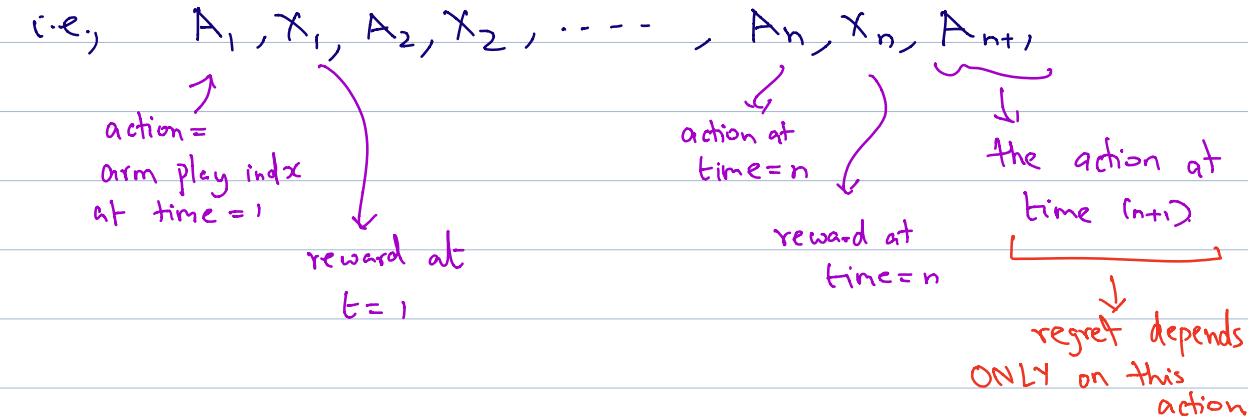


## Algorithms for Pure Exploration

Source: Bandit Algorithms, L&S., Chap 33.

Setting:  $k$ -arms, unstructured environment  $\nu \in \mathcal{E}$ .

Goal: Pure exploration, c.e.,  $\tilde{\pi} = (\pi_t, t=1, 2, \dots, n+1)$  a policy (mapping from history to an arm index).



Simple Regret:  $R_n^{\text{single}}(\pi, \nu) = \mathbb{E}_{\nu \pi} [\Delta_{A_{n+1}}(\nu)]$

time steps provided  
to "experiment"

policy

environment

$\Delta_i(\nu) = \text{gap of } i^{\text{th}} \text{ arm wrt best arm in}$   
 $\text{environment } \nu, \nu \in \mathcal{E}.$

## The Simplest Algorithm: Uniform Explore.

Input:  $n, k$

Iteration: for each  $t=1, 2, \dots, n$ :

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

Output (at time  $n+1$ ):  $A_{n+1} = \arg \max_{1 \leq i \leq k} \hat{\mu}_i(n)$

where  $\hat{\mu}_i(n) = \text{sample average of } i^{\text{th}} \text{ arm rewards}$ .

Thm (33.1 in textbook):  $\mathcal{E} = \mathcal{E}_{SG}^k(1), \forall v \in \mathcal{E}_{SG}^k(1)$ ,  
an environment ( $1\text{-subG}$ ). For all  $n \geq k$ ,

$$R_n^{\text{Simple}}(\pi, v) \leq \min_{\Delta \geq 0} \left\{ \Delta + \sum_{\{i : \Delta_i(v) > \Delta\}} \Delta_i(v) \ell \frac{-\lfloor \frac{n}{k} \rfloor \Delta_i(v)^2}{4} \right\}$$

Proof: follows immediately from  $1\text{-subG}$ . For each arm, we have at least  $\lfloor \frac{n}{k} \rfloor$  samples. Assume (wlog) that arm 1 is optimal (i.e.  $\Delta_1 = 0$ ).

For any  $i$  (sub-optimal arm) with  $\Delta_i > 0$ , we have:

$$P(\hat{\mu}_i(n) \geq \hat{\mu}_1(n)) = P(\underbrace{\hat{\mu}_i(n) - \hat{\mu}_1(n)}_{\text{N}(\mu=0, \sigma^2)} \geq 0)$$

$\sqrt{2}$ -subG with mean  
 $-\Delta_i$

$$\leq e^{-\frac{\sum_i \Delta_i^2}{4}}$$

$$\therefore R_n^{\text{simple}}(\pi, v) = \mathbb{E}_{v \in \Pi} [\Delta_{A_{n+1}(v)}]$$

$$= \sum_{i=2}^k \Delta_i P(A_{n+1} = i) \quad \left( \begin{array}{l} \text{sub optimal} \\ \text{arms} \end{array} \right)$$

Now, choose any  $\Delta > 0$ , and let  $S = \{i \leq k : \Delta_i > \Delta\}$

$$\hookrightarrow = \sum_{i \notin S} \Delta_i P(A_{n+1} = i) + \sum_{i \in S} \Delta_i P(A_{n+1} = i)$$

$$\leq \Delta + \sum_{i \in S} \Delta_i e^{-\frac{\sum_i \Delta_i^2}{4}}$$

②

Note 1: For fixed  $k$ , fixed environment  $v \in \mathcal{E}_{SG}^k(1)$ ,  
the simple regret decays exponentially in time horizon  $n$ . In other words, the prob of not playing the optimal arm at time  $(n+1)$  is exponentially small..

