

Lecture 22: Kullback Leibler-UCB (KL-UCB)

Lecturer: M. K. Hanawal

Scribes: Karan Patel

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

22.1 Recap:

- In previous class mainly we focused on how to compute the expected number of pulls of sub-optimal arm after T rounds.

$$\mathbb{E}[N_{i,(T)}] \leq \frac{4\alpha \log T}{\Delta_i^2} + \frac{\pi^2}{3} + 1$$

- Also we can conclude that $\left(T - (k-1) \left(\frac{4\alpha \log T}{\Delta_i^2} + \left(\frac{\pi^2}{3} + 1\right)\right)\right)$ will be the number of times the optimal arm was pulled after T rounds.
- Also we have seen that pseudo-regret calculation depends on $\mathbb{E}[N_{i,(T)}]$ and it can be calculated as follows:

$$\tilde{R}_T = \sum_{i=1}^K \mathbb{E}[N_i(T)] \Delta_i$$

22.2 UCB

Using previous class results we know that Pseudo regret is given by:

$$\tilde{R}_T \leq \sum_{i \neq i^*} \left(\frac{4\alpha \log T}{\Delta_i^2} + \frac{\pi^2}{3} + 1 \right) \Delta_i$$

$$I_t = \operatorname{argmax}_i \left\{ \hat{\mu}_{i,N_i(t-1)} + \sqrt{\frac{\alpha \log t}{N_i(t-1)}} \right\}$$

$$\text{Let } B_{i,t} = \hat{\mu}_{i,N_i(t-1)} + \sqrt{\frac{\alpha \log t}{N_i(t-1)}}$$

Now by using Hoeffding's inequality, $Pr(\mathbb{E}[x] \geq \epsilon) \leq 2e^{-2n\epsilon^2}$, probability of $B_{i,t}$ can be written as follows:

$$\begin{aligned} Pr \left\{ \hat{\mu}_{i,N_i(t-1)} + \sqrt{\frac{\alpha \log t}{N_i(t-1)}} \geq \mu_i \right\} &\approx \exp \left(-2N_i(t-1) \frac{\alpha \log t}{N_i(t-1)} \right) \\ &= \frac{1}{t^{2\alpha}} \end{aligned}$$

From above simplification of UCB equation we can see that $\log t$ term in UCB ensuring that confidence term $\left(\frac{1}{t^{2\alpha}}\right)$ is not constant but depends on 't'.

22.3 Optimistic Algorithm

22.3.1 Kullback Leibler-UCB (KL-UCB):

Algorithm : KL-UCB

Input: T (Horizons), K (Number of arms)

Initialize: Play each arm once and observe rewards.

for $t = K+1, K+2, \dots, T$ In above equation 'd' represents divergence. Let here

$$B_{i,t} = \max \left\{ q \in [0, 1], d(\hat{\mu}_{i, N_i(t-1)}, q) \leq \frac{\log t + c \log(\log t)}{N_i(t-1)} \right\}$$

Where, 'd' represents the KL-divergence which is given as follows:

$$d(p, q) = p \log \left(\frac{p}{q} \right) + (1 - p) \log \left(\frac{1 - p}{1 - q} \right)$$

We also know that,

$$d(p, q) = \begin{cases} \geq 0, & p \neq q \\ 0, & p = q \end{cases}$$

Further if we fix p and look divergence as function of other variable then for a given p , $d(p, q)$ is strictly convex in q and increasing in the interval $q \in (p, 1)$ as shown in following figure;

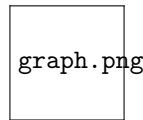


Figure 22.1: Function $d(p, q)$ for a given p

Fixed p in above figure is : $\frac{\log t + \log(\log t)}{N_i(t-1)}$.

Theorem 22.1 Consider a bandit problem with K arms and independent rewards bounded in $[0, 1]$. Let $\epsilon > 0$ and $c=3$. Let i^* be the arm with maximal expected reward μ_{i^*} and i be a sub-optimal arm (i.e. $i \neq i^*$ or $\Delta^* \neq 0$). For any positive integer T , the number of times algorithm KL-UCB chooses arm i is upper-bounded by:

$$\mathbb{E}[N_i(T)] \leq \frac{\log T}{d(\mu_i, \mu_{i^*})} (1 + q) + C_1 \log(\log T) + \frac{C_2(\epsilon)}{T^{\beta(\epsilon)}} \quad \forall \epsilon > 0$$

where $C_1, C_2(\epsilon)$ and $\beta(\epsilon)$ are positive functions of ϵ ($\epsilon > 0$), and T is time horizon.

