

Stochastic bandits: setting and UCB strategy

Stochastic bandits.Finitely many arms.Setting:K arms indexed by $1, 2, \dots, K$ With each arm j is associated a probability distribution ν_j
(over \mathcal{R})
with an expectationAt each round $t = 1, 2, \dots$

- The decision-maker picks $I_t \in \{1, \dots, K\}$, possibly at random
- She gets a reward Y_t drawn at random according to ν_{I_t} (given I_t)
- This is the only feedback she gets / the only observation she has access to.

Aim:We denote by $\mu_j = E(\nu_j)$ the expectation of ν_j (note: operator E vs. expectation E of an expression involving random variables.)Pseudo-regret $\bar{R}_T = T\mu^* - E\left[\sum_{t=1}^T Y_t\right]$ to be controlledwhere $\mu^* = \max_{j \in K} \mu_j$ Useful notation: $\Delta_a = \mu^* - \mu_a$ gap of arm a $\Delta_a = 0$: a is an optimal arm (there can be several of them) $\Delta_a > 0$: a is a suboptimal arm $N_a(T) = \sum_{t=1}^T \mathbb{1}_{I_t = a}$ total number of times that a is pulled.Note:* Pseudo regret \bar{R}_T is a very "expected" notion of regret

$$\bar{R}_T \leq \text{probably } E\left[\max_{a=1, \dots, K} \sum_{t=1}^T \dots - \sum_{t=1}^T Y_t\right]$$

* Can be rewritten (see later) as $\bar{R}_T = \sum_{a=1}^K \Delta_a E[N_a(T)]$

Upper confidence bound [UCB] algorithm:

very popular!

For $t = 1, 2, \dots, K$

- Pull arm $I_t = t$, get a reward Y_t

For $t = K+1, K+2, \dots$

- Pull an arm $I_t \in \operatorname{argmax}_{j \in \{1, \dots, K\}} \left\{ \hat{\mu}_{j|t-1} + \sqrt{\frac{2 \ln t}{N_j(t-1)}} \right\}$

(tie-breaking rule: pick the element with smallest index)

where $N_j(t-1) = \sum_{s=1}^{t-1} \mathbb{1}_{\{I_s=j\}}$

and where $\hat{\mu}_{j|t-1} = \frac{1}{N_j(t-1)} \sum_{s=1}^{t-1} Y_s \mathbb{1}_{\{I_s=j\}}$

always ≥ 1 since each arm was tried sequentially during rounds $1, 2, \dots, K$

- Get a reward Y_t

Theorem:

If the distributions ν_j have supports all included in $[a, b]$, then the pseudo-regret of UCB is smaller than

$$\bar{R}_T \leq \sum_{i: \Delta_i > 0} \left(\frac{8 \ln T}{\Delta_i} + 2 \right)$$

This regret bound is obtained via the following proposition:

Proposition:

If the distributions ν_j have supports all included in $[a, b]$, then

$$\forall i \text{ s.t. } \Delta_i > 0, \quad \mathbb{E}[N_i(T)] \leq \frac{8 \ln T}{\Delta_i^2} + 2.$$

Exercise

The bounds above are called distribution-dependent because they depend heavily on the distributions ν_i at hand (via the gaps $\Delta_i = \mu^* - \mu_i$).

Show the following distribution-free bound (that only

depends on the support $[a, b]$, not on the specific distributions ν_i at hand): for the UCB algorithm,

$$\sup_{\nu_1, \dots, \nu_k \text{ with supports in } [a, b]} \bar{R}_T \leq O(\sqrt{TK \ln T}).$$

Hint: For small values of Δ_i , the bound of the Proposition can be worse than the trivial T bound...

Proof [of the theorem based on the Proposition]:

$$\bar{R}_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T y_t\right]$$

where by definition of the bandit model, \leftarrow Given I_t , y_t is drawn at random according to ν_{I_t}

$$\mathbb{E}[y_t | I_t] = \mu_{I_t}$$

thus (by the tower rule) $\mathbb{E}[y_t] = \mathbb{E}[\mu_{I_t}]$

$$= \sum_j \mu_j \mathbb{E}[\mathbb{1}_{\{I_t=j\}}]$$

Summing over t : $\mathbb{E}\left[\sum_{t=1}^T y_t\right] = \sum_{j=1}^k \mu_j \mathbb{E}[N_j(T)]$

and (in view of $T = \sum_j \mathbb{E}[N_j(T)]$)

$$\begin{aligned} \bar{R}_T &= \sum_j (\mu^* - \mu_j) \mathbb{E}[N_j(T)] = \sum_{j=1}^k \Delta_j \mathbb{E}[N_j(T)] \\ &= \sum_{j: \Delta_j > 0} \Delta_j \mathbb{E}[N_j(T)] \end{aligned}$$

\leftarrow it suffices to consider the suboptimal arms...

We conclude by substituting $\mathbb{E}[N_j(T)] \leq \frac{8 \ln T}{\Delta_j^2} + 2$ and by bounding $2\Delta_j \leq 2$.

NOTE: Keep in mind the rewriting as we will often use it!

$$\begin{aligned} \bar{R}_T &= T\mu^* - \mathbb{E}\left[\sum_{t=1}^T y_t\right] \\ &= \sum_{a=1}^k \Delta_a \mathbb{E}[N_a(T)] \end{aligned}$$

Proof [of the Proposition]:

We fix an optimal arm $a^* \in \{1, \dots, K\}$,
 i.e. s.t. $\mu_{a^*} = \mu^*$.