Note 2: (Worst-case bound over all environments  $v \in \mathcal{E}_{SG}^k(1)$ )

Use B-H inequality to show that  $\exists v \in \mathcal{E}_{SG}^k(1)$  s.t.

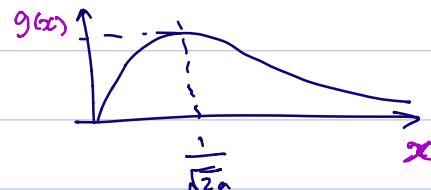
$$R_n^{\text{simple}}(\pi, v) \geq C \sqrt{k/n}$$

With above Uniform Explore policy, recall that

$$R_n^{\text{simple}}(\pi, v) \leq \min_{\Delta \geq 0} \left\{ \Delta + \sum_{\{i : \Delta_i(v) > \Delta\}} \Delta_i(v) e^{-\frac{n}{k} \frac{\Delta_i^2(v)}{4}} \right\}$$

Choose  $\Delta = 2 \sqrt{\frac{\ln(k)}{L^n k}}$ , and assume  $n \geq k$ .

Aside:  $z^* = \max_{x \geq b} x e^{-ax^2}$



$$\therefore b > \frac{1}{\sqrt{2a}} \Rightarrow z^* = g(b)$$

$$R_n^{\text{simple}}(\pi, v) \leq 2 \sqrt{\frac{\ln(k)}{L^n k}} \left( 1 + (k-1) e^{-\ln k} \right)$$

$$\leq 4 \sqrt{\frac{\ln(k)}{L^n k}} \leq 4 \sqrt{\frac{k \ln k}{n}}$$

$\Rightarrow$  Extra  $\sqrt{\ln k}$  factor compared to lower bound.

An alternate policy that (in the worst case) gets rid of the  $\sqrt{\ln(k)}$

Let  $\pi$  be any policy that optimizes for cumulative regret (e.g. UCB). Let  $T_j(t)$  be the number of times that arm  $j \in \{1, 2, \dots, k\}$  is played under this policy  $\pi$ .

Now define the policy  $\hat{\pi} = (\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_n, \hat{\pi}_{n+1})$ ,

with

$$\hat{\pi}_t = \pi_t, \quad t=1, 2, \dots, n$$

$$\text{and } \hat{\pi}_{n+1} \equiv \hat{\pi}_{n+1}(j | a_1, x_1, a_2, x_2, \dots, a_n, x_n)$$

$$= \frac{T_j(n)}{n}, \quad j=1, 2, \dots, k$$

(aka "choose best arm proportional to the relative frequency among played arms").

Prop (33.2 in text): The simple regret of this policy  $\hat{\pi} = (\hat{\pi}_t, t=1, 2, \dots, n, n+1)$  is

$$\text{Given by : } R_n^{\text{simple}}(\hat{\pi}, v) = \underbrace{\frac{R_n(\pi, v)}{n}}$$

cumulative regret  
of policy  $\pi$ .

Proof: Regret decomposition implies

$$R_n(\pi, v) = \sum_{i=1}^k \Delta_i E[T_i(n)]$$

$$= n E \left[ \sum_{i=1}^k \Delta_i \frac{T_i(n)}{n} \right]$$

$$= n E[\Delta_{A_{n+1}}] = n R_n^{\text{simple}}(\hat{\pi}, v). \quad \blacksquare$$

Recall that  $\exists \pi \text{ s.t. } R_n(\pi, v) \leq C\sqrt{nk}$  for

all  $v \in \mathcal{E}_{SG}^k(1)$ , from earlier discussion on worst-case regret bounds.

$$\Rightarrow R_n^{\text{simple}}(\hat{\pi}, v) \leq C\sqrt{\frac{k}{n}}$$

Note 3: The above does NOT mean that  $\hat{\pi}$  is a good policy in general. The reason

is that in a fixed environment (not worst case), the above policy  $\hat{\pi}$  provides only a polynomial decay with  $n$  (time horizon), whereas the Uniform Explore policy provides exponential decay.

Please read Sec. 33.1 for more discussion.

---

Fixed Confidence Level  $\delta \in (0, 1)$

Goal: Identify the best arm, with a prespecified error prob. The number of samples is a free variable, and needs to be minimized.

Formally, let  $\mathcal{F}_t = \sigma\{A_1, X_1, \dots, A_t, X_t\}$ , i.e., the  $\sigma$ -algebra corresponding to the information until time  $t$ .

