

Lecture #8: Pure exploration

All previous Lectures: maximise cumulative reward
→ exploration/exploitation trade-off

In some applications, there is no price for exploring.

Think for example of a researcher testing drugs on mice/artificial human cells or testing products on some people before commercialisation.

Share similarities with regret minimisation, but good algorithms are actually different.

Setting (simple regret)

At each round $t=1, \dots, T$:

- pulls an arm $a_t \in [K]$
- observes $X_{a_t}(t) \sim \nu_{a_t}(\cdot)$

We explore for T rounds
and commit to best action
at time $T+1$.

Goal: minimise simple regret

$$R_T^{\text{exp}} = \mathbb{E} [\mu^* - \mu_{a_{T+1}}]$$

Algorithm: Uniform exploration

For $t=1, \dots, T$:

choose $a_t = 1 + (t \bmod K)$

$$A_{t+1} \in \arg\max_k \hat{\mu}_k(t).$$

Theorem: Uniform-Exploration satisfies for any $\nu \in P([0, 1])^K$

$$R_T^{\text{simple}} \leq \sum_{k, \Delta_k > 0} \Delta_k \exp\left(-\lfloor \frac{T}{K} \rfloor \Delta_k^2\right)$$

Proof

Let b such that $\Delta_b > 0$.

$$\Pr(\hat{\mu}_{k^*}(T) \leq \hat{\mu}_b(T)) = \Pr(\hat{\mu}_{k^*}(T) - \hat{\mu}_b(T) \leq 0)$$

the times where b and k^* are pulled are not random. We can directly apply Hoeffding inequality
 $\geq \lfloor \frac{T}{K} \rfloor$

Assume wlog. $N_{k^*}(T) \leq N_b(T)$

$$\Pr(\hat{\mu}_{k^*}(T) \leq \hat{\mu}_b(T)) \leq \Pr\left(\sum_{r=1}^T X_{k^*}(r) \mathbb{1}_{\{a_r=k^*\}} \leq \frac{N_{k^*}(T)}{N_b(T)} \sum_{r=1}^T X_b(r) \mathbb{1}_{\{a_r=b\}}\right)$$

$$\leq \Pr\left(\underbrace{\sum_{r=1}^T X_{k^*}(r) - \mu_{k^*}^* \mathbb{1}_{\{a_r=k^*\}}}_{N_b(T) \times \text{variable of range}} - \underbrace{\sum_{r=1}^T \frac{N_{k^*}(r)}{N_b(r)} (X_b(r) - \mu_b^*) \mathbb{1}_{\{a_r=b\}}}_{N_b(T) \times \text{variable of range}} \leq -N_b(T) \Delta_b\right)$$

$$\leq \exp\left(-\frac{2 N_b(T)^2 \Delta_b^2}{N_b(T) + \frac{N_{k^*}(T)}{N_b(T)}}\right)$$

$$\leq \exp(-N_b(T) \Delta_b) \leq \exp(-\lfloor \frac{T}{K} \rfloor \Delta_b^2)$$

$$R_T^{\text{simple}} = \sum_{k, \Delta_k > 0} \Delta_k \Pr(a_{t+1} = k)$$

$$\leq \sum_{k, \Delta_k > 0} \Delta_k \Pr(\hat{\mu}_{k^*}(T) \leq \hat{\mu}_b(T))$$

$$\leq \sum_{k, \Delta_k > 0} \Delta_k \exp(-\lfloor \frac{T}{K} \rfloor \Delta_k^2)$$

Theorem UE, distribution free bound

$$\text{for any } v \in \mathbb{R}(0, \Delta), R_T^{\text{simple}} \leq 2 \sqrt{2 \frac{\ln(K)}{T}}$$

Proof

Actually, we could have written for any $\tilde{\delta} \geq 0$

$$R_T^{\text{simple}} \leq \tilde{\delta} + \sum_{k, \Delta_k > \tilde{\delta}} \Delta_k \Pr(a_{t+1} = k)$$

$$\leq \tilde{\delta} + K \tilde{\delta} \exp\left(-\lfloor \frac{T}{K} \rfloor \tilde{\delta}^2\right)$$

for any $\tilde{\delta} \geq 0$.