→ It will show why this algorithm is called UCB:

Because $\hat{\mu}_{j,t-1} + \sqrt{\frac{2 \ln t}{N_j(t-1)}}$ will indeed appear as an upper confidence bound on μ_j

estimate based on the raw performance
 ↔ exploitation of the results

larger for arms not much sampled so far
 ↔ forces some exploration

The UCB algorithm realizes some compromise / trade off between exploitation & exploration.

Later on we compare these statements to the Hoeffding - Azuma inequality

on $\{a_j\}$

LEMMA:

$\forall j, \forall t \geq j$ (so that $N_j(t) \geq 1$)

$\forall \delta \in (0, 1)$,

$$\mathbb{P} \left\{ \mu_j > \hat{\mu}_{j,t} - \sqrt{\frac{\ln(1/\delta)}{2 N_j(t)}} \right\} \geq 1 - \delta.$$

You can by $\mu_j = 1 - \mu_{a^*}$ then $\mu_{a^*} = 1 - \mu_j$ as well and μ_{a^*} supported on $\{a_j\}$

replace

By symmetry: $\forall \delta \in (0, 1)$,

$$\mathbb{P} \left\{ \mu_j < \hat{\mu}_{j,t} + \sqrt{\frac{\ln(1/\delta)}{2 N_j(t)}} \right\} \geq 1 - \delta$$

→ Application:
$$N_i(T) = 1 + \sum_{t=K+1}^T \mathbb{1}_{\left\{ \frac{1}{2} \mathbb{I}_{\mathbb{E}=i} \right\}}$$

We show below that $t \geq K+1$ and $\mathbb{I}_{\mathbb{E}=i}$ entails one of the following:

(i)
$$\hat{\mu}_{i,t-1} > \mu_i + \sqrt{\frac{2 \ln t}{N_i(t-1)}} \quad [\mu_i < \text{lower confidence bound}]$$

(ii)
$$\hat{\mu}_{a^*,t-1} < \mu^* - \sqrt{\frac{2 \ln t}{N_{a^*}(t-1)}} \quad [\mu^* > \text{upper confidence bound}]$$

(iii)
$$N_i(t-1) \leq \frac{8 \ln t}{\Delta^2} \quad [i \text{ not played often yet}]$$

Indeed, we would otherwise have

$$\begin{aligned} \hat{\mu}_{a^*, t-1} + \sqrt{\frac{2 \ln t}{N_{a^*}(t-1)}} &\geq \mu^* && \text{negation of (ii)} \\ &= \mu_i + \Delta_i && \text{definition of } \Delta_i \\ &> \mu_i + 2 \sqrt{\frac{2 \ln t}{N_i(t-1)}} && \left\{ \begin{array}{l} \text{the negation of (iii)} \\ \text{is } \Delta_i^2 > 8 \ln t / N_i(t-1) \end{array} \right. \\ &\geq \hat{\mu}_{i, t-1} + \sqrt{\frac{2 \ln t}{N_i(t-1)}} && \text{negation of (i)} \end{aligned}$$

the inequality between these two quantities would contradict $i \in \arg \max_j \{ \hat{\mu}_{j, t} + \sqrt{2 \ln t / N_j(t-1)} \}$ if $I_t = i$. (that is, $\hat{\mu}_{i, t-1} + \sqrt{2 \ln t / N_i(t-1)} > \hat{\mu}_{a^*, t-1} + \sqrt{2 \ln t / N_{a^*}(t-1)}$)

Thus,
$$E[N_i(T)] \leq 1 + \sum_{t=K+1}^T \mathbb{P} \left(\hat{\mu}_{i, t-1} > \mu_i + \sqrt{\frac{2 \ln t}{N_i(t-1)}} \right) + \sum_{t=K+1}^T \mathbb{P} \left(\hat{\mu}_{a^*, t-1} < \mu^* - \sqrt{\frac{2 \ln t}{N_{a^*}(t-1)}} \right)$$

each $\leq t \delta$ where $\delta = 1/4$

$$\leq 1 + 2 \sum_{t=K+1}^T t^{-3} + E \left[\sum_{t=K+1}^T \mathbb{1}_{\{I_t = i \wedge N_i(t-1) \leq 8 \ln t / \Delta_i^2\}} \right] + E \left[\sum_{t=K+1}^T \mathbb{1}_{\{N_i(t-1) \leq 8 \ln t / \Delta_i^2 \wedge I_t \neq i\}} \right]$$

$\frac{8 \ln t}{\Delta_i^2} \leq 8 \ln T$

$$\leq 1 + 2 \sum_{t=K+1}^T t^{-3} + \left(\frac{8 \ln T}{\Delta_i^2} + 1 \right) - 1$$

deterministically upper bounded by $\left(\frac{8 \ln T}{\Delta_i^2} + 1 \right) - 1$
 as $I_t = i$ only if $N_i(t-1) \leq \frac{8 \ln t}{\Delta_i^2}$
 thus only if $N_i(t) \leq \frac{8 \ln T}{\Delta_i^2} + 1$
 so that the total sum $\sum_{s=1}^t \mathbb{1}_{\{I_s = i\}} = N_i(t)$ is controlled by this number
 -1 because $I_t = i$ is not included in the $\sum_{t=K+1}^T \dots$

Thus:

$$E[N_i(T)] \leq \frac{8 \ln T}{\Delta_i^2} + 2$$

Let's complete the proof of the Lemma:

["Hoeffding-Azuma inequality with a random number of Summands"]