Let  $\tau$  be a stopping time wrt the filtration  $\{\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \dots\}$ . Informally,  $\tau$  is a "causal" r.v., i.e., for each  $t \geq 1$ , there exists a function  $g_t(\cdot)$  such that:

$$\{\tau \leq t\} = \{g_t(A_1, X_1, \dots, A_t, X_t) \leq t\}$$

A Player determines a triplet:  $(\pi, \tau, \psi)$

↓  
 arm play policy ←      ↓  
 $\pi = (\pi_1, \pi_2, \dots)$       Stopping time  
 (when to stop playing)

$\pi_t : (A_1, X_1, \dots, A_{t-1}, X_{t-1}) \mapsto \{1, 2, \dots, k\}$

①  $\pi = (\pi_1, \pi_2, \pi_3, \pi_4, \dots, \pi_t, \dots)$ ,

history of past observations

where  $\pi_t : H_{t-1} \mapsto \{1, 2, \dots, k\}$

②  $\tau$  a stopping time, determines when to stop playing and declare that a specific arm to be the best arm.

③  $\psi \in \{1, 2, \dots, k\}$  is the selection rule, and chooses which arm to select as the best arm.

Sound Triplet  $(\pi, \tau, \psi)$

A triplet is said to be sound at some pre-selected confidence level

$\varsigma \in (0, 1)$  if

$$P_{\pi\tau}(\{\tau < \infty\} \cap \{\Delta_\psi(v) > 0\}) \leq \varsigma$$

In other words, a sound policy that stops in a finite time does not select a sub-optimal arm with probability exceeding  $\delta$ .

Note 4: With multiple optimal arms, a sound policy will not be able to stop in finite time. (impossible to distinguish between noise and arms that are arbitrarily close); see comment in Sec. 33.2.

### Lower Bound

$\Sigma^k$  an arbitrary  $k$ -armed stochastic bandit;  $v \in \Sigma^k$ .

$$i^*(v) = \underset{1 \leq j \leq k}{\operatorname{argmax}} M_j(v)$$

mean reward of  $i^{th}$  arm in environment  $v \in \Sigma^k$ .

$$\Sigma_{\text{alt}}^k(v) = \{v' \in \Sigma : i(v') \cap i(v) = \emptyset\}$$

$\subseteq \Sigma$  s.t. optimal arm in  $v'$  is different from an optimal in  $v$ .

Recall Divergence Decomposition Thm with Stopping Times (HW problem 15.7 in text):

Thm:  $\tau$  a stopping time with respect to filtration

$$\mathcal{F} = \{\mathcal{F}_t, t \geq 1\} \text{ where}$$

$$\mathcal{F}_t = \sigma\{A_1, X_1, \dots, A_t, X_t\}.$$

If  $Z$  is a r.v. that is  $\mathcal{F}_\tau$ -measurable, then

$$D(P_{vZ}, P'_{v'Z}) \leq \sum_{j=1}^k E_v[T_j(\tau)] \cdot D(P_j, P'_j)$$

where  $P_{vZ}$  is the prob law of  $Z$  under the environment  $v$ , and analogous for  $P'_{v'Z}$ .

For comparison, see Set 11 / Lemma 15.1 in text.

There,  $z \rightarrow \pi$ , and  $\tau \rightarrow t$ . The result showed a separation between policy  $\pi$  and arm divergence. A similar result holds (above) even with stopping times.

Note:  $\mathcal{F}_\tau$  is the  $\sigma$ -algebra at a stopping time  $\tau$ ,

defined as follows:  $\mathcal{F} = \{\mathcal{F}_t, t=1, 2, \dots\}$ .

$$\mathcal{F}_\tau = \{A \in \mathcal{F}_\infty : A \cap \{\tau \leq t\} \in \mathcal{F}_t \text{ for } t=1, 2, \dots\}$$

This is indeed a valid  $\sigma$ -algebra, and intuitively, this describes the "information" at the random time  $\tau$ .

Theorem (33.5 in text): Suppose  $(\pi, \tau, \psi)$  is sound for  $\mathcal{E}^k$  at confidence level  $\delta \in (0, 1)$ . Further, let  $v \in \mathcal{E}$ . Then,

$$E_{vv}[\tau] \geq c^*(v) \ln\left(\frac{1}{4\delta}\right), \text{ where}$$

$$c^*(v)^{-1} = \sup_{\alpha \in P^{k-1}} \left( \inf_{v' \in \mathcal{E}_{alt}^k(v)} \sum_{i=1}^k \alpha_i(v) D(v_i, v'_i) \right)$$

Proof: WLOG  $E_{vv}[\tau] < \infty$  (else above immediately true). Thus  $\tau$  is a proper r.v. (i.e. no mass at  $\infty$ ).