Divide both side by $\log T$, we get

$$\frac{\mathbb{E}[N_i(T)]}{\log T} \leq \frac{1}{d(\mu_i, \mu_{i^*})} (1 + q) + \frac{C_1 \log(\log T)}{\log T} + \frac{C_2(\epsilon)}{T^{\beta(\epsilon)} (\log T)}$$

Taking $\limsup_{T \rightarrow \infty}$ both side, we get

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_i(T)]}{\log T} \leq \frac{1}{d(\mu_i, \mu_i^*)}$$

NOTE: Pinsker's inequality, $d(\mu_i, \mu_i^*) > 2(\mu_i^* - \mu_i)^2$.

$$\text{For UCB: } \limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_i(T)]}{\log T} \leq \frac{8\alpha}{2\Delta_i^2}$$

And thus KL-UCB has strictly better upper bound (asymptotically) than UCB, while it has the same range of application.

The regret of the KL-UCB algorithm satisfies :

$$\limsup_{T \rightarrow \infty} \frac{\tilde{R}_T}{\log T} \leq \sum_{i:i \neq i^*} \frac{\mu_{i^*} - \mu_i}{d(\mu_i, \mu_i^*)} = \sum_{i:i \neq i^*} \frac{\Delta_i}{d(\mu_i, \mu_i^*)}$$

$$\text{For UCB: } \limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_i(T)]}{\log T} \leq \frac{8\alpha \sum_{i:i \neq i^*} (\mu_{i^*} - \mu_i)}{2\Delta_i^2} \leq 8\alpha \sum_{i:i \neq i^*} \frac{\Delta_i}{\Delta_i^2}$$

22.3.2 Lower Bound:

Irrespective of algorithm we use, the following holds:

$$\liminf_{T \rightarrow \infty} \frac{\tilde{R}_T}{\log T} \leq \sum_{i \neq i^*} \frac{\Delta_i}{d(\mu_i, \mu_i^*)}$$

Here, $\frac{\tilde{R}_T}{\log T} := \text{Regret is normalized by } \log T$.

i.e. The KL-UCB algorithm thus appears to be (asymptotically) optimal.

Define: KL-UCB relies on the following upper-confidence bound for μ_i :

$$B_{i,t} := u_i(t) = \max \left\{ q > \hat{\mu}_{i, N_i(t-1)} : d(\hat{\mu}_i(t), q) \leq \frac{\log t + 3 \log(\log t)}{N_i(t-1)} \right\}$$

For $x, y \in [0, 1]$, define $d^+(x, y) = d(x, y)$ if $x < y$. Without loss of generality, let $i^* = 1$. The expectation of $N_i(T)$ is upper-bounded by using the following:

$$\begin{aligned} \mathbb{E}[N_i(T)] &= \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{I_t = i\} \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{I_t = i, \mu_1 > u_1(t), \mu_1 \leq u_1(t)\} \right] \end{aligned}$$

(here, μ_1 = true value and $u_1(t)$ = index).

After decomposed above can be written as:

$$\leq \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{\mu_1 > u_1(t)\} \right] + \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{I_t = i, \mu_1 \leq u_1(t)\} \right] \quad (22.1)$$

Lemma 22.1

$$\sum_{t=1}^T \mathbb{1}\{I_t = i, \mu_1 \leq u_1(t)\} \leq \sum_{s=1}^T \mathbb{1}\{sd^+(\hat{\mu}_{i,s}, \mu_1) < \log t + 3 \log(\log t)\}$$

Proof Let we pull some arm i in t_{th} round (i.e. $I_t = i$) such that $\mu_1 < u_1(t)$ which together implies that $u_i(t) \geq u_1(t) > \mu_1$ hence,

$$d^+(\hat{\mu}_i(t), \mu_1) \leq d(\hat{\mu}_i(t), \mu_i(t)) = \frac{\log t + 3 \log(\log t)}{N_i(t)}$$

From above equation we can write

$$\begin{aligned} \sum_{t=1}^T \mathbb{1}\{I_t = i, \mu_1 \leq u_1(t)\} &\leq \sum_{t=1}^T \mathbb{1}\{I_t = i, N_i(t) d^+(\hat{\mu}_i(t), \mu_1) \leq \log t + 3 \log(\log t)\} \\ &= \sum_{t=1}^T \sum_{s=1}^t \mathbb{1}\{N_t(i) = s, I_t = i, sd^+(\hat{\mu}_{i,s}, \mu_1) \leq \log t + 3 \log(\log t)\} \\ &\leq \sum_{t=1}^T \sum_{s=1}^t \mathbb{1}\{N_t(i) = s, I_t = i, \} \mathbb{1}\{sd^+(\hat{\mu}_{i,s}, \mu_1) \leq \log t + 3 \log(\log t)\} \\ &= \sum_{s=1}^T \mathbb{1}\{sd^+(\hat{\mu}_{i,s}, \mu_1) \leq \log t + 3 \log(\log t)\} \sum_{t=s}^T \mathbb{1}\{N_t(i) = s, I_t = i\} \end{aligned}$$

as, for every $s \in \{1, 2, \dots, T\}$, $\sum_{t=s}^T \mathbb{1}\{N_t(i) = s, I_t = i\} = 1$

$$= \sum_{s=1}^T \mathbb{1}\{sd^+(\hat{\mu}_{i,s}, \mu_1) \leq \log t + 3 \log(\log t)\}$$