Setting: Probability distributions ν_1, \dots, ν_K over $[0,1]$
with respective expectations μ_1, \dots, μ_K

At each round, $I_t \in \{1, \dots, K\}$ is picked in a $\sigma(y_1, \dots, y_{t-1})$ -measurable way

then Y_t is drawn independently at random according to ν_{I_t} , given I_t
ie: $Y_t | I_t \sim \nu_{I_t}$

We denote $N_a(t) = \sum_{s=1}^t \mathbb{1}_{I_s=a}$ and assume that each arm a was pulled once in the first K rounds,
so that: $N_a(t) \geq 1 \quad \forall t \geq K$

Then, for $t \geq K$: $\hat{\mu}_{a,t} = \frac{1}{N_a(t)} \sum_{s=1}^t Y_s \mathbb{1}_{I_s=a}$

Lemma: $\forall \delta \in (0,1)$, $\mathbb{P} \left\{ \mu_a > \hat{\mu}_{a,t} - \sqrt{\frac{\ln(1/\delta)}{2N_a(t)}} \right\} \geq 1 - \delta$
(and a similar symmetric statement with $\mu_a < \hat{\mu}_{a,t} + \sqrt{\dots}$)

The proof will be based on the fact that $(Z_t)_{t \geq 0}$, where

$$Z_t = \sum_{s=1}^t (Y_s - \mu_a) \mathbb{1}_{I_s=a}$$

is a martingale w.r.t. $(\mathcal{F}_t)_{t \geq 0} = (\sigma(y_1, \dots, y_t))_{t \geq 0}$

$$\begin{aligned} \mathbb{E} \left[(Y_t - \mu_a) \mathbb{1}_{I_t=a} \mid \mathcal{F}_{t-1} \right] &= \mathbb{E} \left[(Y_t - \mu_a) \mathbb{1}_{I_t=a} \mid I_t, y_1, \dots, y_{t-1} \right] \\ &= \left(\mathbb{E}[Y_t \mid I_t, y_1, \dots, y_{t-1}] - \mu_a \right) \mathbb{1}_{I_t=a} \\ &\stackrel{\text{where we used the bandit model}}{=} \left(\mu_{I_t} - \mu_a \right) \mathbb{1}_{I_t=a} = 0 \text{ a.s.} \end{aligned}$$

Remark: How does this bound compare to what the classical version of the Hoeffding-Azuma says?

Martingale increment $(Y_s - \mu_a) \mathbb{1}_{I_s=a}$ bounded between

$$a_t = -\mu_a \text{ and } b_t = 1 - \mu_a$$

so that (actually in the version I stated, I can have \leq or $<$)

$$(b_t - a_t)^2 = 1$$

$$1 - \epsilon \leq \mathbb{P} \left\{ Z_t < \sqrt{\frac{t}{2}} \sqrt{\frac{\ln \frac{1}{\epsilon}}{t \delta}} \right\} = \mathbb{P} \left\{ N_t(t) (\hat{\mu}_{\text{opt}} - \mu_a) < \sqrt{\frac{t}{2}} \sqrt{\frac{\ln \frac{1}{\epsilon}}{t \delta}} \right\}$$

$$= \mathbb{P} \left\{ \hat{\mu}_{\text{opt}} - \sqrt{\frac{t}{N_t(t)}} \sqrt{\frac{\ln \frac{1}{\epsilon}}{2 N_t(t)}} < \mu_a \right\}$$

versus the bound of our lemma: $1 - \epsilon \leq \mathbb{P} \left\{ \hat{\mu}_{\text{opt}} - \sqrt{\frac{\ln \frac{1}{\epsilon}}{2 N_t(t)}} < \mu_a \right\}$

The proposed deviation essentially differ from a $\sqrt{\frac{t}{N_t(t)}}$ factor, and it is so nice to get rid of it!

Proof: (1) We prove that $\forall x \in \mathbb{R}, \mathbb{E} \left[e^{x Z_t - \frac{x^2}{8} N_t(t)} \right] \leq 1$

We do so by showing that $M_t = \exp \left(x Z_t - \frac{x^2}{8} N_t(t) \right)$ is a supermartingale, so that $\mathbb{E}[M_t] \leq \mathbb{E}[M_0] = 1$.

Indeed, by the conditional version of Hoeffding's lemma,

$$\mathbb{E} \left[e^{x (Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] \leq e^{x^2/8} \quad \text{a.s.} \quad \left. \begin{array}{l} \text{but we} \\ \text{can do} \\ \text{better!} \end{array} \right\}$$

Since \mathcal{I}_t and thus also $\mathbb{1}_{\mathcal{I}_t = a_j}$ are \mathcal{F}_{t-1} -measurable, we get:

$$\mathbb{E} \left[e^{x (Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] = \mathbb{E} \left[e^{x (Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} (\mathbb{1}_{\mathcal{I}_t = a_j} + \mathbb{1}_{\mathcal{I}_t \neq a_j}) \mid \mathcal{F}_{t-1} \right]$$

$$= \mathbb{E} \left[e^{x (Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] \mathbb{1}_{\mathcal{I}_t = a_j} + e^0 \mathbb{1}_{\mathcal{I}_t \neq a_j}$$

given what we had before \downarrow

$$\leq e^{x^2/8} \mathbb{1}_{\mathcal{I}_t = a_j} + \mathbb{1}_{\mathcal{I}_t \neq a_j} = \exp \left(\frac{x^2}{8} \mathbb{1}_{\mathcal{I}_t = a_j} \right)$$