Pick any  $v' \in \mathcal{E}_{alt}^k$ , and  $B = \{\tau < \infty\} \cap \{\psi \notin i^*(v')\}$

i.e.,  $B$  is a bad event in the  $v'$  environment. Further,

$$B^c = \left( \{\tau < \omega\} \cap \{\psi \notin i^*(v')\} \right)^c$$

$$= \{\tau = \omega\} \cup \{\psi \in i^*(v')\}$$

$$= \{\psi \in i^*(v')\}$$

$$\subseteq \{\psi \notin i^*(v)\}$$

$(\because i(v) \cap i(v') = \emptyset)$   
from  $v' \in \Sigma_{alt}^k$

$$= \{\psi \notin i^*(v)\} \cap \{\tau < \omega\} \quad (\because \tau \text{ is a proper r.v.})$$

→ ①

From soundness, we have:

$$P_{v\pi}(\{\tau < \omega\} \cap \{\psi \notin i^*(v)\}) \leq \delta$$

$$P_{v'\pi}(\{\tau < \omega\} \cap \{\psi \notin i^*(v')\}) \leq \delta$$

→ ②

Combining ① and ②, we have:

$$\begin{aligned} 2\delta &\geq P_{v\pi}(\{\tau < \omega\} \cap \{\psi \notin i^*(v)\}) \\ &\quad + P_{v'\pi}(\{\tau < \omega\} \cap \{\psi \notin i^*(v')\}) \end{aligned}$$

$$\stackrel{(2)}{\geq} P_{\nu\pi}(B^c) + P_{\nu'\pi}(B)$$

$$\begin{aligned}
 & \stackrel{\text{B-H inequality}}{\geq} \frac{1}{2} e^{-D(\nu\pi, \nu'\pi)} \\
 (\text{see Set 10; Chap 14}) \quad & \stackrel{\text{intext}}{\geq} \frac{1}{2} e^{-\sum_{j=1}^k E_{\nu\pi}[T_j(\tau)] D(v_j, v'_j)} \\
 & \stackrel{\substack{\text{Divergence} \\ \text{Decomp with} \\ \text{stopping time}}}{\geq} \frac{1}{2} e^{-\sum_{j=1}^k \alpha_j D(v_j, v'_j)} \quad \xrightarrow{\text{L}} \textcircled{3}
 \end{aligned}$$

$\therefore$  From  $\textcircled{3}$ , we have: For any  $v' \in \mathcal{E}_{\text{alt}}^k$ ,

$$\sum_{j=1}^k E_{\nu\pi}[T_j(\tau)] D(v_j, v'_j) \geq \ln\left(\frac{1}{4\delta}\right) \quad \xrightarrow{\text{L}} \textcircled{4}$$

(Aside: this implies  $E_{\nu\pi}[\tau] > 0$ , thus  $0 < \frac{1}{E_{\nu\pi}[\tau]} < \infty$ ).

Finally,

$$\begin{aligned}
 \frac{E_{\nu\pi}[\tau]}{C^*(v)} & \stackrel{\text{defn. of } C^*(v)}{=} E_{\nu\pi}[\tau] \sup_{\alpha \in \mathbb{P}^{k-1}} \inf_{v' \in \mathcal{E}_{\text{alt}}^k} \sum_{j=1}^k \alpha_j D(v_j, v'_j) \\
 & \geq E_{\nu\pi}[\tau] \inf_{v' \in \mathcal{E}_{\text{alt}}^k} \sum_{j=1}^k \left[ \frac{E_{\nu\pi}[T_j(\tau)]}{E_{\nu\pi}[\tau]} \right] D(v_j, v'_j) \quad \begin{array}{l} \text{s.t. } \sum_i \alpha_i = 1 \\ \alpha_i \geq 0 \end{array}
 \end{aligned}$$

$$= \inf_{v' \in \mathcal{E}_{\text{alt}}^k} \sum_{j=1}^k E_{v\pi}[\tau_j(v)] D(v_j, v'_j)$$

$$\geq \ln\left(\frac{1}{48}\right).$$

④



Some intuition on  $c^*(v)$

Fix an environment  $v \in \mathcal{E}^k$ , and let the gap be  $\Delta$  between the best and second best arm. Recall that

$$c^*(v)^{-1} = \sup_{\alpha \in P^{k-1}} \left( \inf_{v' \in \mathcal{E}_{\text{alt}}^k(v)} \sum_{j=1}^k \alpha_j D(v_j, v'_j) \right)$$

Potential problem: for a fixed  $\alpha$ , we might find a sequence of  $v' \in \mathcal{E}_{\text{alt}}^k(v)$  s.t.  $D(v_j, v'_j) \rightarrow 0$ , and thus making  $c^*(v)^{-1} = 0$  for any  $v \in \mathcal{E}^k$ .

The above is not true. Consider as an example

$$\mathcal{E}^k = \mathcal{E}_N^2(1) \quad (\text{i.e., 2-armed 1-Gaussian bandit}).$$

Let's say  $v = (v_1, v_2) = (N(\Delta, 1), N(0, 1))$ ,

$$\text{i.e., } \mu_1 = \Delta, \mu_2 = 0.$$

$$D(\nu_j, \nu'_j) = \frac{1}{2} (\mu_j - \mu'_j)^2 \quad \begin{matrix} \text{(for Gaussian)} \\ \text{v.v.s.} \end{matrix}$$

$\nu' \in \Sigma_{alt}^2 = \{(\mu'_1, \mu'_2) : \mu'_2 > \mu'_1\}$

Suppose we choose  $\mu' = (\mu'_1, \mu'_2)$  s.t.  $\mu'_2 = \mu_2 = 0$   
 (i.e., forcing the divergence for arm 2 to go to zero),  
 the constraint that  $\mu'_1 < \mu'_2 = 0$ , and that  
 $\mu_1 = \Delta > 0 \Rightarrow D(\mu_1, \mu'_1) \geq \frac{\Delta^2}{2} > 0$ .

