

The UCB Algorithm

Upper Confidence Bound

Source: Chap 7, Bandit Algorithms. (textbook)

So far:

ETC : Number of samples per arm predetermined.
 Δ, n known

ETC + Doubling Trick: Δ known, horizon arbitrary.

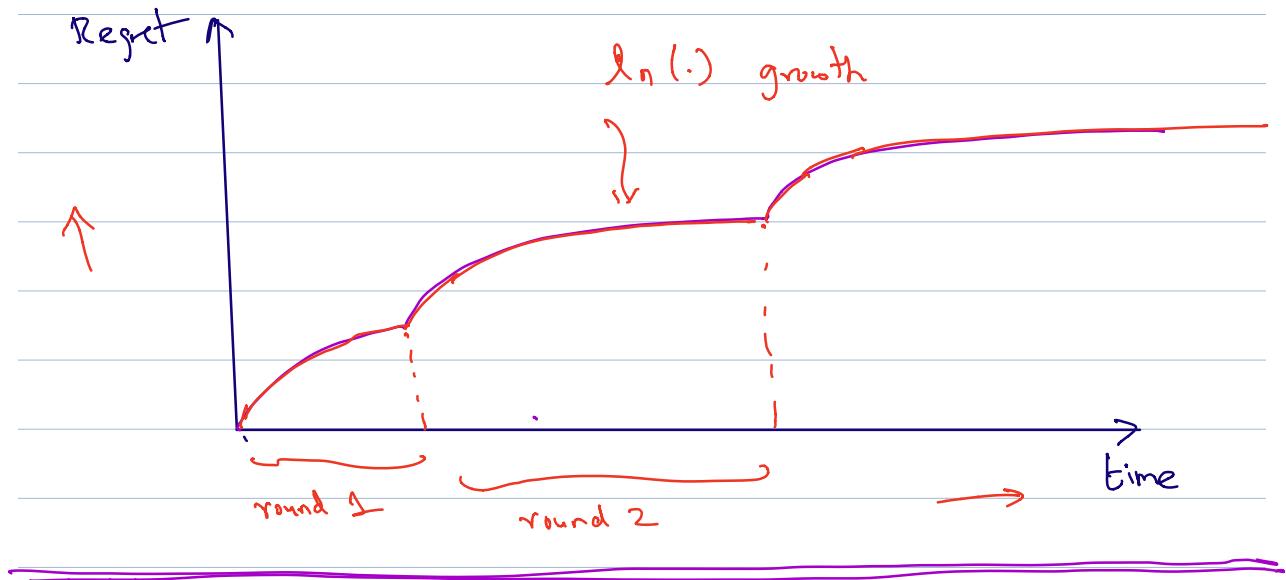
ϵ -Greedy: Samples used for estimation (aka exploration samples) separated from time slots used for regret optimization (aka exploitation samples)
 Δ known, n is arbitrary.

Elimination algo: In each phase, number of samples per remaining arm is predetermined
 Δ unknown, n is fixed

Elimination + Doubling trick: Δ unknown, n arbitrary.

All algos above: Samples for exploration separate from samples from exploitation.

The "problem" with doubling trick is that regret is not "smooth"



VCB Algorithm: Δ unknown,
n arbitrary, regret growth is smooth

Key difference from prev algorithms: Number of samples for exploration for each arm is dynamic, and no separation between exploration and exploitation, i.e., all samples are used to improve estimates.

Technical Challenge: Empirical estimates of arms

are constructed using random number of samples,
so usual sub-Gaussian concentration do not hold.

* Need anytime concentrations:

Let $\{X_i, i=1, 2, \dots\}$ sequence of iid
 Δ -subGaussian r.v.s., and τ a r.v.

Then $P \left(\left| \sum_{i=1}^{\tau} X_i \right| > \sqrt{2\tau(2\log \tau + \log(1/\delta))} \right)$
 $\leq \frac{\pi^2}{3} \delta$.

Compare to usual Δ -subGaussian bound:

deterministic
 $P \left(\sum_{i=1}^n X_i > \sqrt{2n \log(1/\delta)} \right) \leq \delta$

UCB relies on Principle of Optimism: Use the

"most reasonable" optimistic estimate of each arm,
and optimize your decision at each time-step
using this optimistic estimate.

VCB algorithm studied in two steps.

First: $\{\Delta_i\}$ unknown, $n = \text{time-horizon known}$

we will develop the basic algo with
this assumption, and then modify the algo (next class)
to eliminate this requirement.