Put differently,

$$\mathbb{E} \left[e^{x (Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j} - \frac{x^2}{8} \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] \leq 1$$

which entails that

$$\exp \left(x \sum_{s=1}^t (Y_s - \mu_a) \mathbb{1}_{\mathcal{I}_s = a_j} - \frac{x^2}{8} \sum_{s=1}^t \mathbb{1}_{\mathcal{I}_s = a_j} \right)$$

$$= \exp \left(x Z_t - \frac{x^2}{8} N_t(t) \right) = M_t$$

is a supermartingale wrt $\mathcal{F}_t = \sigma(Y_1, \dots, Y_t)$.

$$(2) \text{ We prove that } \forall \varepsilon > 0, \forall \ell \geq 1, \mathbb{P}\{Z_\ell \geq \varepsilon \text{ and } N_\ell(t) = \ell\} \leq \exp(-2\varepsilon^2/\ell)$$

Indeed, by a Markov-Chernoff bounding,

$$\begin{aligned} \forall x > 0, \quad \mathbb{P}\{Z_\ell \geq \varepsilon \text{ and } N_\ell(t) = \ell\} &\leq e^{-x\varepsilon} \mathbb{E}\left[e^{xZ_\ell} \mathbb{1}_{\{N_\ell(t) = \ell\}}\right] \\ &= e^{-x\varepsilon + \frac{x^2\ell}{8}} \mathbb{E}\left[e^{xZ_\ell - \frac{x^2}{8}N_\ell(t)} \mathbb{1}_{\{N_\ell(t) = \ell\}}\right] \\ &\leq e^{-x\varepsilon + \frac{x^2\ell}{8}} \underbrace{\mathbb{E}\left[e^{xZ_\ell - \frac{x^2}{8}N_\ell(t)}\right]}_{\leq 1 \text{ by (1)}} \end{aligned}$$

Optimizing over $x > 0$

(take $x = 4\varepsilon/\ell$) yields the claimed bound.

$$(3) \text{ Conclusion: we prove that } \mathbb{P}\left\{\mu_n \leq \hat{\mu}_{nt} - \sqrt{\frac{\ln(1/\delta)}{2N_n(t)}}\right\} \leq t\delta$$

Indeed, by distinguishing according to the values taken by $N_n(t)$:

$$\begin{aligned} &\mathbb{P}\left\{\mu_n \leq \hat{\mu}_{nt} - \sqrt{\frac{\ln(1/\delta)}{2N_n(t)}}\right\} \\ &= \sum_{\ell=1}^t \mathbb{P}\left\{N_n(t) = \ell \text{ and } \mu_n \leq \hat{\mu}_{nt} - \sqrt{\frac{\ln(1/\delta)}{2\ell}}\right\} \\ &= \sum_{\ell=1}^t \mathbb{P}\left\{N_n(t) = \ell \text{ and } \frac{Z_\ell}{N_n(t)} \geq \sqrt{\frac{\ln(1/\delta)}{2\ell}}\right\} \\ &= \sum_{\ell=1}^t \mathbb{P}\left\{N_n(t) = \ell \text{ and } Z_\ell \geq \sqrt{\ell \ln(1/\delta)/2}\right\} \\ \text{by (2)} \quad &\leq \sum_{\ell=1}^t \exp(-2(\ell \ln(1/\delta)/2)/\ell) = t\delta. \end{aligned}$$

($\sum_{\ell=1}^{t-K+1}$ would be enough)

Two notes on this proof:

- * We noted last week that the conditional version of Hoeffding's lemma could be generalized into

X bounded random variable, U, V two \mathcal{G}_j -measurable random variables
with $U \leq X \leq V$

then $\forall s \in \mathbb{R}$, $\ln \mathbb{E}[e^{sX} | \mathcal{G}_j] \leq s \mathbb{E}[X | \mathcal{G}_j] + \frac{s^2}{8} (V-U)^2$

(Because of the proof technique relying on weights with denominators $\frac{1}{V-U}$, it is safer to assume $V-U > 0$ as, eg, by replacing U and V by $U-\varepsilon$ and $V+\varepsilon$ if needed and letting $\varepsilon \rightarrow 0$; so at the end of the day we may drop the $V-U > 0$ as condition.)

* This extension may be applied to

$$\begin{aligned} X_t &= (Y_t - \mu_a) \mathbb{1}_{\{I_t = a\}} \\ \mathcal{G}_j &= \mathcal{F}_{t-1} \\ U_t &= -\mu_a \mathbb{1}_{\{I_t = a\}} \\ V_t &= (1 - \mu_a) \mathbb{1}_{\{I_t = a\}} \end{aligned}$$

and directly entails

$$\mathbb{E}\left[e^{x(Y_t - \mu_a) \mathbb{1}_{\{I_t = a\}}} \middle| \mathcal{F}_{t-1} \right] \leq \exp\left(\frac{s^2}{8} \mathbb{1}_{\{I_t = a\}} \right)$$

without the need for the

$$1 = \mathbb{1}_{\{I_t = a\}} + \mathbb{1}_{\{I_t \neq a\}}$$

trick used in Step (1).

* The question is:

Don't we have a generalized version of the Hoeffding-Azuma inequality with such predictable ranges $V_t - U_t$?

Yes we do have something in terms of constant upper bounds

$$V_t - U_t \leq \Delta_t \in \mathbb{R} \quad \text{as}$$

But $V_t - U_t = \mathbb{1}_{\{I_t = a\}}$ can only be bounded by $\Delta_t = 1$

So I think that Steps (2) and (3) are needed

Stochastic bandits with a continuum of arms

Stochastic bandits:What about arms indexed by a continuum?