$x \mapsto x \exp(-\lfloor \frac{T}{K} \rfloor x^2)$ is \nearrow on $[\sqrt{\frac{1}{4\lfloor \frac{T}{K} \rfloor}}, +\infty)$

Taking $\tilde{\delta} = \sqrt{\frac{\ln(K) v^{1/2}}{\lfloor T/K \rfloor}}$, we have:

$$R_T^{\text{simple}} \leq \tilde{\delta} + \tilde{\delta} K e^{-\frac{(\ln K) v^{1/2}}{\lfloor T/K \rfloor}} \leq 2\tilde{\delta}$$

$$\leq 2 \sqrt{\frac{\ln K v^{1/2}}{\lfloor T/K \rfloor}} = 2 \sqrt{\frac{\ln K}{\lfloor T/K \rfloor}}$$

if $T \leq K$, $R_T^{\text{simple}} \leq 1$ and the bound holds.

if $T > K$, $2 \lfloor \frac{T}{K} \rfloor \geq \frac{T}{K}$ so

$$R_T^{\text{simple}} \leq 2 \sqrt{2 \frac{K (\ln(K) v^{1/2})}{T}}$$

□

We can show that the minimax lower bound is larger than $c\sqrt{\frac{K}{T}}$, so Uniform-Exploration is nearly optimal in minimax sense.

can we do better than UE? i.e get rid of $\sqrt{ln K}$ term.

Best arm identification

Setting 1: (fixed confidence)

At each round $t = 1, \dots, \infty$:

- agent picks an arm $a_t \in [K]$ (based on previous observations)

- observes $X_{a_t}(t) \sim \nu_{a_t} \in \mathcal{D}$

- decides whether to continue sampling or stop

If stop: return a final choice $\psi \in [K]$

The (random) stopping time is called τ

new

with confidence level $\text{SE}(0, \delta)$

confidence
level

Goal: 1) Have a sound strategy: $P(\tau < \infty \text{ and } \mu_\psi < \mu^*) \leq \delta$ (for any δ)

2) minimize the exploration time $E[\tau]$

Our algorithm will be built on the following lower bound.

Theorem (lower bound)

Let (π, τ, ψ) be a sound strategy for the bandit model \mathcal{D} , with confidence level $\text{SE}(0, \delta)$.

and let $v \in \mathcal{D}^K$. Then:

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

$$c^*(v)^{-1} = \sup_{\alpha \in \mathbb{P}_K} \left(\inf_{v' \in \mathcal{D}_{\text{opt}}(v)} \sum_{k=1}^K \alpha_k \text{KL}(\nu_k, \nu'_{k'}) \right)$$

where $\mathcal{D}_{\text{opt}}(v) = \left\{ v' \in \mathcal{D}^K \mid \arg \max_k E(\nu_k) \cap \arg \max_k E(\nu'_{k'}) = \emptyset \right\}$ i.e. no arm is optimal for both v and v'

Another use of fundamental inequality (with stopping time)

Lemma: (admitted)

For all bandit problems $v = (v_k)_{k \in \mathbb{N}}$ and $v' = (v'_k)_{k \in \mathbb{N}}$ in \mathcal{D}^K with $v_k \ll v'_k$ for all k ,

for all strategies Π , for stopping time τ with respect to the filtration $(\mathcal{F}(H_t))$

and any random variable Z taking values in $[0, 1]$, $\mathcal{F}(H_\tau)$ -measurable,

$$\mathcal{F}(H_\tau) = \left\{ A \in \mathcal{F}(H_\infty) \mid A \cap \{\tau \leq t\} \in \mathcal{F}(H_t) \text{ for all } t \right\}$$

$$\sum_{k=1}^K \mathbb{E}_v[N_k(\tau)] \text{KL}(v_k, v'_k) \geq \text{KL}(\text{Bin}(\mathbb{E}_v[Z]), \text{Bin}(\mathbb{E}_{v'}[Z]))$$

Proof of the Theorem

Assume $\mathbb{E}[\tau] < \infty$ (otherwise the result holds), so that $P(\tau = \infty) = 0$.

