Given $n$ and $\varepsilon$

$$P\left(|\hat{\mu} - \mu| \geq \varepsilon\right) \leq 2 e^{-\left(\frac{n\varepsilon^2}{2\sigma^2}\right)}$$

**Confidence Interval** Given $n$ and $\delta$, with prob $> 1 - \delta$

$$\hat{\mu} \in \left[\mu - \sqrt{\frac{2\sigma^2 \log(2/\delta)}{n}}, \mu + \sqrt{\frac{2\sigma^2 \log(2/\delta)}{n}}\right]$$

**Number of samples** Given $\varepsilon$ and $\delta$, we need

$$n = \frac{2\sigma^2 \log(2/\delta)}{\varepsilon^2}$$

**Assumption:**

From now on, we will assume without loss of generality $\sigma = 1$

\*   at time $t$,

we play $A_t$ and get $x_t \sim P_{A_t}$

\*   $\mu_i = \mathbb{E}_{P_i}[x_t]$ ,   $\mu_* = \max_i \mu_i$

\*   Regret :   $R_n = \mathbb{E}\left[\sum_{t=1}^{n}(\mu_* - x_t)\right]$

Notational
Convention:    $\mu_1 \geq \mu_2 \geq \cdots \geq \mu_k$   ( for purpose
of analysis)

$$\text{Mean Reward} \uparrow$$

$$\mu_1$$
$$\mu_2$$
$$\Delta_2$$
$$\Delta_i \quad \mu_i$$
$$\Delta_k$$
$$\mu_k$$

# Explore - Then- Commit

**Exploration phase**

- play each arm $m$ times and obtain the rewards

**Exploitation phase**
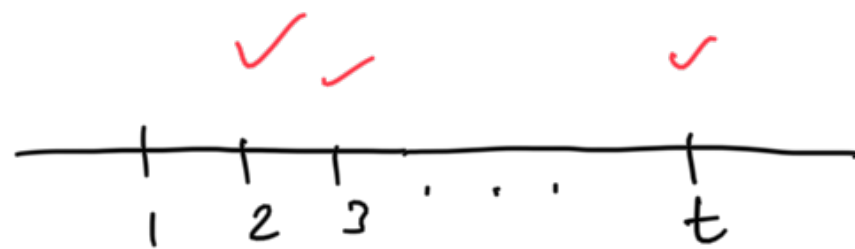
- play the arm with best sample mean

not $t'$ samples

**Sample mean**

$$\widehat{\mu}_i(t) = \frac{\sum_{s=1}^{t} X_s \cdot \mathbb{I}\{A_s = i\}}{T_i(t)}$$

$T_i(t)$ → Total number of times we

Pick
Arm 3

$$\mathbb{I}_{\{A_s = 3\}}$$

✓ ✓            ✓

1  2  3  ·  ·  ·         $t$

* For $1 \leq t \leq mk$

$$A_1 = 1, \quad A_2 = 1, \quad \cdots \quad, A_m = 1$$

Deterministic
exploration

$$A_{m+1} = 2, \quad A_{m+2} = 2, \cdots \quad, A_{2m} = 2$$

⋮

$$A_{m(k-1)+1} = k, \quad \cdots \cdots \quad, A_{mk} = k$$

* for $t > mk$

Exploitation
(or)
commit

$$A_t = \underset{i}{\arg\max} \; \hat{\mu}_i(mk)$$

Regret analysis of Explore - then - Commit

Theorem: Let $n > mk$ ( we have explored for $mk$ rounds and then we are exploiting )

$$R_n \leq m \underbrace{\sum_{i=1}^{k} \Delta_i}_{\text{exploration}} + (n-mk) \underbrace{\sum_{i=1}^{k} \Delta_i \, e^{-\left(\frac{m\Delta_i^2}{4}\right)}}_{\text{exploitation}}$$

Proof: For sake of analysis ( $\mu_1 \geq \mu_2 \geq \cdots \geq \mu_k$ )

we know that $R_n = \sum_{i=1}^{k} \Delta_i \, \mathbb{E}[t_i(n)]$

for first mk (during exploration), we choose

each action deterministically 'm' times

out of mk rounds

$$\mathbb{E}[T_i(n)] = m + (n-mk) \, \mathbb{P}(i^{th} \text{ arm was chosen after mk rounds of exploration})$$

exploitation ↑

arbitrary ties ↘

$$\leq m + (n-mk) \, \mathbb{P}(i^{th} \text{ arm was one of the best arms after mk rounds})$$

↑ arbitrary ties

$$= m + (n-mk) \, \mathbb{P}\left(\hat{\mu}_i(mk) \geq \max_{j \neq i} \hat{\mu}_j(mk)\right)$$

↑ every other arm

$$\mathbb{P}(\text{arm } i \text{ beats all other arms}) \leq \mathbb{P}(\text{arm } i \text{ beats arm } 1)$$

