

Week 7 - Lecture 34

UCB Algorithm

combines exploration & exploitation

Index in UCB for arm i at time t is :

$$UCB_i(t) = \hat{\mu}_i(t-1) + \sqrt{\frac{\alpha \log t}{T_i(t-1)}}$$

$$\text{where } T_i(t-1) = \sum_{s=1}^{t-1} 1_{\{I_s = i\}}$$

Thus :

$$\hat{\mu}_i(t-1) = \frac{\sum_{s=1}^{t-1} 1_{\{I_s = i\}} X_{is}}{T_i(t-1)}$$

\uparrow
random quantity

\leftarrow random quantity

Let :

$$\hat{\mu}_{iu}(t-1) = \frac{\sum_{s=1}^{t-1} X_{is}}{u}$$

\uparrow

fixing no. of
samples observed
till $(t-1)$ to be u

In UCB we have finally :

$$I_t = \operatorname{argmax}_i \text{UCB}_i(t-1)$$

$$\text{Let } \mu^* = \max_i \mu_i \quad \therefore \quad \underset{\substack{\uparrow \\ \text{optimal arm}}}{i^*} = \operatorname{argmax}_i \mu_i$$

$\forall i, i \neq i^* \leftarrow$ sub-optimal arms

$$\Delta_i = \mu^* - \mu_i \quad \leftarrow \text{sub-optimality gap of arm } i$$

$$\Delta = \min_{i \neq i^*} \Delta_i \quad \leftarrow \text{sub-optimality gap}$$

Assuming that optimal arm is unique

$$\therefore \Delta > 0$$

We know :

$$R(\pi, n) = \sum_{i=1}^K \mathbb{E}[T_i(n)] \Delta_i$$

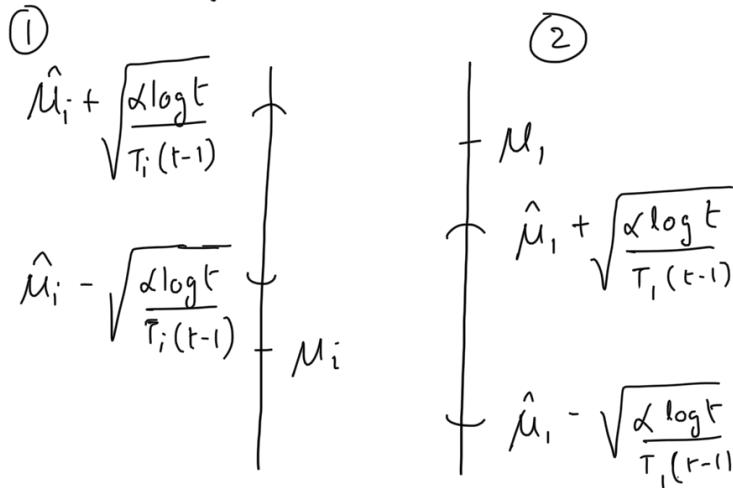
Proof of regret bound for UCB

Assume that sub-optimal arm is played in round t .

$$I_t = i, \quad i \neq i^*$$

[Assuming $i^* = 1$, arm 1 is optimal, without loss of generality.]

How might this have come about:



③ arm is not sampled enough

$$\textcircled{1} \quad \hat{\mu}_i - \sqrt{\frac{\alpha \log t}{T_i(t-1)}} > \mu_i$$

$$\textcircled{2} \quad \hat{\mu}_1 + \sqrt{\frac{\alpha \log t}{T_1(t-1)}} \leq \mu_1$$

$$\textcircled{3} \quad T_i(t-1) \leq \frac{\alpha \log n}{\Delta_i^2/4} \Rightarrow \Delta_i \leq 2 \sqrt{\frac{\alpha \log n}{T_i(t-1)}}$$

Claim: If $I_t = i$ in round t at least one of

①, ② or ③ must hold

To prove this claim:

Suppose none of this holds then $I_t \neq i$

Assume any of ①, ②, ③ holds

Then:

$$\begin{aligned} \hat{\mu}_i + \sqrt{\frac{\alpha \log t}{T_i(t-1)}} &\geq \mu_i \\ \text{// } \text{UCB}_i(t) &= \mu_i + \Delta_i \geq \left(\mu_i + 2 \sqrt{\frac{\alpha \log n}{T_i(t-1)}} \right) \\ &\geq \left(\mu_i + 2 \sqrt{\frac{\alpha \log t}{T_i(t-1)}} \right) \leftarrow \text{playing till round } t, \text{ where } t \leq n \therefore \text{this holds} \\ &> \hat{\mu}_i + \sqrt{\frac{\alpha \log t}{T_i(t-1)}} = \text{UCB}_i(t) \end{aligned}$$

Actually should be
↓

$\Rightarrow I_t \neq i$ a contradiction!

All the conditions have been shown to have been violated, and if this happens then $\text{UCB}_i(t)$ dominates $\text{UCB}_j(t)$, thus the optimal arm should have been played.

Thus if sub-optimal arm is played then one of the three conditions should hold.

$$\begin{aligned}
\Rightarrow \mathbb{E}[T_i(n)] &\leq u + \sum_{t=1}^n \Pr\{I_t = i, \textcircled{1} \text{ holds}\} \\
&\quad + \sum_{t=1}^n \Pr\{I_t = i, \textcircled{2} \text{ holds}\} \\
&\leq u + \sum_{t=1}^n \Pr\{\textcircled{1} \text{ holds}\} + \sum_{t=1}^n \Pr\{\textcircled{2} \text{ holds}\} \\
&\hspace{15em} (\text{upper bounding})
\end{aligned}$$

$$\Pr\{\textcircled{1} \text{ holds}\} = \Pr\left\{\hat{\mu}_i - \sqrt{\frac{\alpha \log t}{T_i(t-1)}} > \mu_i\right\} \text{ --- } \textcircled{i}$$

$$\Pr\{\textcircled{2} \text{ holds}\} = \Pr\left\{\hat{\mu}_i + \sqrt{\frac{\alpha \log t}{T_i(t-1)}} < \mu_i\right\} \text{ --- } \textcircled{ii}$$

We know how to bound \textcircled{i} & \textcircled{ii} ,

will continue with that in the next lecture...