

Online Learning Lecture Notes

Pierre Gaillard

These notes explore fundamental ideas in online learning, where data are processed in real-time, and algorithms are updated dynamically. Topics include online linear and convex optimization, as well as multi-armed bandits. The main algorithms in the field will be introduced, and we will delve into regret minimization concepts for theoretical analysis. Online learning algorithms play a central role in recent advancements in reinforcement learning.

Useful information

- **Pierre Gaillard** (INRIA Grenoble)
- Email: pierre.gaillard@inria.fr
- Relevant references: Cesa-Bianchi and Lugosi [2006], Shalev-Shwartz et al. [2012], Hazan et al. [2016], Lattimore and Szepesvári [2020]
- Content of the class: mostly theoretical (algorithms and proofs), sequential learning with adversarial data, stochastic bandits, adversarial bandits

Contents

1	Introduction	3
2	Online Linear Optimization	6
2.1	The exponentially weighed average forecaster	6
2.2	Application to prediction with expert advice	9
2.3	Linearizing the loss through randomization	12
3	Online Convex Optimization	14
3.1	The Gradient Trick (from linear to convex losses)	14
3.2	Discretized EWA	15
3.3	Online gradient descent	17
3.4	Online Mirror Descent	18
4	Adversarial Bandits	20
4.1	The exponentially weighted average algorithm for bandits	20
4.1.1	Pseudo-regret bound	20
4.1.2	High probability bound on the regret	23
4.2	Adversarial bandits with experts	26
4.3	Adversarial Bandits with side information	28
5	Stochastic multi-armed bandits	30
5.1	Setting: stochastic bandit with finitely many actions	30
5.2	Explore-Then-Commit (ETC)	32
5.3	Upper-Confidence-Bound (UCB)	33
5.4	Other algorithms	36
5.4.1	ϵ -greedy	36
5.4.2	Thompson Sampling	37
5.5	Lower bounds for multi-armed bandit	37
5.5.1	Distribution-free lower bound	37
5.5.2	Distribution-dependent lower bound	37
6	Contextual bandits	39
6.1	Continuous stochastic bandits	39
6.1.1	Contextual bandits through discretization	40
6.1.2	Stochastic Linear bandits	41
6.2	Other possible extensions of bandits	46

1 Introduction

In many applications, the data set is not available from the beginning to learn a model but it is observed sequentially as a flow of data. Furthermore, the environment may be so complex that it is unfeasible to choose a comprehensive model and use classical statistical theory and optimization. A classic example is the spam detection which can be seen as a game between spammer and spam filters. Each trying to fool the other one. Another example, is the prediction of processes that depend on human behaviors such as the electricity consumption. These problems are often not adversarial games but cannot be modeled easily and are surely not i.i.d.

There is a necessity to take a robust approach by using a method that learns as ones goes along, learning from experiences as more aspects of the data and the problem are observed. This is the goal of online learning. The curious reader can know more about online learning in the books Cesa-Bianchi and Lugosi [2006], Hazan et al. [2016], Shalev-Shwartz et al. [2012].

Setting In online learning, a player sequentially makes decisions based on past observations. After committing the decision, the player suffers a loss (or receives a reward depending on the problem). Every possible decision incurs a (possibly different) loss. The losses are unknown to the player beforehand and may be arbitrarily chosen by some adversary. More formally, an online learning problem can be formalized as in Figure 1.

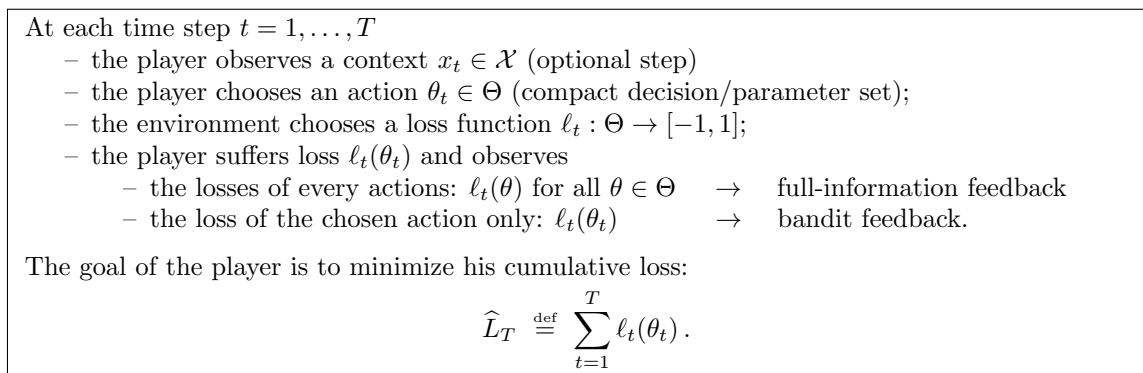


Figure 1: Setting of an online learning problem/online convex optimization

Example 1.1 (Multi-armed bandit). In K -armed bandit, the decision set are K actions (or arms) $\Theta = \{1, \dots, K\}$ and the player only observes the performance of the chosen action (bandit feedback). In this problem, there is an exploration-exploitation trade-off: the player wants to select the best arm as often as possible but he also needs to explore all arms to estimate their performance.

This problem takes his name from slot machines (also known as one-armed bandits because they were originally operated by one lever on the side of the machine) in which some player explores several slot machines and tries to maximize his cumulative gain (or more likely minimize his loss!).

Originally, multi-armed bandit setting was introduced by Thompson in 1933 and motivated by clinical trials. For the t -th patient in some clinical study, one needs to choose the treatment to assign to this patient and observe the response. The goal is to maximize the number of patients healed during the study.

Nowadays, multi-armed bandit is motivated by many applications coming from internet (recommender systems, online advertisements, ...). We will see more on multi-armed bandit in next lectures.

Example 1.2 (Prediction with expert advice). In prediction with expert advice, there is some sequence of observations $y_1, \dots, y_T \in [0, 1]$ to be predicted step by step with the help of expert forecasts. The setting can be formalized as follows: at each time step $t \geq 1$

- the environment reveals experts forecasts $x_t(k)$ for $k = 1, \dots, K$
- the player chooses a weight vector $p_t \in \Delta_K \stackrel{\text{def}}{=} \{p \in [0, 1]^K : \sum_{k=1}^K p_k = 1\}$
(here θ_t is denoted p_t and $\Theta = \Delta_K$)
- the player forecasts $\hat{y}_t = \sum_{k=1}^K p_t(k)x_t(k)$
- the environment reveals $y_t \in [0, 1]$ and the player suffers loss $\ell_t(p_t) = \ell(\hat{y}_t, y_t)$ where $\ell : [0, 1]^2 \rightarrow [0, 1]$ is a loss function.

Considering $\Theta := \Delta_K$ and $\theta_t := p_t$, this setting can be recovered by the online learning setting of Figure 1. The inputs correspond to the expert advice $x_t(k)$ that are often revealed before the learner makes his decision p_t .

Player's performance is then measured via a loss function $\ell_t(p_t) = \ell(\hat{y}_t, y_t)$ which measures the distance between the prediction \hat{y}_t and the output y_t . Typical loss functions are the squared loss $\ell(\hat{y}_t, y_t) = (\hat{y}_t - y_t)^2$, the absolute loss $\ell(\hat{y}_t, y_t) = |\hat{y}_t - y_t|$ or the absolute percentage of error $\ell(\hat{y}_t, y_t) = |\hat{y}_t - y_t|/|y_t|$. All these loss functions are convex, which will play an important role in the analysis.

How to measure the performance: the regret Of course, if the environment chooses large losses $\ell_t(x)$ for all decisions $\theta \in \Theta$, it is impossible for the player to ensure small cumulative loss. Therefore, one needs a relative criterion: the regret of the player is the difference between the cumulative loss he incurred and that of the best fixed decision in hindsight.

Definition 1 (Regret). *The regret of the player with respect to a fixed parameter $\theta^* \in \Theta$ after T time steps is*

$$R_T(\theta^*) \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(\theta^*).$$

The regret (or uniform regret) is defined as $R_T \stackrel{\text{def}}{=} \sup_{\theta^* \in \Theta} R_T(\theta^*)$.

We have some bias-variance decomposition:

$$\sum_{t=1}^T \ell_t(\theta_t) = \underbrace{\inf_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)}_{\text{Approximation error = how good the possible actions are.}} + \underbrace{R_T}_{\text{Sequential estimation error of the best action}}$$

We will focus on the regret in these lectures. The goal of the player is to ensure a sublinear regret $R_T = o(T)$ as $T \rightarrow \infty$ and this for any possible sequence of losses ℓ_1, \dots, ℓ_T . In this case, the average performance of the player will approach on the long term the one of the best decision.

Remarks Let us makes some remarks:

- Except in the stochastic bandit part, we will not make any random assumption on the process generating the losses ℓ_t . The latter are deterministic and may be chosen by some adversary. Typically, the problem can be seen as a game between the player who aims at optimizing with respect to $\theta_1, \dots, \theta_T$ against an environment who aims at maximizing with respect to ℓ_t, \dots, ℓ_T and θ^* . Player's goal is to approach the quantity:

$$\inf_{\theta_1} \sup_{\ell_1} \inf_{\theta_2} \sup_{\ell_2} \dots \inf_{\theta_T} \sup_{\ell_T} \sup_{\theta^* \in \Theta} R_T(\theta^*).$$

- Note that the loss functions ℓ_t depend on the round t . This may be caused by many phenomena. We provide here some possible reasons. This may be because
 - of some observation to be predicted if $\ell_t(x) = \ell(x, y_t)$. For instance, if the goal is to predict the evolution of the temperature y_1, \dots, y_T , the latter changes over time and a prediction x is evaluated with $\ell_t(x) = (x - y_t)^2$.
 - the environment is stochastic and the variation over time t models some noise effect.

