

Lecture #3: Stochastic bandits (part 1)

Full Information Setting

At each round $t = 1, \dots, T$:

- agent picks an arm $a \in \{1, \dots, K\}$ (possibly at random)

- observes reward vector $X(t) \in [0, 1]^K$

gets reward $X_{a_t}(t)$.

$$R_T = \max_{k \in [K]} \sum_{t=1}^T X_k(t) - \sum_{t=1}^T X_{a_t}(t)$$

| a_t is $\tau(U, X(1), \dots, X(t-1))$
 | randomization | measurable

As learning with experts, but:

- rewards instead of loss ($l_t \leftrightarrow 1 - X(t)$)
- choose pure actions (K -simplex $\leftrightarrow \{1, \dots, K\}$)
but can randomize over actions.

The $X(t)$ were chosen adversarially (worst case) in 1st lecture.

What if instead they are stochastic?

Assume

- $(X_a)_t$ are iid.

- $X_a(t) \sim \nu_a$ with $E[X_a(t)] = \mu_a$.

Problem should be easier?

→ not really: we proved the lower bound in this setting:

for any algorithm, with $X_R(t) \sim \text{Ber}\left(\frac{1}{2}\right)$

$$\mathbb{E}[R_T] \geq \sqrt{\frac{T}{2} \ln K}$$

However, we can have much better results with the pseudo-regret:

$$\bar{R}_T = \max_{k \in [K]} \sum_{t=1}^T \mu_k = \sum_{t=1}^T \mu_{a_t}$$

↳ expectation w.r.t. the realizations of $X(t)$
but still a random variable!

Previous example yields $\bar{R}_T = 0$. Makes sense: we cannot guess in advance heads or tails.

⚠ Warning: $E[\bar{R}_T] \neq E[R_T]$

Actually, $E[R_T] > E[\bar{R}_T]$. Why?

$$\bar{R}_T = T \max_k \mu_k - \sum_{t=1}^T \mu_{a_t}$$

→ from now on, we will write R_T for the pseudo-regret.

Notations:

- $\mu^* = \max_k \mu_k$

- $\Delta_k = \mu^* - \mu_k \begin{cases} > 0 & \text{for sub-optimal arms} \\ = 0 & \text{for optimal arms} \end{cases}$

- $\Delta = \min_{k, \Delta_k > 0} \Delta_k$

$$N_k(t) = \sum_{s=1}^t \mathbb{1}_{a_s=k}$$

number of pulls on arm k .

Lemma:

$$\text{For any policy, } LR_T = \sum_{k=1}^K \Delta_k N_k(T)$$

Proof:

$$\begin{aligned} R_T &= \sum_{t=1}^T \mu^* - \mu_A \\ &= \sum_{t=1}^T \mu^* - \sum_{k=1}^K \mathbb{1}_{a_t=k} \mu_k \\ &= \sum_{k=1}^K \sum_{t=1}^T (\mu^* - \mu_k) \mathbb{1}_{a_t=k} \\ &= \sum_{k=1}^K \Delta_k \sum_{t=1}^T \mathbb{1}_{a_t=k} \\ &= \sum_{k=1}^K \Delta_k N_k(T) \end{aligned}$$

□.

Greedy algorithm (or Follow The Leader)

choose a_1 arbitrarily

For $t \geq 2$:

$$a_t \in \operatorname{argmax}_{a \in [K]} \left(\sum_{\tau=1}^{t-1} X_a(\tau) \right)$$

Theorem

For any $(\mu_1, \dots, \mu_K) \in [0, 1]^K$ and $T \in \mathbb{N}$, Greedy satisfies in the Full Information setting:

$$\mathbb{E}[R_T] \leq \sum_{k: \Delta_k > 0} \frac{1}{\Delta_k}$$

Proof: $\mathbb{E}[R_T] = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(T)]$

Let us bound $\mathbb{E}[N_k(T)]$ for any k with $\Delta_k > 0$.