Setting 1: Arms indexed by $x \in A$, where A is some possibly large set
 With each arm $x \in A$ is associated a probability distribution \mathcal{J}_x over \mathbb{R} s.t. $E(\mathcal{J}_x)$ exists
 At each round, the decision-maker picks $I_t \in A$,
 gets a reward Y_t drawn at random according to \mathcal{J}_{I_t}
 (given I_t); and this is the only feedback she gets.

Definition: $f: x \in A \mapsto E(\mathcal{J}_x)$ is the mean-payoff function
Regret: $\bar{R}_T = T \sup_{x \in A} f(x) - E\left[\sum_{t=1}^T Y_t\right]$

Setting 2: [special case] \rightarrow Noisy optimization of a function.

We fix $f: A \rightarrow \mathbb{R}$

The noise is given by a sequence of iid random variables

$\varepsilon_1, \varepsilon_2, \dots$

When $I_t \in A$ is picked, $Y_t = f(I_t) + \varepsilon_t$

\hookrightarrow Special case of setting #1 where \mathcal{J}_x is the distribution of $f(x) + \varepsilon_1$ (all these distributions have the same shape given by the common distribution of the ε_j)

We of course need conditions for the regret to be minimized.

Definition: Let \mathcal{F} be a set of possible bandit problems $\mathcal{J} = (\mathcal{J}_x)_{x \in A}$

The regret can be controlled (in a non-uniform way) against \mathcal{F} if:

we also say that (A, \mathcal{F}) is tractable

there exists a strategy s.t. $\forall \mathcal{J} \in \mathcal{F}, \bar{R}_T = o(T)$.

Ex: $A = \{1, \dots, K\}$ and $\mathcal{F} = (\mathcal{P}([0,1]))^K$, the set of all K -tuples of probability distributions over $[0,1]$
 the case of finitely many arms with bounded distributions
 \rightarrow UCB does the job.

Counter-example: $A = [0,1]$ and $\mathcal{F} = (\mathcal{P}([0,1]))^{[0,1]}$
 illustrating that continuity is a minimal requirement.
 all bandit problems $(\nu_x)_{x \in [0,1]}$ with distributions ν_x having support $[0,1]$.

Indeed: Consider $(\delta_{x_0})_{x \in [0,1]}$ the bandit problem in which each arm x is associated with the Dirac mass on 0

Fix any strategy: it gets $Y_t = 0$ $\forall t$ and uses a sequence of (possibly) random choices $I_t, t \geq 1$
 Since probability distributions can only have at most countably many atoms,
 $\mathcal{Y} = \{x \in [0,1] : \exists t \mid \mathbb{P}\{I_t = x\} > 0 \text{ under } (\delta_{x_0})_{x \in [0,1]}\}$
 is countable. In particular, $[0,1] \setminus \mathcal{Y}$ is non empty.

But the strategy behaves the same under the problem $(\nu'_x)_{x \in [0,1]}$ in which $\begin{cases} \nu'_x = \delta_0 & \forall x \neq x_0 \\ \nu'_{x_0} = \delta_{x_0} & \text{for one fixed } x_0 \in [0,1] \setminus \mathcal{Y} \end{cases}$

With probability 1, it thus never hits x_0 .
 Therefore, $Y_t = 0$ a.s. $\forall t$ and $\bar{R}_T = \frac{T}{T} - \mathbb{E}\left[\sum_{t=1}^T Y_t\right] = \frac{T}{T}$.

Actually, continuity is sufficient for the regret to be controlled, as long as A is not too large.

Theorem: Let A be a ^{cont} metric space and let \mathcal{F} be the set of bandit problems $(\nu_x)_{x \in A}$
 with : $\rightarrow \forall x, \nu_x$ is a distribution over $[0,1]$
 \rightarrow a continuous mean-payoff function $f, x \mapsto \mathbb{E}(\nu_x)$

The regret can be controlled against $\mathcal{F}^{\text{cont}}$ if and only if A is separable.

Corollary → Let \mathcal{F}^{all} be the family of all bandit models $(\mu_x)_{x \in A}$ with distributions ν_x over $[0,1]$. Then the regret against \mathcal{F}^{all} can be controlled if and only if A is at most countable.
 let A be any set.

Before we prove these facts, consider the following more concrete example, in which, by strengthening the regularity requirement on the mean-payoff function, we can even get rates.

Exercise: (Lipschitz bandits) Let $A = [0,1]$ and let \mathcal{F}^{Lip} be the family of bandit models $(\mu_x)_{x \in [0,1]}$ with distributions ν_x over $[0,1]$ and with mean-payoff functions that are Lipschitz. Exhibit a strategy based on UCB + a sequence of discretizations of $[0,1]$ into K bins (to be refined over time) such that:

Hint:

First, prove a performance bound by splitting $[0,1]$ into $[i/K, (i+1)/K]$ with $i=1, \dots, K$ for a fixed K , where each bin $[i/K, (i+1)/K]$ plays the role of an i in a bandit problem with finitely many arms. Then discuss how to pick K over the time, as we do in the next proof.

$\forall \epsilon \in \mathcal{F}^{\text{Lip}}$,

$$\bar{R}_T \leq (3L + 6\sqrt{8 \ln T} + 2)T^{2/3} + 2$$

where L is the Lipschitz constant of the mean-payoff function of ϵ .

