots, K\}$$

→ Play all active actions once

Deactivate all a : $UCB_t(a) \leq LCB_t(a')$ for any $a' \in \mathcal{A}$

$$\hat{\mu}_t(a) + \sigma_t(a) \leq \hat{\mu}_t(a') - \sigma_t(a')$$

Note: a^* is always active.

whp. $UCB_t(a^*) = \hat{\mu}_t(a^*) + \sigma_t(a^*)$

$$\geq \mu(a^*)$$

$$\geq \mu(a') \geq \hat{\mu}_t(a') - \sigma_t(a') = LCB_t(a')$$

If a active $\underbrace{\mu(a^*) - \mu(a)}_{=: \Delta(a)} \leq 2\sigma_t(a^*) + 2\sigma_t(a) = 4\sigma_t(a)$

$$\because n_t(a^*) = n_t(a) \Rightarrow \text{active arm @ } t$$

$$\Delta(a) = O\left(\sqrt{\frac{\log T}{n_T(a)}}\right) \leftarrow \star$$

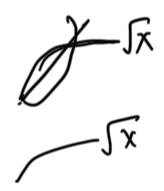
Let t be last round @ which a is active

$$n_T(a) \leq n_t(a) + 1$$

$$R(t) = \sum_{a=1}^K n_t(a) \Delta(a) = \sum_{a=1}^K n_t(a) O\left(\sqrt{\frac{\log T}{n_t(a)}}\right)$$

$$= O(K\sqrt{\log T}) + \sum_{a=1}^K \sqrt{n_t(a)}$$

\sqrt{x} is concave



$$\leq O(K\sqrt{\log T}) \sqrt{\frac{1}{K} \sum_{a=1}^K n_t(a)}$$

Jensen's ineq.

$$= O(\sqrt{Kt \log T}) \leftarrow \text{vs. } T^{2/3} (K \log T)^{1/3}$$

Gap-dependent bounds or Instance-based bounds

$$\star \Rightarrow n_T(a) = O\left(\frac{\log T}{(\Delta(a))^2}\right) \leftarrow$$

$$R(T) = \sum_a n_T(a) \Delta(a) \leq \sum_a \Delta(a) O\left(\frac{\log T}{(\Delta(a))^2}\right)$$

$$= O\left(\sum_a \frac{1}{\Delta(a)} \cdot \log T\right) \leftarrow$$

(4) Upper Confidence Band (UCB) simplify (Optimism under uncertainty)

pick arm which $\max_a UCB_t(a)$

$$a_t = \arg \max_a UCB_t(a)$$

$$\hat{\mu}_t(a) + \sigma_t(a)$$

UCB is high either 1) reward is high (exploit)

2) uncertainty is high (explore)

$$\begin{aligned} \Delta(a_t) = \mu(a^*) - \mu(a_t) &\leq \text{UCB}_t(a^*) - (\hat{\mu}(a_t) - \sigma_t(a_t)) \\ &\leq \text{UCB}_t(a_t) - \hat{\mu}(a_t) + \sigma_t(a_t) \leq 2\sigma_t(a_t) \\ &\stackrel{\text{from UCB sampling}}{=} O\left(\sqrt{\frac{\log T}{n_t(a_t)}}\right) \text{ why} \end{aligned}$$

Same Regret bounde
 instance relative
 instance independent.

Lower Bounds

1. For any bandit algo \exists a problem instance $E[R(T)] = \Omega(kT) \leftarrow$

2. For any bandit algo with non-adaptive exp \exists a problem instance
 s.t. $E[R(T)] = \Omega(k^{1/3} T^{2/3})$

3. For any bandit algo with non-adaptive exp $\forall \lambda \in (0, 1)$ $E[R(T)] = O(T^\lambda)$
 $\forall \lambda \in [2/3, 1)$ for all problem instances then for any problem instance
 a random permutation of arms yields

$$E[R(T)] = \Omega\left(T^\lambda \sum_a \Delta(a)\right) \quad \lambda = \underbrace{2(1-\gamma)}_{\forall \gamma \in [2/3, 1)}$$

$$= \lambda \log T$$

4. No algorithm can achieve $E[R(T)] = O(C_I \log t)$ for all problem instances, where C_I - depends on instance I but not on t .

Bandits with prior information

Constrained means - linear, Lipschitz $\mu \in \mathcal{F}$
 holds for continuous arms

Bayesian prior $P(\mu)$

Lipschitz bandits $x \in \mathcal{X}$ continuous arms/actions

$$|\mu(x) - \mu(x')| \leq L |x - x'| \quad \forall x, x' \in \mathcal{X} = [0, 1]$$

$P(x, x')$



$$E[R(T)] = T\mu(x^*) - \sum_{t=1}^T \mu(x_t)$$

$$= \underbrace{T\mu(x^*) - T\mu_k^*}_{\text{discretization error}} + \underbrace{\left(T\mu_k^* - \sum_{t=1}^T \mu(x_t) \right)}_{\text{K-armed bandit regret.}}$$

best amongst K arms (centers of bins)

$$\leq O\left(T \frac{L}{K} + \sqrt{KT \log T}\right) \leftarrow \text{adaptive}$$

$$= O\left((L \log T)^{1/3} T^{2/3}\right) \quad K?$$