

UCB over arbitrary time
horizons.

Source: Chap 8, Bandit Algorithms, L.K.S, 2019.

So far: UCB(s) with $s = \frac{1}{n^2}$, regret
bounds.

→ We assumed knowledge of time horizon n , and set $s = \frac{1}{n^2}$.

→ $\{\mathbb{D}_i\}$ are unknown.

UCB with knowledge of n :

Prev: $(1/s) = n^2$.

Instead: $f(t) = 1 + t \ln^2(t)$: Time varying
parameter.

i.e.,

$$A_t = \arg \max_{1 \leq j \leq k} \hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln f(t)}{T_j(t-1)}}$$

$$f(t) = 1 + t \ln^2(t).$$

Note: As t increases $f(t)$ grows slightly faster than linear in time.

Setting: As before with k arms, unstructured environment, 1-subGaussian noise in rewards.

Thm (8.1 in text):

$$R_n \leq \sum_{\{j: \Delta_j > 0\}} R_{n,j}$$

$$R_{n,j} = \inf_{\varepsilon_j \in (0, \Delta_j)} \Delta_j \left(1 + \frac{5}{\varepsilon_j^2} + \frac{2 \left(\ln(f(n)) + \sqrt{\pi \ln(f(n))} + 1 \right)}{(\Delta_j - \varepsilon_j)^2} \right)$$

$$\lim_{n \rightarrow \infty} \frac{R_n}{\ln(n)} \leq \sum_{\{j: \Delta_j > 0\}} \left(\frac{2}{\Delta_j} \right)$$

If $\varepsilon_j = \frac{\Delta_j}{2}$,

$$R_n \leq \sum_{\{j: \Delta_j > 0\}} \left(\Delta_j + \frac{8 \left(\ln f(n) + \sqrt{\pi \ln f(n)} \right) + 28}{\Delta_j} \right)$$

$$\text{Recall } f(n) = 1 + n \ln^2(n)$$

$$\therefore R_n \leq C \sum_{\{j : \Delta_j > 0\}} \left(\Delta_j + \frac{\ln(n)}{\Delta_j} \right)$$

fixed constant.

Aside: We first need the following bound.

Counting the expected number of sub-optimal arm plays:

$\{x_i, i=1, 2, \dots, n\}$ iid, with

$$E[x_i] = 0, x_i \sim 1\text{-subGaussian}.$$

$$\hat{m}_t = \frac{1}{t} \sum_{s=1}^t x_s.$$

$\varepsilon > 0, \alpha > 0$ fixed parameters.

$$C_n = \sum_{t=1}^n \chi_{\left\{ \hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \varepsilon \right\}}$$

interpretation of $\chi_j(\cdot)$ with parameter $\ln(f(t)) = a$ and at some time such that t samples has been acquired by that time, and this index function exceeds the true mean by $\varepsilon > 0$.

Also note that for $u = (2a/\varepsilon^2)$, we have:

$$C_n \leq C'_n \equiv u + \sum_{t=\lceil u \rceil}^n \chi_{\left\{ \hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \varepsilon \right\}}$$

Claim (Lemma 8.2): Setting as above. Then,

$$\begin{aligned} E[C_n] &\leq E[C'_n] \\ &\leq 1 + \frac{2}{\varepsilon^2} (a + \sqrt{\pi a} + 1). \end{aligned}$$

Proof:

$$E[C_n] \leq E[C_n']$$

$$\begin{aligned} &= u + \sum_{t=1}^n P\left(\hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \varepsilon\right) \\ \frac{2a}{\varepsilon^2} &\quad \downarrow \begin{cases} = u + \\ 1 - \text{Pr}[\cdot] \end{cases} \\ &\leq u + \sum_{t=1}^n e^{-\frac{t}{2}(\varepsilon - \sqrt{\frac{2a}{t}})^2} \end{aligned}$$

$$\leq (1+u) + \int_u^\infty e^{-\frac{t}{2}(\varepsilon - \sqrt{\frac{2a}{t}})^2} dt$$

Aside:

$$= \int_0^\infty e^{-x^2/2} \cdot 2 \left(\frac{x + \sqrt{2a}}{\varepsilon} \right) dx$$

$$x = \varepsilon \sqrt{t} - \sqrt{2a}$$

Use Gaussian Integral formula:

(Ref: Wikipedia

$$\int_0^\infty x^{m+1} e^{-ax^2} dx = \frac{m!}{2a^{\frac{m+1}{2}}} \quad \text{on Gaussian Integrals.)}$$

$$\leq 1 + \frac{2}{\varepsilon^2} (a + \sqrt{\pi a} + 1)$$



Back to Main Regret Result:

Proof of Thm :

$$R_n = \sum_{\{j : \Delta_j > 0\}} \Delta_j E[\tau_j(n)] \quad (\text{Regret Decomp.})$$

Pick any arm $j \in \{1, \dots, k\}$ with $\Delta_j > 0$.