Proof of the Corollary:

We endow A with the discrete topology, i.e., choose the distance

$$d(x,y) = \mathbb{1}_{\{x \neq y\}}. \quad \text{Then:}$$

1. All applications $f: A \rightarrow \mathbb{R}$ are continuous
2. A is separable if and only if A is at most countable.

Proof of the Theorem: It relies on the possibility or impossibility of uniform exploration of the arms

1) If A is separable: let $(z_n)_{n \in \mathbb{N}}$ be a collection of points in A that is dense

We pick actions in a triangular fashion:

Regime 1: UCB based on z_1, z_2 (fresh start) : $I_1^{(1)} \dots I_4^{(1)}$
 \vdots

Regime r : UCB based on z_1, z_2, \dots, z_r (fresh start) $I_1^{(r)} \dots I_{(r+1)^2}^{(r)}$
 \vdots

In Regime r : $(r+1)^2 \max_{s \leq r} f(z_s) - E \left[\sum_{t=S_r+1}^{S_r+(r+1)^2} Y_t \right]$ (*)

Starts at time $S_r + 1 = 2^2 + \dots + r^2 + 1$

$\leq c \sqrt{r^3 \ln r}$
 well-chosen numerical constant

order of magnitude of the distribution free regret bound for UCB on $(r+1)^2$ steps with r arms (we saw this bound as an exercise)

Now, let $\epsilon > 0$ and let $\tilde{z}_\epsilon \in \mathbb{N}^*$ s.t. $\mu_{\tilde{z}_\epsilon} \geq \sup_A f - \epsilon$
 (\tilde{z}_ϵ exists by separability of A and continuity of f)

In particular, $\max_{s \leq \tilde{z}_\epsilon} f(z_s) \geq \sup_A f - \epsilon$ (**)

We denote by r_T the index of the regime where T lies:

we have that S_{r_T} is of the order of r^3

thus that r_T is of the order of $T^{1/3}$

in particular, $r_T = O(T^{1/3})$

The regret can be decomposed (for T large enough) as

$$\begin{aligned}
 \bar{R}_T &= T \sup_A f - \mathbb{E} \left[\sum_{t=1}^T y_t \right] = \text{sum of the regrets of each regime} \\
 &= \underbrace{\sum_{r=1}^{\tilde{r}_\varepsilon - 1} (r+1)^2}_{\substack{\text{initial regimes,} \\ \text{regret bounded} \\ \text{by their lengths} \\ = O(1)}} + \underbrace{\sum_{r=\tilde{r}_\varepsilon}^{r_T-1} \left((r+1)^2 \varepsilon + c \sqrt{r^3 \ln r} \right)}_{\substack{\text{cf. bounds (*)} \\ \text{and (**)}}} + \underbrace{(r_T+1)^2}_{\substack{\text{regime } r_T \\ \text{may be} \\ \text{incomplete}}} \\
 &\leq T\varepsilon + c \sum_{r=\tilde{r}_\varepsilon}^{r_T-1} r^{3/2} \sqrt{\ln r} = O(T^{2/3}) \\
 &\leq T\varepsilon + O(r_T^{5/2} \sqrt{\ln r_T}) \\
 &= T\varepsilon + O(T^{5/6} \sqrt{\ln T})
 \end{aligned}$$

All in all, $\limsup \frac{\bar{R}_T}{T} \leq \varepsilon$, which is true $\forall \varepsilon > 0$,

that is, $\lim \frac{\bar{R}_T}{T} = 0$ as requested.

2) If A is not separable:

* We use the following characterization of separability (which relies on Zorn's lemma):

A metric space X is separable if and only if it contains no uncountable subset \mathcal{D} s.t. $\rho = \inf \{d(x, y) : x, y \in \mathcal{D}\} > 0$.

In particular, if A is not separable, there exist an uncountable subset $\mathcal{D} \subset A$ and $\rho > 0$ such that the balls $B(a, \rho/2)$, with $a \in \mathcal{D}$, are all disjoint.

\Rightarrow No probability distribution over A can give a positive mass to all these balls.

* We consider the bandit models $\nu^{(a)}$ inducing mean-payoff functions $f^{(a)} : x \in A \mapsto (1 - \frac{d(x, a)}{\rho/2})^+$; in particular, $\nu^{(a)} = \delta_a$ for $x \notin B(a, \rho/2)$.
 \uparrow $f^{(a)}$ is indeed continuous.

We proceed as in the example showing the necessity of continuity when $A = [0, 1]$ and consider the bandit model $(\delta_a)_{a \in A}$, as well as any strategy and the laws induced by the \mathcal{I}_t under this model: let ν_t

be the law of \mathcal{I}_t under $(\delta_a)_{a \in A}$ and let $\nu = \sum_{t \geq 1} \frac{1}{2^t} \nu_t$.

[as only countably many balls can have a positive mass under ν

\hookrightarrow There exists $a \in A$ s.t. $\nu(B(a, \rho/2)) > 0$, that is, s.t., $\forall t \geq 1$, $\mathbb{P}(\mathcal{I}_t \in B(a, \rho/2) \text{ under } (\delta_a)_{a \in A}) = 0$.

The considered strategy is therefore such that the \mathcal{I}_t have the same distribution under $(\delta_a)_{a \in A}$ and $\nu^{(a)}$. In particular,

$\mathbb{E}[\sum_{t=1}^T Y_t] = 0$ in both cases, but in the latter case,

$\sup f^{(a)} = 1$, so that $\bar{R}_T = T$ against $\nu^{(a)}$. The regret is not controlled against $\nu^{(a)} \in \mathcal{F}^{\text{cont}}$.