Let $k^* \in \operatorname{argmax}_k \mu_k$.

$$\mathbb{E}[N_k(T)] \leq \sum_{t=1}^T \mathbb{P}\left(\frac{1}{T} \sum_{s=1}^t X_k(s) - X_{k^*}(s) \geq 0\right)$$

$$\leq \sum_{t=1}^T \mathbb{P}\left(\frac{1}{T} \sum_{s=1}^t (X_k(s) - \mu_k) - \frac{1}{T} \sum_{s=1}^t (X_{k^*}(s) - \mu_{k^*}) \geq t\Delta_k\right)$$

$$\leq \sum_{t=1}^T e^{-t\Delta_k^2} \quad \leftarrow \frac{e^{-t\Delta_k^2}}{1 - e^{-\Delta_k^2}} = \frac{1}{e^{\Delta_k^2} - 1}$$

Hoeffding inequality.

$$\leq \frac{1}{\Delta_k^2}$$

$$e^{-1} \geq \Delta_k^2$$

$$\therefore \mathbb{E}[R_T] = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(T)]$$

$$\leq \sum_{\substack{k \\ \mu_k < \mu^{*}}} \Delta_k \cdot \frac{1}{\Delta_k^2}$$

□.

Bandit Setting (random table model)

At each round $t = 1, \dots, T$:

- agent picks an arm $a_t \in \{1, \dots, K\}$ (possibly at random)
- observes and gets reward $X_{a_t}^{(t)} \in [0, 1]$

| a_t is $r(U_0, X_{a_1}^{(1)}, U_1, \dots, X_{a_{t-1}}^{(t-1)}, U_{t-1})$ -measurable

$$R_T = \max_{k \in [K]} \sum_{t=1}^T \mu_k - \sum_{t=1}^T \mu_{a_t}$$

→ only observe the reward of the pulled arm

→ exploration vs exploitation trade-off

estimate optimal arm by pulling all arms

maximize reward by pulling arm which seems the best

This setting is sometimes called random table model and is known to be equivalent (from a probabilistic point of view) to the following stack of rewards model.

Bandit Setting (stack of rewards model)

At each round $t = 1, \dots, T$:

- agent picks an arm $a_t \in \{1, \dots, K\}$ (possibly at random)

- observes and gets reward $X_{a_t}^{(N_{a_t}(t))} \in [0, 1]$

| a_t is $r(U_0, X_{a_1}^{(1)}, U_1, \dots, X_{a_{t-1}}^{(N_{a_{t-1}}(t-1))}, U_{t-1})$ -measurable

Same definition of N_{a_t}

Stack of rewards model allows easier proofs, but heavier in notations.

Unless specified otherwise, we will consider the random table model in the following.

Notation

$$\cdot \hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^t X_k(s) \frac{1}{\{a_s=k\}} \quad (\text{empirical mean})$$

$$\text{Random stack} \rightarrow \hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^{N_k(t)} X_k(s)$$

Greedy algorithm (Bandit setting)

For $t=1, \dots, k$:

$$a_t = t$$

For $t \geq K+1$:

$$a_t \in \operatorname{argmax}_{k \in [K]} \hat{\mu}_k(t-1)$$

Theorem For $v_1 = \operatorname{Ber}\left(\frac{3}{4}\right)$, $v_2 = \operatorname{Ber}\left(\frac{1}{4}\right)$, Greedy satisfies

in the bandit setting:

$$\mathbb{E}[R_T] \geq \frac{T-1}{32}$$

Proof:

$$\Pr(X_1(1) = 0, X_2(2) = 1) = \left(\frac{1}{4}\right)^2 = \frac{1}{16}$$

If $X_1(1) = 0$ and $X_2(2) = 1$, Greedy will keep pulling the arm 2

until T , so that

$$\mathbb{E}[N_2(T)] \geq \frac{T-1}{16}, \quad \square$$

Greedy does not explore enough. It can underestimate the optimal arm and never pull it again.

Explore-then-Commit algorithm

parameter $n \in \mathbb{N}^*$

For $t=1, \dots, nk$: explore by drawing each arm n times.

