

Sequential Learning, Sequential Optimization

M2 course for tracks:
 1- Probability & Statistics
 2- Optimization
 3- Data Science

Gilles Stoltz: CNRS researcher at HEC Paris
 stoltz@hec.fr

Exam: written exam in unlimited time (9h40 - [...])
 on Wednesday March 29 at HEC Paris

! No other option! (no project!) → Only 3 A4
 2-sided cheat sheets

Beware: 5 ECTS for 18h in Data Science track:

A good deal but highly technical course:

definition / theorem / proof (sometimes: long and technical proofs)

6 × 3h
 from Feb. 1st
 to March 15
 (except: Feb. 15)

Pre-requisites:

- Martingale theory
- Measure theory / Integration theory
- But no real knowledge of statistics is required

Lecture notes:

On my website, "Enseignements" page
<http://stoltz.perso.math.cnrs.fr>

We will study:
 - algorithms with associated upper bounds of performance
 - lower bounds on the performance of any algorithm

Read them in details, meditate the proof techniques and summarize them on your cheat sheets

Also put FYI the statement of last year's exam

(was a disaster in the Data Science track: beware, work is needed!)

! Send me an email for registration!

Course 1 / February 1st, 2017

1. Framework and setting of stochastic bandits
2. A popular algorithm: UCB

↑ almost: up to a concentration lemma
(Hoeffding-Azuma for a random number of summands)

basic unconditional one

conditional version

0. Reminder on the Hoeffding-Azuma inequality

↳ There are three exercises for you!

Lemma (Hoeffding): X random variable s.t. $X \in [a, b]$ a.s.

Then $\forall s \in \mathbb{R}$,

$$\ln \mathbb{E}[e^{s(X-\mathbb{E}X)}] = \ln \mathbb{E}[e^{sX}] - s \mathbb{E}X \leq \frac{s^2(b-a)^2}{8}$$

Proof (most elegant one I know of):

$$\Psi(s) = \ln \mathbb{E}[e^{sX}] \text{ defined for all } s \in \mathbb{R}$$

Ψ is differentiable at each $s \in \mathbb{R}$: cf. X bounded, thus $\eta \mapsto X e^{\eta X}$ locally dominated around s by an integrable r.v. independent of η .
 thus $\eta \mapsto \mathbb{E}[e^{\eta X}]$ differentiable at s with derivative $\mathbb{E}[X e^{sX}]$

with

$$\Psi'(s) = \frac{\mathbb{E}[X e^{sX}]}{\mathbb{E}[e^{sX}]}$$

Similarly, Ψ is twice differentiable at each $s \in \mathbb{R}$, with:

$$\Psi''(s) = \frac{\mathbb{E}[X^2 e^{sX}] \mathbb{E}[e^{sX}] - (\mathbb{E}[X e^{sX}])^2}{(\mathbb{E}[e^{sX}])^2} = \text{Var}_{\mathbb{Q}}(X)$$

under the probability \mathbb{Q} defined by

$$\frac{d\mathbb{Q}}{d\mathbb{P}}(\omega) = \frac{e^{sX(\omega)}}{\mathbb{E}[e^{sX}]}$$

$$X \in [a, b]: \quad \text{Var}_{\mathbb{Q}}(X) = \inf_{\mu \in \mathbb{R}} \mathbb{E}_{\mathbb{Q}}[(X-\mu)^2] \leq \mathbb{E}_{\mathbb{Q}}\left[\left(X - \frac{a+b}{2}\right)^2\right] = \frac{(b-a)^2}{4}$$

Taylor: $\exists x$ s.t. $\Psi(s) = \underbrace{\Psi(0)}_{=0} + \underbrace{s \Psi'(0)}_{=0} + \frac{s^2}{2} \underbrace{\Psi''(x)}_{\leq (b-a)^2/4}$

ie, $\ln \mathbb{E}[e^{sX}] \leq \frac{s^2}{8} (b-a)^2$

The Hoeffding-Azuma inequality

Theorem: Let $(\mathcal{F}_t)_{t \geq 0}$ be a filtration and let $(X_t)_{t \geq 1}$ be a sequence of adapted random variables (ie, $\forall t \geq 1$, X_t is \mathcal{F}_t -measurable), that are bounded: $\forall t$, $a_t \leq X_t \leq b_t$ a.s., where $a_t, b_t \in \mathbb{R}$.