Tools for regret lower bounds for stochastic bandits

[Back to K-armed bandits]

Overview of the next topic: Fix a model \mathcal{D} , known to the decision-maker, i.e. a collection of probability distributions over \mathcal{R} with an expectation.

Assume that y_1, \dots, y_k are unknown but that the decision-maker knows $y_j \in \mathcal{D}$ $\forall j$.

What are the best bounds on $\bar{R}_T = T\mu^* - \mathbb{E} \left[\sum_{t=1}^T y_t \right]$?

One can show matching upper and lower bounds (with associated strategies):

$$\bar{R}_T \text{ is at best of the order of } \left(\sum_{a: \Delta_a > 0} \frac{\Delta_a}{K_{\text{inf}}(\bar{y}_a, \mu^*, \mathcal{D})} \right) \ln T$$

where

$$K_{\text{inf}}(\bar{y}_a, \mu^*, \mathcal{D}) = \inf \left\{ \text{KL}(\bar{y}_a, y_a^0) : \begin{array}{l} y_a^0 \in \mathcal{D} \\ \mathbb{E}(y_a^0) > \mu^* \end{array} \right\}$$

Kullback-Leibler divergence

expectation of y_a^0

We will only prove the lower bound part

but not

discuss the strategy, called KL-UCB, to achieve the bound.

} it's just a matter of time we would need about 3h to discuss this strategy

|| * But * before we do that, I guess that some reminder of basic and non-basic results about KL divergences would be needed!

The Kullback-Leibler divergence: definition and basic properties.

Definition (intrinsic): Let P, Q be two probability measures over (Ω, \mathcal{F})

$$KL(P, Q) = \begin{cases} +\infty & \text{if } P \text{ is not absolutely continuous w.r.t } Q \\ \int_{\Omega} \left(\frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ = \int_{\Omega} \left(\ln \frac{dP}{dQ} \right) dP & \text{if } P \ll Q \end{cases}$$

Basic facts:

- Existence of the defining integral when $P \ll Q$: because $\psi: x \mapsto x \ln x$ is bounded from below on $[0, +\infty)$
- $KL(P, Q) \geq 0$ and $0 = KL(P, Q)$ if and only if $P = Q$:

It suffices to consider the case $P \ll Q$: because ψ is strictly convex, Jensen's inequality indicates that

$$KL(P, Q) = \int_{\Omega} \psi\left(\frac{dP}{dQ}\right) dQ \geq \psi\left(\underbrace{\int_{\Omega} \frac{dP}{dQ} dQ}_{=1}\right) = 0$$

with equality if and only if $\frac{dP}{dQ}$ is Q -a.s constant, i.e., $P = Q$

Exercise: A useful rewriting. Prove the following result:

Assume $P \ll Q$ and let ν be any probability measure over (Ω, \mathcal{F})

such that $P \ll \nu$ and $Q \ll \nu$. Denote $f = \frac{dP}{d\nu}$ and $g = \frac{dQ}{d\nu}$.

Then:

$$\begin{aligned} KL(P, Q) &= \int_{\Omega} \frac{f}{g} \ln\left(\frac{f}{g}\right) g d\nu \\ &= \int_{\Omega} \ln\left(\frac{f}{g}\right) f d\nu \end{aligned}$$

Beware: with the usual measure-theoretic conventions, if $x \neq 0$ and $y = 0$, then $x \neq y \frac{x}{y}$ \hookrightarrow you therefore need to proceed with care!

Lemma (contraction of entropy; also known as data-processing inequality):

Let P, Q be two probability measures over (Ω, \mathcal{F})

Let $X: (\Omega, \mathcal{F}) \rightarrow (\Omega', \mathcal{F}')$ be any random variable

Denote by P^X and Q^X the laws of X under P and Q .

Then:

$$KL(P^X, Q^X) \leq KL(P, Q)$$

Proof:

We may assume that $P \ll Q$, otherwise $KL(P, Q) = +\infty$ and the inequality is true. We show that we then have

$$P^X \ll Q^X, \quad \text{with} \quad \frac{dP^X}{dQ^X} = \mathbb{E}_Q \left[\frac{dP}{dQ} \mid X = \cdot \right] \stackrel{\text{not.}}{=} \gamma$$

$$\text{ie, } \gamma(x) = \mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right].$$

Indeed, for all $B \in \mathcal{F}'$:

$$\begin{aligned} P^X(B) &= P\{X \in B\} = \int_{\Omega} \mathbb{1}_B(X) \frac{dP}{dQ} dQ \stackrel{\text{tower rule}}{=} \int_{\Omega} \mathbb{1}_B(X) \mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right] dQ \\ &\stackrel{\text{not.}}{=} \int_{\Omega} \mathbb{1}_B(X) \gamma(x) dQ \stackrel{\text{by definition of } Q^X}{=} \int_{\Omega'} \mathbb{1}_B \gamma dQ^X. \end{aligned}$$

Therefore,

$$\begin{aligned} KL(P^X, Q^X) &= \int_{\Omega'} \gamma \ln \gamma dQ^X = \int_{\Omega} \gamma(x) \ln \gamma(x) dQ \\ &= \int_{\Omega} \left(\mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right] \ln \mathbb{E}_Q \left[\frac{dP}{dQ} \mid X \right] \right) dQ \quad \left. \begin{array}{l} \text{definition} \\ \text{of } \gamma \end{array} \right\} \\ &\leq \int_{\Omega} \mathbb{E}_Q \left[\frac{dP}{dQ} \ln \frac{dP}{dQ} \mid X \right] dQ \quad \left. \begin{array}{l} \text{conditional} \\ \text{version of} \\ \text{Jensen's inequality} \end{array} \right\} \\ &\stackrel{\text{tower rule.}}{\leq} \int_{\Omega} \left(\frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ = KL(P, Q) \end{aligned}$$

References:

- The proof above is due to Ali and Silvey (1966), but it's far from being well-known!

