

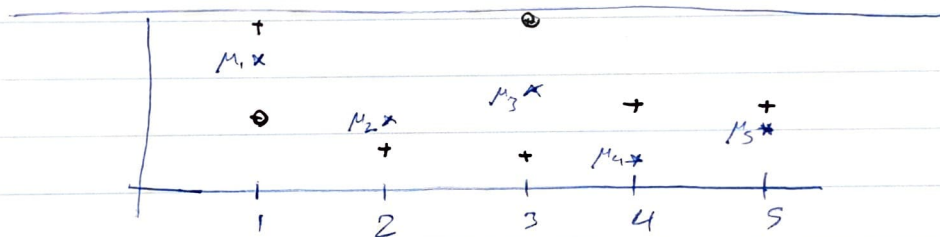
## CS 59000 RL

## Exploration and exploitation.

From last lecture we prove that:

$$R_n(\pi, \nu) = R_n = \sum_i^K \underbrace{\Delta_i}_{\text{suboptimality gap}} \underbrace{E[T_i(n)]}_{\text{the number of time arm } i \text{ is pulled.}}$$

$$\Delta_i = \mu_1 - \mu_i$$



**Explore - Then - Commit.**

Choose each arm/action sequentially  $m$  times.

Then follow the best one estimated.

Two phases

- Pure exploration phase:  $1 \leq t \leq mK$

$$A_t = t \bmod K + 1$$

- Pure exploitation phase: for each arm

$$\text{estimate } \mu_i \text{ in } mK \text{ : } \hat{\mu}_i = \frac{1}{m} \sum_{t=1}^{mK} X_t I\{A_t = i\}$$

$$\text{when } mK \gg n : \hat{\mu}_i = \frac{1}{m \lfloor \frac{n}{K} \rfloor} \sum_{t=1}^n X_t I\{A_t = i\}$$

2)

Then find  $a \in \arg \max \hat{\mu}_i$   
and follow  $a$  for the remaining part  
i.e.  $A_t = a \quad t > mK$

Theorem: For  $\mathcal{V}$ , with  $R$ -sub Gaussian arms,  
the regret of ETC (Explore then Commit) satisfies:

$$R_n \leq m \wedge \left\lceil \frac{n}{K} \right\rceil \sum_{i=1}^K \Delta_i + (n - mK)^+ \sum \Delta_i \exp\left(-\frac{m \Delta_i^2}{4R}\right)$$

$$\text{when } mK < n \rightarrow R_n \leq m \sum_{i=1}^K \Delta_i + (n - mK) \sum_{i=1}^K \Delta_i \exp\left(-\frac{m \Delta_i^2}{4R}\right)$$

$$a \wedge b : \min\{a, b\}, \quad a \vee b : \max\{a, b\}$$

Proof:  $R_n = \sum_{i=1}^K \Delta_i E[T_i(n)]$

$$E[T_i(n)] \leq m \wedge \left\lceil \frac{n}{K} \right\rceil + (n - mK)^+ \mathbb{P}(A_{mK+1} = i)$$

$$\leq m \wedge \left\lceil \frac{n}{K} \right\rceil + (n - mK)^+ \mathbb{P}(\hat{\mu}_i(mK) \geq \max_{j \neq i} \hat{\mu}_j(mK))$$

$$\leq m \wedge \left\lceil \frac{n}{K} \right\rceil + (n - mK)^+ \mathbb{P}(\hat{\mu}_i(mK) \geq \hat{\mu}_1(mK))$$

$$= m \wedge \left\lceil \frac{n}{K} \right\rceil + (n - mK)^+ \mathbb{P}(\hat{\mu}_i - \mu_i - (\hat{\mu}_1(mK) - \mu_1) \geq \Delta_i)$$

$\mu_1 - \mu_i$

3)

Note the  $\underbrace{\hat{\mu}_i^{(mk)} - \mu_i}_{\frac{R}{m} \text{ - sub Gaussian}} - \underbrace{(\mu_i^{(mk)} - \mu_i)}_{\frac{R}{m} \text{ - sub Gaussian}}$

$\Rightarrow$  the whole thing is  $\frac{2R}{m}$  - sub-Gaussian.

$\hookrightarrow$  application of Hoeffding's inequality:

$$\mathbb{P}(\hat{\mu}_i^{(mk)} - \mu_i - (\mu_i^{(mk)} - \mu_i) \geq \Delta_i) \leq \exp\left(-\frac{m \Delta_i^2}{4R}\right)$$

$$\hookrightarrow R_n \leq m \mathbb{1}\left[\frac{n}{K} \geq 1\right] \sum_{i=1}^K \Delta_i + (n - mk)^+ \sum_{i=1}^K \Delta_i \exp\left(-\frac{m \Delta_i^2}{4R}\right)$$

If the learner explores a lot ( $m \gg 1$ ), the second term vanishes. But, the first term blows up. we just explored too much. If  $m$  is small, we explore a little, the first term is small but we most likely will choose a wrong arm for the exploitation and the second term is over kill.

