

Lecture #4: Stochastic bandits (part 2)

Explore-then-Commit

(input $N \in \mathbb{N}$)

Pull each arm N times.

Let $k^* \in \operatorname{argmax}_{k \in \{1, \dots, K\}} \hat{\mu}_k(NK)$

Play k^* until time T

Simple algorithm clearly separating exploration from exploitation.
Easy analysis.

Theorem ETC with $N = \frac{\ln T}{\Delta^2}$ satisfies

$$R_T \leq \sum_{k=1}^K \frac{\Delta_k}{\Delta^2} \ln T + \Delta_k.$$

Proof: Similarly to Greedy in the full info setting:

For $n=1, \dots, N$ let $t_k(n)$ be the deterministic time where k is pulled for the n -th time.

By Hoeffding inequality:

$$\mathbb{P}\left(\frac{1}{N} \sum_{n=1}^N X_{k^*}(t_k(n)) - X_{k^*}(k^*) \geq 0\right) \leq e^{-N \Delta_k^2}$$

$$\text{So: } R_T \leq \sum_{k=1}^K N \Delta_k + \sum_{k=1}^K \Delta_k (T-N) e^{-N \Delta_k^2}$$

Plugging the value of N :

$$R_T \leq \sum_{k=1}^K \frac{\Delta_k}{\Delta_k^2} \ln T + \Delta_k. \quad \square$$

Remark: • Again we can get a distribution-free bound scaling as $O((K \ln T)^{1/3} T^{2/3})$, and the instance dependent version requires knowledge of Δ

Two main drawbacks of these methods:

- they require knowledge of Δ .
- they scale in $\frac{1}{\Delta^2} (\ln T)^{2/3}$ in distribution-free bounds)

This is because they use a uniform exploration: each arm is explored the

same amount of time.

exploration rounds depend on past observations.

A better strategy is to use an adaptive exploration: better arms are explored more often. The idea is that a very bad arm is quicker to detect as sub-optimal.

Successive Eliminations

Let $K = [K]$

While $\text{Card}(K) > 1$:

 Pull each arm in K once

 For $k \in K$:

$$\text{if } \hat{\mu}_k(t) + \sqrt{\frac{2 \ln T}{N_k(t)}} \leq \max_{k' \in K} \hat{\mu}_{k'}(t) - \sqrt{\frac{2 \ln T}{N_{k'}(t)}} \text{ then } K \leftarrow K \setminus \{k\}$$

 Pull the only arm in K until the end

Theorem: For SE, the regret satisfies for any $T \in \mathbb{N}$:

$$R_T \leq \sum_{k, \Delta_k} \left(\frac{32 \ln T}{\Delta_k} + 1 \right) + \frac{K}{T}$$

Proof: Define the clean event

$$\mathcal{E} = \left\{ \begin{array}{l} \forall k \neq k^*, \forall t \in [T], \hat{\mu}_k(t) - \mu_k \leq \sqrt{\frac{2 \ln T}{N_k(t)}} \\ \forall t \in [T], \hat{\mu}_{k^*}(t) - \mu_{k^*} \geq -\sqrt{\frac{2 \ln T}{N_{k^*}(t)}} \end{array} \right\}$$

Thanks to our concentration lemma on $\hat{\mu}_k$:

$$P(\mathcal{E}) \geq 1 - K \sum_{t=1}^T \frac{1}{T^4} \geq 1 - \frac{K}{T^3}$$

We now bound $\mathbb{E}[N_k(T) \mathbb{1}_{\mathcal{E}}]$.

Note that when \mathcal{E} holds, we always have:

$$\hat{\mu}_k^*(t) + \sqrt{\frac{2 \ln T}{N_k(t)}} \geq \mu_k^* \geq \mu_k \geq \hat{\mu}_k(t) - \sqrt{\frac{2 \ln T}{N_k(t)}}$$

So k^* is never eliminated from \mathcal{K} .

For a suboptimal arm k , let N_k be the smallest integer such that:

$$4 \sqrt{\frac{2 \ln T}{N_k(t)}} \leq \Delta_k$$

$$\text{i.e. } N_k = \left\lceil \frac{32 \ln T}{\Delta_k^2} \right\rceil.$$

Then once all arms in \mathcal{K} have been pulled N_k times, we have if \mathcal{E} holds

$$\hat{\mu}_k(t) + \sqrt{\frac{2 \ln T}{N_k}} \leq \mu_k + 2 \sqrt{\frac{\ln T}{N_k}} \leq \mu_k^* - 2 \sqrt{\frac{\ln T}{N_k}} \leq \hat{\mu}_k^*(t) - \sqrt{\frac{\ln T}{N_k}}$$

So k is eliminated after at most N_k pulls if \mathcal{E} holds:

$$\mathbb{E}[N_k(T) \mathbb{1}_{\mathcal{E}}] \leq \left\lceil \frac{32 \ln T}{\Delta_k^2} \right\rceil$$

