· Let $x \in [0,1]$ with mean $\mu$;

then, $\forall \lambda \in \mathbb{R}$

$$E\left[e^{\lambda(x-\mu)}\right] \leq e^{\lambda^2/8};$$

and $E\left[e^{\lambda(\mu-x)}\right] \leq e^{\lambda^2/8}.$ $\Bigg\} \leftarrow$

In general:

Def: A rv is $\sigma$-subgaussian if $\forall \lambda \in \mathbb{R}$

$$E\left[\exp(\lambda x)\right] \leq \exp\left(\lambda^2 \sigma^2/2\right)$$

Cor: Let $x_i - \mu$ be independent $\sigma$-sub gaussian r.v. Set $x = \dfrac{x_1 + \cdots + x_n}{n}$;

$\boxed{\dfrac{\sigma}{\sqrt{n}}}$

$c X$
$\uparrow$
$|c| \sigma$

Then

ⓐ $E[x] = \mu$;

ⓑ Set $\hat{\mu} = \dfrac{x_1 + \cdots + x_n}{n}$;

Then $\boxed{\mathbb{P}\left[\hat{\mu} \geq \mu + \varepsilon\right] \leq \exp\left(\dfrac{-n\varepsilon^2}{2\sigma^2}\right)}$

and $\Pr\left[\hat{\mu} \leq \mu - \varepsilon\right] \leq \exp\left(\dfrac{-n\varepsilon^2}{2\sigma^2}\right)$ ✓

$\exp\left(-\log(1/\delta)\right)$

In particular: $\forall \delta \in [0,1]$

$$\Pr\left[\mu \leq \hat{\mu} + \sqrt{\dfrac{2\sigma^2 \log(1/\delta)}{n}}\right] \geq 1 - \delta$$

$\& \Pr\left[\mu \geq \hat{\mu} - \sqrt{\dfrac{2\sigma^2 \log(1/\delta)}{n}}\right] \geq 1 - \delta$

• Ex of subGaussian r.v:

ⓐ If $X$ is Gaussian with mean $0$ & variance $\sigma^2$, $X$ is $\sigma$-subgaussian

ⓑ If $X$ has mean zero & $|X| \leq B$, $B \geq 0$, $X$ is $B$-subgaussian. ✓

ⓒ $X$ has mean zero & $X \in [a, b]$ then $X$ is $\dfrac{b-a}{2}$ subgaussian

· Explore then commit bandit:
_____

· Assume all reward distributions for all arms
  is 1-Subgaussian; (ie) $\sigma = 1$ ✓  ⟵

· Algorithms make use of this knowledge $f$ ⊙

There are $k$ actions; The algorithm will
explore for $\underline{mk}$ rounds & choose a single
action for $\overline{\text{remaining}}$ rounds.  }

Let $\hat{\mu}_i(t) =$ avg reward obtained from

arm $i$, after round $t$:

(ie) $\left| \hat{\mu}_i(t) = \frac{1}{T_i(t)} \sum_{s=1}^{t} \mathbb{1}\{A_s = i\} x_s. \right\|$

Here $T_i(t) = \sum_{s=1}^{t} \mathbb{1}\{A_s = i\} = $ # action $i$ has

been played till round $t$;

<u>Explore then commit;</u>

$(\circ, \circ, ', ', ', ', ')$

Input $m$

$$\underbrace{1 \ 2 \ 3 \ -- \ k \ \ 1 \ 2 \ --- k \ \ 1 \ 2 \ -- \ k \ \ 1 \ 2 \ --- k}_{mk}$$

In round $t$ choose

$$A_t = \left\{ \begin{array}{ll} t \bmod k & \text{if } \underline{t \le mk} \\ \underset{i}{\operatorname{argmax}} \ \hat{\mu}_i(mk), & \underline{t > mk} \end{array} \right\}$$

<u>Recall:</u>

$$\boxed{\Delta_i(\nu) = \mu^*(\nu) - \mu_i(\nu);}$$

Suppose $\boxed{\mu_1 = \mu^*}$ then $\Delta_i = \underline{\mu_1 - \mu_i}$ ✓

• The above is deterministic till round $mk$ and chooses each action $m$ times;