For $t \geq nk+1$:

pull the best empirical arm until the end, i.e.

$$a_t = \arg\max_k \hat{\mu}_k(nk)$$

Simple algorithm clearly separating exploration from exploitation.
Easy analysis

Theorem:

For any $1 \leq n \leq T/K$, ETC

has expected regret

$$\mathbb{E}[R_T] \leq n \sum_{k=1}^K \Delta_k + (T-nk) \sum_{k=1}^K \Delta_k \exp(-n\Delta_k^2)$$

Proof:

$$R_T = \sum_{k=1}^K \Delta_k N_k(T).$$

$$\text{if } n \leq T/K, \quad N_k(T) = \begin{cases} n & \text{if } k \neq \arg\max \hat{\mu}_k(nk) \\ n + t - nk & \text{if } k = \arg\max \hat{\mu}_k(nk) \end{cases}$$

$$R_T \leq n \sum_{k=1}^K \Delta_k + (T-nk) \sum_{k=1}^K \Delta_k \quad \text{if } k = \arg\max \hat{\mu}_k(nk).$$

$$\text{Let } k^* = \arg\max_k \mu_k \quad (\mu_k = \hat{\mu}_k)$$

$$\mathbb{E}[R_T] \leq n \sum_{k=1}^K \Delta_k + (T-nk) \sum_{k=1}^K \Delta_k P(k = \arg\max \hat{\mu}_k(nk))$$

$$\leq n \sum_{k=1}^K \Delta_k + (T-nk) \sum_{k=1}^K \Delta_k P(\hat{\mu}_k(nk) \geq \hat{\mu}_{k^*}(nk))$$

$$\Pr(\hat{\mu}_k(nK) \geq \hat{\mu}_{k+}(nK)) = \Pr\left(\sum_{s=1}^n X_k(s) - \sum_{s=1}^n X_{k+}(s) \geq 0\right)$$

$$= \Pr\left(\sum_{s=1}^n (X_k(s) - \mu_k) - \sum_{s=1}^n (X_{k+}(s) - \mu_{k+}) \geq n\Delta_k\right)$$

Hoeffding: $\leq e^{-n\Delta_k^2}$ \square .

- n too large \rightarrow explore too much
- n too small \rightarrow not enough exploration, might pull suboptimal arm for T steps.

what n should we choose?

for $\Delta = \min_{k, \Delta_k > 0} \Delta_k$ and $n = \lceil \frac{\ln(T)}{\Delta^2} \rceil$

$$\mathbb{E}[R_T] \leq \sum_{k=1}^K \frac{\Delta_k \ln T}{\Delta^2} + \sum_{k=1}^K \Delta_k$$

- Actually, we even showed a **high probability regret bound**, i.e. with $n = \lceil \frac{\ln(K/\delta)}{\Delta^2} \rceil$, $E[T]$ satisfies with probability at least $1-\delta$,

$$R_T \leq \lceil \frac{\ln(K/\delta)}{\Delta^2} \rceil \sum_{k=1}^K \Delta_k$$

- ETC is easy to analyze \rightarrow direct application of Hoeffding inequality

Yet, this use of Hoeffding ineq. is not always possible.

Instead, we use the following concentration lemma.

Lemma: (bandit concentration)

For any bandit algorithm, any $k \in [K]$, $t \in \mathbb{N}$, $\delta \in (0, 1)$:

$$\Pr(\mu_k - \hat{\mu}_k(t) \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(t)}}) \leq t\delta.$$

$$\Pr(\hat{\mu}_k(t) - \mu_k \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(t)}}) \leq t\delta.$$

1) This is not a trivial consequence of Hoeffding inequality,

$N_k(t)$ is a random variable and $\hat{\mu}_k(t), N_k(t)$ are not independent!