Let $v' \in \mathcal{D}_{\text{act}}(v)$. We define the $\mathcal{F}(H_\tau)$ -measurable r.v.

$\bar{\tau} = \mathbb{1}_{\{\tau < \infty \text{ and } v \notin \arg\max_k \mathbb{E}(v'_k)\}}$. Then the fundamental inequality (with stopping time)

yields:

$$\begin{aligned} \sum_{k=1}^K \mathbb{E}_v[N_k(\tau)] \text{KL}(v_k, v'_k) &\geq \text{KL}(\text{Bin}(\mathbb{E}_v[Z]), \text{Bin}(\mathbb{E}_{v'}[Z])) \\ &> 1-\delta & \leq \delta & \text{as } \Pi \text{ is sound with confidence level } \delta \\ &\geq (1-\delta) \ln \left(\frac{1-\delta}{\delta} \right) + \delta \ln \left(\frac{\delta}{1-\delta} \right). \end{aligned}$$

$$= (1-\delta) \ln \left(\frac{1-\delta}{\delta} \right) \geq \ln \left(\frac{1}{4\delta} \right)$$

Let $\alpha_k = \frac{\mathbb{E}[N_{\alpha}(\tau)]}{\mathbb{E}_v[\tau]}$ $\alpha \in P_K$ and we have shown:

$$\mathbb{E}_v[\tau] \sum_{k=1}^K \alpha_k \text{KL}(v_k, v'_k) \geq \ln \left(\frac{1}{4\delta} \right).$$

for any $v' \in D_{\text{alt}}(v)$
and a specific $\alpha \in P_K$.
That is independent of v' .

i.e.: $\mathbb{E}_v[\tau]$ $\sup_{\alpha \in P_K} \inf_{v' \in D_{\text{alt}}(v)} \sum_{k=1}^K \alpha_k \text{KL}(v_k, v'_k) \geq \ln \left(\frac{1}{4\delta} \right)$

$c(v)$

□

Track-and-stop algorithm

Idea of the algorithm is to track the lower bound and stop when τ is larger than the estimated lower bound.

$\alpha_k^*(v)$ corresponds to the proportion of pulls on k .

i.e. we should stop when

$$\inf_{v' \in D_{\text{alt}}(v)} \sum_{k=1}^K N_k(\tau) \text{KL}(v_k, v'_k) \geq \ln \left(\frac{1}{\delta} \right).$$

z_r

Problem: v is unknown, but can be estimated.

Assume in the following $\mathcal{D} = \{N(\mu, 1) | \mu \in \mathbb{R}\}$.

In that case, we can approximate Z_t by

$$Z_t := \inf_{\mu' \in M_{\text{alt}}(\hat{\mu})} \sum_{k=1}^K N_k(t) (\hat{\mu}_k(t) - \mu')^2.$$

$$M_{\text{alt}}(\hat{\mu}) = \left\{ \mu' \mid \arg \max_{\mu'} \hat{\mu}_k \text{ marginal } \hat{\mu}_k = \mu' \right\}$$

Track-and-stop algorithm Input δ and $\beta_t(\delta)$

For $t=1, \dots, K$:

Pull $a_t = h$.

While $Z_t < \beta_t(\delta)$

| if $\min_k N_k(t) \leq \sqrt{\delta}$ then pull $a_{t+1} \in \arg \min_k N_k(t)$ forced exploration

| else choose $a_{t+1} \in \arg \max_{\mu'} \hat{F}_{\mu'}(t) - N_k(t)$ track

| stop and return $\psi \in \arg \max_{\mu'} \hat{\mu}_k(t)$ stop.

$$\hat{\alpha}(t) \in \arg \max_{\alpha \in \mathcal{P}_K} \inf_{\nu' \in D_{\text{alt}}(\hat{\mu}(t))} \sum_{k=1}^K \alpha_k \text{KL}(\hat{\nu}(t), \nu')$$

For our Gaussian setting,

Theorem There exists a choice $\beta_t(\delta)$ such that track-and-stop

is sound for the Gaussian setting and for any ν with a unique optimal arm:

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau]}{\ln(1/\delta)} = c(v)$$

Lemma: let $f: [K, +\infty) \rightarrow \mathbb{R}$ be given by

$$f(x) = \exp(K-x) \left(\frac{x}{K}\right)^K \text{ and } \beta_f(\delta) = K \ln(\delta t^2 + 1) + f^{-1}(\delta).$$

Then for $\bar{\tau} = \min\{t \mid \exists r \geq \beta_f(\delta)\}$, it holds $P(\arg\max_k \hat{\mu}_k(t) \neq \arg\max_k \mu_k) \leq \delta$.

Proof of the lemma:

If $|\arg\max_k \hat{\mu}_k(t)| > 1$, then $\bar{\tau} = 0$. So assume in the following $\bar{\tau} < \infty$ and

$$|\arg\max_k \hat{\mu}_k(\bar{\tau})| = 1.$$

denote ψ

By definition of $M_{\text{alt}}(\hat{\mu})$, for any $b^* \in \arg\max_k \mu_k$

$$P(b^* \neq \psi \text{ and } \bar{\tau} < \infty) = P(\mu_b \in M_{\text{alt}}(\hat{\mu}(\bar{\tau})) \text{ and } \bar{\tau} < \infty)$$

By def'n of $\bar{\tau}$, for any $\mu' \notin M_{\text{alt}}(\hat{\mu}(\bar{\tau}))$: $\frac{1}{2} \sum_{a=1}^K N_a(\bar{\tau}) (\hat{\mu}_a(\bar{\tau}) - \mu'_a)^2 \geq \beta_c(\delta)$, so

$$P(b^* \neq \psi \text{ and } \bar{\tau} < \infty) \leq P\left(\frac{1}{2} \sum_{a=1}^K N_a(\bar{\tau}) (\hat{\mu}_a(\bar{\tau}) - \mu_a)^2 \geq \beta_c(\delta)\right).$$

$$(\beta_f(\delta) = K \ln(\delta t^2 + 1) + f^{-1}(\delta))$$

We can show the two following concentration inequalities, which allow to conclude (see exercises)

$$P(\exists t \geq 1, \sum_{a=1}^K \frac{N_a(t)}{2} (\hat{\mu}_a(t) - \mu_a)^2 \geq K \ln(\delta t^2 + 1) + \alpha) \leq \left(\frac{\alpha}{K}\right)^K \exp(K-\alpha) = f(\delta)$$

]

Proof of the theorem

Define for any $\nu' \in \Delta^K$, when $\Delta = \{N(\mu, t) | \mu \in \mathbb{R}\}$.

$$\alpha^*(\nu') = \operatorname{argmax}_{\alpha \in P_K} \inf_{\tilde{\nu}' \in \Delta^K(\nu')} \sum_{k=1}^K \alpha_k \text{KL}(\nu'_k, \tilde{\nu}'_k).$$

We admit the following lemma:

Lemma

If ν' admits a unique optimal arm:

1) $\alpha^*(\nu')$ is unique

2) α^* is continuous at ν' .

Define the distance in Δ : $d(\nu, \nu') = \max_k (\mathbb{E}[\nu'_k] - \mathbb{E}[\nu_k])$,

$$\text{and } \phi(\nu, \alpha) = \frac{1}{2} \min_{k \neq k'} \frac{\alpha_k - \alpha_{k'}}{\alpha_k + \alpha_{k'}} \Delta_k^2.$$

Admit that $\alpha^*(\nu) = \operatorname{argmax}_{\alpha \in P_K} \Phi(\nu, \alpha)$ and Φ is concave in α .

Let $\varepsilon > 0$ be a small constant and define the random times

$$\tau_\nu(\varepsilon) = 1 + \max\{t \mid d(\hat{\nu}^t, \nu) \geq \varepsilon\}$$

$$\tau_\alpha(\varepsilon) = 1 + \max\{t \mid \|\alpha^*(\nu) - \alpha^*(\hat{\nu}^t)\|_\infty \geq \varepsilon\}$$

$$\tau_T(\varepsilon) = 1 + \max\{t \mid \|\alpha^*(\nu) - \frac{N(t)}{t}\|_\infty \geq \varepsilon\}.$$

Note these are not stopping times.