More generally, we cannot construct an alternative environment  $\nu'$  s.t.  $D(\mu_j, \mu'_j) \approx 0$  for all arms. Specifically for the case where

$$\nu = (N(\mu_1, 1) \ N(\mu_2, 1)); \quad \nu \in \Sigma_N^2 (1)$$

with  $\mu_1 > \mu_2$ , we have:

$$\Sigma_{alt}^2 = \{ (N(\mu'_1, 1) \ N(\mu'_2, 1)) : \mu'_2 > \mu'_1 \}$$

$$\text{Then, } C^*(\nu) = \sup_{\alpha \in P^{k-1}} \left( \inf_{\nu' \in \Sigma_{alt}^k(\nu)} \sum_{j=1}^k \alpha_j D(\nu_j, \nu'_j) \right)$$

$$= \sup_{\alpha \in [0,1]} \inf_{\nu' \in \Sigma_{alt}^2} \alpha \frac{(\mu_1 - \mu'_1)^2}{2} + (1-\alpha) \frac{(\mu_2 - \mu'_2)^2}{2}$$

$$= \sup_{\alpha \in [0,1]} \frac{1}{2} \alpha (1-\alpha) (\mu_1 - \mu_2)^2 \quad (\text{check!})$$

$$= \frac{1}{8} (\mu_1 - \mu_2)^2, \text{ with } \alpha^* = \frac{1}{2}$$

(See textbook, sec. 33.2.1 for a more extensive discussion of  $c^*(\gamma)$ )

### The Track and Stop Algorithm

Goal: Algorithm that asymptotically meets the lower bound:  $(\pi, \tau, \psi)$  s.t. lower bound met as  $\delta \rightarrow 0$ .

Setting:  $\Sigma = \sum_N^k (1)$ , i.e., rewards are Gaussian with  $\sigma^2 = 1$  for each arm.

Intuition: We need an algorithm  $(\pi, \tau, \psi)$  s.t.

$$E_{\pi}[\tau] \approx c^*(\gamma) \cdot \ln\left(\frac{4}{\delta}\right).$$

$\downarrow$   
as  $\delta \rightarrow 0$

Now recall the lower bound proof. Towards the end,

we had the following chain:

$$\frac{E_{\nu\pi}[\tau]}{C^*(\nu)} \geq \inf_{\nu' \in \mathcal{E}_{alt}^k} \sum_{j=1}^k E_{\nu\pi}[\tau_j(\tau)] D(v_j, v'_j)$$

If we hope that this is tight (at least in an asymptotic sense), then, we should expect that a "good"  $(\pi, \tau, \psi)$  triplet should satisfy:  $\ln\left(\frac{1}{\epsilon}\right) + \ln\left(\frac{1}{\delta}\right)$

$$\inf_{v' \in E_{\text{alt}}^k} \sum_{j=1}^k E_{v''} [T_j(z)] D(v_j, v'_j) \approx \underbrace{\ln\left(\frac{1}{4\delta}\right)}_{\hookrightarrow \textcircled{1}}$$

In the regime where  $S \rightarrow 0$ , we expect  $T \rightarrow \infty$ , and  $T_j(\tau) \rightarrow \infty$ . To get asymptotic optimality, we thus have some slack: Instead of hoping to satisfy ①, we can try to satisfy:

$$\inf_{v^* \in E_{\text{alt}}^k} \sum_{j=1}^k E_{\nu\pi} [T_j(z)] D(v_j, v'_j) \approx \underbrace{h(t)}_{\substack{\text{some slowly growing function of } t \\ \text{s.t. this term is not dominant, but}}} + \underbrace{\ln\left(\frac{1}{s}\right)}_{\substack{\text{nevertheless gives flexibility to} \\ \text{ensure sufficient concentrations.}}}$$

In a Gaussian setting ( $\mathcal{E}_N^k(\cdot)$ ),  $D(v_i, v'_i) = \frac{(u_i - u'_i)^2}{2}$

(a) Thus, we are looking for a criteria that uses the above quadratic.

(b) Further, since we do not have access to  $\{u_j, j=1, 2, \dots, k\}$ , and we are on a single sample path (i.e. no  $E_{\pi\tau}[\cdot]$ ), we can hope to replace the appropriate parameters by a sample average and/or single sample.

(c) Recall in the lower bound that  $\alpha_j(\tau) = \frac{E_{\pi\tau}[T_j(\tau)]}{E_{\pi\tau}[\tau]}$ ,

$j=1, 2, \dots, k$ . Thus, the intuition is that this is the fraction of time to spend sampling arm  $j$ .