Hoeffding inequality indeed gives

$$\Pr\left(\frac{1}{n} \sum_{s=1}^n X_k(s) - \mu_k \geq \sqrt{\frac{\ln(1/\delta)}{2n}}\right) \leq e^{-\ln(1/\delta)} = \delta.$$

But here, n is a random variable and is not independent from $\hat{\mu}_k(t)$

• What if instead we used Azuma-Hoeffding on $(X_k(s) - \mu_k) \mathbf{1}_{\{a_s=k\}}$?

martingale increment bounded between

$-\mu_k$ and $1 - \mu_k$.

$$\Pr\left(\sum_{s=1}^t (X_k(s) - \mu_k) \mathbf{1}_{\{a_s=k\}} \geq \sqrt{\frac{t}{2} \ln(1/\delta)}\right) \leq \delta t$$

$$\Pr(\hat{\mu}_k(t) - \mu_k \geq \sqrt{\frac{t}{N_k(t)} \frac{\ln(1/\delta)}{2N_k(t)}}) \leq \delta t$$

differences with our Lemma

getting rid of this $\sqrt{\frac{t}{N_k(t)}}$ factor is a big deal!

Proof 1 (stack of rewards model)

$$Z_n = \sum_{k=1}^n (X_k(\omega) - \mu_k)$$

$$\hat{\mu}_k(r) - \mu_k = \frac{1}{N_k(r)} Z_{N_k(r)}$$

Hoeffding inequality yields for any n .

$$P\left(\frac{1}{n} Z_n \geq \sqrt{\frac{\ln(1/\delta)}{2n}}\right) \leq \delta$$

Moreover, $N_k(t) \in [0, H]$ a.s., so we can do a union bound:

$$P\left(\exists n \in [1, H], \frac{Z_n}{n} \geq \sqrt{\frac{\ln(1/\delta)}{n}}\right) \leq \delta H. \quad (\text{bound is automatic for } n=0)$$

so that in particular

$$P\left(\frac{Z_{N_k(t)}}{N_k(t)} \geq \sqrt{\frac{\ln(1/\delta)}{n}}\right) \leq \delta t. \quad \text{by inclusion of the event.}$$

Symmetric arguments for the second inequality. \square

Proof 2 (random table model)

$$\text{Let } Z_t = \sum_{k=1}^t (X_k(\omega) - \mu_k) \mathbb{1}_{a_k=k}. \quad \hat{\mu}_k(r) - \mu_k = \frac{1}{N_k(r)} Z_t$$

Here using a union bound on all the possible values of $(1_{a_k=k})_{k=1 \dots t}$ is a bad idea...

1) We first prove that $\forall x \in \mathbb{R}, \mathbb{E}[e^{x Z_t - \frac{x^2}{8} N_k(t)}] \leq 1$.

For that, we show that $M_t = \exp(x Z_t - \frac{x^2}{8} N_k(t))$ is a supermartingale, so that

$$\mathbb{E}[M_t] \leq \mathbb{E}[M_0] = 1.$$

$$\text{Let } \mathcal{F}_{t-1} = \sigma(X_{a_1}(1), \dots, X_{a_{t-1}}(t-1))$$

at t is \mathcal{F}_{t-1} measurable so that:

$$\begin{aligned} \mathbb{E}[M_r | F_{r-1}] &= \mathbb{E}\left[e^{(X_{\alpha(t)} - \mu_{\alpha}) \cdot \frac{x^2}{2}} \mathbf{1}_{a_r=k} | F_{r-1}\right] M_{t-1} \\ &= \left(\mathbb{E}\left[e^{(X_{\alpha(t)} - \mu_{\alpha}) \cdot \frac{x^2}{2}} \mathbf{1}_{a_r=k}\right] \mathbf{1}_{a_r=k} + \mathbb{E}\left[e^{(X_{\alpha(t)} - \mu_{\alpha}) \cdot \frac{x^2}{2}} \mathbf{1}_{a_r \neq k}\right] \mathbf{1}_{a_r \neq k}\right) M_{t-1} \end{aligned}$$

Hoeffding's lemma (conditional) gives $\ln(\mathbb{E}[e^{(X_{\alpha(t)} - \mu_{\alpha}) \cdot \frac{x^2}{2}} | F_{r-1}]) \leq \frac{x^2}{8}$

$$\Rightarrow \mathbb{E}[e^{(X_{\alpha(t)} - \mu_{\alpha}) \cdot \frac{x^2}{2}} | F_{r-1}] \leq 1.$$

$$\therefore \mathbb{E}[M_r | F_{r-1}] \leq (\mathbf{1}_{a_r=k} + \mathbf{1}_{a_r \neq k}) M_{t-1}$$

$$\leq M_{t-1}.$$

$$\text{So we showed } \mathbb{E}[e^{x z_r - \frac{x^2}{8} N_k(t)}] \leq 1.$$

2) We now prove that $\forall \varepsilon > 0, \forall n \geq 1$,

$$P(z_r \geq \varepsilon \text{ and } N_k(t) = n) \leq e^{-\frac{2\varepsilon^2}{n}}$$