Recall: $\boxed{\overline{R}_n = \sum_{i=1}^{k} \Delta_i \, \mathbb{E}[T_i(n)];}$

**Clearly:**

$$\mathbb{E}\left[T_i(n)\right] = \underbrace{m} + \underbrace{(n - mk)\, \Pr\left[A_{mk+1} = i\right]}$$

$$\leq m + (n - mk)\, \Pr\left[\hat{\mu}_i(mk) \geq \max_{\hat{j} \neq i} \hat{\mu}_j(mk)\right]$$

**Now**

$$\boxed{\Pr\left(\hat{\mu}_i(mk) \geq \max_{\hat{j} \neq i} \hat{\mu}_j(mk)\right)} \quad \Longleftarrow$$

$$\leq \Pr\left(\hat{\mu}_i(mk) \geq \hat{\mu}_1(mk)\right)$$

$$= \Pr\left[\hat{\mu}_i(mk) - \hat{\mu}_1(mk) \geq 0\right]$$

$$= \Pr\left[\hat{\mu}_i(mk) - \hat{\mu}_1(mk) \geq \Delta_i - \Delta_i\right]$$

$$= \Pr\left[\hat{\mu}_i(mk) - \hat{\mu}_1(mk) + \Delta_i \geq \Delta_i\right]$$

$$= \Pr\left[\underbrace{\left(\hat{\mu}_i(mk) - \mu_i\right)}_{\substack{\uparrow \text{ subgaussian} \\ \sigma = 1}} - \underbrace{\left(\hat{\mu}_1(mk) - \mu_1\right)}_{\substack{\uparrow \text{ subgaussian} \\ \sigma = 1}} \geq \Delta_i\right]$$

$$\boxed{\Delta_i = \mu_1 - \mu_i}$$

**H.W.** $X_1$ subgaussian $\underbrace{\phantom{\qquad}}_{\sigma_1}$  $X_2$ subgaussian then $\underbrace{\phantom{\qquad}}_{\sigma_2}$

$$\dfrac{X_1 + X_2 \leq \sqrt{\sigma_1^2 + \sigma_2^2}\ \text{subgaussian};}{}$$

first arm is pulled $m$ times till round $m$;

$$\hat{\mu}_i = \dfrac{X_i + X_{i+k} + \cdots + X_{i+(m-1)k}}{m}$$

$$X_i + X_{i+k} + \cdots + X_{i+(m-1)k} \leq \sqrt{\underbrace{1 + 1 + \cdots + 1}_{m}}$$

$$= \sqrt{m}\ \text{subgaussian};$$

$$\therefore \dfrac{X_i + X_{i+k} + \cdots + X_{i+(m-1)k}}{m} \leq \dfrac{1}{\sqrt{m}}\ \text{subgaussian};$$

$$\therefore \left(\hat{\mu}_i - \mu_i\right) - \left(\hat{\mu}_1 - \mu_1\right)$$

$\boxed{\dfrac{1}{\sqrt{m}}}$ $\quad \boxed{\dfrac{1}{\sqrt{m}}}$

$$\therefore \quad \sqrt{\frac{1}{m} + \frac{1}{m}} = \boxed{\sqrt{\frac{2}{m}}} \quad \text{Gaussian}$$

$$\therefore \quad \Pr\left[ \hat{\mu_i}(mk) - \mu_i - \left( \hat{\mu_1}(mk) - \mu_1 \right) \geq \underline{\underline{\Delta_i}} \right]$$

$$\leq \exp\left( \frac{-\Delta_i^2}{2\left(\sqrt{\frac{2}{m}}\right)^2} \right)$$

$$= \exp\left( \frac{-\Delta_i^2 m}{4} \right) \checkmark \qquad \overset{(m\tau-)}{R_n} = \sum_{a=1}^{h} \Delta_a E\left( T_a(n) \right)$$

$$\therefore \quad \boxed{R_n \leq m \sum_{i=1}^{k} \Delta_i + (n - mk) \sum_{i=1}^{h} \Delta_i \exp\left( \frac{-\Delta_i^2 m}{4} \right)}$$

- If $m$ large first term dominates
- If $m$ small, prob of choosing wrong arm ↑ & second dominates