Then (« probabilistic version »)

$$\forall \varepsilon > 0, \quad \mathbb{P}\left\{ \sum_{t=1}^T X_t - \sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}] \geq \varepsilon \right\} \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{t=1}^T (b_t - a_t)^2}\right)$$

or (« statistical version », totally equivalent)

$$\forall \delta \in (0, 1), \quad \text{with probability at least } 1 - \delta, \quad \sum_{t=1}^T X_t - \sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}] \leq \sqrt{\frac{\sum_{t=1}^T (b_t - a_t)^2}{2} \ln \frac{1}{\delta}}$$

Note: Hoeffding's inequality is the special case when all X_t are independent and $\mathcal{F}_{t-1} = \sigma(X_1, \dots, X_{t-1})$, so that $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = \mathbb{E}[X_t]$.

Basic ingredient of the proof: extension of Hoeffding's lemma to conditional expectations

Lemma: X random variable s.t. $X \in [a, b]$ a.s.

Then, for all σ -algebras \mathcal{G} , for all $s \in \mathbb{R}$,

$$\ln \mathbb{E}\left[e^{s(X - \mathbb{E}[X | \mathcal{G}])} | \mathcal{G}\right] = \ln \left(\mathbb{E}\left[e^{sX} | \mathcal{G}\right] \right) - s \mathbb{E}[X | \mathcal{G}] \leq \frac{s^2}{8} (b-a)^2$$

(We will discuss the proof in a minute... let's first prove the theorem based on this lemma.)

Proof (of the theorem):

Markov-Chernoff bounding (= Markov's inequality after taking exponents):

$$\text{We denote } S_T = \sum_{t=1}^T X_t - \sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}]$$

\uparrow martingale = sum of \uparrow martingale increments or martingale differences

The « probabilistic version » is about upper bounding $\mathbb{P}\{S_T > \varepsilon\}$:

$$\mathbb{P}\{S_T > \varepsilon\} = \mathbb{P}\{e^{\lambda S_T} > e^{\lambda \varepsilon}\} \stackrel{\text{Markov's inequality}}{\leq} e^{-\lambda \varepsilon} \mathbb{E}[e^{\lambda S_T}]$$

\uparrow
 $\forall \lambda > 0$

We show by induction that $\mathbb{E}[e^{\lambda S_T}] \leq \exp\left(\frac{\lambda^2}{8} \sum_{t=1}^T (b_t - a_t)^2\right)$

- For $T=1$, true by the conditional version of Hoeffding's lemma and the fact that $S_1 = X_1 - \mathbb{E}[X_1 | \mathcal{F}_0]$ with $X_1 \in [a_1, b_1]$ + taking expectations
- For $T-1 \rightarrow T$, where $T \geq 2$:

The extension of Hoeffding's lemma ensures that

$$\mathbb{E}[e^{\lambda(X_T - \mathbb{E}[X_T | \mathcal{F}_{T-1}])} | \mathcal{F}_{T-1}] \leq e^{\lambda^2(b_T - a_T)^2/8}$$

so that

$$\begin{aligned} \mathbb{E}[e^{\lambda S_T}] &= \mathbb{E}[\mathbb{E}[e^{\lambda S_T} | \mathcal{F}_{T-1}]] \\ &= \mathbb{E}[e^{\lambda S_{T-1}} \mathbb{E}[e^{\lambda(X_T - \mathbb{E}[X_T | \mathcal{F}_{T-1}])} | \mathcal{F}_{T-1}]] \\ &\stackrel{\text{by the induction hypothesis}}{\leq} e^{\lambda^2(b_T - a_T)^2/8} \times \mathbb{E}[e^{\lambda S_{T-1}}] \\ &\leq \exp\left(\lambda^2 \sum_{t=1}^T (b_t - a_t)^2/8\right) \end{aligned}$$