Trade off Between Exploration and Exploitation

4)

Consider a 2 armed bandit: ( $mK > n$ )

$$R_n \leq m \Delta_2 + n \Delta_2 \exp\left(-\frac{m \Delta_2^2}{4R}\right)$$

$$\hookrightarrow \text{choose } m = \max\left\{1, \frac{4R}{\Delta_2^2} \log\left(\frac{n \Delta_2^2}{4}\right)\right\}$$

$$\min\{n \Delta_2, \dots\}$$

$$\Rightarrow R_n \leq \Delta_2 + \frac{4}{\Delta_2} \left(1 + \max\left\{0, \log\left(\frac{n \Delta_2^2}{4R}\right)\right\}\right)$$

Important note: the bound is logarithmic.

in  $n$ , nice

algorithm.

How come? we were able to do so, because we are assuming we know  $\Delta_2$ , what a bummer.



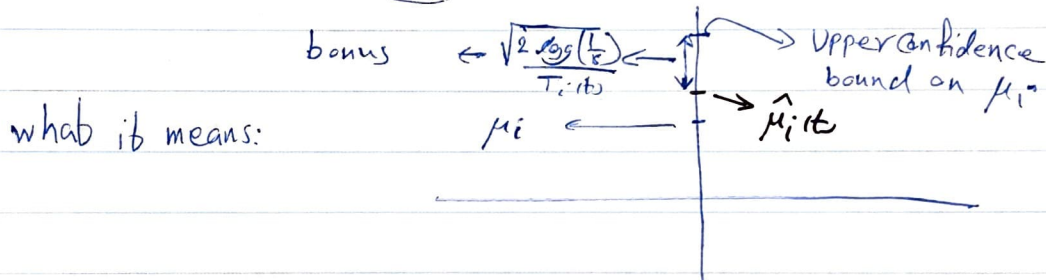
5)

# The Upper Confidence bound (UCB)

Intuition: Let pull arm  $i$  for  $T_i(t)$  times.  
Loosely speaking (using Hoeffding's)

$$\underbrace{\mu_i - \underbrace{\hat{\mu}_i(t)}_{\text{empirical estimate}}} \leq \underbrace{\sqrt{\frac{2 \log(\frac{1}{\delta})}{T_i(t)}}}_{\text{the number of sample.}} \rightarrow \text{at least } 1 - \delta$$

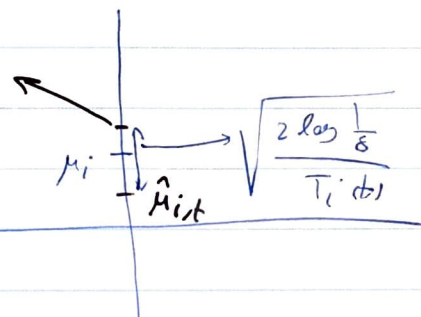
$\hat{\mu}_i(t) = \mu_i / T_i(t)$



$$\mu_i \leq \mu_{i,t} + \sqrt{\frac{2 \log \frac{1}{\delta}}{T_i(t)}}$$

assuming  
 $R=1$   
↳ arms are  
1-sub-Gaussian

Upper Confidence  
bound on  $\mu_i$

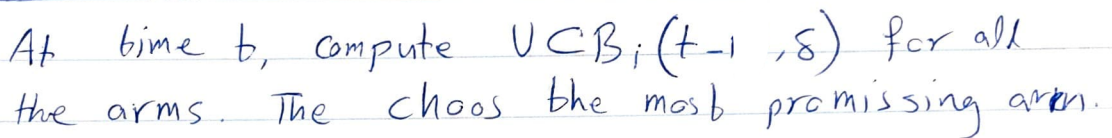


$\mu_i(t) + \sqrt{\frac{2 \log(\frac{1}{\delta})}{T_i(t)}}$  represent upper confidence

bound!

$$UCB(t-1, \delta) = \hat{\mu}_i(t) + \sqrt{\frac{2 \log \frac{1}{\delta}}{T_i(t)}}$$

How a UCB algorithm work?



It is known as Optimism in the Face of Uncertainty (OFU)

$$R_n \leq 3 \sum_{i=1}^K \Delta_i + \sum_{i: \Delta_i > 0} \frac{16 \log(n)}{\Delta_i}$$

when  $\delta = \frac{1}{n^2}$