

Lecture #9: 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 $X_{a_t}(t) \sim \nu_{a_t} \in \mathcal{D}$

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

Goal: minimise simple regret

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

Algorithm: Uniform exploration

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

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

$A_{T+1} \in \operatorname{argmax}_k \hat{\mu}_k(T)$.

Theorem: Uniform-Exploration satisfies for any $v \in \mathcal{P}([0, 1])^k$

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

Proof Let k such that $\Delta_k > 0$.

$$\mathbb{P}(\hat{\mu}_{a^*}(T) \leq \hat{\mu}_k(T)) = \mathbb{P}(\hat{\mu}_{a^*}(T) - \hat{\mu}_k(T) \leq 0)$$

$N_{a^*}(T)$ and $N_k(T)$ are not random. We can directly apply Hoeffding inequality $\triangleright \lfloor \frac{T}{k} \rfloor$

$$\begin{aligned} \mathbb{P}(\hat{\mu}_{a^*}(T) \leq \hat{\mu}_k(T)) &\leq \exp(-2(N_k(T) + N_{a^*}(T)) \Delta_k^2) \\ &\leq \exp(-2 \lfloor \frac{T}{k} \rfloor \Delta_k^2). \end{aligned}$$

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

$$\leq \sum_{k, \Delta_k > 0} \Delta_k \mathbb{P}(\hat{\mu}_{a^*}(T) \leq \hat{\mu}_k(T))$$

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

Theorem UE, distribution free bound

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

Proof

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

$$R_T^{\text{simple}} \leq \tilde{\Delta} + \sum_{k, \Delta_k > \tilde{\Delta}} \Delta_k \mathbb{P}(a_{T+1} = k)$$

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

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

Taking $\tilde{\Delta} = \sqrt{\frac{\ln(K)+v}{2 \lfloor \frac{T}{K} \rfloor}}$, we have:

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

$$\leq \sqrt{2 \frac{\log K + v}{\lfloor \frac{T}{K} \rfloor}}$$

if $T \leq 2K$, $R_T^{\text{simple}} \leq 1$ and the bound holds.

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

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

□

(see exercise section # 6)

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.

Reduction from cumulative regret.

For any strategy $(\pi_t)_{t=1, \dots, T}$; we can define $\tilde{\pi}$ s.t.

$$\tilde{\pi}_t = \pi_t \text{ for } t=1, \dots, T$$

$$P_{\tilde{\pi}}(a_{t+1}=k | \mathcal{F}_T) = \frac{N_k(T)}{T}.$$

Proposition for any instance v , $R_T^{\text{simple}}(v, \tilde{\pi}) = \frac{R_T(v, \tilde{\pi})}{T}$.

Proof

$$R_T^{\text{simple}}(v, \tilde{\pi}) = \sum_{k=1}^K P_{\tilde{\pi}}(a_{T+1}=k) \Delta_k$$

$$= \sum_{k=1}^K E[P_{\tilde{\pi}}(a_{T+1}=k) | \mathcal{F}_T] \Delta_k = \frac{1}{T} \sum_{k=1}^K \Delta_k E[N_k(T)] = \frac{R_T(v, \tilde{\pi})}{T} \quad \square$$

Corollary: there exists a strategy with a simple regret $R_T^{\text{simple}} \leq c \sqrt{\frac{K}{T}}$ for a universal constant c . (Reduction 11.6.8 with $\tilde{\pi}$)

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

- Goal:
- 1) Have a **sound** strategy: $\mathbb{P}(\tau < \infty \text{ and } \mu_\Psi < \mu^*) \leq \delta$ (for any $\delta \in (0, 1)$)
with confidence level $\delta \in (0, 1)$ confidence level
 - 2) Minimize the exploration time $\mathbb{E}[\tau]$

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 $\delta \in (0, 1)$, and let $\nu \in \mathcal{D}^K$. Then:

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

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

$$\text{where } \mathcal{D}_{\text{det}}(\nu) = \left\{ \nu' \in \mathcal{D}^K \mid \operatorname{argmax}_k \mathbb{E}[\nu_k] \cap \operatorname{argmax}_k \mathbb{E}[\nu'_k] = \emptyset \right\}$$

i.e. no arm is optimal for both ν and ν'

Another use of fundamental inequality (with stopping time)

Lemma: (admitted)

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

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

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

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

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

Proof of the Theorem

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

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

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

yields:

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

$\geq 1-\delta$

$\leq \delta$

as Π is sound with confidence level δ

$$\geq (1-\delta) \ln\left(\frac{1-\delta}{\delta}\right) + \delta \ln\left(\frac{\delta}{1-\delta}\right).$$

$$= (1-2\delta) \ln\left(\frac{1-\delta}{\delta}\right) \geq \ln\left(\frac{1}{4\delta}\right)$$