$k = 2$ case: say $\Delta_1 = 0$, $\Delta = \Delta_2$

$$\boxed{R_n \leq m\Delta + (n - 2m)\Delta \exp\left( \frac{-m\Delta^2}{4} \right)} \leftarrow$$

$$\leq \; m\Delta + n\Delta \exp\left(-m\frac{\Delta^2}{4}\right)$$

$$\Delta\!\!\!/ + n\!\!\!\!\!\diagup e^{\left(-m\frac{\Delta^2}{4}\right)} \cdot \left(-\frac{\Delta^2}{4}\right) = 0$$

$$1 + n e^{-m\frac{\Delta^2}{4}} \cdot \left(-\frac{\Delta^2}{4}\right) = 0$$

$$\sim \frac{1}{n} = \frac{\Delta^2}{4} \cdot \frac{1}{e^{m\Delta^2/4}}$$

$$e^{m\Delta^2/4} = \frac{\Delta^2 n}{4}$$

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

$$\boxed{\;\therefore m = \frac{4}{\Delta^2}\log\left(\frac{n\Delta^2}{4}\right)\;} \quad\leftarrow$$

$$\boxed{\text{let } m = \max\left\{1, \; \left\lceil \frac{4}{\Delta^2}\log\left(\frac{n\Delta^2}{4}\right)\right\rceil\right\}}$$

$$\boxed{R_n \leq \Delta + c\sqrt{n}}$$ ✓ $\bar{R}_n \leq \boxed{n^{2/3}}$

Sublinear in n.

$$\boxed{\text{OPTIMISM in the face of uncertainty.}}$$

. Use the observed data to assign to each arm an "upper confidence bound", which with high probability is an overestimate of unknown mean.

. So the other arm can be played only if its UCB is greater than that of optimal arm. But that is larger than the mean of the optimal arm.

This can't happen too often — as that arm is played more its data will force us to reduce the value we have assigned to its UCB. Eventually it will fall below UCB of optimal arm.

$\underline{\text{Recall}}$: $(X_t)_{t=1}^n$ a seq of $1$-subgaussian

r.v , mean $= \mu$ & set $\hat{\mu} = \frac{1}{n} \sum X_t$ ;

$$\mathbb{P}_\delta \left[ \mu > \hat{\mu} + \sqrt{\frac{2\log(1/\delta)}{n}} \right] \leq \boxed{\delta} \quad \forall \delta \in (0,1)$$

• In round $t$, learner has seen $T_i(t-1)$

Samples of arm $i$ ; empirical mean $\bar{\mu}_i(t-1)$

Then $\underline{UCB_i(t-1,\delta)} = \begin{cases} \infty & \text{if } \underline{T_i(t-1) = 0} \\ \hat{\mu}_i(t-1) + \sqrt{\dfrac{2\log(1/\delta)}{T_i(t-1)}} \end{cases}$

is a reasonable assignment of a
UCB for the unknown mean of arm $i$

$\hat{\mu}_i \rightarrow$ empirical mean of rewards of arm $i$

$\sqrt{\dfrac{2\log(1/\delta)}{T_i(t-1)}}$   confidence width

## Algorithm:

Input $k, \delta$

for $t \in 1, \ldots, n$ do

      Choose $A_t = \text{argmax}_i \; UCB_i(t-1, \delta)$

      observe $X_t$ and update $UCB$

end

- $(1/\delta)$ - called confidence level ✓

$= \underline{\mathbf{\delta}}$: Can width of confidence interval for given confidence level, be significantly decreased? Turns out we cannot improve it significantly. →

- $\delta$ - should be small for ensuring optimism with high probability!

- But then the confidence interval becomes larger and suboptimal arms may be

played excessively.

- <u>Idea</u>: If the confidence interval fails &
$UCB_{i^*}$ drops below $\mu_{i^*}$ — we may not play
$i^*$ anymore & could get linear regret;
<u>Select</u> ( $\delta = 1/n$ ) in expectation the contribution
to regret of this case is small.

<span style="color:red">PROBLEM</span>: <span style="color:red">$T_i(t-1)$ is a r.v.</span>

So $\delta$ should be smaller — $1/n^2$

THM: Stochastic $k$-armed 1-subgaussian
bandit. For any horizon $n$, if $\delta = 1/n^2$

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

$R_n = O(\sqrt{n})$