Setting: Unstructured environment, k arms,

$P_{A_i} \sim (\mu_i + 1\text{-subGaussian})$ reward.

$$V_j(t-1, s) = \hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{T_j(t-1)}} \quad \begin{matrix} \text{empirical mean} \\ \text{exploration bonus} \end{matrix} \quad \begin{matrix} \dots \rightarrow (\text{known for now}) \end{matrix}$$

$$\text{where } \hat{\mu}_j(t-1) = \frac{1}{T_j(t-1)} \sum_{s=1}^{t-1} X_s \chi_{\{A_s=j\}}$$

Also, let $V_j(t-1, s) = \infty$ if $T_j(t-1) = 0$.

(this forces the player to play each arm at least once initially).

UCB(δ) algorithm:

At each time t : $A_t = \underset{1 \leq j \leq k}{\operatorname{argmax}} V_j(t-1, \delta)$,

Intuition

* Consider the best arm, (arm 1 with $\mu_1 > \mu_2 \geq \mu_3 \dots$)

for t large, hopefully the player has played this arm a large number of times, c.e;

$|T_1(t-1) - (t-1)|$ is small.

$$\therefore V_1(t-1) = \hat{\mu}_1(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{T_1(t-1)}}$$

number of samples $T_1(t-1)$ large enough, so $\hat{\mu}_1(t-1) \approx \boxed{\mu_1}$

vanishingly small

$$\therefore V_1(t-1) \approx \mu_1$$

* Consider any other arm $j \geq 2$.

$T_j(t-1)$ is small, but notice that

$$\sqrt{\frac{2 \ln(1/\delta)}{T_j(t-1)}} \text{ is a } \underline{\text{confidence bound}}$$

for the estimator $\hat{\mu}_j(t-1)$.

Specifically, if $\underbrace{T_j(t-1)}_{\text{estimated with } n_j \text{ samples}} = n_j$ (say)

Then, $(\hat{\mu}_j^{(n_j)} - \mu_j) \sim \frac{1}{\sqrt{n_j}} - \text{sub Gaussian.}$

$$P\left(\hat{\mu}_j^{(n_j)} > \mu_j + \sqrt{\frac{2 \ln(1/\delta)}{n_j}}\right) \leq \delta$$

This is true for any $n_j \geq 1$!

Choose $\delta = 1/n^2$, and we have that

$$V_j(t-1) = \hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{T_j(t-1)}} \text{ is}$$

a high probability upper bound on μ_j .

Thus, these choices of $V_i(t-1)$, $V_j(t-1)$ gives the player the confidence that the TRUE MEAN μ_j is smaller than the TRUE MEAN μ_i of arm 1 by only comparing $V_i(t-1)$ and $V_j(t-1)$.

In symbols:

Suppose you observe this
↓

$$\mu_i \approx V_i(t-1) > V_j(t-1)$$

$$= \hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{T_j(t-1)}}$$

$$\stackrel{\text{w.p. } (1-\delta)^2}{\geq} \mu_j$$

∴ The algorithm choice $A_t = \operatorname{argmax}_j V_j(t-1)$
"ensures" that with high prob:

$$\text{IF } V_i(t-1) > V_j(t-1)$$

we are indeed

$$\Rightarrow \mu_i > \mu_j \Rightarrow \text{playing best arm}$$

The problem in making this rigorous: $T_j(t-1)$

is random, and depends on past samples.

However, the anytime concentration result suggests that the above intuition is approximately correct.

Formal Analysis of UCB(δ)

(Thm 7.1 in text).

Thm: With setting as above, and $\delta = \frac{1}{n^2}$,

$$R_n(\pi, v) \leq 3 \sum_{i=1}^k D_i + \sum_{\{i : \Delta_i > 0\}} \left(\frac{16 \ln(n)}{\Delta_i} \right)$$

Proof:

Construct the sequences of r.v.s:

→ sample

$$\begin{array}{ccccccc} X_{11} & X_{21} & X_{31} & \cdots & X_{r1} & \cdots & X_{n1} \\ X_{12} & X_{22} & X_{32} & \cdots & X_{r2} & \cdots & X_{n2} \\ \vdots & \vdots & \vdots & & \vdots & & \vdots \\ X_{1k} & X_{2k} & \cdots & X_{rk} & \cdots & X_{nk} \end{array}$$

where $X_{rj} \sim P_{aj} = M_j + 1\text{-subGaussian}$.

There are indep. r.v.s. that will be sampled as needed. For instance, the third time we play arm 5, the r.v. X_{35} is used. In general,

$$X_t = X_{T_{A_t}^{(t)} A_t}$$

reward at t time t The cumulative number of times that arm A_t has been played at time t . the arm index played at time t

$$R_n = \sum_{j=1}^k \Delta_j E[T_j(n)] \quad \begin{matrix} \text{(Regret} \\ \text{Decomp. Lemma)} \end{matrix}$$

Fix any arm $j \neq 1$. Define

$$G_j = \left\{ M_1 < \min_{1 \leq s \leq n} U_1(s) \right\} \cap \left\{ \hat{M}_{j|u_j} + \sqrt{\frac{2}{u_j} \ln(\frac{n}{s})} < M_1 \right\}$$