Substituting above:

$$\mathbb{P}\{S_T > \varepsilon\} \leq \inf_{\lambda > 0} \exp\left(-\lambda \varepsilon + \frac{\lambda^2}{8} \sum_{t=1}^T (b_t - a_t)^2\right)$$

strictly convex function to minimize in the exponent:
minimum achieved at λ^*
such that $\lambda^* \sum_{t=1}^T (b_t - a_t)^2/4 = \varepsilon$ (gradient vanishes)
i.e. $\lambda^* = 4\varepsilon / \sum_{t=1}^T (b_t - a_t)^2$

$$\uparrow = \exp\left(-\frac{2\varepsilon^2}{\sum_{t=1}^T (b_t - a_t)^2}\right)$$

→ It only remains to prove the extension of Hoeffding's lemma to conditional expectations.

Back to Hoeffding's lemma with conditional expectations:

Proof 1? Can we take the proof of Hoeffding's lemma we just saw and replace all E by $E[\cdot | \mathcal{G}_j]$?

$\Psi(s) = \ln E[e^{sx} | \mathcal{G}_j]$ → The theorem of differentiation under $E[\cdot]$ only requires dominated convergence, which holds true for $E[\cdot | \mathcal{G}_j]$ as well. Thus, we also have a theorem of differentiation under $E[\cdot | \mathcal{G}_j]$:

$$\text{a.s., } \Psi''(s) \text{ exists and equals } \Psi''(s) = \frac{E[X^2 e^{sx} | \mathcal{G}_j] E[e^{sx} | \mathcal{G}_j] - (E[X e^{sx} | \mathcal{G}_j])^2}{(E[e^{sx} | \mathcal{G}_j])^2}$$

= some conditional variance under a different probability measure?

Yes, using the notion of « regular conditional probability », which always exists in our case we could perform a change of measure again and identify a conditional variance.

But PLEASE! this is extremely heavy maths... I can't inflict that to you....

Proof 2 Too bad for elegance, let's get back to the original proof of Hoeffding's (unconditional) lemma, which only relies on calculus:

$$Y = X - E[X | \mathcal{G}_j] \in [A, B] \quad \text{where } A = a - E[X | \mathcal{G}_j] \leq 0 \\ B = b - E[X | \mathcal{G}_j] \geq 0 \\ \text{and both } \mathcal{G}_j\text{-measurable and } B-A = b-a > 0$$

$$Y = \frac{B-Y}{B-A} A + \frac{Y-A}{B-A} B$$

↑ convex weights

Since $y \mapsto e^{sy}$ is convex:

$$e^{sY} \leq \frac{B-Y}{B-A} e^{sA} + \frac{Y-A}{B-A} e^{sB}$$

Taking $E[\cdot | \mathcal{G}_j]$: using $E[Y | \mathcal{G}_j] = 0$ and A, B \mathcal{G}_j -measurable:

$$E[e^{sY} | \mathcal{G}_j] \leq \frac{B}{B-A} e^{sA} - \frac{A}{B-A} e^{sB}$$

← note that $B/(B-A)$ and $-A/(B-A)$ are convex weights

Now, by a function study (the very same as the one we performed in the proof of the unconditional version of Hoeffding's lemma) — or even by the latter lemma itself:

$$\forall u, v \in \mathbb{R}, \quad \forall p \in (0, 1), \quad \ln(p e^{su} + (1-p)e^{sv}) \leftarrow \ln \text{ of expected value of } e^{sz} \text{ where } z = \begin{cases} u & \text{w.p. } p \\ v & \text{w.p. } 1-p \end{cases}$$

$$\leq s(pu + (1-p)v) + \frac{s^2}{8}(v-u)^2$$

↖ expected value of Z
↙ range is $[u, v]$

In particular,

$$\frac{B}{B-A} e^{sA} - \frac{A}{B-A} e^{sB} \leq \exp\left(s\left(\frac{BA}{B-A} - \frac{AB}{B-A}\right) + \frac{s^2}{8}(B-A)^2\right)$$

$$= \exp\left(\frac{s^2}{8}(b-a)^2\right)$$

Summarizing:

$$\mathbb{E}[e^{sY} | \mathcal{G}] \leq \exp\left(\frac{s^2}{8}(b-a)^2\right)$$

$$= \mathbb{E}[e^{sX} | \mathcal{G}] \times \exp(-s \mathbb{E}[X | \mathcal{G}])$$

