

Lecture #3: Stochastic bandits (part 1)

Full Information Setting

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

- agent picks an arm $a_t \in \{1, \dots, K\}$ (possibly at random)
 - observed reward vector $X(t) \in [0, 1]^K$
- gets reward $X_{a_t}(t)$.

a_t is $\sigma(U, X(1), \dots, X(t-1))$ measurable
(\uparrow possible randomization)

$$R_T = \max_{k \in [K]} \sum_{t=1}^T X_k(t) - \sum_{t=1}^T X_{a_t}(t)$$

As learning with experts, but:

- rewards instead of loss ($l_t \leftrightarrow 1 - X(t)$)
- choose pure actions (K -simplex $\leftrightarrow \{1, \dots, K\}$)
but can randomize over actions.

The $X(t)$ were chosen adversarially (worst case) in 1st lecture.

What if instead they are stochastic?

Assume

• $(X_{a_t})_t$ are iid.

• $X_k(t) \sim \nu_k$ with $E[X_k(t)] = \mu_k$.

Problem should be easier?

\hookrightarrow not really: we proved the lower bound in this setting:

for any algorithm, with $X_k(t) \sim \text{Ber}(\frac{1}{2})$

$$\mathbb{E}[R_T] \gtrsim \sqrt{\frac{T}{2} \ln K}$$

However, we can have much better results with the pseudo-regret:

$$\bar{R}_T = \max_{k \in [K]} \mathbb{E} \left[\sum_{t=1}^T X_{k^*}(t) \right] - \mathbb{E} \left[\sum_{t=1}^T X_{a_t}(t) \right]$$

↳ expectation w.r.t the realizations of $X(t)$

and eventually the agent policy

Previous example yields $\bar{R}_T = 0$. Makes sense: we cannot guess in advance heads or tails.

⚠ Warning: $\mathbb{E}[R_T] \neq \mathbb{E}[\bar{R}_T]$

Actually, $\mathbb{E}[R_T] > \mathbb{E}[\bar{R}_T]$. Why?

$$\bar{R}_T = T \max_k \mu_k - \mathbb{E} \left[\sum_{t=1}^T X_{a_t}(t) \right]$$

→ from now on, we will write R_T for the pseudo-regret.

Notations: • $\mu^* = \max_k \mu_k$

• $\Delta_k = \mu^* - \mu_k$ $\begin{cases} > 0 & \text{for sub-optimal arms} \\ = 0 & \text{for optimal arms} \end{cases}$

• $\Delta = \min_{k, \Delta_k > 0} \Delta_k$

$$N_k(t) = \sum_{s=1}^t \mathbb{1}_{(a_s=k)}$$

number of pulls on arm k .

Lemma:

$$\text{For any policy, } R_T = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(T)]$$

Proof:

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \mu^* - X_{a_t}(t) \right]$$

$$= \mathbb{E} \left[\sum_{t=1}^T \mu^* - \sum_{k=1}^K \mathbb{1}_{a_t=k} X_k(t) \right]$$

$$= \sum_{t=1}^T \sum_{k=1}^K \mathbb{E} \left[(\mu^* - X_k(t)) \mathbb{1}_{a_t=k} \right]$$

$$= \sum_{k=1}^K \sum_{t=1}^T (\mu^* - \mu_k) \mathbb{E} \left[\mathbb{1}_{a_t=k} \right]$$

$X_k(t) \perp a_t$

$$= \sum_{k=1}^K \Delta_k \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}_{a_t=k} \right]$$

$$= \sum_{k=1}^K \Delta_k \mathbb{E} [N_k(T)]$$

□

Greedy algorithm (or Follow The Leader)

choose a_1 arbitrarily

For $t \geq 2$:

$$a_t \in \operatorname{argmax}_{k \in [K]} \sum_{s=1}^{t-1} X_k(s)$$

Theorem

For any $(\mu_1, \dots, \mu_K) \in [0, 1]^K$ and $T \in \mathbb{N}$, Greedy satisfies in the Full Information setting:

$$R_T \leq \sum_{k: \Delta_k > 0} \frac{1}{\Delta_k}$$

Proof: $R_T = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(T)]$

Let us bound $\mathbb{E}[N_k(T)]$ for any k with $\Delta_k > 0$.

