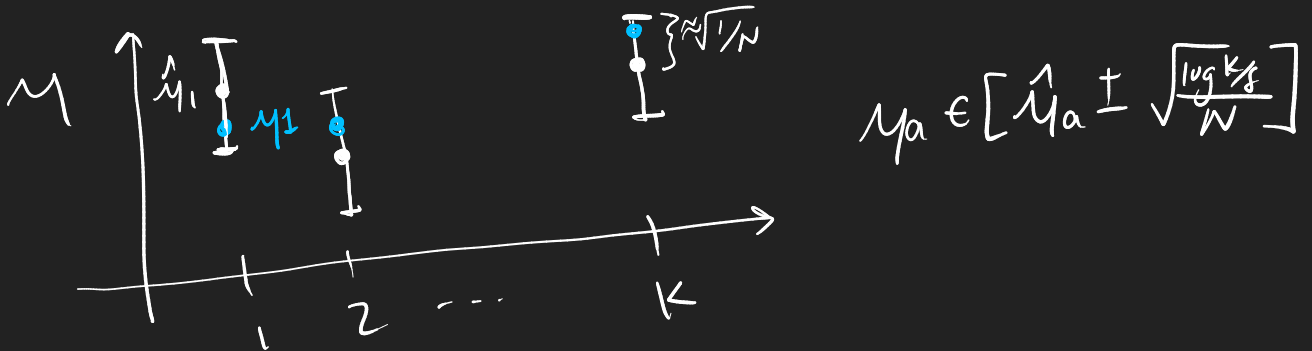


# 1) Explore-then-Commit

$$R(T) = \sum_{t=1}^T \mu^* - \mu_{a_t} = \underbrace{\sum_{t=1}^{NK} (\mu^* - \mu_{a_t})}_{R_1 \leq NK} + \underbrace{\sum_{t=1+NK}^T (\mu^* - \mu_{a_t})}_{R_2}$$

assume  $r_t \in [0, 1]$

Consider the difference  $\hat{\mu}_a - \mu_a$



$$\mu_a \in \left[ \hat{\mu}_a \pm \sqrt{\frac{\log(k/s)}{N}} \right]$$

Lemma: After exploration

$$|\hat{\mu}_a - \mu_a| \lesssim \sqrt{\frac{\log(k/s)}{N}} \quad \text{w.p. } 1-\delta$$

Proof: Hoeffding & union bound  $P(A \cap B) \leq P(A) + P(B)$

Hoeffding's Bound: if  $r_i \in [0, 1]$  and  $\mathbb{E}[r_i] = \mu$   
then  $r_1, \dots, r_N$  iid

$$\left| \frac{1}{N} \sum_{i=1}^N r_i - \mu \right| \lesssim \sqrt{\frac{\log(k/s)}{N}}$$

$$R_2 = \sum_{t=NK+1}^T \mu^* - \mu_{\hat{a}^*} = (T-NK) (\mu^* - \mu_{\hat{a}^*})$$

$$\leq (T-NK) \left( \mu_{\hat{a}^*} + \sqrt{\frac{\log(k/s)}{N}} \right)$$

$$\text{w.p. } 1-\delta \quad \rightarrow \quad - \left( \hat{\mu}_{\hat{a}^*} - \sqrt{\frac{\log(k/s)}{N}} \right)$$

$$\leq (T-NK) \left[ \underbrace{\hat{\mu}_{\hat{a}^*} - \mu_{\hat{a}^*}}_{\leq 0} + 2\sqrt{\frac{\log(k/s)}{N}} \right]$$

$$R(T) = R_1 + R_2 \leq \underline{NK} + \underline{2T \sqrt{\frac{\log(K/s)}{N}}} \quad \text{w.p. } 1-\delta$$

minimize w.r.t.  $N$

$$N = \left( \frac{T}{2K} \sqrt{\log(K/s)} \right)^{2/3}$$

$$R(T) \lesssim T^{2/3} K^{1/3} (\log \frac{K}{s})^{1/3}$$

$$\frac{R(T)}{T} = T^{-1/3} \rightarrow 0 \quad \text{as } T \rightarrow \infty \quad \checkmark$$

## 2) UCB Algorithm

$$\frac{\delta}{KT} \cdot K \cdot T = \delta$$

Alg 4: UCB

Initialize  $\hat{U}_0^a = \infty$

for  $t=1, \dots, T$

$a_t = \operatorname{argmax}_a \hat{U}_t^a$

update  $\hat{\mu}_{t+1}^{a_t}$  and  $N_{t+1}^{a_t}$

$$\hat{\mu}_t^a + \sqrt{\frac{\log(KT/s)}{N_t^a}} \quad \text{w.p. } 1-\delta$$

↓ Exploration bonus

# times arm  $a$  pulled:  $N_t^a = \sum_{k=1}^t \mathbb{1}\{a_k = a\}$

Empirical mean:  $\hat{\mu}_t^a = \sum_{k=1}^t r_k \mathbb{1}\{a_k = a\} / N_t^a$

## UCB Analysis:

Intuition: case 1)  $a_t$  has a large CI (high uncertainty) → exploring

case 2)  $a_t$  has a small CI (good arm) → exploiting

Regret-at- $t$ :  $\mu^* - \mu_{a_t} \leq \hat{U}_t^{a^*} - \mu_{a_t}$   $\mu^*$  in CI

$$\begin{aligned} &\leq \hat{u}_t^{a_t} - \mu_{a_t} \quad a_t \text{ argmax} \\ &= \hat{\mu}_t^{a_t} + \sqrt{\frac{\log(KT/s)}{N_t^{a_t}}} - \mu_{a_t} \\ &\leq 2 \sqrt{\frac{\log(KT/s)}{N_t^{a_t}}} \end{aligned}$$

$$R(T) = \sum_{t=1}^T \mu^* - \mu_{a_t} \leq 2 \sqrt{\log(KT/s)} \underbrace{\sum_{t=1}^T \sqrt{1/N_t^{a_t}}}$$

Claim:  $\sum_{t=1}^T \sqrt{1/N_t^{a_t}} \leq \sqrt{KT}$

$$R(T) \leq 2 \sqrt{KT \log(KT/s)} \quad O(T^{1/2})$$

sublinear!  $\epsilon$ - $\bar{F}$ - $C$   $O(T^{2/3})$

Proof of claim:

$$\begin{aligned} \sum_{t=1}^T \sqrt{1/N_t^{a_t}} &= \sum_{t=1}^T \sum_{a=1}^k \mathbb{1}\{a_t = a\} \sqrt{\frac{1}{N_t^a}} \\ &= \sum_{a=1}^k \left( \sum_{t=1}^T \mathbb{1}\{a_t = a\} \sqrt{\frac{1}{N_t^a}} \right) \\ &= \sum_{a=1}^k \underbrace{\sum_{i=1}^{N_T^a} \sqrt{1/i}} \\ &\leq \sum_{a=1}^k \sqrt{N_T^a} \leq \sqrt{KT} \end{aligned}$$

Notice:  $\sum_{a=1}^k N_T^a = T$  one arm per timestep

$$\left(\frac{1}{K}\right) \sum_{a=1}^K \sqrt{N_T^a} \stackrel{\text{Jensen's}}{\leq} \sqrt{\frac{1}{K} \sum_{a=1}^K N_T^a} = \sqrt{T/K}$$

Principle: "Optimism in the Face of Uncertainty"