- Typical proofs in the more recent literature:

- either focus on the discrete case (Cover and Thomas, 2006)

- or use the duality / variational formula for the KL (Massart, 2007; Boucheron, Lugosi, Massart, 2013)

- The joint convexity of KL, which we discuss below, is typically proved in a tedious way, relying on the rewriting of Exercise 1 and the joint convexity of $(x, y) \in [0, +\infty)^2 \mapsto \left(\frac{x}{y} \ln \frac{x}{y}\right)_+$

We may see it instead as a consequence of the data-processing inequality:

Corollary (joint convexity of KL): For all probability distributions $\mathbb{P}_1, \mathbb{P}_2$ and $\mathbb{Q}_1, \mathbb{Q}_2$ over the same measurable space (Ω, \mathcal{F}) , and all $d \in (0, 1)$,

$$KL((1-d)\mathbb{P}_1 + d\mathbb{P}_2, (1-d)\mathbb{Q}_1 + d\mathbb{Q}_2) \leq (1-d)KL(\mathbb{P}_1, \mathbb{Q}_1) + dKL(\mathbb{P}_2, \mathbb{Q}_2)$$

Proof: We augment (Ω, \mathcal{F}) into $(\Omega \times \{1, 2\}, \mathcal{F}')$ where

$$\mathcal{F}' = \mathcal{F} \otimes \{\emptyset, \{1\}, \{2\}, \{1, 2\}\}$$

We define the random pair (X, J) by the projections

$$X: (\omega, j) \mapsto \omega \quad \text{and} \quad J: (\omega, j) \mapsto j$$

Let \mathbb{P} be a probability measure on $(\Omega \times \{1, 2\}, \mathcal{F}')$ such that

$$\begin{cases} J \sim 1 + \text{Ber}(d) \\ X | J=j \sim \mathbb{P}_j \end{cases} \quad (\text{and a similar definition for } \mathbb{Q} \text{ based on } \mathbb{Q}_1, \mathbb{Q}_2)$$

that is, $\forall j \in \{1, 2\} \quad \forall A \in \mathcal{F} \quad \mathbb{P}(A \times \{j\}) = \left((1-d) \mathbb{1}_{\{j=1\}} + d \mathbb{1}_{\{j=2\}} \right) \mathbb{P}_j(A)$

$$\text{Now, } P^X = (1-d)P_1 + dP_2$$

$$Q^X = (1-d)Q_1 + dQ_2$$

and (as we prove below) $KL(P^X, Q^X) = (1-d)KL(P_1, Q_1) + dKL(P_2, Q_2)$
 so that the result follows from the data-processing inequality.

Indeed: we may assume with no loss of generality, given $d \in (0,1)$, that $P_1 \ll Q_1$ and $P_2 \ll Q_2$, so that $P^X \ll Q^X$ with

$$\frac{dP}{dQ}(w, j) = \mathbb{1}_{\{j=1\}} \frac{dP_1}{dQ_1}(w) + \mathbb{1}_{\{j=2\}} \frac{dP_2}{dQ_2}(w)$$

This entails that

$$\begin{aligned} KL(P, Q) &= \int_{\Omega \times \{1,2\}} \left(\frac{dP}{dQ}(w, j) \ln \frac{dP}{dQ}(w, j) \right) dQ(w, j) \\ &= \int_{\Omega \times \{1\}} \left(\frac{dP}{dQ}(w, 1) \ln \frac{dP}{dQ}(w, 1) \right) \mathbb{1}_{\Omega \times \{1\}}(w, j) dQ(w, j) \\ &\quad + \int_{\Omega \times \{2\}} \left(\frac{dP}{dQ}(w, 2) \ln \frac{dP}{dQ}(w, 2) \right) \mathbb{1}_{\Omega \times \{2\}}(w, j) dQ(w, j) \\ &= \int_{\Omega} \left(\frac{dP_1}{dQ_1}(w) \ln \frac{dP_1}{dQ_1}(w) \right) (1-d) dQ_1(w) + \int \dots \\ &= (1-d) KL(P_1, Q_1) + d KL(P_2, Q_2) \end{aligned}$$

we just use that for $f \geq a$ constant, $\int f d\mu = \int f \mathbb{1}_a d\mu + \int f \mathbb{1}_{f > a} d\mu$ whether f is integrable or not

on $\Omega \times \{1\}$, dQ is $(1-d)dQ_1$

KL for product measures. (\Leftrightarrow The independent case)

Proposition: Let (Ω, \mathcal{F}) and (Ω', \mathcal{F}') be two measurable spaces,
 let P, Q be two probability measures over (Ω, \mathcal{F})
 P', Q' over (Ω', \mathcal{F}')

and denote by $P \otimes P'$ and $Q \otimes Q'$ the product distributions over $(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}')$. Then:

$$KL(P \otimes P', Q \otimes Q') = KL(P, Q) + KL(P', Q')$$