- of a changing environment. For instance, if the player is playing a game against some adversary that evolves and adapts to its strategy. A typical example is the case of spam detections. If the player tries to detect spams, while some spammers (the environment) try at the same time to fool the player with new spam strategies.

Exercise 1.1. Instead considering the regret with respect to a fixed $\theta^* \in \Theta$, one would be tempted to minimize the quantity

$$R_T^* \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \inf_{\theta \in \Theta} \ell_t(\theta)$$

where the infimum is inside the sum. Show that the environment can ensure R_T^* to be linear in T by choosing properly the loss functions ℓ_t .

2 Online Linear Optimization

We will start with the simple case where the decision set Θ is the K -dimensional simplex

$$\Delta_K \stackrel{\text{def}}{=} \{p \in [0, 1]^K : \sum_{k=1}^K p_k = 1\}. \quad (\text{decision set})$$

Since the decisions θ_t are probability distributions in $\Theta = \Delta_K$, in this part we will denote them by p_t instead of θ_t . We assume that the loss functions ℓ_t are linear

$$\forall p \in \Theta, \quad \ell_t(p) = \sum_{k=1}^K p(k)g_t(k) \in [-1, 1] \quad (\text{linear loss})$$

where $g_t = (g_t(1), \dots, g_t(K)) \in [-1, 1]^K$ is a loss vector chosen by the environment at round t .

2.1 The exponentially weighed average forecaster

How to choose the weights p_t ? At round t the player needs to choose a weight vector $p_t \in \Delta_K$. The question is how to choose it? The idea is to give more weight to actions that performed well in the past. But we should not give all the weight to the current best action, otherwise it would not work (see exercises). The exponentially weighted average forecaster (EWA) also called Hedge performs this trade-off by choosing a weight that decreases exponentially fast with the past errors.

The Exponentially weighted average forecaster (EWA)
 Parameter: $\eta > 0$
 Initialize: $p_1 = (\frac{1}{K}, \dots, \frac{1}{K})$
 For $t = 1, \dots, T$
 - select p_t ; incur loss $\ell_t(p_t) = p_t^\top g_t$ and observe $g_t \in [-1, 1]^K$;
 - update for all $k \in \{1, \dots, K\}$

$$p_{t+1}(k) = \frac{e^{-\eta \sum_{s=1}^t g_s(k)}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^t g_s(j)}}.$$

Exercise 2.1. Consider the strategy, called “Follow The Leader” (FTL) that puts all the mass on the best action so far:

$$p_t \in \arg \min_{p \in \Theta} \sum_{s=1}^{t-1} \ell_s(p). \quad (\text{FTL})$$

1. Show that $p_t(k) > 0$ implies that $k \in \arg \min_j \sum_{s=1}^{t-1} g_s(j)$
2. Show that the regret of FTL might be linear: i.e., there exists a sequence $g_1, \dots, g_T \in [-1, 1]^K$ such that $R_T \geq (1 - 1/K)T$.

The following theorem proves that EWA, which is a smoothed version of FTL, achieves sublinear regret.

Theorem 1. Let $T \geq 1$. For all sequences of loss vectors $g_1, \dots, g_T \in [-1, 1]^K$, EWA achieves the bound

$$R_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t) - \min_{p \in \Delta_K} \sum_{t=1}^T \ell_t(p) \leq \eta \sum_{t=1}^T \sum_{k=1}^K p_t(k)g_t(k)^2 + \frac{\log K}{\eta}, \quad (1)$$

where we recall $\ell_t : p \in \Delta_K \mapsto p^\top g_t$. Therefore, for the choice $\eta = \sqrt{\frac{\log K}{T}}$, EWA satisfies the regret bound $R_T \leq 2\sqrt{T \log K}$.

This regret bound is optimal (see Cesa-Bianchi and Lugosi [2006]).

Exercise 2.2. Generalize the above theorem when the losses $\ell_1, \dots, \ell_T \in [-B, B]$ for some $B > 0$.

Proof. First, we remark that by definition of $\ell_t : p \mapsto p \cdot g_t$ we have

$$\begin{aligned} R_T &\stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t) - \min_{p \in \Delta_K} \sum_{t=1}^T \ell_t(p) \\ &= \sum_{t=1}^T p_t \cdot g_t - \min_{p \in \Delta_K} \sum_{t=1}^T p \cdot g_t \\ &= \sum_{t=1}^T p_t \cdot g_t - \min_{p \in \Delta_K} \sum_{k=1}^K \sum_{t=1}^T p(k) g_t(k). \end{aligned}$$

Now, we can see that the minimum over $p \in \Delta_K$ is reached on a corner of the simplex. Therefore

$$R_T = \sum_{t=1}^T p_t \cdot g_t - \min_{1 \leq k \leq K} \sum_{t=1}^T g_t(k).$$

We denote $W_t(j) = e^{-\eta \sum_{s=1}^t g_s(j)}$ and $W_t = \sum_{j=1}^K W_t(j)$. The proof will consist in upper-bounding and lower-bounding W_T . We have

$$\begin{aligned} W_t &= \sum_{j=1}^K W_{t-1}(j) e^{-\eta g_t(j)} && \leftarrow W_t^{(j)} = W_{t-1}(j) e^{-\eta g_t(j)} \\ &= W_{t-1} \sum_{j=1}^K \frac{W_{t-1}(j)}{W_{t-1}} e^{-\eta g_t(j)} \\ &= W_{t-1} \sum_{j=1}^K p_t(j) e^{-\eta g_t(j)} && \leftarrow p_t(j) = \frac{e^{-\eta \sum_{s=1}^{t-1} g_s(j)}}{\sum_{k=1}^K e^{-\eta \sum_{s=1}^{t-1} g_s(k)}} = \frac{W_{t-1}(j)}{W_{t-1}} \\ &\leq W_{t-1} \sum_{j=1}^K p_t(j) (1 - \eta g_t(j) + \eta^2 g_t(j)^2) && \leftarrow e^x \leq 1 + x + x^2 \text{ for } x \leq 1 \\ &= W_{t-1} (1 - \eta p_t \cdot g_t + \eta^2 p_t \cdot g_t^2), \end{aligned}$$

where we assumed in the inequality $-\eta g_t(j) \leq 1$ and where we denote $g_t = (g_t(1), \dots, g_t(K))$, $g_t^2 = (g_t(1)^2, \dots, g_t(K)^2)$ and $p_t = (p_t(1), \dots, p_t(K))$. Now, using $1 + x \leq e^x$, we get:

$$W_t \leq W_{t-1} \exp(-\eta p_t \cdot g_t + \eta^2 p_t \cdot g_t^2).$$

By induction on $t = 1, \dots, T$, this yields using $W_0 = K$

$$W_T \leq K \exp\left(-\eta \sum_{t=1}^T p_t \cdot g_t + \eta^2 \sum_{t=1}^T p_t \cdot g_t^2\right). \quad (2)$$

On the other hand, upper-bounding the maximum with the sum,

$$\exp\left(-\eta \min_{j \in [K]} \sum_{t=1}^T g_t(j)\right) \leq \sum_{j=1}^K \exp\left(-\eta \sum_{t=1}^T g_t(j)\right) \leq W_T.$$

Combining the above inequality with Inequality (2) and taking the log, we get

$$-\eta \min_{j \in [K]} \sum_{t=1}^T g_t(j) \leq -\eta \sum_{t=1}^T p_t \cdot g_t + \eta^2 \sum_{t=1}^T p_t \cdot g_t^2 + \log K. \quad (3)$$

Dividing by η and reorganizing the terms proves the first inequality:

$$R_T = \sum_{t=1}^T p_t \cdot g_t - \min_{1 \leq j \leq K} \sum_{t=1}^T g_t(j) \leq \eta \sum_{t=1}^T p_t \cdot g_t^2 + \frac{\log K}{\eta}$$

Optimizing η and upper-bounding $p_t \cdot g_t^2 \leq 1$ concludes the second inequality. \square

Anytime algorithm (the doubling trick) The previous algorithm EWA depends on a parameter $\eta > 0$ that needs to be optimized according to K and T . For instance, for EWA using the value

$$\eta = \sqrt{\frac{\log K}{T}}.$$

the bound of Theorem 1 is only valid for horizon T . However, the learner might not know the time horizon in advance and one might want an algorithm with guarantees valid simultaneously for all $T \geq 1$. We can avoid the assumption that T is known in advance, at the cost of a constant factor, by using the so-called *doubling trick*. The general idea is the following. Whenever we reach a time step t which is a power of 2, we restart the algorithm (forgetting all the information gained in the past) setting η to $\sqrt{\log K/t}$. Let us denote EWA-doubling this algorithm.

Theorem 2 (Anytime bound on the regret). *For all $T \geq 1$, the pseudo-regret of EWA-doubling is then upper-bounded as:*

$$R_T \leq 7\sqrt{T \log K}.$$

The same trick can be used to turn most online algorithms into anytime algorithms (even in more general settings: bandits, general loss, ...). We can use the *doubling trick* whenever we have an algorithm with a regret of order $\mathcal{O}(T^\alpha)$ for some $\alpha > 0$ with a known horizon T to turn it into an algorithm with a regret $\mathcal{O}(T^\alpha)$ for all $T \geq 1$.

Another solution is to use time-varying parameters η_t replacing T with the current value of t . The analysis is however less straightforward.

Exercise 2.3. Prove a regret bound for the time-varying choice $\eta_t = \sqrt{\log K/t}$ in EWA.