Let $\alpha_k = \frac{\mathbb{E}_v[N_k(\tau)]}{\mathbb{E}_v[\tau]}$ $\alpha \in \mathcal{P}_K$ 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 $v' \in \mathcal{D}_{\text{alt}}(v)$
and a specific $\alpha \in \mathcal{P}_K$
that is independent of v'

$$\text{so: } \mathbb{E}_v[\tau] \sup_{\alpha \in \mathcal{P}_K} \inf_{v' \in \mathcal{D}_{\text{alt}}(v)} \sum_{k=1}^K \alpha_k \text{KL}(v_k, v'_k) \geq \ln\left(\frac{1}{4\delta}\right)$$

$$\underline{c(v)}$$

□

Stop-and-track 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

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

τ_r

Problem: v is unknown, but can be estimated.

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

In that case, we can approximate Z_t by

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

$$\mathcal{M}_{\text{alt}}(\hat{\mu}) = \left\{ \mu' \mid \arg\max_k \hat{\mu}_k \neq \arg\max_k \mu'_k = \beta \right\}$$

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

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

Pull $a_t = t$.

While $Z_t < \beta_T(\delta)$

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

else choose $a_{t+1} \in \arg\max_k \hat{\mu}_k(t) - N_k(t)$ track

stop and return $\psi \in \arg\max_k \hat{\mu}_k(t)$ stop.

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

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:

$$\lim_{\delta \rightarrow 0} \frac{E[\tau]}{\ln(1/\delta)} = c^*(\nu).$$

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

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

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

Sketch of proof of the Theorem

(detailed proof in
exercise session #6)

Taking $\beta_r(\delta)$ as in the lemma, the strategy is sound.

The idea of the proof is to show that we have: $\hat{\mu}_k(t) \approx \mu_k$ and $\frac{N_k(t)}{t} \approx \hat{\alpha}_k(t) \approx \alpha_k^*(\nu)$

From there:

$$Z_t = \frac{1}{2} \inf_{\mu \in \mathcal{H}_t(\hat{\mu}(t))} \sum_{k=1}^K N_k(t) (\hat{\mu}_k(t) - \mu_k)^2$$

$$\approx t \inf_{\mu \in \mathcal{H}_t(\mu)} \sum_{k=1}^K \frac{\alpha_k^*(\nu)}{2} (\mu_k - \mu_k)^2$$

$$= \frac{t}{c^*(\nu)}$$

The algorithm should then stop when $\frac{t}{c^*(\nu)} \geq \beta_r(\delta) = (1 + o(1)) \ln(1/\delta)$.

$$\text{ie } \tau \sim \ln(1/\delta)$$

Proof of the Lemma:

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

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

define Ψ

By definition of $\mathcal{M}_{\text{alt}}(\hat{\mu})$, for any $k^* \in \arg\max_k \mu_k$

$$\mathbb{P}(k^* \neq \Psi \text{ and } \tau < \infty) = \mathbb{P}(\mu \in \mathcal{M}_{\text{alt}}(\hat{\mu}(\tau)) \text{ and } \tau < \infty)$$

By defn of Z_t , for any $\mu' \in \mathcal{M}_{\text{alt}}(\hat{\mu}(\tau))$: $\frac{1}{2} \sum_{k=1}^K N_k(\tau) (\hat{\mu}_k(\tau) - \mu'_k)^2 \geq \beta_c(\delta)$, so

$$\mathbb{P}(k^* \neq \Psi \text{ and } \tau < \infty) \leq \mathbb{P}\left(\frac{1}{2} \sum_{k=1}^K N_k(\tau) (\hat{\mu}_k(\tau) - \mu_k)^2 \geq \beta_c(\delta)\right). \quad (\beta_c(\delta) = K \ln(t^2 + t) \psi^{-1}(\delta))$$

We can show the two following concentration inequalities, which allow to conclude

$$(1) \mathbb{P}(\exists t \geq 1, \frac{N_k(t)}{2} (\hat{\mu}_k(t) - \mu_k)^2 \geq \ln(\frac{1}{\delta}) + \ln(t^2 + t)) \leq \delta$$

$$(2) \mathbb{P}(\exists t \geq 1, \sum_{k=1}^K \frac{N_k(t)}{2} (\hat{\mu}_k(t) - \mu_k)^2 \geq K \ln(t^2 + t) + \alpha) \leq \left(\frac{\alpha}{K}\right)^K \exp(K - \alpha)$$

Proofs are omitted
Easier to prove in the alternative
model where $Y(t) = X_{a_t}(N_{a_t}(t))$.

□

Setting 2: (fixed budget)

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

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

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

Goal: maximize $\mathbb{P}(\Psi \in \arg\max_k \mu_k)$

→ much border problem