Let $k^* \in \arg \max_k \mu_k$.

$$\mathbb{E}[N_k(T)] \leq \sum_{t=1}^T \mathbb{P}\left(\frac{1}{t} \sum_{s=1}^t X_k(s) - X_{k^*}(s) \geq 0\right)$$

$$\leq \sum_{t=1}^T \mathbb{P}\left(\sum_{s=1}^t (X_k(s) - \mu_k) - \sum_{s=1}^t (X_{k^*}(s) - \mu_{k^*}) \geq t\Delta_k\right)$$

$$\leq \sum_{t=1}^T e^{-t\Delta_k^2} \leq \frac{e^{-\Delta_k^2}}{1 - e^{-\Delta_k^2}} = \frac{1}{e^{\Delta_k^2} - 1}$$

Hoeffding inequality

$$e^{\Delta_k^2} - 1 \geq \Delta_k^2$$

$$\leq \frac{1}{\Delta_k^2}$$

$$\text{So } R_T = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(T)]$$

$$\leq \sum_{\substack{k \\ \mu_k < \mu^*}} \Delta_k \cdot \frac{1}{\Delta_k^2} \quad \square$$

Bandit Setting

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

- agent picks an arm $a_t \in \{1, \dots, K\}$ (possibly at random)
- observed and gets reward $X_{a_t}(t) \in [0, 1]$

a_t is $\sigma(V_1, X_1(1), V_2, \dots, X_{t-1}(1), V_{t-1})$ -measurable

$$R_T = \max_{k \in [K]} \sum_{t=1}^T X_k(t) - \sum_{t=1}^T X_{a_t}(t)$$

→ only observe the reward of the pulled arm

→ exploration vs exploitation trade-off

estimate optimal arm by pulling all arms

maximize reward by pulling arm which seems the best

Notation

$$\hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^t X_k(s) \mathbb{1}_{\{a_s=k\}} \quad (\text{empirical mean})$$

Greedy algorithm (Bandit setting)

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

$$a_t = t$$

For $t \geq K+1$:

$$a_t \in \operatorname{argmax}_{k \in [K]} \hat{\mu}_k(t-1)$$

Theorem

For $v_1 = \operatorname{Ber}\left(\frac{3}{4}\right)$, $v_2 = \operatorname{Ber}\left(\frac{1}{4}\right)$, Greedy satisfies

in the bandit setting:

$$R_T \geq \frac{T-1}{32}$$

Proof:

$$\mathbb{P}(X_1(1) = 0, X_2(2) = 1) = \left(\frac{1}{4}\right)^2 = \frac{1}{16}$$

If $X_1(1) = 0$ and $X_2(2) = 1$, Greedy will keep pulling the arm 2 until T , so that: $\mathbb{E}[N_2(T)] \geq \frac{T-1}{16}$ \square

Greedy does not explore enough. It can underestimate the optimal arm and never pull it again.

Lemma: (bandit concentration)

For any bandit algorithm, any $k \in [K]$, $T \in \mathbb{N}$, $\delta \in (0, 1)$.

$$\mathbb{P}\left(\mu_k - \hat{\mu}_k(T) \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(T)}}\right) \leq \delta$$

$$\mathbb{P}\left(\hat{\mu}_k(T) - \mu_k \geq \sqrt{\frac{\ln(1/\delta)}{2N_k(T)}}\right) \leq \delta$$

1) This is not a trivial consequence of Hoeffding inequality,

$N_k(t)$ is a random variable and $\hat{\mu}_k(t), N_k(t)$ are not independent!