Proof of Theorem 2. For simplicity we assume $T = 2^{M+1} - 1$. The regret of EWA-doubling is then upper-bounded as:

$$\begin{aligned} R_T &= \sum_{t=1}^T \ell_t(p_t) - \min_{p \in \Delta_K} \sum_{t=1}^T \ell_t(p) \\ &\leq \sum_{t=1}^T \ell_t(p_t) - \sum_{m=0}^M \min_{p \in \Delta_K} \sum_{t=2^m}^{2^{m+1}-1} \ell_t(p) \\ &= \sum_{m=0}^M \underbrace{\sum_{t=2^m}^{2^{m+1}-1} \ell_t(p_t) - \min_{p \in \Delta_K} \sum_{t=2^m}^{2^{m+1}-1} \ell_t(p)}_{R_m}. \end{aligned}$$

Now, we remark that each term R_m corresponds to the expected regret of an instance of EWA over the 2^m rounds $t = 2^m, \dots, 2^{m+1} - 1$ and run with the optimal parameter $\eta = \sqrt{\log K / 2^m}$. Therefore, using Theorem 1, we get $R_m \leq 2\sqrt{2^m \log K}$, which yields:

$$R_T \leq \sum_{m=0}^M 2\sqrt{2^m \log K} \leq 2(1 + \sqrt{2})\sqrt{2^{M+1} \log K} \leq 7\sqrt{T \log K}.$$

□

Improvement for small losses The first inequality in Theorem 1 is sometimes called improvement for small losses when losses take values in $[0, 1]$. Let's define $\widehat{L}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t)$ the loss of the algorithm and $L_T^* \stackrel{\text{def}}{=} \min_{p \in \Delta_K} \sum_{t=1}^T \ell_t(p)$. Then, the regret is upper-bounded by

$$\begin{aligned} R_T \stackrel{\text{def}}{=} \widehat{L}_T - L_T^* &\leq \frac{\log K}{\eta} + \eta \sum_{t=1}^T p_t \cdot g_t^2 \\ &\leq \frac{\log K}{\eta} + \eta \sum_{t=1}^T p_t \cdot g_t = \frac{\log K}{\eta} + \eta \widehat{L}_T. \end{aligned}$$

Therefore, rearranging the terms

$$(1 - \eta)\widehat{L}_T - (1 - \eta)L_T^* \leq \frac{\log K}{\eta} + \eta L_T^*,$$

which implies

$$R_T \leq \frac{\log K}{\eta(1 - \eta)} + \frac{\eta}{1 - \eta} L_T^*.$$

Optimising in $\eta \approx \sqrt{(\log K) / L_T^*}$ we get $R_T \lesssim \sqrt{(\log K) L_T^*}$ which is small whenever some parameter achieves a small cumulative loss.

2.2 Application to prediction with expert advice

The preceding section considers linear loss functions. Yet, it can yield non-trivial regret bounds for general convex losses. We consider here an application to the setting of prediction with expert advice detailed in Example 1.2. The goal is to minimize the regret with respect to the best expert

$$R_T^{\text{expert}} \stackrel{\text{def}}{=} \sum_{t=1}^T \ell(\widehat{y}_t, y_t) - \min_{1 \leq k \leq K} \sum_{t=1}^T \ell(x_t(k), y_t),$$

where $\widehat{y}_t = p_t \cdot x_t$ are the prediction of the algorithm and y_t the observations to be predicted sequentially.

Convex loss function ℓ . We state below a corollary to Theorem 1 when the loss functions $\ell(\cdot, \cdot)$ are convex in their first argument.

Corollary 1 (Regret of EWA for prediction with expert advice and convex loss). *Let $T \geq 1$. Assume that the loss function $\ell : (x, y) \in \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$ is convex and takes values in $[-1, 1]$. Then, EWA applied with the vector vectors $g_t = (\ell(x_t(1), y_t), \dots, \ell(x_t(K), y_t)) \in [-1, 1]^K$ has a regret upper-bounded by*

$$R_T^{\text{expert}} \leq 2\sqrt{T \log K}$$

where $\widehat{y}_t = p_t \cdot x_t$ and were $\eta > 0$ is well-tuned.

Therefore, the average error of the algorithm will converge to the average error of the best expert. This is the case for the square loss, the absolute loss or the absolute percentage of error.

Proof. It suffices to remark that by convexity of $\ell(\cdot, \cdot)$ in its first argument

$$\begin{aligned} R_T^{\text{expert}} &= \sum_{t=1}^T \ell(p_t \cdot x_t, y_t) - \min_{1 \leq k \leq K} \sum_{t=1}^T \ell(x_t(k), y_t) \\ &\leq \sum_{t=1}^T p_t \cdot g_t - \min_{1 \leq k \leq K} \sum_{t=1}^T g_t(k) \stackrel{\text{def}}{=} R_T. \end{aligned}$$

The result is then obtained by Theorem 1. □

Exp-concave loss function Here, we show that a faster rate can be obtained (with EWA) if the loss function are exp-concave.

Definition 2 (η -exp-concavity). For $\eta \in \mathbb{R}$, a function f is said to be η -exp-concave if $x \mapsto e^{-\eta f(x)}$ is concave.

Exp-concavity is stronger than convexity but weaker than strong convexity. Indeed, exp-concave functions are convex because $-\log$ is convex and decreasing. Furthermore, any η -exp-concave function is also η' -exp-concave for $0 \leq \eta' \leq \eta$.

In prediction with expert advice, if the loss are generated from a fixed loss function $\ell_t(p) = \ell(p \cdot \ell_t, y_t)$, then ℓ_t are η -exp-concave if $\hat{y} \mapsto \ell(\hat{y}, y_t)$ are η -exp-concave for all y_t . We can compute η for some common loss functions:

- the squared loss: $\ell : (\hat{y}, y) \in [0, 1]^2 \mapsto (\hat{y} - y)^2$, then ℓ_t are $1/2$ -exp-concave. Indeed, let $y \in [0, 1]$ and denote $G : \hat{y} \mapsto \exp(-\eta(\hat{y} - y)^2)$. Then, $G''(\hat{y}) = G(\hat{y})(4\eta^2(\hat{y} - y)^2 - 2\eta)$. Thus G is concave if and only if $(\hat{y} - y)^2 \leq 1/(2\eta)$ which is satisfied for $\eta = 1/2$. This is also the case in higher dimensions with $\ell(\hat{y}, y) = \|\hat{y} - y\|^2$. If the observations and prediction $\hat{y}, y \in [0, B]$, then the ℓ_t are $1/(2B^2)$ -exp-concave
- the relative entropy (or Kullback–Leibler divergence): $\ell : (\hat{y}, y) \in [0, 1]^2 \mapsto y \log(y/\hat{y}) - (1 - y) \log((1 - y)/(1 - \hat{y}))$. Then the functions ℓ_t are 1 -exp-concave. This loss can for instance used for density estimation of the sequence y_1, \dots, y_T .
- the linear loss $\ell(\hat{y}, y) = \hat{y} \cdot y$, the absolute loss $\ell(\hat{y}, y) = |\hat{y} - y|$ or the absolute percentage of error are however not η -exp-concave for any $\eta > 0$.

Corollary 2 (Regret of EWA for prediction with expert advice and exp-concave loss). In the setting of prediction with expert advice, if the loss functions $\ell(\cdot, y_t)$ are η -exp-concave for all y_t , then EWA run with vectors $g_t = (\ell(x_t(1), y_t), \dots, \ell(x_t(K), y_t)) \in \mathbb{R}^K$ with parameter $\eta > 0$ satisfies

$$R_T^{\text{expert}} \leq \frac{\log K}{\eta},$$

for all $T \geq 1$.

The worst-case regret does not increase with T but grows logarithmically in the dimension K .

Proof. The proof is similar to the original proof of EWA. We define $W_t(i) = e^{-\eta \sum_{s=1}^t g_s(i)}$ and $W_t =$

$\sum_{i=1}^K W_t(i)$. We have

$$\begin{aligned}
W_t &= \sum_{j=1}^N W_{t-1}(j) e^{-\eta g_t(j)} && \leftarrow W_t(j) = W_{t-1}(j) e^{-\eta g_t(j)} \\
&= W_{t-1} \sum_{j=1}^N \frac{W_{t-1}(j)}{W_{t-1}} e^{-\eta g_t(j)} \\
&= W_{t-1} \sum_{j=1}^N p_t(j) e^{-\eta g_t(j)} && \leftarrow p_t(j) = \frac{e^{-\eta \sum_{s=1}^{t-1} g_s(j)}}{\sum_{k=1}^N e^{-\eta \sum_{s=1}^{t-1} g_s(k)}} = \frac{W_{t-1}(j)}{W_{t-1}} \\
&\leq W_{t-1} \exp(-\eta \ell(p_t \cdot x_t, y_t)) && \leftarrow \eta\text{-exp-concavity}
\end{aligned}$$

Now, by induction on $t = 1, \dots, T$, this yields using $W_0 = K$

$$W_T \leq K \exp\left(-\eta \sum_{t=1}^T \ell(\hat{y}_t, y_t)\right). \quad (4)$$

On the other hand, upper-bounding the maximum with the sum,

$$\exp\left(-\eta \min_{j \in [K]} \sum_{t=1}^T g_t(j)\right) \leq \sum_{j=1}^K \exp\left(-\eta \sum_{t=1}^T g_t(j)\right) \leq W_T.$$

Combining the above inequality with Inequality (4) and taking the log concludes the proof. \square

Continuous EWA Similarly as for Section 2.2, Theorem 2 does not really control the true regret R_T since it controls the regret with respect to Dirac mass $\min_{1 \leq k \leq K} \sum_{t=1}^T \ell_t(\delta_k)$ instead of the one with respect to all convex combinations $\min_{p \in \mathcal{X}} \sum_{t=1}^T \ell_t(p)$. A true upper-bound on the regret can be obtained by using a continuous version of EWA:

$$p_t = \frac{\int_{\mathcal{X}} p e^{-\eta \sum_{s=1}^{t-1} \ell_s(p)} d\mu(p)}{\int_{\mathcal{X}} e^{-\eta \sum_{s=1}^{t-1} \ell_s(p)} d\mu(p)},$$

where μ is the uniform (Lebesgue) measure on $\mathcal{X} = \Delta_K$.