Indeed, by Markov-Chernoff bounding for any $x > 0$:

$$\begin{aligned} P(z_r \geq \varepsilon \text{ and } N_k(t) = n) &\leq e^{-nx\varepsilon} \mathbb{E}\left[e^{xz_r} \mathbf{1}_{\{N_k(t) = n\}}\right] \\ &= e^{-x\varepsilon + \frac{x^2}{8}n} \mathbb{E}\left[e^{xz_r - \frac{x^2}{8}N_k(t)} \mathbf{1}_{\{N_k(t) = n\}}\right] \\ &\leq e^{-x\varepsilon + \frac{x^2}{8}n} \mathbb{E}\left[e^{xz_r - \frac{x^2}{8}N_k(t)}\right] \\ &\leq e^{-x\varepsilon + \frac{x^2}{8}n} \quad \text{(} \leq 1 \text{ thanks to 1)} \end{aligned}$$

Taking $x = \frac{4\varepsilon}{n}$ finally yields

$$P(z_r \geq \varepsilon \text{ and } N_k(t) = n) \leq e^{-\frac{2\varepsilon^2}{n}}. \quad (\star)$$

3) We conclude using a union bound: (as in the stock of rewards model)

$$\begin{aligned}
 \Pr\left(\hat{\mu}_k(t) - \mu_k \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(t)}}\right) &= \sum_{n=1}^t \Pr\left(\hat{\mu}_k(t) - \mu_k \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(t)}} \text{ and } N_k(t)=n\right) \\
 &= \sum_{n=1}^t \Pr\left(\frac{Z_t}{N_k(t)} \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(t)}} \text{ and } N_k(t)=n\right) \\
 &= \sum_{n=1}^t \Pr\left(Z_t \geq \sqrt{\frac{n \ln(1/\delta)}{2}} \text{ and } N_k(t)=n\right) \\
 &\leq \sum_{n=1}^t e^{-\ln(1/\delta)} = t\delta. \quad \square
 \end{aligned}$$

Notes on this proof:

- We saw last week that the conditional version of Hoeffding's lemma could be generalized into

X bounded random variable, U, V two \mathcal{G} -measurable random variables with $U \leq X \leq V$ a.s.

then $\forall y \in \mathbb{R}$,

$$\ln \mathbb{E}[e^{yX} | \mathcal{G}] \leq y \mathbb{E}[X | \mathcal{G}] + \frac{y^2}{2}(V-U)^2$$

This can be applied to

$$Z_t = (X_k(t) - \mu_k) \mathbf{1}_{\{a_t=k\}}$$

$$g = F_{t-1}$$

$$U_t = -\mu_k \mathbf{1}_{\{a_t=k\}}$$

$$V_t = (1-\mu_k) \mathbf{1}_{\{a_t=k\}}$$

and directly entails $\mathbb{E}\left[e^{\eta(X_k(t) - \mu_k)} \mathbf{1}_{\{a_t=k\}} \mid F_{t-1}\right] \leq \exp\left(\frac{\eta^2}{2} \mathbf{1}_{\{a_t=k\}}\right)$ without the need for the

$1 = \mathbf{1}_{\{a_t=k\}} + \mathbf{1}_{\{a_t \neq k\}}$ trick used in step 1).

The question is: || Don't we have a generalized version of the Hoeffding-Azuma inequality with such predictable ranges $V_T - U_T$?

Yes, we do have something in terms of constant upper bounds $V_T - U_T \leq \Delta_T \in \mathbb{R}$ a.s.

but $V_T - U_T = \mathbb{1}_{\text{for } k \geq 1}$ can only be bounded by $\Delta_T = 1$ here, so steps 2) and 3) are still needed.