Event A                    Event B

$$A \subseteq B$$

$$\mathbb{P}(\text{arm } i \text{ beats arm } 1) = \mathbb{P}(\hat{\mu}_i(mk) \geqslant \hat{\mu}_1(mk))$$

$$= \mathbb{P}(\hat{\mu}_i(mk) - \hat{\mu}_1(mk) \geqslant 0)$$

$$= \mathbb{P}(\underbrace{\hat{\mu}_i(mk) - \mu_i}_{\substack{\text{Sample mean} \\ - \\ \text{True Mean} \\ \text{for arm } i}} + \underbrace{\mu_i - \mu_1}_{-\Delta_i} + \underbrace{\mu_1 - \hat{\mu}_1(mk)}_{\substack{\text{Sample mean} \\ - \\ \text{True Mean} \\ \text{for arm } 1}} \geqslant 0)$$

$$Y = \hat{\mu}_i(mk) - \mu_i \qquad\qquad Y' = \mu_1 - \hat{\mu}_1(mk)$$

$Y$ and $Y'$ are $\dfrac{1}{\sqrt{m}}$ subgaussian

$Y + Y'$ to be $\sqrt{\dfrac{2}{m}}$ subgaussian

$$= \mathbb{P}\left( \underbrace{\hat{\mu}_i(mk) - \mu_i}_{Y} + \underbrace{\mu_i - \mu_1}_{-\Delta_i} + \underbrace{\mu_1 - \hat{\mu}_1(mk)}_{Y'} \geq 0 \right)$$

$$= \mathbb{P}\left( Y + Y' \geq \Delta_i \right)$$

$$\leq e^{-\left( m \frac{\Delta_i^2}{4} \right)}$$

$$R_n \leq m \underbrace{\sum_{i=1}^{k} \Delta_i}_{\text{exploration}} + (n - mk) \underbrace{\sum_{i=1}^{k} \Delta_i \, e^{-\left( \frac{m\Delta_i^2}{4} \right)}}_{\text{exploitation}}$$

Gap Dependent Bound: Consider $k = 2, \ \Delta_1 = 0, \ \Delta_2 = \Delta$  ← Case

Let use assume we know the gap $\Delta$

$$R_n \leq m\Delta + (n - 2m)\Delta e^{-\left(\frac{m\Delta^2}{4}\right)} \quad —$$

$$\left( \text{let us also say } n \text{ is much larger than } 2m \right)$$

$$R_n \leq m\Delta + n\Delta e^{-\left(\frac{m\Delta^2}{4}\right)} \quad —$$

to minimise the R.H.S, diff w.r.t $m$

$$\Delta + n\Delta\left(-\frac{\Delta^2}{4}\right)e^{-\left(\frac{m_* \Delta^2}{4}\right)} = 0$$

$$\cancel{\Delta} = n\cancel{\Delta}\left(\frac{\Delta^2}{4}\right)e^{-\left(\frac{m_* \Delta^2}{4}\right)}$$

$$e^{-\left(\frac{m_* \Delta^2}{4}\right)} \quad m \Delta^2$$

$$\frac{m_* \Delta^2}{4} = \log\left(\frac{n\Delta^2}{4}\right)$$

**Optimal exploration**

$$m_* = \frac{4}{\Delta^2} \log\left(\frac{n\Delta^2}{4}\right)$$

↑

notice the $\Delta^2$

**Arm 1**



**Arm 2**

$$m_* = \max\left\{ 1, \left\lceil \frac{4}{\Delta^2} \log\left(\frac{n\Delta^2}{4}\right) \right\rceil \right\}$$

$$- \left(\frac{m\Delta^2}{4}\right)$$

$$R_n \leq m\Delta + n\Delta e$$

$$R_n \leq \Delta \max\left\{1, \left\lceil \sqrt{\frac{4}{\Delta^2} \log\left(n\frac{\Delta^2}{4}\right)} \right\rceil \right\}$$

$$\underbrace{\phantom{\Delta \max\left\{1, \left\lceil \sqrt{\frac{4}{\Delta^2}}\right\rceil\right\}}}_{\text{Term I}}$$

$$e^{-5} \leq e^{-\frac{5}{2}}$$

$$+$$

$$n\Delta \, e^{-\left(\frac{\Delta^2}{4} \max\left\{1, \left\lceil \frac{4}{\Delta^2} \log\left(n\frac{\Delta^2}{4}\right)\right\rceil\right\}\right)}$$

$$\underbrace{\phantom{n\Delta \, e^{-\left(\frac{\Delta^2}{4}\right)}}}_{\text{Term II}}$$

$$\lceil x \rceil \leq 1 + x$$

**Term I**

$$\Delta \quad \text{or} \quad \Delta + \frac{4}{\Delta} \log\left(\frac{n\Delta^2}{4}\right)$$

$$\Delta \lceil x \rceil \leq \Delta + \Delta x$$

**Term II**

$$n\Delta \, e^{-\left(\frac{\cancel{\Delta^2}}{4} \cdot \frac{\cancel{4}}{\cancel{\Delta^2}} \log\left(\frac{n\Delta^2}{4}\right)\right)}$$

$$= n\Delta \, \frac{4}{n\Delta^2} \qquad = \frac{4}{\Delta}$$

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

as $\Delta \to 0$, above bound belongs to $\infty$

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