Hoeffding inequality indeed gives

$$\mathbb{P}\left(\frac{1}{n} \sum_{s=1}^n X_k(s) - \mu_k \geq \sqrt{\frac{\ln(1/\delta)}{2n}}\right) \leq e^{-\ln(1/\delta)} = \delta$$

But here, n is a random variable and is not independent from $\hat{\mu}_k(t)$

• What if instead we used Azuma-Hoeffding on $(X_k(s) - \mu_k) \mathbb{1}_{\{a_s = k\}}$?

martingale increment bounded between $-\mu_k$ and $1 - \mu_k$.

$$\mathbb{P}\left(\sum_{s=1}^t (X_k(s) - \mu_k) \mathbb{1}_{\{a_s = k\}} \geq \sqrt{\frac{t}{2} \ln(1/\delta)}\right) \leq \delta$$

$$\mathbb{P}\left(\hat{\mu}_k(t) - \mu_k \geq \sqrt{\frac{t}{N_k(t)} \frac{\ln(1/\delta)}{2 N_k(t)}}\right) \leq \delta$$

differences with our Lemma
getting rid of this $\sqrt{\frac{t}{N_k(t)}}$ factor is a big deal!

Proof: Let $Z_t = \sum_{s=1}^t (X_k(s) - \mu_k) \mathbb{1}_{\{a_s = k\}}$.

1) We first prove that $\forall z \in \mathbb{R}, \mathbb{E}\left[e^{z Z_t - \frac{z^2}{8} N_k(t)}\right] \leq 1$.

For that, we show that $M_t = \exp\left(z Z_t - \frac{z^2}{8} N_k(t)\right)$ is a supermartingale, so that

$$\mathbb{E}[M_t] \leq \mathbb{E}[M_0] = 1.$$

$$\text{Let } \mathcal{F}_{t-1} = \sigma(U, X_{a_1}(1), \dots, X_{a_{t-1}}(t-1))$$

a_t is \mathcal{F}_{t-1} measurable, so that:

$$\begin{aligned} \mathbb{E}[M_t | \mathcal{F}_{t-1}] &= \mathbb{E}\left[e^{\frac{x(X_k(t) - \mu_k) - \frac{x^2}{2}}{1}} \mathbb{1}_{a_t=k} | \mathcal{F}_{t-1} \right] M_{t-1} \\ &= \left(\mathbb{E}\left[e^{\frac{x(X_k(t) - \mu_k) - \frac{x^2}{2}}{1}} | \mathcal{F}_{t-1} \right] \mathbb{1}_{a_t=k} + \mathbb{1}_{a_t \neq k} \right) M_{t-1} \end{aligned}$$

Hoeffding's lemma (conditional) gives $\mathbb{E}\left[e^{\frac{x(X_k(t) - \mu_k)}{1}} | \mathcal{F}_{t-1} \right] \leq \frac{x^2}{8}$

$$\text{so } \mathbb{E}\left[e^{\frac{x(X_k(t) - \mu_k) - \frac{x^2}{2}}{1}} | \mathcal{F}_{t-1} \right] \leq 1.$$

$$\begin{aligned} \text{ie } \mathbb{E}[M_t | \mathcal{F}_{t-1}] &\leq (\mathbb{1}_{a_t=k} + \mathbb{1}_{a_t \neq k}) M_{t-1} \\ &\leq M_{t-1}. \end{aligned}$$

$$\text{so we showed } \mathbb{E}\left[e^{xZ_t - \frac{x^2}{2} N_k(t)} \right] \leq 1.$$

2) We now prove that $\forall \epsilon > 0, \forall n \geq 1,$

$$\mathbb{P}(Z_t \geq \epsilon \text{ and } N_k(t) = n) \leq e^{-\frac{2\epsilon^2}{n}}$$