Theorem 3. *Let $T \geq 1$. For all sequences of η -exp-concave losses ℓ_1, \dots, ℓ_t the continuous EWA forecaster satisfies*

$$R_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t) - \inf_{p \in \mathcal{X}} \sum_{t=1}^T \ell_t(p) \leq \frac{1 + (K-1) \log(T+1)}{\eta}$$

Proof. The proof starts similarly to the one of Theorem 2. Let us denote $W_t(p) = e^{-\eta \sum_{s=1}^t \ell_s(p)}$, $W_t = \int_{\mathcal{X}} W_t(p) d\mu(p)$ and $d\hat{\mu}_t(p) = W_t(p) d\mu(p) / W_t$. Then,

$$\begin{aligned}
W_T &= \int_{\mathcal{X}} e^{-\eta \sum_{t=1}^T \ell_t(p)} d\mu(p) \\
&= W_{T-1} \int_{\mathcal{X}} \frac{W_{T-1}(p)}{W_{T-1}} e^{-\eta \ell_T(p)} d\mu(p) \\
&= W_{T-1} \int_{\mathcal{X}} e^{-\eta \ell_T(p)} d\hat{\mu}_{T-1}(p) && \leftarrow p_T = \int_{\mathcal{X}} p d\hat{\mu}_{T-1}(p) \\
&\leq W_{T-1} \exp(-\eta \ell_T(p_T)) && \leftarrow \eta\text{-exp-concavity} \\
&\leq \exp\left(-\eta \sum_{t=1}^T \ell_t(p_t)\right), && \leftarrow \text{induction}
\end{aligned} \quad (5)$$

The second part of the proof to lower-bound W_T is however less straightforward. For simplicity, let us assume that ℓ_t are continuous on \mathcal{X} (do the general case as exercise). Therefore the infimum is a minimum and let $p^* \in \arg \min_{p \in \mathcal{X}} \sum_{t=1}^T \ell_t(p)$ and define

$$\mathcal{X}_\varepsilon \stackrel{\text{def}}{=} \left\{ (1-\varepsilon)p^* + \varepsilon q, \quad q \in \mathcal{X} \right\}, \quad \varepsilon \in (0, 1).$$

By exconcaavity of ℓ_t , we have for all t and all $p = (1-\varepsilon)p^* + \varepsilon q$

$$e^{-\eta \ell_t(p)} \geq (1-\varepsilon)e^{-\eta \ell_t(p^*)} + \varepsilon e^{-\eta \ell_t(q)} \geq (1-\varepsilon)e^{-\eta \ell_t(p^*)}$$

Therefore, for all $p \in \mathcal{X}_\varepsilon$

$$e^{-\eta \sum_{t=1}^T \ell_t(p)} \geq (1-\varepsilon)^T e^{-\eta \sum_{t=1}^T \ell_t(p^*)}$$

Integrating both parts over \mathcal{X}_ε and using $\mu(\mathcal{X}_\varepsilon) = \varepsilon^{K-1} \mu(\mathcal{X})$ (exercise) we get

$$W_T \geq \int_{\mathcal{X}_\varepsilon} e^{-\eta \sum_{t=1}^T \ell_t(p)} d\mu(p) \geq \mu(\mathcal{X}) \varepsilon^{K-1} (1-\varepsilon)^T e^{-\eta \sum_{t=1}^T \ell_t(p^*)}.$$

Combining with (5), using $W_0 = \mu(\mathcal{X})$, taking the log and reorganizing the terms yields

$$R_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t) - \sum_{t=1}^T \ell_t(p^*) \leq \frac{(K-1) \log \frac{1}{\varepsilon} + T \log \frac{1}{1-\varepsilon}}{\eta}.$$

Optimizing $\varepsilon = 1/(T+1)$ concludes the proof since

$$T \log \frac{1}{1-\varepsilon} = T \log \left(1 + \frac{1}{T} \right) \leq 1.$$

□

Though the nice theoretical result, this algorithm is complicated to implement because of the integral. In practice, p_t can be computed by using $(1/T)$ -discretization grid of \mathcal{X} (bad complexity of order T^K !) or by using Monte-Carlo methods to approximate the integral. We will see in next lectures efficient algorithms with similar guarantees.

2.3 Linearizing the loss through randomization

Setting: Θ finite, non-convex loss functions $\ell_t : \Theta \rightarrow [-1, 1]$.

In this section, we consider a *finite set of decision* $\Theta = \{1, \dots, K\}$ and we assume that the player is restricted to play an action in Θ . In other words, the player cannot play convex combinations of the actions as it was done for prediction with expert advice. For instance, we may want to build a recommender system to recommend movies to customers. The loss function are *arbitrary bounded loss functions* $\ell_t : \Theta \rightarrow [-1, 1]$.

Need of a random strategy The following proposition shows that the choice θ_t cannot be deterministic in this setting. Otherwise, the adversary may fool the player by taking ℓ_t depending on θ_t .

Proposition 1. *Any deterministic algorithm may incur a linear regret. In other words, we can find some sequence of losses ℓ_t such that $R_T \gtrsim T$.*

Proof. Since θ_t is deterministic, the loss function ℓ_t can depend on θ_t . We then choose $\ell_t(\theta_t) = 1$ and $\ell_t(\theta) = 0$ for $\theta \neq \theta_t$. Then one of the chosen actions was picked less than T/K times so that $\max_{1 \leq k \leq K} \ell_t(k) \leq T/K$. Therefore, $R_T \geq (1 - 1/K)T$. □

From the above proposition, we see that the strategy of the learner needs to be random. Therefore, instead of choosing an action in $\{1, \dots, K\}$, the player chooses a probability distribution $p_t \in \Delta_K := \{p \in [0, 1]^K : \sum_k p_k = 1\}$ and draws $\theta_t \sim p_t$. And we recover the setting with actions played in the simplex Δ_K .

A random regret The regret R_T will be here a random quantity that depends on the randomness of the algorithm (and eventually of the data). We will thus focus on upper-bounding the regret:

- with high-probability: $R_T \leq \varepsilon$ with probability at least $1 - \delta$;
- in expectation: $\mathbb{E}[R_T] \leq \varepsilon$.

From high-probability bound to expected bound. Note that since the losses are bounded in $[0, 1]$ a bound in high probability entails a bound in expectation. If $R_T \leq \varepsilon$ with probability at least $1 - \delta$, then

$$\mathbb{E}[R_T] \leq \mathbb{E}[R_T | R_T \leq \varepsilon] \mathbb{P}(R_T \leq \varepsilon) + \mathbb{E}[R_T | R_T \geq \varepsilon] \mathbb{P}(R_T \geq \varepsilon) \leq \varepsilon + T\delta. \quad (6)$$

Another useful (and often better) tool to transform a high-probability bound into a bound in expectation is the inequality $\mathbb{E}[X] = \int_0^\infty \mathbb{P}(X \geq \varepsilon) d\varepsilon$ for nonnegative random variable X .

From expected bound to high-probability bound. On the other hand, since the losses are bounded, using Hoeffding's inequality a bound in expectation entails a bound in high probability at the cost of an additive term of order $\sqrt{T \log(1/\delta)}$ in the regret.

Proposition 2. *Drawing $\theta_t \sim p_t$ where p_t are chosen by EWA satisfies the expected regret*

$$\mathbb{E}[R_T] = \mathbb{E} \left[\sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \right] \leq 2\sqrt{T \log K}$$

for η well tuned.

Exercise 2.4. Using Hoeffding's inequality, provide a bound on the regret R_T with probability $1 - \delta$.

Proof. Using $g_t = (\ell_t(1), \dots, \ell_t(K)) \in [-1, 1]^K$, from Theorem 1 of last class, we have

$$\sum_{t=1}^T p_t \cdot g_t - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \leq 2\sqrt{T \log K}.$$

It suffices then to take the expectation and remark that

$$\mathbb{E}[\ell_t(\theta_t)] = \mathbb{E}[\mathbb{E}[\ell_t(\theta_t) | p_t]] = \mathbb{E}[p_t \cdot g_t].$$

□

It is worth pointing out that we did not make any assumption on the loss function ℓ_t beside boundedness. In particular, it can be non-convex.

Example 2.1 (Online classification). Assume that you may want to predict a sequence of labels $y_1, \dots, y_T \in \{0, 1\}$ (such as spams) based on expert advice $x_t(k) \in \{0, 1\}$ (such as different spam detectors). Then, using the losses $\ell_t(k) = \mathbf{1}_{x_t(k) \neq y_t}$, EWA ensures

$$\mathbb{E} \left[\sum_{t=1}^T \mathbf{1}_{\theta_t \neq y_t} - \min_{1 \leq k \leq K} \sum_{t=1}^T \mathbf{1}_{x_t(k) \neq y_t} \right] \leq 2\sqrt{T \log K}.$$

Hence, the expected number of mistakes of the algorithms will not be much larger than the one of the best expert. This is valid though the loss function is nonconvex.

3 Online Convex Optimization

In this section, we consider convex loss functions ℓ_t and introduce several well-known algorithms.