Motivated by above, three ideas to develop an algorithm, characterized by  $(\pi, \tau, \psi)$ :

① Choosing  $\tau$ : In ② above, replace  $E_{\pi\tau}[\cdot]$  by the sample (i.e., get rid of  $E[\cdot]$ ), and replace  $u_j$  by  $\hat{u}_j(t)$  = empirical estimate of arm  $j$ 's mean at time  $t$ . Finally, estimate best current

arm using the empirical averages (i.e., no UCB or any variance correction). Then define a surrogate for LHS in ② as  $Z_t$ , where:

$$Z_t = \inf_{v' \in E_{alt}^k(\hat{v}(t))} \sum_{j=1}^k T_j(t) \frac{(\hat{\mu}_j(t) - \mu'_j)^2}{2}$$

i.e., avoid the currently estimated best arm.

use sample instead of  $E_{\pi^*}[v]$

use empirical estimate to replace unknown  $\mu_i$

- Here, there is a unique best arm,
- defined by  $i^*(t) = \operatorname{argmax}_j \hat{\mu}_j(t)$ .
- If several  $\hat{\mu}_j(t)$  are equal, pick
- any one of them to be the (sole)
- declared best arm

Original Expression in ②

$$\inf_{v' \in E_{alt}^k} \sum_{j=1}^k E_{\pi^*}[T_j(v)] D(v_j, v'_j)$$

Stopping criteria motivated by ②: Stop at first

$$t \geq k : Z_t > \beta_t(s)$$

$$\text{where } \beta_t(s) \approx \underbrace{2k(\ln t) + \ln(1/s)}_{h(t) \text{ is slow wrt } t} + \underbrace{o(\ln t + \ln \frac{1}{s})}_{\text{exact correction to be specified later.}}$$

② Choosing  $\pi = (\pi_t, t=1, 2, 3, \dots)$ , the sampling strategy:

Since LB uses  $d_j(\nu) = \frac{E_{\nu\nu}[T_j(\nu)]}{E_{\nu\nu}[\nu]}$ , we will

try to satisfy for any  $t \geq k$ :  $T_j(t) \approx \hat{\alpha}_j^*(\hat{\nu}) \cdot t$ .

The problem is that we do not know  $\nu$ . However, we do have an approximation  $\nu \approx \hat{\nu}(t) = (\hat{\mu}_1(t), \dots, \hat{\mu}_k(t))$

$\therefore$  Choose action  $A_{t+1} = \arg \max_{1 \leq j \leq k} (t \hat{\alpha}_j^*(t) - T_j(t))$

where  $\hat{\alpha}_j^*(t) = \alpha^*(\hat{\nu}(t))$

i.e.,

$$\hat{\alpha}^*(\hat{\nu}(t)) = \arg \sup_{\alpha \in \mathbb{P}^{k-1}} \inf_{\nu' \in \mathcal{E}_{alt}^k(\hat{\nu}(t))} \sum_{j=1}^k (t \alpha_j) \frac{(\hat{\mu}_j(t) - \mu'_j)^2}{2}$$

As above, if there are more than one best arm, pick one of them to be the best (breaking ties arbitrarily).

③ Ensure sufficient number of samples for each arm for sufficient concentrations, so that above steps of "replace expectations by sample-average" work!

At any time  $t$  s.t.  $T_j(t) \leq \sqrt{t}$ , for some  $j \in \{1, 2, \dots, k\}$ , choose

$$A_{t+1} = \underset{1 \leq j \leq k}{\operatorname{arg\,min}} T_j(t).$$

④ Decision  $\Psi$ : Stop when ① above is violated.

$$\psi = i^*(v(t)) = \arg \max_{1 \leq j \leq k} \hat{\mu}_j(t)$$

Return:

$$\tau = t$$

## Track-and-Stop Algorithm

(Garrivier & Kaufmann 2016).

Input:  $\delta$ ,  $\beta_t(\delta)$

Initial: Play each arm once ( $K$  rounds total).

Iteration: If  $z_t < \beta_t(s)$ , then:

Case (i)  $\underset{1 \leq j \leq k}{\operatorname{argmin}} T_j(t) \leq \sqrt{t}$  :  $A_{t+1} = \underset{j}{\operatorname{argmin}} T_j(t)$

Case (ii)  $\underset{1 \leq j \leq k}{\operatorname{argmin}} T_j(t) > \sqrt{t}$  :  $A_{t+1} = \underset{j}{\operatorname{argmax}} \left( t \hat{\alpha}_j^*(t) - T_j(t) \right)$

Observation: Reward sample  $X_{A_{t+1}(t)}$ , update arm statistics

$t \rightarrow t+1$  — — — — —

Return:  $\psi = \underset{j}{\operatorname{argmax}} \hat{\mu}_j(t)$ ,  $T = t$

Choice of  $\beta_t(s)$ : We earlier discussed that  $\beta_t(s)$  should scale as  $2k \ln t + \ln(\frac{1}{\delta}) + o(\dots)$ .