Note: $\hat{M}_{j|u_j}$ = empirical estimate for arm j with u_j samples

Side note: I am closely following notation from text, despite minor inconsistencies. E.g:

$$x_{i,j} \quad \begin{matrix} \text{Sample index} \\ \text{arm index} \end{matrix}$$
$$\hat{\mu}_{j,u} \quad \begin{matrix} \text{# of samples index} \\ \text{i.e., estimate} \\ \text{constructed from } u \text{ samples} \end{matrix}$$

Also recall : $\hat{\mu}_{j,s}$ (i.e., estimate at time s)

G_{ij} is the good event:

* The UCB index for arm 1 always exceeds the true mean for arm 1

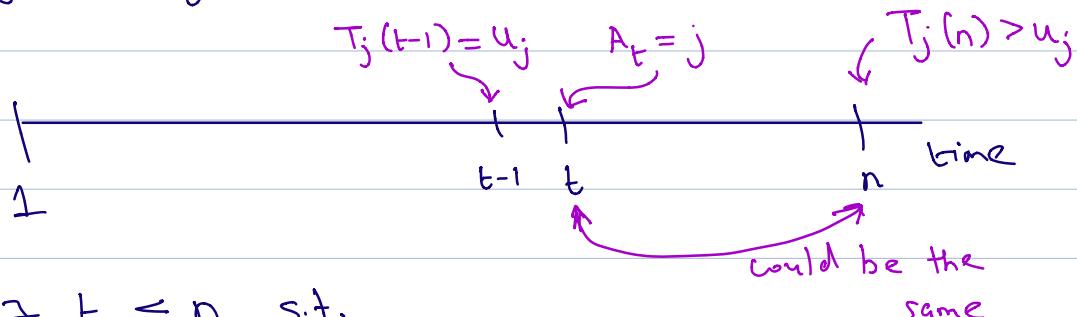
AND

* If the UCB index for arm j was constructed using u_j samples, then this index would be smaller than the true mean for arm 1

Two claims for G_j :

a) If G_j occurs, then $T_j(n) \leq u_j$

Suppose not, i.e., suppose that G_j holds but $T_j(n) > u_j$.



Then $\exists t \leq n$ s.t.

$$A_{\text{not } t} = j \quad \text{and} \quad T_j(t-1) = u_j.$$

$$V_j(t-1) = \hat{\mu}_j(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{T_j(t-1)}}$$

$$= \hat{\mu}_{ju_j} + \sqrt{\frac{2 \ln(1/\delta)}{u_j}}$$

$$< \mu,$$

(G_j defn)

$$< V_1(t-1)$$

Since $U_i(t-1) > U_j(t-1)$, $A_t \neq j$.

This is a contradiction.

$$(b) P(G_j^c) \leq n\delta + \epsilon \stackrel{-u_j \Delta_j^2 / 8}{\leq} (n+1)\delta$$

$$\text{for } u_j \geq \frac{8 \ln(1/\delta)}{\Delta_j^2}$$

Recall:

$$G_j^c = \left\{ M_1 \geq \min_{1 \leq t \leq n} U_1(t) \right\} \cup$$

$$\left\{ \hat{M}_j u_j + \sqrt{\frac{2 \ln(1/\delta)}{u_j}} \geq M_1 \right\}$$

$$P(M_1 \geq \min_{1 \leq t \leq n} U_1(t))$$

$$= P(M_1 \geq \min_{1 \leq t \leq n} \hat{M}_1(t-1) + \sqrt{\frac{2 \ln(1/\delta)}{T_1(t-1)}})$$

time
 ↓
 $\hat{M}_1(t-1)$
 constructed using $T_1(t-1)$ samples
 $\subseteq \{1, 2, \dots, N\}$

t: time
 r: samples
 note: switching from time to number of samples.

$$\leq P \left(\bigcup_{r=1}^n \left\{ M_1 \geq \hat{M}_{1,r} + \sqrt{\frac{2 \ln(1/\delta)}{r}} \right\} \right)$$

Union bound
 # of samples

$$\leq \sum_{r=1}^n P\left(\mu_i \geq \hat{\mu}_{ir} + \sqrt{\frac{2 \ln(1/\delta)}{r}}\right)$$

\downarrow

1-subG concentration $\leq \sum_{r=1}^n s$

$$= n s.$$

i.e., $P\left(\mu_i \geq \min_{1 \leq t \leq n} U_i(t)\right) \leq n s$

Next consider the other term in G_{ij}^c :

$$\left\{ \hat{\mu}_{ju_j} + \sqrt{\frac{2 \ln(1/\delta)}{u_j}} \geq \mu_i \right\}$$