3.1 The Gradient Trick (from linear to convex losses)

Setting: simplex decision set $\Theta = \Delta_K$, convex and differentiable loss functions

In this section, we consider the simplex decision set $\Theta = \Delta_K$ and thus we will denote by p_t (instead of θ_t) the actions played by the player. Moreover, we assume the losses to be convex and Lipschitz

$$\forall p_t \in \Theta, \quad \|\nabla \ell_t(p_t)\|_\infty \leq G.$$

We will see a simple trick, so-called *the gradient trick* that allows to extend the results we saw for linear losses to convex losses. The resulted algorithm is called the Exponentiated Gradient forecaster (EG). It consists in playing EWA with the gradients $g_t = \nabla \ell_t \in [-G, G]^K$ as loss vectors.

Theorem 4. *Let $T \geq 1$. For all sequences of convex differentiable losses $\ell_1, \dots, \ell_T : \Theta \rightarrow \mathbb{R}$ with bounded gradient $\max_{p \in \Theta} \|\nabla \ell_t(p)\|_\infty \leq G$, EWA applied with $g_t = \nabla \ell_t$ achieves the regret bound*

$$R_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t) - \min_{p \in \Theta} \sum_{t=1}^T \ell_t(p) \leq \eta G^2 T + \frac{\log K}{\eta}. \quad (7)$$

Therefore, for the choice $\eta = \frac{1}{G} \sqrt{\frac{\log K}{T}}$, EWA satisfies the regret bound $R_T \leq 2G\sqrt{T \log K}$.

Proof. Therefore, applying the regret bound of EWA (see Theorem 1 of last class) we get

$$\sum_{t=1}^T p_t \cdot g_t - \min_{p \in \Delta_K} \sum_{t=1}^T p \cdot g_t \leq \eta \sum_{t=1}^T \sum_{k=1}^K p_t(k) g_t(k)^2 + \frac{\log K}{\eta}.$$

Remark that the theorem also holds for loss vectors $g_t \in [-G, G]^K$ as soon as $\eta \leq 1/G$. Upper-bounding $g_t(j)^2 \leq \|\nabla \ell_t(p_t)\|_\infty^2 \leq G^2$, substituting $g_t = \nabla \ell_t(p_t)$, this yields for all $p \in \Delta_K$

$$\sum_{t=1}^T p_t \cdot \nabla \ell_t(p_t) - p \cdot \nabla \ell_t(p_t) \leq \eta T G^2 + \frac{\log K}{\eta}.$$

But by convexity of the losses, we have the gradient inequality

$$\ell_t(p_t) - \ell_t(p) \leq (p_t - p) \cdot \nabla \ell_t(p_t),$$

which yields

$$\sum_{t=1}^T \ell_t(p_t) - \ell_t(p) \leq \eta T G^2 + \frac{\log K}{\eta}.$$

The proof is concluded by optimizing $\eta = \frac{1}{G} \sqrt{\frac{\log K}{T}}$. □

Example 3.1 (Prediction with expert advice (continued)). In prediction with expert advice, a sequence of observations $y_1, \dots, y_T \in [0, 1]$ is to be predicted with the help of K expert advice $x_t(k) \in [0, 1]$ for $1 \leq k \leq K$. The learner predict $\hat{y}_t = \sum_{k=1}^K p_t(k) x_t(k)$ and suffers a loss $\ell(\hat{y}_t, y_t)$. If the loss function is convex and Lipschitz in its first argument we can apply Theorem 4 with $\ell_t : p \mapsto \ell(p \cdot x_t, y_t)$. For instance, with the absolute loss, $G = 1$ and EG satisfies a bounded regret with respect to any fixed convex combination of experts:

$$\sum_{t=1}^T |\hat{y}_t - y_t| - \min_{p \in \Theta} \sum_{t=1}^T |p \cdot x_t - y_t| \leq 2\sqrt{T \log K}.$$

Hence, on the long run we perform as good as the best convex combination of the experts which may outperform the best expert. This may leads to much better performance than a simple EWA on the experts if

$$\min_{p \in \Theta} \sum_{t=1}^T |p \cdot x_t - y_t| \ll \min_{k \in [K]} \sum_{t=1}^T |x_t(k) - y_t|.$$

Convex hull of finite point set It is worth pointing out that the simplex decision set Δ_K can be generalized with any convex hull of a finite point set $S = \{\theta(1), \dots, \theta(K)\}$:

$$\text{Conv}(S) = \left\{ \sum_{i=1}^K p_i \theta(i) : \forall i, p_i > 0 \text{ and } \sum_{i=1}^K p_i = 1 \right\}.$$

Transforming the loss functions, EWA can be applied to compete with such sets as shown by the theorem below.

Theorem 5. *Let $T \geq 1$. Let $\Theta \subset \mathbb{R}^d$ be a convex set and $S = \{\theta(1), \dots, \theta(K)\} \in \Theta^K$ with diameter $D \geq \max_{i,j} \|\theta(i) - \theta(j)\|_1$. Let $\ell_1, \dots, \ell_T : \Theta \rightarrow \mathbb{R}$ be an arbitrary sequence of convex differentiable losses with bounded gradient $\max_{\theta \in \Theta} \|\nabla \ell_t(\theta)\|_\infty \leq G$. Then, EWA applied with $g_t = \nabla \tilde{\ell}_t$ where $\tilde{\ell}_t : p \mapsto \ell_t \left(\sum_{i=1}^K p(i) \theta(i) \right)$ achieves the regret bound*

$$R_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \text{Conv}(S)} \sum_{t=1}^T \ell_t(\theta) \leq 2GD\sqrt{T \log K},$$

where $\theta_t = \sum_{k=1}^K p_t(k) \theta(k)$

Such a trick can be used for instance to compete with the ℓ_1 -balls using $S = \{\theta \in \mathbb{R}^d : \|\theta\|_1 = R, \|x\|_0 = 1\}$. Since ℓ_p -balls are contained into the ℓ_1 -ball (of possibly larger radius depending on p) this can also be used to compete against any ℓ_p -ball for $p \geq 1$. This trick was introduced by Kivinen and Warmuth [1997] for the EG \pm forecaster.

3.2 Discretized EWA

Setting: general compact decision set, β -Hölder loss functions

In this section, we aim at designing a procedure for general compact decision set Θ . We will assume for simplicity that $\Theta \subset \mathbb{R}^d$ with $\max_{\theta, \theta' \in \Theta} \|\theta - \theta'\| \leq D$, where $\|\cdot\|$ denotes the Euclidean norm. If the loss functions ℓ_t are β -Hölder, i.e.,

$$|\ell_t(\theta) - \ell_t(\theta')| \leq c \|\theta - \theta'\|^\beta$$

there exists a simple solution: approximate Θ with a finite discretization grid Θ_ε and apply EWA on Θ_ε . If Θ or the losses are non-convex, one needs to use the random EWA (see Section 2.3) and bound the regret with high-probability. For convenience, we will assume Θ and the loss functions ℓ_t to be convex so that the algorithm can play convex combinations of points in Θ_ε and all quantities are deterministic.

Lemma 1. *Let $\varepsilon > 0$. Let $\Theta \subset \mathbb{R}^d$ such that $\max_{\theta, \theta' \in \Theta} \|\theta - \theta'\| \leq D$. Then, there exists $\Theta_\varepsilon \subset \Theta$ such that*

$$\text{Card}(\Theta_\varepsilon) \lesssim \left(\frac{D}{\varepsilon}\right)^d \quad \text{and} \quad \forall x \in \Theta, \exists x' \in \Theta_\varepsilon \quad \|\theta - \theta'\| \leq \varepsilon,$$

where \lesssim denotes a rough inequality (up to multiplicative constants and logarithmic terms).

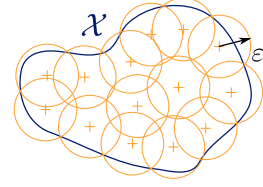
Remark. Remark that a set finite Θ_ε which approximate Θ at radius ε , is called an ε -covering of Θ . The cardinal of the smallest ε -covering is called the *covering number* of Θ . This cardinal is heavily used in theory to analyze the complexity of general spaces Θ . It heavily differentiates parametric spaces with covering number of order $(1/\varepsilon)^d$ with nonparametric spaces (spaces of functions) for which the logarithm of the covering number (or metric entropy) is of order $(1/\varepsilon)^d$.

Proof sketch. We only provide the high-level idea of the proof. First, by properties of the Lebesgue measure in d -dimension, denoting $\mathcal{B}_2(r)$ is the ℓ_2 -ball of radius $r > 0$, we have

$$\text{Vol}(\mathcal{B}_2(r)) = \frac{\pi^{d/2}}{\Gamma(n/2 + 1)} r^d,$$

where Γ is the Euler's gamma function. Therefore,

$$\text{Vol}(\Theta) \leq \text{Vol}(\mathcal{B}_2(D/2)) = \left(\frac{D}{2\varepsilon}\right)^d \text{Vol}(\mathcal{B}_2(\varepsilon)),$$



and thus approximatively $\left(\frac{D}{2\varepsilon}\right)^d$ balls of radius ε are sufficient to cover Θ . □

Theorem 6 (Discretized EWA). *Let $T \geq 1$, $\varepsilon, D > 0$. Let Θ be a compact convex subset of \mathbb{R}^d such that $\max_{\theta, \theta' \in \Theta} \|\theta - \theta'\| \leq D$. Let Θ_ε be an ε -covering of Θ with smallest cardinal. Then, for all sequences of β -Hölder convex losses $\ell_1, \dots, \ell_T : \Theta \rightarrow [0, 1]$, EWA played on the finite set of action Θ_ε with optimized η satisfies the regret bound*

$$R_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \lesssim \sqrt{Td \left(\log D + \frac{1}{\beta} \log(cT) \right)}.$$

Exercise 3.1. Provide a bound on the expected regret for random EWA when the losses and the decision set are non-convex.

