

Week 8 - Lecture 36

Divergence

Let p & q be any distribution, then divergence between them is defined as:

$$d(p, q) = \sum_{i=1}^N p_i \log \frac{p_i}{q_i} \geq 0$$

Proof: of $d(p, q) \geq 0$

$$d(p, q) = - \sum_{i=1}^N p_i \log \frac{q_i}{p_i}$$

Apply Jensen's inequality

KL-UCB Algorithm

Input : No. of arms = K , constant c
Horizon = T

Initialize : Play each arm once

For $t = K+1, K+2, \dots, T$

$$I_t \leftarrow \arg \max_{a \in [K]} \left[\max \left\{ q \in \Theta : N_a^{(t-1)} d \left(\frac{S_a^{(t-1)}}{N_a^{(t-1)}}, q \right) \leq \log t + c \log(\log t) \right\} \right]$$

↖ index for each arm

$r \leftarrow$ observed reward from I_t

$$N_{I_t}(t) \leftarrow N_{I_t}(t-1) + 1 \quad : \text{No. of pulls of arms}$$

$$S_{I_t}(t) \leftarrow S_{I_t}(t-1) + r \quad : \text{Cumulative reward}$$

$\frac{S_a^{(t-1)}}{N_a^{(t-1)}}$ & q are Bernoulli distributions

Divergence for Bernoulli:

$$p, (1-p)$$

$$q, (1-q)$$

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

If $q \in [p, 1]$, \therefore all q 's that are larger than p

then $d(p, q) = 0$ for $q = p$

and increasing for all $q \in [p, 1]$

$$\therefore \text{ If } q \in \left[\frac{S_a(t-1)}{N_a(t-1)}, 1 \right]$$

Then the divergence of the index will increase till it reaches a point where the inequality is violated.

Thus at the point the inequality is violated, that would be $\max q$, & that q will be the upper confidence bound.

Theorem: Consider a K -armed stochastic bandit with support in $[0, 1]$.

Let $\epsilon > 0$ & set $C = 3$.

For any T , the number of times KL-UCB chooses a sub-optimal arm $K (\neq K^*)$ is bounded as

$$\mathbb{E}[N_K(T)] \leq \underbrace{\frac{\log T (1+\epsilon)}{d(\mu_K, \mu_{K^*})}}_{\substack{\uparrow \text{largest term} \\ \text{compare with} \\ \text{UCB } \log T \text{ term}}} + \underbrace{C_1 \log(\log T)}_{\substack{\leftarrow \log(\log T) \\ \text{will be} \\ \text{very small} \\ \text{thus can} \\ \text{ignore}}} + \underbrace{C_2(\epsilon)}_{\substack{\frac{1}{T^{\beta(\epsilon)}} \leftarrow T \text{ is already in} \\ \text{denominator } \therefore \text{ignore}}}$$

where C_1 is constant

C_2 and β are positive functions of ϵ

Lemma: Pinsker Inequality

[Proof is found in introduction]

$$\forall p, q \in [0, 1]$$

... for me ...
Theoretic analysis
or in the book]

$$d(p, q) \geq 2(p - q)^2$$

$$\begin{aligned} \therefore d(\mu_k, \mu_{k^*}) &\geq 2(\mu_k - \mu_{k^*})^2 \\ &= 2\Delta_k^2 \end{aligned}$$

$$\mathbb{E}[T_k(n)] \leq \frac{6 \log n}{\Delta_k^2} + \pi^2/3 + 1$$

for UCB

$$\mathbb{E}[N_k(T)] \leq \frac{\log T}{d(\mu_k, \mu_{k^*})} (1 + \varepsilon) + C_1 \log(\log T) + \frac{C_2(\varepsilon)}{T^{\beta(\varepsilon)}}$$

for KL-UCB

$$\leq \frac{\log T}{2\Delta_k^2} (1 + \varepsilon) + C_1 \log(\log T) + \frac{C_2(\varepsilon)}{T^{\beta(\varepsilon)}}$$

Thus we have a better bound on the expected number of pulls \therefore would have a better regret bound.