For unbounded, but σ -sub-Gaussian Variables $X_k(t)$, we still have, thanks to the stack of rewards model: $\Pr(\mu_k - \hat{\mu}_k(t) \geq \sqrt{\frac{2\ln(1/\delta)}{\sigma^2 N_k(t)}}) \leq \delta$.

ϵ -Greedy sequence of probabilities ϵ_t .

For $t=1, \dots, K$:

$$a_t = t$$

For $t \geq K+1$:

$\left\{ \begin{array}{l} \text{with proba } \epsilon_t, \quad a_t \sim \mathcal{U}([K]) \quad \text{explore uniformly at random} \\ \text{with proba } 1-\epsilon_t, \quad a_t \in \operatorname{argmax}_{k \in [K]} \hat{\mu}_k(t-1) \end{array} \right.$

Theorem

For $\epsilon_t = \min\left\{1, \frac{cK}{t\Delta^2}\right\}$ where c is a large enough universal constant, ϵ -greedy satisfies for a large enough universal constant c'

$$R_T \leq \frac{c'}{\Delta^2} \sum_{t=1}^T (\Delta \ln T + 1)$$

Proof: For any δ with $\Delta_\delta > 0$,

$$P(a_r = b) \leq \frac{\varepsilon_r}{K} + P(\hat{\mu}_\alpha(t-1) \geq \hat{\mu}_{\alpha^*}(t-1))$$

$$\leq \frac{\varepsilon_r}{K} + P(\hat{\mu}_\alpha(t-1) - \mu_\alpha \geq \frac{\Delta_\delta}{2}) + P(\mu_{\alpha^*} - \hat{\mu}_{\alpha^*}(t-1) \geq \frac{\Delta_\delta}{2}).$$

$$\begin{aligned} P(\hat{\mu}_\alpha(t-1) - \mu_\alpha \geq \frac{\Delta_\delta}{2}) &\leq \sum_{n=1}^{L(x_r)} P(N_\alpha(t-1) \leq x_r) + \sum_{n=L(x_r)+1}^{T-1} P(\hat{\mu}_\alpha(t-1) - \mu_\alpha \geq \frac{\Delta_\delta}{2} \text{ and } N_\alpha(t-1) = n) \\ &\leq P(N_\alpha(T-1) \leq x_r) + \sum_{n=L(x_r)+1}^{T-1} e^{-\frac{\Delta_\delta^2}{2}n} \quad (*) \\ &\leq \underbrace{P(N_\alpha^R(T-1) \leq x_r)}_{\substack{\text{number of times } \alpha \text{ is pulled at random} \\ (\text{in following the } E \text{ path event})}} + 2 \frac{e^{-\frac{\Delta_\delta^2}{2}x_r}}{\Delta_\delta^2} \end{aligned}$$

where $0 \leq x_r \leq s$

$$E[N_\alpha^R(T-1)] = 1 + \frac{1}{K} \sum_{\alpha=K+1}^{T-1} \varepsilon_\alpha = \frac{1}{K} \sum_{\alpha=1}^{T-1} \varepsilon_\alpha$$

$$Var(N_\alpha^R(T-1)) = \sum_{\alpha=K+1}^{T-1} \frac{\varepsilon_\alpha(1-\frac{\varepsilon_\alpha}{K})}{K} \leq \frac{1}{K} \sum_{\alpha=1}^{T-1} \varepsilon_\alpha$$

Recall

Bernstein Inequality

Lecture #2

Let X_1, \dots, X_T be random variables in $[0, 1]$ s.t. $Var[X_\alpha | X_1, \dots, X_{\alpha-1}] = \sigma_\alpha^2$