Proof. Let $K = \text{Card}(\Theta_\varepsilon)$. Let us order the elements of $\Theta_\varepsilon = \{\theta(1), \dots, \theta(K)\}$. Therefore, at time $t \geq 1$, EWA chooses a weight vector $p_t \in \Delta_K$ and predict the weighted average $\theta_t = \sum_{k=1}^K p_t(k)\theta(k) \in \Theta$. Applying the regret bound of EWA, we get

$$\sum_{t=1}^T \sum_{k=1}^K p_t(k) \ell_t(\theta(k)) - \min_{1 \leq j \leq K} \sum_{t=1}^T \ell_t(\theta(j)) \leq 2\sqrt{T \log K}. \quad (8)$$

Let $\theta^* \in \Theta$ and $\theta(k^*) \in \Theta_\varepsilon$ such that $\|\theta^* - \theta(k^*)\| \leq \varepsilon$. Because the losses are β -Hölder and convex, we have

$$\begin{aligned} R_T &= \sum_{t=1}^T \ell_t(\theta_t) - \ell_t(\theta^*) \\ &\stackrel{\text{Convexity}}{\leq} \sum_{t=1}^T \sum_{k=1}^K p_t(k) \ell_t(\theta(k)) - \ell_t(\theta^*) \quad \leftarrow \theta_t = \sum_{k=1}^K p_t(k)\theta(k) \\ &\leq \sum_{t=1}^T \sum_{k=1}^K p_t(k) \ell_t(\theta(k^*)) - \ell_t(\theta(k^*)) + \sum_{t=1}^T |\ell_t(\theta(k^*)) - \ell_t(\theta^*)| \\ &\stackrel{(8)}{\leq} 2\sqrt{T \log K} + T \max_{1 \leq t \leq T} |\ell_t(\theta^*) - \ell_t(\theta(k^*))| \\ &\stackrel{\beta\text{-Hölder}}{\leq} 2\sqrt{T \log K} + cT\varepsilon^\beta \\ &\stackrel{\text{Lem. 1}}{\lesssim} \sqrt{Td \log \left(\frac{D}{\varepsilon} \right)} + cT\varepsilon^\beta. \end{aligned}$$

Optimizing $\varepsilon^\beta = 1/cT$, hence $\varepsilon = (cT)^{-1/\beta}$, we get

$$R_T \lesssim \sqrt{Td \left(\log D + \frac{1}{\beta} \log(cT) \right)}.$$

□

Though this algorithm is theoretically convenient since it can deal with general compact sets Θ and general loss functions (which can be non-convex and non-differentiable). It suffers two considerable drawbacks:

- *computational complexity*: the algorithm needs to consider a discretization space of cardinal $(X/\varepsilon)^d$ which is of order $O(T^{d/\beta})$. This is prohibitive in practice.
- *bad regret dependence on the dimension*: the regret bound is of order $O(\sqrt{dT \log T})$. We will see how to have no dependence on d when Θ is bounded in ℓ_2 -norm.

3.3 Online gradient descent

Setting: convex differentiable Lipschitz loss function, convex and compact decision set Θ

In this section, we provide another algorithm that solves the drawbacks of the discretized EWA seen in the previous section. This algorithm is called, Online Gradient Descent, and is due to Zinkevich [2003]. It is an online variant of the well-known Gradient Descent algorithm in optimization.

Online Gradient Descent (OGD)

Parameter: $\eta > 0$

Initialize: $\theta_1 \in \Theta$ arbitrarily chosen

For $t = 1, \dots, T$

- select θ_t ; incur loss $\ell_t(\theta_t)$ and observe $\ell_t : \Theta \rightarrow [0, 1]$;
- compute the gradient $\nabla \ell_t(\theta_t)$
- update

$$\theta_{t+1} = \Pi_{\Theta}(\theta_t - \eta \nabla \ell_t(\theta_t)).$$

where Π_{Θ} is the Euclidean projection onto Θ .

Theorem 7. Let $D, G, \eta > 0$. Assume that $\max_{\theta, \theta' \in \Theta} \|\theta - \theta'\| \leq D$ and. Then for any sequence ℓ_1, \dots, ℓ_T of convex differentiable loss functions such that $\max_{\theta \in \Theta} \|\nabla \ell_t(\theta)\| \leq G$, the regret of OGD satisfies

$$\sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \leq \frac{D^2}{2\eta} + \frac{\eta}{2} G^2 T.$$

In particular, for $\eta = \frac{D}{G\sqrt{T}}$, we have $R_T \leq DG\sqrt{T}$.

Exercise 3.2. Prove an upper-bound on the regret of OGD

- a) when η is calibrated with a doubling trick.
- b) when η is calibrated using a time-varying parameter η_t

Exercise 3.3. Prove an upper-bound on the regret of OGD with respect to any sequence of points $\theta_1^*, \dots, \theta_t^* \in \Theta$ such that $\sum_{t=2}^T \|\theta_t^* - \theta_{t-1}^*\| \leq X$

$$\sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(\theta_t^*) \leq \dots$$

Remark. Assume that $\Theta = \Delta_K$ is the simplex and the loss functions are sub-differentiable convex functions with $\|\nabla \ell_t\|_{\infty} \leq G_{\infty}$. Then both EG and OGD are possible algorithms (see Theorems 4 and 7). We saw in Theorem 4 that EG has a regret bound $R_T \leq 2G_{\infty}\sqrt{T \log K}$. In this case, for all $p, p' \in \Delta_K$

$$\|p - p'\| = \sum_{k=1}^K (p(i) - p'(i))^2 \leq \sum_{i=1}^K |p(i) - p'(i)| \leq \sum_{i=1}^K p(i) + p'(i) = 2,$$

and $\|\nabla\ell_t(p)\| \leq \sqrt{K}\|\nabla\ell_t(p)\|_\infty \leq \sqrt{K}G_\infty$. Therefore, the regret of OGD is upper-bounded by $R_t \leq G_\infty\sqrt{2KT}$. To summarize

$$\text{EG: } R_T \leq 2G_\infty\sqrt{T \log K} \quad \text{and} \quad \text{OGD: } R_T \leq \sqrt{2KT}.$$

The dependence on K of OGD is suboptimal in this case. This is solved by OMD, a generalization of both algorithms.

Proof of Theorem 7. Let $\theta^* \in \arg \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)$ and denote $z_{t+1} = \theta_t - \eta \nabla \ell_t(\theta_t)$ so that by definition of θ_{t+1} in the algorithm, we have $\theta_{t+1} = \Pi_\Theta(z_{t+1})$. By convexity, the regret can be upper-bounded as

$$\begin{aligned} R_T &\stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \ell_t(\theta^*) \leq \sum_{t=1}^T \nabla \ell_t(\theta_t) \cdot (\theta_t - \theta^*) \\ &= \frac{1}{\eta} \sum_{t=1}^T (z_{t+1} - \theta_t) \cdot (\theta_t - \theta^*) \quad \leftarrow \nabla \ell_t(\theta_t) = \frac{z_{t+1} - \theta_t}{\eta}. \end{aligned}$$

Then, we use the equality $\|x - y\|^2 = \|x\|^2 + \|y\|^2 - 2x \cdot y$ for all $x, y \in \Theta$ so that

$$x \cdot y = \frac{\|x\|^2 + \|y\|^2 - \|x - y\|^2}{2}.$$

Applying it with $x = z_{t+1} - \theta_t$ and $y = \theta_t - \theta^*$ en substituting into the above regret bound, this yields

$$R_T \leq \frac{1}{2\eta} \sum_{t=1}^T (\|z_{t+1} - \theta_t\|^2 + \|\theta^* - \theta_t\|^2 - \|z_{t+1} - \theta^*\|^2)$$

Then, using $\|z_{t+1} - \theta_t\| = \eta \|\nabla \ell_t(\theta_t)\| \leq \eta G$ and $\|\theta^* - \theta_t\| \leq \|\theta^* - z_t\|$ because Θ is convex and $\theta_t = \Pi_\Theta(z_t)$, we get

$$R_T \leq \frac{1}{2\eta} \sum_{t=1}^T (\eta^2 G^2 + \|\theta^* - z_t\|^2 - \|z_{t+1} - \theta^*\|^2).$$

The last terms telescope, therefore summing over t concludes the proof

$$R_T \leq \frac{\eta G^2 T}{2} + \frac{\|\theta^* - \theta_1\|^2}{2\eta} \leq \frac{\eta G^2 T}{2} + \frac{D^2}{2\eta}.$$

□

3.4 Online Mirror Descent

Online Mirror Descent (OMD) is a generalization of OGD to better exploit the geometry of the decision space Θ . OMD is the online counterpart of the *Mirror Descent* algorithm from convex optimization. The generality of OMD comes from the updates being performed into a dual space which is defined by a convex differentiable regularization function $R : \Theta \rightarrow \mathbb{R}$.

Before stating the algorithm, we need to define the Bregman divergence.

Definition 3 (Bregman divergence). *For any continuously differentiable convex function R , the Bregman divergence with respect to R is defined as*

$$D_R(x||y) \leq R(x) - R(y) - \nabla R(y) \cdot (x - y) \quad \forall x, y \in \Theta.$$

It is the difference between the value of the regularization function at x and the value of its first order Taylor approximation. It is nonnegative but not symmetric. Online Mirror Descent is then defined as follows.