Suppose at time t , $\{A_t = j\}$. There are four possibilities:

(a) $v_i(t-1) \leq \mu_i - \varepsilon, v_j(t-1) < \mu_j - \varepsilon$

(b) $v_i(t-1) \leq \mu_i - \varepsilon, v_j(t-1) \geq \mu_j - \varepsilon$

(c) $v_i(t-1) > \mu_i - \varepsilon, v_j(t-1) < \mu_j - \varepsilon$

(d) $v_i(t-1) > \mu_i - \varepsilon, v_j(t-1) \geq \mu_j - \varepsilon$

$$\therefore \{A_t = j\} \subseteq \{A_t = j, v_i(t-1) \leq \mu_i - \varepsilon\}$$

covers (d)

covers (a), (b)

V $\{A_t = j, v_j(t-1) \geq \mu_j - \varepsilon\}$

$$\subseteq \{v_i(t-1) \leq \mu_i - \varepsilon\} \cup \{v_j(t-1) \geq \mu_j - \varepsilon, A_t = j\}$$

$$\therefore T_j(n) = \sum_{t=1}^n \chi_{\{A_t = j\}}$$

$$\leq \sum_{t=1}^n \chi_{\{v_i(t-1) \leq \mu_i - \varepsilon\}} + \sum_{t=1}^n \chi_{\{v_j(t-1) \geq \mu_j - \varepsilon, A_t = j\}}$$

arm 1 is atypically small
 arm j is atypically large whenever played.

(I)

(II)

$$(I) : E \left[\sum_{t=1}^n \chi_{\{v_i(t-1) \leq \mu_i - \varepsilon\}} \right]$$

$$= E \left[\sum_{t=1}^n \chi_{\left\{ \hat{\mu}_i(t-1) + \sqrt{\frac{2 \ln f(t)}{T_i(t-1)}} \leq \mu_i - \varepsilon \right\}} \right]$$

of samples $\in \{1, 2, \dots, n\}$
 (say r)

union bound (over r) + linearity of expectation (over t)

$$\leq \sum_{t=1}^n \sum_{r=1}^{\tau} P\left(\hat{\mu}_{ir} + \sqrt{\frac{2 \ln f(t)}{r}} \leq M_i - \varepsilon\right)$$

↓ time ↓ samples

1-subG conc.

$$\leq \sum_{t=1}^n \sum_{r=1}^{\tau} e^{-\left(\sqrt{\frac{2 \ln f(t)}{r}} + \varepsilon\right)^2 / 2}$$

check

these
yourself

$$\leq \sum_{t=1}^n \frac{1}{f(t)} \sum_{r=1}^{\tau} e^{-\tau \varepsilon^2 / 2}$$

$$\leq \frac{5}{\varepsilon^2}$$

$$\text{(II)} : E \left[\sum_{t=1}^n \chi_{\{U_j(t-1) \geq M_i - \varepsilon, A_t = j\}} \right]$$

$$= E \left[\sum_{t=1}^n \chi_{\{\hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln f(t)}{T_j(t-1)}} \geq M_i - \varepsilon, A_t = j\}} \right]$$

$$\leq E \left[\sum_{t=1}^n \chi_{\{\hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln f(t)}{T_j(t-1)}} \geq M_i - \varepsilon, A_t = j\}} \right]$$

↓ over time

Prev. in part (I), we used a union bound to deal with $T_j(t-1) \in \{1, 2, \dots, n\}$. Here, we have more information. We observe that the indicator function $\chi_{\{A_t=j\}}$ equals 1 only when $A_t=j$. But this INCREMENTS $T_j(t) \leftarrow T_j(t-1) + 1$.

Thus, $T_j(t-1)$ can equal $r \in \{1, 2, \dots, n\}$ at most once over $t \in \{1, 2, \dots, n\}$.

$$\leq E \left[\sum_{r=1}^n \chi_{\{\hat{\mu}_{jr} + \sqrt{\frac{2 \ln f(n)}{r}} \geq \mu_j - \varepsilon\}} \right]$$

note that this
is over samples

$$= \sum_{r=1}^n P \left(\hat{\mu}_{jr} - \mu_j + \sqrt{\frac{2 \ln f(n)}{r}} \geq \Delta_j - \varepsilon \right)$$

$$\leq 1 + \frac{2}{(\Delta_j - \varepsilon)^2} \left(\ln f(n) + \sqrt{n \ln f(n)} + 1 \right)$$

Substitute in Regret Decomp to get

the main result. Choose $\varepsilon = \frac{1}{(\ln(n))^{1/4}}$

and let $n \rightarrow \infty$ to get the asymptotic upper bound.

3

Note: Please read Sec 8.2 for discussion on other choices of $f(t)$.