Proof 3 My preferred (not only because I found by myself):
Hoeffding's lemma in its unconditional version ENTAILS the conditional version! This is because Hoeffding's lemma holds for all probability distributions — we should play with this fact.

For all $A \in \mathcal{G}$ s.t. $\mathbb{P}(A) > 0$, let $\mathbb{P}_A = \mathbb{P}(\cdot | A)$, the conditional distribution given the event A .

The unconditional version of Hoeffding's lemma ensures that

$$\forall A \in \mathcal{G} \text{ s.t. } \mathbb{P}(A) > 0, \quad \forall s \in \mathbb{R}, \quad \ln \mathbb{E}_A[e^{sX}] \leq s \mathbb{E}_A[X] + \frac{s^2}{8}(b-a)^2$$

Why do we consider the \mathbb{E}_A ? Because $\mathbb{E}[X | \mathcal{G}]$ is the unique \mathcal{G} -measurable random variable such that

$$\forall A \in \mathcal{G}, \quad \mathbb{E}[X \mathbb{1}_A] = \mathbb{E}[\mathbb{E}[X | \mathcal{G}] \mathbb{1}_A]$$

or, equivalently,

$$\forall A \in \mathcal{G} \text{ s.t. } \mathbb{P}(A) > 0, \quad \mathbb{E}_A[X] = \mathbb{E}_A[\mathbb{E}[X | \mathcal{G}]].$$

Now, consider the event $H = \left\{ \mathbb{E}[e^{sX} | \mathcal{G}] - e^{s \mathbb{E}[X | \mathcal{G}]} e^{\frac{s^2}{8}(b-a)^2} > 0 \right\} \in \mathcal{G}$

We want to prove that $P(H) = 0$. We proceed by contradiction:
 if we had $P(H) > 0$, then P_H would be defined and

$$\begin{aligned}
 E_H[e^{\lambda X}] &= E_H[E[e^{\lambda X} | \mathcal{G}]] > E_H[e^{\lambda E[X | \mathcal{G}]}] e^{\lambda^2(b-a)^2/8} \\
 &\quad \uparrow \text{a property of conditional expectations} \quad \uparrow \text{by definition of } H \text{ and because a function } > 0 \text{ on a set with probability } > 0 \text{ has an expectation } > 0 \\
 &\quad \geq e^{\lambda E_H[E[X | \mathcal{G}]]} e^{\lambda^2(b-a)^2/8} \\
 &\quad \quad \quad \uparrow \text{(unconditional) Jensen's inequality} \\
 &= e^{\lambda E_H[X]} e^{\lambda^2(b-a)^2/8}
 \end{aligned}$$

which would be in contradiction with the unconditional Hoeffding's lemma:

$$\ln E_H[e^{\lambda X}] \leq \lambda E_H[X] + \frac{\lambda^2}{8} (b-a)^2.$$

We conclude that $P(H) = 0$.

A final remark:

A better version of Hoeffding's lemma in its conditional form is the following one:

$$\begin{aligned}
 \text{Hoeffding's lemma:} \quad & X \text{ random variable s.t. there exists } G \text{ } \mathcal{G}\text{-measurable} \\
 & \text{and } a, b \in \mathbb{R} \text{ with: } G+a \leq X \leq G+b \text{ a.s.} \\
 \text{Then } \forall s \in \mathbb{R}, \quad & \ln E[e^{sX} | \mathcal{G}] - s E[X | \mathcal{G}] \leq \frac{s^2}{8} (b-a)^2
 \end{aligned}$$

→ In my first statement, I only considered $G \equiv 0$.