Online Mirror Descent (OMD)Parameters: $\eta > 0$, regularization function R Initialize: $z_1 \in \mathbb{R}^d$ such that $\nabla R(z_1) = 0$ and $\theta_1 = \arg \min_{\theta \in \Theta} B_R(\theta || y_1)$ For $t = 1, \dots, T$

- select θ_t ; incur loss $\ell_t(\theta_t)$ and observe $\ell_t : \Theta \rightarrow [0, 1]$;
- compute the gradient $\nabla \ell_t(\theta_t)$
- update z_t such that

$$\nabla R(z_{t+1}) = \nabla R(\theta_t) - \eta \nabla \ell_t(\theta_t).$$

- project according to the Bregman divergence

$$\theta_{t+1} \in \arg \min_{\theta \in \Theta} D_R(\theta || z_{t+1}).$$

Theorem 8. Let $t \geq 1$. Let Θ be a compact and convex set. Then, for all sequences of convex subdifferentiable loss functions $\ell_1, \dots, \ell_T : \Theta \rightarrow [0, 1]$, the regret of OMD is upper-bounded as

$$\sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \leq \frac{D}{\eta} + \frac{1}{\eta} \sum_{t=1}^T D_{R^*}(\nabla R(\theta_t) - \eta \nabla \ell_t(\theta_t) || \nabla R(\theta_t))$$

where $D \geq \max_{\theta \in \Theta} |R(\theta)|$ and R^* is the Fenchel conjugate of R defined as $R^*(z) \stackrel{\text{def}}{=} \max_{\theta \in \Theta} \{\theta \cdot z - R(\theta)\}$.

The proof can be found for instance in Bubeck et al. [2012]. EG and OGD are two particular cases of Online Mirror Descent.

Example 3.2 (Balls in $\mathbb{R}^d = \text{OGD}$). If $\Theta \subset \mathbb{R}^d$, we can choose $R(x) = \frac{1}{2} \|x\|^2$. Then $\nabla R(x) = x$ and $D_R(x || y) = \frac{1}{2} \|x - y\|^2$. Therefore, the update of OMD becomes $y_{t+1} = \theta_t - \eta \nabla \ell_t(\theta_t)$ and $\theta_{t+1} = \Pi_{\Theta}(y_{t+1})$. We recover the online gradient descent algorithm.

Example 3.3 (Simplex = EG). If $\Theta = \Delta_K$. We can choose the negative entropy

$$R(x) = \sum_{i=1}^K x(i) \log x(i).$$

In this case, $\nabla R(x)_i = 1 + \log x(i)$ and the Bregman Divergence is $D_R(x || y) = \sum_{i=1}^K x(i) \log(x(i)/y(i))$ also known as the Kullback-Leibler divergence. The update of OMD is then

$$1 + \log(y_{t+1}(i)) = 1 + \log \theta_t(i) - \eta g_t(i),$$

where $g_t = \nabla \ell_t(\theta_t) \in \mathbb{R}^K$. This can be rewritten

$$y_{t+1}(i) = \theta_t(i) e^{-\eta [\nabla \ell_t(\theta_t)]_i}.$$

The projection to the simplex is a simple renormalization (left as exercise), we thus get

$$\theta_{t+1}(i) = \frac{\theta_t(i) e^{-\eta g_t(i)}}{\sum_{k=1}^K \theta_t(k) e^{-\eta g_t(k)}},$$

and we recover the update of EG (i.e., EWA applied with the gradient trick $g_t = \nabla \ell_t(\theta_t)$).

4 Adversarial Bandits

In previous chapters, we considered the full-information feedback and the bandit feedback with stochastic loss functions. In *full information with finite decision set* $\Theta = [K] \stackrel{\text{def}}{=} \{1, \dots, K\}$, we saw the Random Exponentially Weighted Average (EWA) forecaster. It is defined as

$$p_t(k) = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s(k)}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s(j)}}. \quad (\text{EWA})$$

and draws $\theta_t = k$ with probability $p_t(k)$. If $-\eta \ell_t(j) \leq 1$ (see the proof of EWA in first lecture), it satisfies the upper-bound:

$$\sum_{t=1}^T p_t \cdot \ell_t - \min_{1 \leq j \leq K} \sum_{t=1}^T \ell_t(j) \leq \eta \sum_{t=1}^T \sum_{k=1}^K p_t(k) \ell_t(k)^2 + \frac{\log K}{\eta}. \quad (*)$$

Since the decision θ_t is random, we assume that ℓ_t cannot depend on θ_t but may depend on past information $\sigma(p_1, \ell_1, x_1, \dots, x_{t-1}, p_t)$. The above bound can be converted into a bound on the expected regret for well-calibrated learning rate η

$$\mathbb{E}[R_T] = \mathbb{E} \left[\sum_{t=1}^T \ell_t(\theta_t) - \min_{k \in [K]} \sum_{t=1}^T \ell_t(k) \right] \leq 2\sqrt{T \log K}.$$

In this chapter, we will see adversarial bandits: that is bandit feedback (only $\ell_t(\theta_t)$ is observed at the end of round t by the player) with an adversarial sequence of loss function ℓ_t (i.e., no stochastic assumptions). Note that we turn back to losses instead of rewards but we will come back to rewards whenever it makes the proof easier. Remember that the lower-bound on the regret in the worst-case is of order $O(\sqrt{TK})$.

4.1 The exponentially weighted average algorithm for bandits

We consider Setting 1 with bandit feedback, finite decision space $\Theta = [K] \stackrel{\text{def}}{=} \{1, \dots, K\}$ and adversarial losses. To emphasize that the action is in $[K]$, we denote by k_t the action chosen by the player (instead of θ_t). We do not assume the loss functions ℓ_t to be linear nor convex (the decision space is not). Similarly to Random EWA the chosen action $k_t \in [K]$ is sampled randomly from a distribution p_t chosen at round t by the player. We will provide an algorithm called Exp3 inspired by EWA.

4.1.1 Pseudo-regret bound

Let us denote the regret with respect to action $k \in [K]$ by

$$R_T(k) \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(k_t) - \sum_{t=1}^T \ell_t(k).$$

Instead of minimizing the *expected regret* $\mathbb{E}[R_T] = \mathbb{E}[\max_k R_T(k)]$, we will start with an easier objective, the *pseudo-regret* defined as

$$\bar{R}_T \stackrel{\text{def}}{=} \max_{k \in [K]} \mathbb{E}[R_T(k)] = \max_{k \in [K]} \mathbb{E} \left[\sum_{t=1}^T \ell_t(k_t) - \sum_{t=1}^T \ell_t(k) \right]. \quad (\text{pseudo regret})$$

It is worth pointing out that the expectations are taken with respect to the randomness of the algorithm: the decisions k_t are random. We can distinguish two types of adversaries:

- *oblivious adversary*: all the loss functions ℓ_1, \dots, ℓ_t are chosen in advance before the game starts and do not depend on the past player decisions k_1, \dots, k_T . In this case, the losses $\ell_t(k)$ are deterministic and there is thus equality: $\bar{R}_T = \mathbb{E}[R_T]$.

- *adaptive adversary*: the loss function ℓ_t at round $t \geq 1$ may depend on past information $\sigma(k_1, \dots, k_{t-1})$. It is thus random. By Jensen's inequality $\max_{k \in [K]} \mathbb{E}[R_T(k)] \leq \mathbb{E}[\max_{k \in [K]} R_T(k)]$ and thus $\bar{R}_T \leq \mathbb{E}[R_T]$.

The EXP3 algorithm Ideally, we would like to reuse our algorithm EWA that assigned weights

$$\forall k \in [K], \quad p_t(k) = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s(k)}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s(j)}}. \quad (\text{EWA})$$

Unfortunately this is not possible since the player does not observe $\ell_t(k)$ for $k \neq k_t$. The high-level idea of Exp3 is to replace $\ell_t(k)$ with an unbiased estimate that is observed by the player. A first idea would be to use $\ell_t(k)$ if we observe it and 0 otherwise:

$$g_t(k) = \begin{cases} \ell_t(k) & \text{if } k = k_t \quad \leftarrow \text{i.e., decision } k \text{ is observed} \\ 0 & \text{otherwise} \end{cases}.$$

However, this estimate is biased:

$$\mathbb{E}_{k_t \sim p_t} [g_t(k_t)] = p_t(k) \ell_t(k) \neq \ell_t(k).$$

In other words, the actions that are less likely to be chosen by the algorithm (small weight $p_t(k)$) are more likely to be unobserved and incur 0 loss. We need to correct this phenomenon. Therefore we choose

$$g_t(k) = \frac{\ell_t(k)}{p_t(k)} \mathbb{1}_{k = k_t}, \quad (9)$$

which leads to the algorithm EXP3 detailed below.

EXP3

Parameter: $\eta > 0$

Initialize: $p_1 = (\frac{1}{K}, \dots, \frac{1}{K})$

For $t = 1, \dots, T$

- draw $k_t \sim p_t$; incur loss $\ell_t(k_t)$ and observe $\ell_t(k_t) \in [0, 1]$;
- update for all $k \in \{1, \dots, K\}$

$$p_{t+1}(k) = \frac{e^{-\eta \sum_{s=1}^t g_s(k)}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^t g_s(j)}}, \quad \text{where } g_s(k) = \frac{\ell_s(k)}{p_s(k)} \mathbb{1}_{k = k_s}$$

Then applying the Inequality (*) for EWA with the substituted losses g_t , we get the following theorem.