More precisely, let  $f(x) = e^{k-x} \cdot \left(\frac{x}{k}\right)^k$ , and

$$\beta_t(s) = k \ln(t^2 + t) + \underbrace{f'(s)}_{\approx (1+o(1)) \ln(\frac{1}{\delta})}$$

Lemma below  $\Rightarrow$  Track-and-Stop is sound.

Lemma (33.7 in text):  $\beta_t(s)$  as above, and

$$\tau = \min \{ t \geq k : Z_t \geq \beta_t(s) \}. \text{ Then}$$

$$P\left(i^*(\hat{\gamma}(\tau)) \neq i^*(\gamma)\right) \leq \delta$$

$\underbrace{\downarrow}_{\text{estimated best arm}}$        $\underbrace{\rightarrow}_{\text{true best arm in environment } \gamma}$

This also formally justifies intuition ① on the selection of the stopping time. To prove this, we need some concentration results.

Lemma or concentration (33.8 in text):  $\{X_t, t=1, 2, 3, \dots\}$   
iid  $N(\mu, 1)$  vs.  $\hat{\mu}_n = \frac{1}{n} \sum_{t=1}^n X_t$ . Then:

$$P\left(\exists n \geq 1 \text{ s.t. } (\hat{\mu}_n - \mu)^2 \geq \frac{2}{n} \left( \ln(\frac{1}{\delta}) + \ln(n^2) \right) \right)$$

$$\leq \delta.$$

Pf: Gaussian tail bound + Union bound. □

Prop. (33.9 in text):  $g: \mathbb{N} \rightarrow \mathbb{R}$  an increasing function.

further, define  $\{S_{i1}, S_{i2}, \dots\}$ ,  $i=1, 2, \dots, k$ , sequences of independent r.v.s. s.t. for any  $s \in (0, 1)$ , this sequence satisfies:

$$P\left(\exists r \in \mathbb{N} : S_{ir} \geq g(s) + \ln(\frac{1}{s})\right) \leq s$$

Then, for any  $x \geq 0$ , we have:

$$\begin{aligned} P\left(\exists \{\tau_1, \tau_2, \dots, \tau_k\} \in \mathbb{N}^k : \sum_{i=1}^k S_{ir_i} \geq kg\left(\sum_{i=1}^k \tau_i\right) + x\right) \\ \leq \left(\frac{x}{k}\right)^k e^{(k-x)} \end{aligned}$$

Proof:

Use the fact: If  $\{X_t, t \geq 1\}$  indep with  $P(X_t \leq a) \leq a$   $\forall a \in [0, 1]$ , then, for any  $\varepsilon > 0$ ,

$$P\left(\sum_{t=1}^n \ln\left(\frac{1}{X_t}\right) \geq \varepsilon\right) \leq \left(\frac{\varepsilon}{n}\right)^n e^{(n-\varepsilon)}$$

(Problem 5.16 in text).

Now, let  $W_i = \max\{w \in [0, 1] : S_{ir} \leq g(s) + \ln(\frac{1}{s}w) \text{ } \forall r = 1, 2, \dots\}$ ,

for  $i = 1, 2, \dots, k$ . Then for any  $\{r_1, r_2, \dots, r_k\}$  positive integers, we have

$$\sum_{i=1}^k S_i r_i \leq \sum_{i=1}^k g(r_i) + \sum_{i=1}^k \ln(1/w_i)$$

$$\leq k g\left(\sum_{i=1}^k r_i\right) + \sum_{i=1}^k \ln(1/w_i)$$

trivially, as  $g(\cdot)$  is increasing

Now,  $\{w_i, i=1, 2, \dots, k\}$  indep. and  $P(w_i \leq a) \leq a$   
 $\forall a \in [0, 1]$ . Thus, use fact above  $\square$

Back to Lemma (33.7); pp. 23 in Notes above:

Need to show that:  $P(i^*(\hat{\gamma}(\tau)) \neq i^*(\gamma)) \leq \delta$

First note that the Track-and-Stop algorithm will output a unique  $i^*(\hat{\gamma}(\tau))$ . Specifically, if there exists two arms that have  $\hat{\mu}_i(t) = \hat{\mu}_j(t)$  at some time  $t \geq 1$ , then  $Z_t = 0$ , and thus  $0 = Z_t \leq \beta_t(s)$ . Thus, we will not stop at this time.