By using above lemma, equation (22.1) can be written as follows:

$$\begin{aligned} &\leq \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{\mu_1 > u_1(t)\} \right] + \mathbb{E} \left[\sum_{s=1}^T \mathbb{1}\{sd^+(\hat{\mu}_{i,s}, \mu_1) < \log t + 3 \log(\log t)\} \right] \\ &\leq \sum_{t=1}^T \mathbb{P}(\mu_1 > u_1(t)) + \mathbb{E} \left[\sum_{s=1}^T \mathbb{1}\{sd^+(\hat{\mu}_{i,s}, \mu_1) < \log t + 3 \log(\log t)\} \right] \end{aligned}$$

The first summand is upper-bounded as follows:

$$\mathbb{P}(\mu_1 > u_1(t)) \leq \frac{e[\log(t)^3 + 3 \log(\log(t))]}{t \log(t)^3}$$

For the second summand let

$$K_T = \lfloor \frac{1 + \epsilon}{d^+(\mu_i, \mu_1)} (\log T + 3 \log(\log T)) \rfloor \quad (22.2)$$

then

$$\begin{aligned}
& \sum_{s=1}^T \mathbb{P} (sd^+(\hat{\mu}_{i,s}, \mu_1) < \log T + 3 \log(\log T)) \\
& \leq K_T + \sum_{s=K_T+1}^{\infty} \mathbb{P} (sd^+(\hat{\mu}_{i,s}, \mu_1) < \log t + 3 \log(\log t)) \\
& \leq K_T + \sum_{s=K_T+1}^{\infty} \mathbb{P} (K_T d^+(\hat{\mu}_{i,s}, \mu_1) < \log t + 3 \log(\log t))
\end{aligned}$$

by using equation(22.2) above inequality can be written as:

$$\leq K_T + \sum_{s=K_T+1}^{\infty} \mathbb{P} \left(d^+(\hat{\mu}_{i,s}, \mu_1) < \frac{d(\mu_i, \mu_1)}{1 + \epsilon} \right) \quad (22.3)$$

Lemma 22.2 For each $\epsilon > 0$, there exist $C_2(\epsilon) > 0$ and $\beta(\epsilon) > 0$ such that

$$\sum_{s=K_T+1}^{\infty} \mathbb{P} \left(d^+(\hat{\mu}_{i,s}, \mu_1) < \frac{d(\mu_i, \mu_1)}{1 + \epsilon} \right) \leq \frac{C_2(\epsilon)}{T^{\beta(\epsilon)}}$$

Proof If $d^+(\hat{\mu}_{i,s}, \mu_1) < \frac{d(\mu_i, \mu_1)}{1 + \epsilon}$, then $\hat{\mu}_{i,s} > r(\epsilon)$ where $r(\epsilon) \in (\mu_i, \mu_1)$ such that $d(r(\epsilon), \mu_1) = \frac{d(\mu_i, \mu_1)}{1 + \epsilon}$. Hence,

$$\begin{aligned}
\mathbb{P} \left(d^+(\hat{\mu}_{i,s}, \mu_1) < \frac{d(\mu_i, \mu_1)}{1 + \epsilon} \right) & \leq \mathbb{P} (d(\hat{\mu}_{i,s}, \mu_i) > d(r(\epsilon), \mu_i), \mu_{i,s} > \mu_i) \\
& \leq \mathbb{P} (\mu_{i,s} > r(\epsilon)) \leq \exp(-sd(r(\epsilon), \mu_i))
\end{aligned}$$

and

$$\sum_{s=K_T+1}^{\infty} \mathbb{P} \left(d^+(\hat{\mu}_{i,s}, \mu_1) < \frac{d(\mu_i, \mu_1)}{1 + \epsilon} \right) \leq \frac{\exp(-d(r(\epsilon), \mu_i)K_T)}{1 - \exp(-d(r(\epsilon), \mu_i))} \leq \frac{C_2(\epsilon)}{T^{\beta(\epsilon)}}$$

$$\text{with } C_2(\epsilon) = \frac{1}{1 - \exp(-d(r(\epsilon), \mu_i))} \approx O(\epsilon^{-2})$$

$$\beta(\epsilon) = \frac{(1 + \epsilon)d(r(\epsilon), \mu_i)}{d(\mu_i, \mu_1)} \approx O(\epsilon^2)$$

By using result of above lemma, equation(22.3) can be written as:

$$\leq \frac{1 + \epsilon}{d^+(\mu_i, \mu_1)} (\log T + 3 \log(\log T)) + \frac{C_2(\epsilon)}{T^{\beta(\epsilon)}}$$

References

- [1] Garivier, Aurlen, and Olivier Capp. "The KL-UCB algorithm for bounded stochastic bandits and beyond." In Proceedings of the 24th annual Conference On Learning Theory, pp. 359-376. 2011.