Indeed, by Markov-Chernoff bounding for any $x > 0$:

$$\begin{aligned} \mathbb{P}(Z_t \geq \epsilon \text{ and } N_k(t) = n) &\leq e^{-x\epsilon} \mathbb{E}\left[e^{xZ_t} \mathbb{1}_{\{N_k(t) = n\}} \right] \\ &= e^{-x\epsilon + \frac{x^2}{2}n} \mathbb{E}\left[e^{xZ_t - \frac{x^2}{2}N_k(t)} \mathbb{1}_{\{N_k(t) = n\}} \right] \\ &\leq e^{-x\epsilon + \frac{x^2}{2}n} \mathbb{E}\left[e^{xZ_t - \frac{x^2}{2}N_k(t)} \right] \\ &\leq e^{-x\epsilon + \frac{x^2}{2}n} \quad \left(\leq 1 \text{ thanks to 1) } \right) \end{aligned}$$

Taking $x = \frac{4\epsilon}{n}$ finally yields

$$\mathbb{P}(Z_t \geq \epsilon \text{ and } N_k(t) = n) \leq e^{-\frac{2\epsilon^2}{n}}.$$

3) We conclude using a union bound:

$$\begin{aligned}
 \mathbb{P}(\hat{\mu}_R(t) - \mu_R \geq \sqrt{\frac{\ln(1/\delta)}{2N_R(t)}}) &= \sum_{n=1}^r \mathbb{P}(\hat{\mu}_n(t) - \mu_n \geq \sqrt{\frac{\ln(1/\delta)}{2N_n(t)}} \text{ and } N_n(t) = n) \\
 &= \sum_{n=1}^r \mathbb{P}\left(\frac{Z_t}{N_n(t)} \geq \sqrt{\frac{\ln(1/\delta)}{2N_n(t)}} \text{ and } N_n(t) = n\right) \\
 &= \sum_{n=1}^r \mathbb{P}(Z_t \geq \sqrt{\frac{n \ln(1/\delta)}{2}} \text{ and } N_n(t) = n) \\
 &\leq \sum_{n=1}^r e^{-\ln(1/\delta)} = t\delta. \quad \square
 \end{aligned}$$

Notes on this proof:

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

X bounded random variable, U, V two g -measurable random variables with $U \leq X \leq V$ a.s.

then $\forall \gamma \in \mathbb{R}$,

$$\ln \mathbb{E}[e^{\gamma X} | g] \leq \gamma \mathbb{E}[X | g] + \frac{\gamma^2}{8} (V - U)^2$$

This can be applied to

$$Z_t = (X_n(t) - \mu_n) \mathbb{1}_{\{t \leq R\}}$$

$$g = \mathcal{F}_{t-1}$$

$$U_t = -\mu_n \mathbb{1}_{\{t \leq R\}}$$

$$V_t = (1 - \mu_n) \mathbb{1}_{\{t \leq R\}}$$

and directly entails

$$\mathbb{E}\left[e^{\gamma (X_n(t) - \mu_n) \mathbb{1}_{\{t \leq R\}}} \mid \mathcal{F}_{t-1}\right] \leq \exp\left(\frac{\gamma^2}{8} \mathbb{1}_{\{t \leq R\}}\right)$$

without the

need for the

$$1 = \mathbb{1}_{\{t \leq R\}} + \mathbb{1}_{\{t > R\}} \text{ trick used in step 1.}$$

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

but $V_t - U_t = \mathbb{1}_{\{a_t = 1\}}$ can only be bounded by $\Delta_t = 1$ here, so steps 2) and 3) are still needed.

For unbounded, but σ -sub-Gaussian variables $X_k(t)$, we still have:
$$P(\mu_k - \hat{\mu}_k(t) \geq \sqrt{\frac{2 \ln(1/\delta)}{\sum_{s=1}^t N_k(s)}}) \leq \delta.$$

ϵ -Greedy

sequence of probabilities ϵ_t .