Since $u_j \geq \frac{8 \ln(1/\delta)}{\Delta_j^2}$, we have

$$\sqrt{\frac{2 \ln(1/\delta)}{u_j}} \leq \frac{\Delta_j}{2}$$

$$\therefore P\left(\hat{\mu}_{ju_j} + \sqrt{\frac{2 \ln(1/\delta)}{u_j}} \geq \mu_1\right)$$

$\underbrace{\quad}_{\leq \frac{\Delta_j}{2}}$ = $\mu_j + \Delta_j$

$$\leq P\left(\hat{\mu}_{ju_j} - \mu_j \geq \frac{\Delta_j}{2}\right)$$

$$\leq e^{-u_j \Delta_j^2 / 8}$$

(1 - SubG concentration)

Thus, for any $u_j \geq \frac{8 \ln(1/\delta)}{\Delta_j^2}$, we have

$$P(G_j^c) \leq n\delta + e^{-u_j \Delta_j^2 / 8} \leq (n+1)\delta$$

Summary: G_j s.t. $T_j(n) \leq u_j$

and $P(G_j^c) \leq n\delta + e^{-u_j \Delta_j^2 / 8}$

for $u_j \geq \frac{8 \ln(1/\delta)}{\Delta_j^2}$ $\leq (n+1)\delta$

$$\therefore E[\tau_j(n)] = E\left[\underbrace{\tau_j(n)}_{\leq u_j} \chi_{G_j}\right]$$

$$+ E\left[\underbrace{\tau_j(n)}_{\leq n} \chi_{G_j^c}\right]$$

$$\leq u_j + n P(G_j^c)$$

$$\text{Choose } u_j = \sqrt{\frac{8 \ln(1/\delta)}{\Delta_j^2}}$$

\Rightarrow

$$E[\tau_j(n)] \leq \sqrt{\frac{8 \ln(1/\delta)}{\Delta_j^2}} + n(n+1)\delta$$

Choose $\delta = 1/n^2$:

$$E[\tau_j(n)] \leq \frac{16 \ln(n)}{\Delta_j^2} + 3$$

\downarrow
~~✓~~

$$\therefore \Delta_j E[\tau_j(n)] \leq S \Delta_j + \frac{16 \ln(n)}{\Delta_j}$$

$$\Rightarrow R_n \leq 3 \sum_{j=1}^k \Delta_j + \sum_{j: \Delta_j > 0} \frac{16 \ln(n)}{\Delta_j}$$



So far: Instance-dependent bound on regret, i.e., the bound depends on $\{\Delta_j, j=1, 2, \dots, k\}$.

* Bound decays as $\ln(n)$ over the time horizon.

Alternative: Instance-independent bound that does not depend on $\{\Delta_j, j=1, 2, \dots, k\}$ but grows as \sqrt{n} over time.

Thm (7.2 in text): Suppose $\delta = 1/n^2$. Then, for any $v \in \mathcal{E}_{SG}^k(1)$,

$$R_n \leq 3 \sum_{j=1}^k \Delta_j + 8 \sqrt{nk \ln(n)}$$

Proof:

$$R_n = \sum_{j=1}^k \Delta_j E[\tau_j(n)] \leq 3 + \frac{16 \ln(n)}{\Delta_j^2}$$
$$= \sum_{j=1}^k \Delta_j E[\tau_j(n)] \quad \text{from } \oplus$$

$\left\{ j : \Delta_j \geq \sqrt{\frac{16 k \ln(n)}{n}} \right\}$ two pages back.

$$+ \sum_{\left\{ j : \Delta_j < \sqrt{\frac{16 k \ln(n)}{n}} \right\}} \Delta_j E[\tau_j(n)] \leq \sqrt{\frac{16 k \ln(n)}{n}}$$

$$\leq \sum_{\left\{ j : \Delta_j \geq \sqrt{\frac{16 k \ln(n)}{n}} \right\}} \left(3\Delta_j + \frac{16 \ln(n)}{\Delta_j} \right)$$

$$+ \sqrt{\frac{16 k \ln(n)}{n}} \sum_{\left\{ j : \Delta_j < \sqrt{\frac{16 k \ln(n)}{n}} \right\}} E[\tau_j(n)] \leq n$$

$$\leq 3 \sum_{j=1}^k \Delta_j + \sum_{\{j : \Delta_j \geq \sqrt{\frac{16kn \ln(n)}{n}}\}} \left(\frac{16 \ln(n)}{\Delta_j} \right)$$

$\approx \sqrt{\frac{16kn \ln(n)}{n}}$

$| \cdot | \leq k$

$$\leq 3 \sum_{j=1}^k \Delta_j + k \cdot \sqrt{n \frac{16 \ln(n)}{k}}$$

$$+ \sqrt{16 kn \ln(n)}$$

$$= 3 \sum_{j=1}^k \Delta_j + 8 \sqrt{nk \ln(n)}$$