Then for all $\varepsilon > 0$:

$$P\left(\sum_{\alpha=1}^T X_\alpha - E[X_\alpha | X_1, \dots, X_{\alpha-1}] \leq -\varepsilon\right) \leq \exp\left(\frac{-\varepsilon^2/2}{\sum_{\alpha=1}^T \sigma_\alpha^2 + \frac{\varepsilon}{2}}\right)$$

so here for $x_r = \frac{1}{2K} \sum_{\alpha=1}^{T-1} \varepsilon_\alpha$

$$\mathbb{P}(N_a^R(t-1) \leq x_t) = \mathbb{P}\left(N_a^R(t-1) - \mathbb{E}[N_a^R(t-1)] \leq -x_t\right)$$

$$\leq \exp\left(-\frac{x_t^2/2}{\frac{5}{2}x_t}\right) = e^{-\frac{x_t}{5}}$$

Moreover: $x_t = \frac{1}{2K} \sum_{j=1}^{t-1} \min\left(1, \frac{cK}{\Delta_j^2}\right)$

$$T_a t \geq \left\lfloor \frac{\Delta^2}{cK} \right\rfloor + 1, \quad x_t = \left\lfloor \frac{cK}{\Delta^2} \right\rfloor \cdot \frac{1}{2K} + \frac{1}{2K} \sum_{j=\left\lfloor \frac{cK}{\Delta^2} \right\rfloor + 1}^{t-1} \frac{cK}{\Delta_j^2}$$

$$\geq \left(\frac{c}{2\Delta^2} - \frac{\Delta K}{2\Delta^2} \right) + \frac{c}{2\Delta^2} \ln\left(\frac{t-1}{\left\lfloor \frac{cK}{\Delta^2} \right\rfloor}\right)$$

$$\frac{b}{2} \geq \int_{a-1}^b \frac{1}{s} ds$$

$$x_t \geq \frac{c-1}{2\Delta^2} \ln\left(\frac{e(T)\Delta^2}{cK}\right)$$

Recap: $\mathbb{P}(a_t = k) \leq \frac{\varepsilon_k}{K} + \mathbb{P}(\hat{\mu}_a(t-1) - \mu_k \geq \frac{\Delta_k}{2}) + \mathbb{P}(\mu_a^* - \hat{\mu}_a(t-1) \geq \frac{\Delta_k}{2})$.

With

$$\mathbb{P}(\hat{\mu}_a(t-1) - \mu_k \geq \frac{\Delta_k}{2}) \leq e^{-\frac{x_t}{5}} + 2 \frac{e^{-\frac{\Delta_k}{2}x_t}}{\Delta_k}$$

$$\text{so } \mathbb{P}(a_t = k) \leq \frac{c}{\Delta^2 T} + 2 e^{-x_t/5} + \frac{4}{\Delta_k^2} e^{-\Delta_k^2 x_t/2}$$

With $x_t \geq \frac{c-1}{2\Delta^2} \ln\left(\frac{e(T)\Delta^2}{cK}\right)$

thus yields for large enough c and

$$T \geq \left\lfloor \frac{D^2}{cK} \right\rfloor + 1,$$

$$\Pr(a_T = k) = O\left(\frac{1}{\Delta^2 T}\right) \text{ so that for large enough } c', c:$$

so that:

$$\sum_{t=\left\lfloor \frac{D^2}{cK} \right\rfloor + 1}^T \Pr(a_t = k) \leq \frac{c'}{\Delta^2} \ln(T)$$

hence regret is bounded as:

$$R_T \leq \frac{cK}{\Delta^2} + \sum_{k=1}^K \frac{c' \Delta_k}{\Delta^2} \ln(T)$$

$$\leq \frac{c'}{\Delta^2} \sum_{k=1}^K (\ln(T) + 1) \quad \square$$

does not depend on any parameters ($\mu, K, \Delta, T \dots$)

Remarks • The bound above is called instance dependent as it heavily relies on parameters of the instance Δ_k .

A different choice of ϵ_t can lead to the following distribution-free bound for ϵ -Greedy:

$$R_T \leq O\left((K \ln T)^{1/3} T^{2/3}\right)$$

Two main drawbacks of these methods: (ETC and E-greedy)

- they require knowledge of Δ .
- they scale in $\frac{1}{\Delta^2}$ ($nT^{2/3}$ 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.

Moreover, ETC needs knowledge of T .

→ if T is unknown, we can use the doubling trick (for any algo)