1) We are gonna use the first concentration inequality admitted in the proof of the Lemma that guarantees soundness of Track-and-Stop.

(a) Define the random variable:

$$\Lambda = \min\{\lambda \geq 1 \mid d(\hat{\nu}^t, \nu) \leq \sqrt{\frac{2 \ln(2\lambda K t(t+1))}{\min_k N_k(t)}} \text{ for all } t\}.$$

Show that $\mathbb{E}[\ln(\Lambda)^2] < \infty$.

(b) Prove that $\mathbb{E}[\tau_\nu(\varepsilon)] < \infty$ for all $\varepsilon > 0$.

2) Prove that $\mathbb{E}[\tau_\alpha(\varepsilon)] < \infty$ for all $\varepsilon > 0$.

3) Prove that $\mathbb{E}[\tau_T(\varepsilon)] < \infty$ for all $\varepsilon > 0$.

4)

(a) Define for any $\varepsilon > 0$

$$\tau_\beta(\varepsilon, \delta) = 1 + \max\{t \mid t\Phi(\nu, \alpha^*(\nu)) < \beta_t(\delta) + \varepsilon t\}$$

$$\text{and } u(\varepsilon) = \sup_{\nu', \alpha} \{\Phi(\nu, \alpha^*(\nu) - \Phi(\nu', \alpha)) \mid d(\nu', \nu) \leq \varepsilon, \|\alpha - \alpha^*(\nu)\|_\infty \leq \varepsilon\}.$$

Show that $\mathbb{E}[\tau] \leq \mathbb{E}[\tau_\nu(\varepsilon)] + \mathbb{E}[\tau_T(\varepsilon)] + \mathbb{E}[\tau_\beta(u(\varepsilon), \delta)]$.

(b) Conclude that $\lim_{\delta \rightarrow 0^+} \frac{\mathbb{E}[\tau]}{\ln(1/\delta)} \leq c^*(\nu)$.

Solution: 1) a) Define the random variable:

$$\Lambda = \min\{\lambda \geq 1 \mid d(\hat{\nu}^t, \nu) \leq \sqrt{\frac{2 \ln(2\lambda K t(t+1))}{\min_k N_k(t)}} \text{ for all } t\}.$$

The mentioned concentration inequality (in the hint) implies that $\mathbb{P}(\Lambda \geq x) \leq 1/x$. So that

$$\begin{aligned} \mathbb{E}[\ln(\Lambda)^2] &= \int_0^\infty \mathbb{P}(\Lambda \geq \exp(\sqrt{x})) dx \\ &\leq \int_0^\infty \exp(-\sqrt{x}) dx \\ &= 2 \int_0^\infty ue^{-u} du = 2. \end{aligned}$$

b) By definition of Λ ,

$$\tau_\nu(\varepsilon) \leq 1 + \max\{t \mid \sqrt{\frac{2 \ln(\Lambda) + 2 \ln(2Kt(t+1))}{\min_k N_k(t)}} > \varepsilon\}.$$

The forced exploration in the algorithm means that $N_k(t) \geq \frac{\sqrt{t}}{2}$ almost surely for t large enough. Hence,

$$\tau_\nu(\varepsilon) = \mathcal{O}(\ln(\Lambda)^2).$$

So finally $\mathbb{E}[\tau_\nu(\varepsilon)] < \infty$.

2) This is a consequence of 1) and the continuity of α^* at ν .

3) For t large enough ($t = \Omega(\tau_\alpha(\varepsilon/2K))$), the algorithm (tracking part dominates the forced exploration at some point) implies that

$$N_k(t) \leq t(\alpha_k^*(\nu) + \frac{\varepsilon}{2K}) + 1 \leq t(\alpha_k^*(\nu) + \frac{\varepsilon}{K}).$$

But since $\sum_k N_k(t) = t$, this also implies

$$N_k(t) \geq t(\alpha_k^*(\nu) - \varepsilon).$$

This thus allows to conclude as $\mathbb{E}[\tau_T(\varepsilon)] \leq \mathcal{O}(\mathbb{E}[\tau_\nu(\varepsilon/2K)])$.