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

$a_t = t$

For $t \geq K+1$:

with proba ϵ_t , $a_t \sim \mathcal{U}([K])$

explore uniformly at random

with proba $1 - \epsilon_t$, $a_t \in \arg \max_{k \in [K]} \hat{\mu}_k(t-1)$

Theorem

For $\epsilon_t = \min \left\{ 1, \frac{cK}{t\Delta^2} \right\}$ where c is a large enough universal constant, ϵ -greedy satisfies for a large enough universal constant c'

$$R_T \leq \frac{c'}{\Delta^2} \sum_{k=1}^K (\Delta_k \ln T + 1)$$

Proof: For any k with $\Delta_k > 0$,

$$P(a_t = k) \leq \frac{\epsilon_t}{K} + P(\hat{\mu}_k(t-1) \geq \hat{\mu}_k^*(t-1))$$

$$\leq \frac{\epsilon_t}{K} + P\left(\hat{\mu}_k(t-1) - \mu_k \geq \frac{\Delta_k}{2}\right) + P\left(\mu_k^* - \hat{\mu}_k^*(t-1) \geq \frac{\Delta_k}{2}\right)$$

$$P\left(\hat{\mu}_k(t-1) - \mu_k \geq \frac{\Delta_k}{2}\right) = \sum_{n=1}^{t-1} P\left(\hat{\mu}_k(t-1) - \mu_k \geq \frac{\Delta_k}{2} \text{ and } N_k(t-1) = n\right)$$

$$\leq \sum_{n=1}^{L(t)} P(N_k(t-1) = n) + \sum_{n=L(t)+1}^{t-1} P\left(\hat{\mu}_k(t-1) - \mu_k \geq \frac{\Delta_k}{2} \text{ and } N_k(t-1) = n\right)$$

where $0 \leq x \leq t-1$

$$\leq P(N_k(t-1) \leq x_t) + \sum_{n=L(t)+1}^{t-1} e^{-\frac{x \Delta_k}{2} n + \frac{x^2}{2} n} \mathbb{E}\left[e^{x Z_{t-1} - \frac{x^2}{2} N_k(t-1)} \mathbb{1}_{\{N_k(t-1) = n\}}\right]$$

similarly to previous proof

with $Z_{t-1} = \sum_{s=2}^{t-1} (X_k(s) - \mu_k) \mathbb{1}_{\{a_s = k\}}$

for any $x > 0$

$$\leq P(N_k(t-1) \leq x_t) + \sum_{n=L(t)+1}^{t-1} e^{-\frac{\Delta_k^2}{2} n}$$

$\mathbb{E}\left[e^{x Z_{t-1} - \frac{x^2}{2} N_k(t-1)}\right] \leq 1$
and taking $x = \frac{\Delta_k}{2}$

$$\leq P\left(N_k^R(t-1) \leq x_t\right) + \frac{e^{-\frac{\Delta_k^2}{2} x_t}}{\Delta_k^2}$$

number of times k is pulled at random
(ie following the ϵ -greedy)

$$\mathbb{E}[N_k^R(t-1)] = 1 + \frac{1}{K} \sum_{s=k+1}^{t-1} \epsilon_s = \frac{1}{K} \sum_{s=1}^{t-1} \epsilon_s$$

$$\text{Var}(N_k^R(t-1)) = \sum_{s=k+1}^{t-1} \frac{\epsilon_s(1-\frac{\epsilon_s}{K})}{K} \leq \frac{1}{K} \sum_{s=1}^{t-1} \epsilon_s$$

Recall

Bernstein Inequality

Let X_1, \dots, X_T be random variables in $[0, 1]$ s.t. $\text{Var}[X_s | X_1, \dots, X_{s-1}] = \sigma_s^2$.

Then for all $\varepsilon > 0$:
$$\mathbb{P}\left(\sum_{s=1}^T X_s - \mathbb{E}[X_s | X_1, \dots, X_{s-1}] \leq -\varepsilon\right) \leq \exp\left(\frac{-\varepsilon^2/2}{\sum_{s=1}^T \sigma_s^2 + \frac{\varepsilon}{2}}\right)$$

so here for $x_t = \frac{1}{2K} \sum_{s=1}^{t-1} \varepsilon_s$

$$\begin{aligned} \mathbb{P}(N_2^R(t-1) \leq x_t) &= \mathbb{P}\left(N_2^R(t-1) - \mathbb{E}[N_2^R(t-1)] \leq -x_t\right) \\ &\leq \exp\left(-\frac{x_t^2/2}{\frac{5}{2} x_t}\right) = e^{-\frac{x_t}{5}} \end{aligned}$$