Why can we consider general G 's?

Proof 2: $Y = X - E[X | \mathcal{G}] = (X - G) - E[(X - G) | \mathcal{G}]$ ^{still} $\in [a, b]$
 and everything goes through

Proof 3: Some adaptation is needed as X may be unbounded! Only $X - G$ is bounded between a and b .

So, using the same proof scheme, we can conclude that

$$\ln E[e^{s(X-G)} | \mathcal{G}] \leq s E[X-G | \mathcal{G}] + \frac{s^2}{8} (b-a)^2$$

But I am not sure we may write e.g.

$$\mathbb{E}[X-G | \mathcal{G}] = \mathbb{E}[X | \mathcal{G}] - G$$

We may, under an additional assumption that X is integrable.

Some issues with the $\mathbb{E}[e^{\lambda(X-G)} | \mathcal{G}]$ term: I would be more comfortable with adding the assumption that $e^{\lambda X} \in L^1$ $\forall \lambda$.

Hence:

Exercise 0: Can anyone make his proof #3 work without these extra assumptions?

Now, we get back to the Hoeffding-Azuma inequality:

[Hoeffding-Azuma: holds therefore for adapted processes such that for all t , there exist $(a_t, b_t) \in \mathbb{R}^2$ with $G_t + a_t \leq X_t \leq G_t + b_t$ a.s. G_t \mathcal{F}_{t-} -measurable]

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 \mathbb{R})
with an expectation.At 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.

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

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

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

- Get a reward Y_t

always ≥ 1 since each arm was tried sequentially during rounds 1, 2, ..., K

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

Theorem:

If the distributions ν_j have supports all included in $[0, 1]$, then the pseudo-regret of UCB is smaller than

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

where $\Delta_i = \mu^* - \mu_i$ is the gap of arm i .

This regret bound is obtained via the following proposition:

Proposition:

If the distributions ν_j have supports all included in $[0, 1]$, 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 1.

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 $[a, 1]$, not on the specific distributions ν_i at hand): for the UCB algorithm,

$$\sup_{\nu_1, \dots, \nu_K \text{ with supports in } [a, 1]} \bar{R}_T \leq O(\sqrt{TK \ln T}).$$

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 $T = \sum_j \mathbb{E}[N_j(T)]$)

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

) 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 $2\Delta_j \leq 2$.

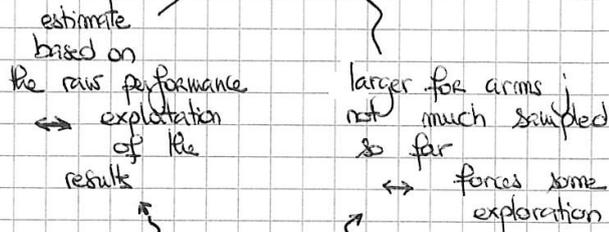
NOTE: Keep in mind the rewriting $\bar{R}_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T Y_t\right]$ as we will often use it!

$$= \sum_{a=1}^K \Delta_a \mathbb{E}[N_a(T)]$$

Proof [of the Proposition]: We fix an optimal arm $a^* \in \{1, \dots, K\}$, $i \neq a^*$ s.t. $\mu_{a^*} = \mu_i^*$.

→ It will show why this algorithm is called UCB:

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



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

We compare these statements to the Hoeffding - Azuma inequality in the next lecture

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

$\forall \delta \in (0, 1)$, $\mathbb{P}\left\{ \mu_j > \hat{\mu}_{j,t} - \sqrt{\frac{\ln(1/\delta)}{2 N_j(t)}} \right\} \geq 1 - \delta$

or

By symmetry: $\forall \delta \in (0, 1)$, $\mathbb{P}\left\{ \mu_j < \hat{\mu}_{j,t} + \sqrt{\frac{\ln(1/\delta)}{2 N_j(t)}} \right\} \geq 1 - \delta$