4) a) When $t \geq \max\{\tau_\nu(\varepsilon), \tau_t(\varepsilon), \tau_\beta(u(\varepsilon), \delta)\}$,

$$\begin{aligned} tZ_t &= t\Phi(\hat{\nu}^t, \frac{N_k(t)}{t}) \\ &\geq t(\Phi(\nu, \alpha^*(\nu)) - u(\varepsilon)) \geq \beta_t(\delta). \end{aligned}$$

So $\tau \leq \max\{\tau_\nu(\varepsilon), \tau_t(\varepsilon), \tau_\beta(u(\varepsilon), \delta)\}$.

b) Taking the limit $\delta \rightarrow 0$, there only remains the last term (the pull choices do not depend on δ), i.e.

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau]}{\ln(1/\delta)} \leq \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\beta(u(\varepsilon), \delta)]}{\ln(1/\delta)}.$$

Since $\beta_t(\delta) = \frac{\ln(1/\delta)}{t} + o(\frac{\ln(1/\delta)}{t})$, we have $\tau_\beta(u(\varepsilon), \delta) = \frac{\ln(1/\delta)}{\Phi(\nu, \alpha^*(\nu)) - u(\varepsilon)} + o(\ln(1/\delta))$. Continuity of Φ at $(\nu, \alpha^*(\nu))$ ensures that $\lim_{\varepsilon \rightarrow 0} u(\varepsilon) = 0$ and the result then follows taking $\varepsilon \rightarrow 0$.

Setting 2: (fixed budget)

At each round $t=1, \dots, T$:

- agent picks an arm $a_t \in [K]$ (based on previous observations)
- observes $X_{a_t}(t) \sim \nu_{a_t} \in \mathcal{D}$

After round T , return $\psi \in [K]$.

Goal: minimize $P(\mu^* > \mu_\psi)$

→ much harder problem

Complexities: $H_1 = \sum_{k, \Delta_k > 0} \frac{1}{\Delta_k}$, $H_2 = \max_{k, \Delta_k > 0} \frac{K}{\Delta_k}$

standard values: $\Delta_{(1)} \leq \Delta_{(2)} \leq \dots \leq \Delta_{(K)}$

$$H_2 \leq H_1 \leq \ln(2K) H_2$$

Lower bound of order $\begin{cases} \exp(-\frac{T}{\ln(K) H_2}) & \text{if } H_2 \text{ unknown} \\ \exp(-\frac{T}{H_2}) & \text{if } H_2 \text{ known} \end{cases}$

First approach: uniform exploration of the arms. Good baseline, but not very good when arms have very different means.

Sequential Halving:

Set $L = \lceil \log_2(K) \rceil$ and $A_1 = [K]$

For $l = 1, \dots, L$:

Pull each arm in A_l $T_l = \lfloor \frac{T}{|A_l|} \rfloor$ times

Let $\hat{\mu}_i^l$ be the empirical mean of arm i based only on these last T_l samples

Let A_{l+1} contain the top $\lceil \frac{|A_l|}{2} \rceil$ arms in A_l

Return ψ as the only arm in A_{L+1}

Then:

If the distributions are 1-sub-Gaussian, then Sequential Halving satisfies

$$P(\mu^+ > \mu_\psi) \leq 3 \log_2(K) \exp\left(-\frac{T}{16H_2 \log(K)}\right)$$

Remarks

- close to lower bound
- For uniform exploration, we can bound this probability by
$$\sum_{k: \Delta_k > 0} \exp\left(-\frac{\mathbb{E}[k] \Delta_k^2}{4}\right).$$
- VE slightly better than SH when $\Delta_a = \Delta$ for any a ($H_2 = \frac{K}{\Delta^2}$)
but SH much better than VE when $\Delta_2 = \Delta \ll 1$ ($H_2 = \frac{1}{\Delta^2}$)
 $\Delta_a = 1$ for $a \geq 2$