WLOG, assume  $\Delta_1 = 0$  (i.e., arm 1 is the best arm),

Then, with  $\mu'_j = \mu_j'(\nu')$ , (i.e., the mean of arm  $j$   
in environment  $\nu'$ )

By the defn of  $(\bar{z}_t, z)$ , we have:

$$\bar{z}_t = \inf_{\nu' \in \Sigma_{alt}^k(\hat{\nu}(z))} \sum_{j=1}^k T_j(z) \frac{(\hat{\mu}_j(z) - \mu'_j)^2}{2} \geq \beta_t(s).$$

$$\Rightarrow \forall \nu' \in \Sigma_{alt}^k(\hat{\nu}(z)), \sum_{j=1}^k T_j(z) \frac{(\hat{\mu}_j(z) - \mu'_j)^2}{2} \geq \beta_t(s).$$

$\Rightarrow$

$$\left\{ \nu' \in \Sigma_{alt}^k(\hat{\nu}(z)) \right\} \subseteq \left\{ \nu \in \Sigma : \frac{1}{2} \sum_{j=1}^k T_j(z) (\hat{\mu}_j(z) - \mu'_j)^2 \geq \beta_t(s) \right\}$$

Next, observe that the error event:

$$\left\{ i \notin i^*(\hat{\nu}(z)) \right\} = \left\{ \nu \in \Sigma_{alt}^k(\hat{\nu}(z)) \right\}$$

true best arm  $\downarrow$   
 true environment

Combining above:

$$P(i \notin i^*(\hat{\nu}(z))) \leq P \left( \frac{1}{2} \sum_{j=1}^k T_j(z) (\hat{\mu}_j(z) - \mu'_j)^2 \geq \beta_t(s) \right)$$

true mean  $\downarrow$   
 $\tilde{\phantom{x}}^2$

Combining with the concentration results above, the lemma now follows.  $\square$

Theorem (33.6 in text):  $(\pi, \tau, \psi)$  according to the Track-and-Stop policy.  $\beta_t(s)$  chosen as above. Further  $\mathcal{E} = \mathcal{E}_N^k(\iota)$ , with  $|\iota^*(v)| = 1$ . (wLOG,  $\iota^*(v) = 1$ ). Then:

- i) Track-and-Stop is sound.
- ii)  $\lim_{\delta \rightarrow 0} \frac{\text{Ev}_{\pi}[\tau]}{\ln(1/\delta)} = c^*(v)$  (asympt. optimal).

Proof:

- i) Soundness proved in Lemma 33.7 (above).
- ii) (Sketch of proof):

$$Z_t = \inf_{v \in \mathcal{E}_{alt}^k(\hat{v}(t))} \frac{1}{2} \sum_{j=1}^k T_j(t) (\hat{u}_j(t) - u'_j)^2$$

$\approx (\alpha_j^*(v) \cdot t)$ , assuming that

we have sufficient samples s.t.  $\hat{v}(t) \approx v$

$$\approx t \inf_{v' \in \mathcal{E}_{alt}^k(\hat{v}(t))} \frac{1}{2} \sum_{j=1}^k \alpha_j^*(v) (\underbrace{u_j(v)}_{\hat{v}(t) \approx v} - \underbrace{u'_j}_{\text{assuming } \hat{u}_j(t) \approx u_j(v)})^2$$

$$= t \cdot (c^*(\tau))^{-1}$$

$\therefore$  At  $t = \tau$ ,

$$\frac{\tau}{c^*(s)} \approx \mathbb{E}_\tau \geq \beta_t(s) \approx 2k \ln \tau + \ln(\frac{1}{\delta})$$

from above intuition      ↑  
 Stopping Criterion      ↑

ignore as small wrt  $\tau$   
 +  $\delta (\ln(\tau/\delta))$   
 ignore as small wrt anything else

i.e., as  $s \rightarrow 0$ ,  $\frac{\tau}{c^*(s)} \approx \ln(\frac{1}{\delta})$

□

### Fixed Horizon Bandit: Best Arm Identification.

Setting: Similar as before, but stop at a fixed pre-determined (known) time and declare one of the arms as the best one.

Algorithm: Elimination  $\equiv$  Successive Halving Algorithm.

Time horizon ( $n$ ) divided into  $L$  rounds, where  $L = \lceil \log_2 k \rceil$ . In each round, eliminate half the arms.

Input:  $n, L = \lceil \log_2 k \rceil$ .

Initial:  $A_1 = \{1, 2, \dots, k\}$ .

Iteration: For each  $l = 1, 2, \dots, L$

a.  $T_l = \left\lfloor \frac{n}{L \cdot |A_l|} \right\rfloor$

b. Draw  $T_l$  samples for every arm  $j \in A_l$ .

c. Compute  $\hat{\mu}_j^l$  for each  $j \in l$  using the  $T_l$  samples from this round.

d.  $A_{l+1} \subseteq A_l : \text{top } \left\lceil \frac{|A_l|}{2} \right\rceil \text{ arms}$

Termination: Return sole surviving arm.

Theorem:  $\forall \epsilon \in \mathcal{E}_{SG}^k(1)$ , with  $M_1 \geq M_2 \geq \dots \geq M_k$ ,  
and  $\pi$  the algorithm above. Then:

$$P_{\pi}(\Delta_{A_{l+1}} > 0) \leq 3 \log_2 k \epsilon - n / (6 H_2(\mu) \log_2 k)$$

i.e. finds  $\nwarrow$  a best arm

$$\text{where } H_2(u) = \max_{\{i : \Delta_i > 0\}} \frac{i}{\Delta_i^2} .$$

①