Notes: You can log $1-x \sim x^a$ then $\mu_{a^*} = 1 - \mu_i$ as well and μ_{a^*} supported on $[0, 1]$. replace

Application: $N_i(T) = 1 + \sum_{t=K+1}^T \mathbb{1}_{\{I_t = i\}}$

We show below that $t \geq K+1$ and $I_t = i$ entails one of the following:

- (i) $\hat{\mu}_{i,t-1} > \mu_i + \sqrt{\frac{2 \ln t}{N_i(t-1)}}$ [$\mu_i < \text{lower confidence bound}$]
- (ii) $\hat{\mu}_{a^*,t-1} < \mu_{a^*} - \sqrt{\frac{2 \ln t}{N_{a^*}(t-1)}}$ [$\mu_{a^*} > \text{upper confidence bound}$]
- (iii) $N_i(t-1) \leq \frac{8 \ln t}{\Delta^2}$ [i 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 $\mathbb{I}_t = i$, that is, $i \in \arg \max_j \left\{ \hat{\mu}_{j, t-1} + \sqrt{\frac{2 \ln t}{N_j(t-1)}} \right\}$

Thus,

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

each $\leq t \delta$ where $\delta = \frac{1}{t^4}$

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

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

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

deterministically upper bounded by $\left(\frac{8 \ln T}{\Delta_i^2} + 1 \right) - 1$

as $\mathbb{I}_t = i$ only if $N_i(t-1) \leq \frac{8 \ln T}{\Delta_i^2} + 1$
 thus only if $N_i(t) \leq \frac{8 \ln T}{\Delta_i^2} + 1$

so that the total sum $\sum_{s=1}^t \mathbb{1}_{\mathbb{I}_s = i} = N_i(t)$ is controlled by this number

-1 because $\mathbb{I}_1 = i$ is not included in the $\sum_{t=K+1}^T \dots$

Thus:

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

Proof of the lemma (Hoeffding-Azuma inequality with a random number of summands):

We only start (we will finish it next week!)

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

(0) $(Z_t)_{t \geq 0}$ is a martingale wrt. $(\mathcal{F}_t)_{t \geq 0} = (\sigma(Y_1, \dots, Y_t))$
 where $\mathcal{F}_0 = \{\emptyset, \Omega\}$ trivial σ -algebra

Indeed: each \mathcal{I}_t is \mathcal{F}_{t-1} -measurable (picked based only on past payoffs)

thus Z_t is \mathcal{F}_t -adapted

Showing that it is a martingale amounts to showing

$$E[(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a} \mid Y_1, \dots, Y_{t-1}] \stackrel{?}{=} 0 \text{ a.s.}$$

but since \mathcal{I}_t is \mathcal{F}_{t-1} -measurable, this quantity equals

$$E[(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a} \mid Y_1, \dots, Y_{t-1}, \mathcal{I}_t]$$

$$= (E[Y_t \mid \mathcal{I}_t, Y_1, \dots, Y_{t-1}] - \mu_a) \mathbb{1}_{\mathcal{I}_t = a}$$

$$= (\mu_{\mathcal{I}_t} - \mu_a) \mathbb{1}_{\mathcal{I}_t = a} = 0 \text{ as desired}$$

by the bandit model, Y_t is drawn independently at random given \mathcal{I}_t . Plus by the very bandit model, this conditional expectation equals $\mu_{\mathcal{I}_t}$

Exercise 2: To conclude the proof, show that

(1) For all $x \in \mathbb{R}$, $(M_t) = \left(\exp(x Z_t - \frac{x^2}{8} N_t(t)) \right)_{t \geq 0}$ is an adapted supermartingale
 \hookrightarrow in particular $E[M_t] \leq 1$ for all t

(2) $\forall \varepsilon > 0, \forall \ell \geq 1, \mathbb{P}\{Z_t \geq \varepsilon \text{ and } N_t(t) = \ell\} \leq e^{-2\varepsilon^2/\ell}$

and then conclude.