Moreover:
$$x_t = \frac{1}{2K} \sum_{s=1}^{t-1} \min\left(1, \frac{cK}{s\Delta^2}\right)$$

$\forall a \geq \lfloor \frac{\Delta^2}{cK} \rfloor + 1$:
$$x_t = \lfloor \frac{cK}{\Delta^2} \rfloor \cdot \frac{1}{2K} + \frac{1}{2K} \sum_{s=\lfloor \frac{cK}{\Delta^2} \rfloor + 1}^{t-1} \frac{cK}{s\Delta^2}$$

$$\geq \left(\frac{c}{2\Delta^2} - \frac{\delta/K}{2\Delta^2}\right) + \frac{c}{2\Delta^2} \ln\left(\frac{t-1}{\lfloor \frac{cK}{\Delta^2} \rfloor}\right)$$

$$\int_{a-1}^b \frac{1}{s} ds \geq \int_{a-1}^b \frac{1}{s} ds$$

$$x_t \geq \frac{c-1}{2\Delta^2} \ln\left(\frac{e(t-1)\Delta^2}{cK}\right)$$

Recap:
$$\mathbb{P}(a_t = b) \leq \frac{\varepsilon_t}{K} + \mathbb{P}\left(\hat{\mu}_a(t-1) - \mu_a \geq \frac{\Delta_k}{2}\right) + \mathbb{P}\left(\mu_a - \hat{\mu}_a(t-1) \geq \frac{\Delta_k}{2}\right)$$

with

$$IP(\hat{\mu}_n(t-1) - \mu_n \geq \frac{\Delta_k}{2}) \leq e^{-\frac{\gamma t}{5}} + \frac{2e^{-\frac{\Delta_k^2 \gamma t}{2}}}{\Delta_k^2}$$

$$\text{and } \gamma t \geq \frac{c-1}{2\Delta^2} \ln\left(\frac{e(t-1)\Delta^2}{cK}\right)$$

$$\text{So } P(a_t = k) \leq \frac{c}{\Delta^2 t} + 2e^{-\frac{\gamma t}{5}} + \frac{4}{\Delta_k^2} e^{-\frac{\Delta_k^2 \gamma t}{2}}$$

$$\leq \frac{c}{\Delta^2 t} + 2 \left(\frac{cK}{e(t-1)\Delta^2} \right)^{\frac{c-1}{40}}$$

$$+ \frac{4}{\Delta_k^2} \left(\frac{cK}{e(t-1)\Delta^2} \right)^{\frac{c-1}{4}} \left(\frac{\Delta_k^2}{\Delta^2} \right)$$

for $c \geq 11$: $P(a_t = k) \leq \frac{c}{\Delta^2 t} + \frac{c'}{\Delta^2} \frac{\ln(a(t-1))}{(a(t-1))^2}$ $a = \frac{e\Delta^2}{cK}$
for some constant c' .

$$\sum_{t=\lfloor \frac{T}{a} \rfloor}^T P(a_t = k) \leq \frac{C_1}{\Delta^2} \ln(T) + \frac{C_2}{\Delta^2}$$

So that for a universal constant C'' large enough

↳ does not depend on any parameter $K, T, \mu, \Delta, \text{etc.}$

$$R_T \leq \left\lceil \frac{1}{a} \right\rceil + \sum_{k=2}^K \frac{C''}{\Delta^2} \Delta_k \ln(T)$$

$$\leq \frac{C''}{\Delta^2} \sum_{k=1}^K (\Delta_k \ln(T) + 1)$$

□

Remarks • the bound above is called instance dependent as it heavily relies on parameters of the instance Δ_k

A different choice of ϵ_t can lead to the following distribution-free bound for ϵ -greedy:

$$R_T \leq O\left((K \ln T)^{1/3} T^{2/3}\right)$$

see exercise session #2

• the instance dependent bound requires a priori knowledge of Δ , which is usually unknown.