Finally:

$$R_T \leq \sum_{k, \Delta_k > 0} \Delta_k \left(\mathbb{E}[N_k(T) \mathbb{1}_\epsilon] + \mathbb{E}[N_k(T) \mathbb{1}_{\text{not } \epsilon}] \right)$$

$$\leq \sum_{k, \Delta_k > 0} \Delta_k \sqrt{\frac{32 \ln T}{\Delta_k^2}} + T(1 - \mathbb{P}(\epsilon))$$

$$\leq \sum_{k, \Delta_k > 0} \left(32 \frac{\ln T}{\Delta_k} + 1 \right) + \frac{K}{T} \quad \square$$

Remarks

• SE assumes a prior knowledge of T .
 assuming T is not too restrictive in practice, as we can use the doubling trick
 see exercise session #2

• We can easily get a better constant than 32

• This instance dependent bound also implies a distribution free bound $O(\sqrt{TK \ln T})$

see exercise session #2.

Upper Confidence Bound (UCB)

Pull each arm once

For $t \geq K+1$:

$$a_t \in \arg \max_{k \in [K]} \underbrace{\hat{\mu}_k(t-1) + \sqrt{\frac{2 \ln(t)}{N_k(t-1)}}}_{\text{UCB score}}$$

- Greedy, but with UCB scores
 \rightarrow no underestimation of μ_{k^*} (with high probability)
- No prior knowledge of T .

Theorem

For any TCIN, the regret of UCB satisfies

$$R_T \leq \sum_{k, \Delta_k > 0} \left(8 \frac{\ln T}{\Delta_k} + 2 \right)$$

Proof:

For $t \geq K+1$ and $k \neq k^*$, let

$$E_{k,t} = \left\{ \begin{array}{l} \hat{\mu}_k(t) - \mu_k \leq \sqrt{\frac{2 \ln t}{N_k(t)}} \\ \hat{\mu}_{k^*}(t) - \mu_{k^*} \geq -\sqrt{\frac{2 \ln t}{N_{k^*}(t)}} \end{array} \right\}$$

$$P(E_t) \geq 1 - \frac{2}{t^3}$$

If $E_{k,t}$ holds and $k \neq k^*$ is pulled at time t , then:

$$\hat{\mu}_k(t) + \sqrt{\frac{2 \ln t}{N_k(t-1)}} \geq \hat{\mu}_{k^*}(t) + \sqrt{\frac{2 \ln t}{N_{k^*}(t-1)}}$$

$$E_{k,t} \text{ holds, so } \mu_k + 2\sqrt{\frac{2 \ln t}{N_k(t-1)}} \geq \hat{\mu}_k(t) + \sqrt{\frac{2 \ln t}{N_k(t-1)}}$$

$$\text{and } \hat{\mu}_{k^*}(t) + \sqrt{\frac{2 \ln t}{N_{k^*}(t-1)}} \geq \mu_{k^*}$$

In particular:

$$\mu_k + 2\sqrt{\frac{8 \ln T}{N_k(t-1)}} \geq \mu_k^*$$

$$\text{so } (\epsilon_{k,t} \text{ and } a_t = k) \Rightarrow N_k(t-1) \leq \frac{8 \ln T}{\Delta_k^2}$$

From here for $k \neq k^*$

$$\mathbb{E}[N_k(T)] = 1 + \mathbb{E}\left[\sum_{t=k+1}^T \mathbb{1}(a_t = k \text{ and } \epsilon_{k,t}) + \mathbb{1}(a_t = k \text{ and not } (\epsilon_{k,t}))\right]$$

$$\leq 1 + \mathbb{E}\left[\sum_{t=k+1}^T \mathbb{1}(a_t = k \text{ and } N_k(t-1) \leq \frac{8 \ln T}{\Delta_k^2})\right] + 2 \sum_{t=k+1}^T \frac{1}{t^3}$$

$$\leq 1 + \mathbb{E}\left[\sum_{t=k+1}^T \mathbb{1}(a_t = k \text{ and } N_k(t-1) \leq \frac{8 \ln T}{\Delta_k^2})\right] + 2 \int_1^{\infty} \frac{1}{s^3} ds$$

$$\leq 1 + \mathbb{E}\left[\left(\frac{8 \ln T}{\Delta_k^2} + 1\right) - 1\right] + \mathbb{E}[T^{-2}]_{1}^{\infty}$$

$$\leq 2 + \frac{8 \ln T}{\Delta_k^2} \quad \square$$

- The $\frac{8 \sum \ln T}{k \Delta_k^2 \Delta_k}$ instance dependent bound is nearly optimal. We'll see in exercise session that UCB can

be made optimal with respect to the lower bound we are going to prove next week.

- UCB is said to use the optimism in the face of uncertainty principle: aiming at the best statistically possible scenario is a good strategy here.

- Previous algorithms/results hold for independent bounded rewards $X_k(t) \in [0, 1]$

They can be easily extended to independent σ sub-gaussian rewards, as similar concentration bounds hold.

eg UCB scores become

$$\hat{\mu}_k(t-1) + \sqrt{\frac{\sigma^2 \ln(t)}{2N_k(t-1)}} \rightarrow \text{same regret bounds, rescaled by } \sigma$$

What if σ is unknown? can be estimated see exercise session #3.