Theorem 9. *Let $T \geq 1$. The pseudo-regret of EXP3 run with $\eta = \sqrt{\frac{\log K}{KT}}$ is upper-bounded as:*

$$\bar{R}_T \leq 2\sqrt{KT \log K}.$$

Proof. Apply EWA to the estimated losses $g_t(j)$ that are completely observed (nonnegative but not bounded), we get from Inequality (*) and taking the expectation:

$$\mathbb{E} \left[\sum_{t=1}^T p_t \cdot g_t - \min_{j \in [K]} \sum_{t=1}^T g_t(j) \right] \leq \frac{\log K}{\eta} + \eta \sum_{t=1}^T \mathbb{E} [p_t \cdot g_t^2]. \quad (10)$$

Now we compute the expectations. Denote by $\mathcal{F}_{t-1} \stackrel{\text{def}}{=} \sigma(p_1, \ell_1, k_1, \dots, k_{t-1}, p_t, \ell_t)$ the past information available at round t for the adversary (which cannot use the randomness of k_t but can use p_t). Note that ℓ_t and p_t are \mathcal{F}_{t-1} -measurable by assumption. We have

$$\forall j \in [K] \quad \mathbb{E} [g_t(j) | \mathcal{F}_{t-1}] = \mathbb{E} \left[\frac{\ell_t(j)}{p_t(j)} \mathbb{1}_{j = k_t} | \mathcal{F}_{t-1} \right] = \sum_{k=1}^K p_t(k) \frac{\ell_s(j)}{p_t(j)} \mathbb{1}_{j = k} = \ell_t(j)$$

thus the estimated losses are unbiased $\mathbb{E}[g_t(j)] = \mathbb{E}[\ell_t(j)]$ and

$$\begin{aligned}\mathbb{E}[p_t \cdot g_t] &= \mathbb{E}\left[\sum_{j=1}^K p_t(j)g_t(j)\right] = \mathbb{E}\left[\sum_{j=1}^K p_t(j)\mathbb{E}[g_t(j)|\mathcal{F}_{t-1}]\right] \\ &= \mathbb{E}\left[\sum_{j=1}^K p_t(j)\ell_t(j)\right] = \mathbb{E}\left[\mathbb{E}[\ell_t(k_t)|\mathcal{F}_{t-1}]\right] = \mathbb{E}[\ell_t(k_t)].\end{aligned}$$

Therefore, we can lower-bound the left-hand side:

$$\begin{aligned}\mathbb{E}\left[\sum_{t=1}^T p_t \cdot g_t - \min_{j \in [K]} \sum_{t=1}^T g_t(j)\right] &\geq \max_{j \in [K]} \mathbb{E}\left[\sum_{t=1}^T p_t \cdot g_t - \sum_{t=1}^T g_t(j)\right] \\ &= \max_{j \in [K]} \mathbb{E}\left[\sum_{t=1}^T \ell_t(k_t) - \sum_{t=1}^T \ell_t(j)\right] = \bar{R}_T.\end{aligned}$$

On the other hand, the expectation of the right-hand side satisfies

$$\begin{aligned}\mathbb{E}[p_t \cdot g_t^2] &= \mathbb{E}\left[\sum_{j=1}^K p_t(j)g_t(j)^2\right] = \mathbb{E}\left[\sum_{j=1}^K p_t(j)\mathbb{E}[g_t(j)^2 | \mathcal{F}_{t-1}]\right] \\ &= \mathbb{E}\left[\sum_{j=1}^K \sum_{k=1}^K p_t(j)p_t(k)\left(\frac{\ell_t(j)}{p_t(j)}\mathbb{1}_{j=k}\right)^2\right] \\ &= \mathbb{E}\left[\sum_{j=1}^K \sum_{k=1}^K p_t(k)\frac{\ell_t(j)^2}{p_t(j)}\mathbb{1}_{j=k}\right] \\ &= \mathbb{E}\left[\sum_{j=1}^K \ell_t(j)^2\right] \leq K.\end{aligned}$$

Substituting into Inequality (10) yields

$$\bar{R}_T \leq \frac{\log K}{\eta} + \eta KT.$$

and optimizing $\eta = \sqrt{KT/(\log K)}$ concludes. \square

The issue with the above regret bound is that it bounds the pseudo-regret and not the expected regret. This is because we have

$$\mathbb{E}\left[\min_j \sum_{t=1}^T g_t(j)\right] \leq \min_j \mathbb{E}\left[\sum_{t=1}^T g_t(j)\right] = \min_{j \in [K]} \mathbb{E}\left[\sum_{t=1}^T \ell_t(j)\right]$$

but not

$$\mathbb{E}\left[\min_j \sum_{t=1}^T g_t(j)\right] \not\leq \mathbb{E}\left[\min_j \sum_{t=1}^T \ell_t(j)\right]. \quad (11)$$

Hence, controlling the cumulative loss against the best estimated action only controls the pseudo regret and not the true regret.

4.1.2 High probability bound on the regret

Gains versus losses In this part, we will switch the analysis from losses $\ell_t(k)$ to gains $g_t(k) = 1 - \ell_t(k) \in [0, 1]$ because the core idea of the next algorithm is easier to see with gains. Remark that the loss and gain versions are symmetric via the transformation $g_t(k) = 1 - \ell_t(k)$. The regret in terms of gains is defined as

$$R_T \stackrel{\text{def}}{=} \max_{k \in [K]} \sum_{t=1}^T g_t(k) - \sum_{t=1}^T g_t(k_t).$$

Using EWA with full information from (*), if $\eta g_t(k) \leq 1$, we also have for gains the inequality

$$\max_{1 \leq j \leq K} \sum_{t=1}^T g_t(j) - \sum_{t=1}^T p_t \cdot g_t \leq \eta \sum_{t=1}^T p_t \cdot g_t^2 + \frac{\log K}{\eta}, \quad \text{where } p_t(k) = \frac{e^{\eta \sum_{s=1}^{t-1} g_s(k)}}{\sum_{j=1}^K e^{\eta \sum_{s=1}^{t-1} g_s(j)}}. \quad (12)$$

High-level idea of EXP3.P The high-level idea of the next algorithm is to ensure that the estimators $\hat{g}_t(k)$ of the gains satisfy

$$\mathbb{E} \left[\max_j \sum_{t=1}^T \hat{g}_t(j) \right] \geq \mathbb{E} \left[\max_j \sum_{t=1}^T g_t(j) \right] \quad (13)$$

so that controlling the performance with respect to the estimated gains (left-hand side) also controls the performance with respect to the true gains (right-hand side). This was not the case of the estimators used for EXP3 (see (11)). To ensure (13), we add a bias term β to the estimators $\hat{g}_t(k)$ as follows:

$$\hat{g}_t(k) \stackrel{\text{def}}{=} \frac{g_t(k) \mathbf{1}_{k = k_t} + \beta}{p_t(k)} \quad (14)$$

In contrary to (9), the estimator is indeed biased

$$\mathbb{E}[\hat{g}_t(k) | \mathcal{F}_{t-1}] = g_t(k) + \frac{\beta}{p_t(k)}, \quad (15)$$

where we recall that $\mathcal{F}_{t-1} \stackrel{\text{def}}{=} \sigma(p_1, k_1, g_1, \dots, k_{t-1}, p_t, g_t)$ contains the information up to time t available to the environment. We have the following Lemma:

Lemma 2. For any $\delta > 0$, with probability $1 - \delta$ and $\beta \in (0, 1)$,

$$\sum_{t=1}^T \hat{g}_t(j) \geq \sum_{t=1}^T g_t(j) - \frac{\log(1/\delta)}{\beta}.$$

Proof. Let $\beta \in (0, 1)$, from Markov's inequality, we have

$$\begin{aligned} \mathbb{P} \left(\sum_{t=1}^T \hat{g}_t(j) \geq \sum_{t=1}^T g_t(j) - \frac{\log(1/\delta)}{\beta} \right) &= \mathbb{P} \left(\exp \left(\beta \sum_{t=1}^T (g_t(j) - \hat{g}_t(j)) \right) \geq \delta^{-1} \right) \\ &\leq \delta \mathbb{E} \left[\exp \left(\beta \sum_{t=1}^T (g_t(j) - \hat{g}_t(j)) \right) \right]. \end{aligned}$$

It only remains to upper-bound the expectation in the right-hand side by 1, which we do now. Since $\beta \in (0, 1)$ and $\hat{g}_t(j) \geq \beta/p_t(j)$, we have $\beta(g_t(j) - \hat{g}_t(j) + \beta/p_t(j)) \leq 1$. Therefore, we can use the inequality

$e^x \leq 1 + x + x^2$ for $x \leq 1$, which entails

$$\begin{aligned} \mathbb{E} \left[\exp \left(\beta (g_t(j) - \hat{g}_t(j)) \right) \middle| \mathcal{F}_{t-1} \right] &= \mathbb{E} \left[\exp \left(\beta \left(g_t(j) - \hat{g}_t(j) + \frac{\beta}{p_t(j)} \right) \right) \middle| \mathcal{F}_{t-1} \right] \exp \left(-\frac{\beta^2}{p_t(j)} \right) \\ &\leq \mathbb{E} \left[\left(1 + \beta \left(g_t(j) - \hat{g}_t(j) + \frac{\beta}{p_t(j)} \right) + \beta^2 \left(g_t(j) - \hat{g}_t(j) + \frac{\beta}{p_t(j)} \right)^2 \right) \middle| \mathcal{F}_{t-1} \right] e^{-\frac{\beta^2}{p_t(j)}} \\ &\stackrel{(15)}{=} \left(1 + \beta^2 \mathbb{E} \left[\left(g_t(j) - \hat{g}_t(j) + \frac{\beta}{p_t(j)} \right)^2 \middle| \mathcal{F}_{t-1} \right] \right) e^{-\frac{\beta^2}{p_t(j)}} \end{aligned}$$

where the last equality is by (15) and because $p_t(j)$ is \mathcal{F}_{t-1} -measurable. Now,

$$\begin{aligned} \mathbb{E} \left[\left(g_t(j) - \hat{g}_t(j) + \frac{\beta}{p_t(j)} \right)^2 \middle| \mathcal{F}_{t-1} \right] &= \text{Var} \left(\hat{g}_t(j) \middle| \mathcal{F}_{t-1} \right) = \text{Var} \left(\frac{g_t(j) \mathbb{1}_{j=k_t}}{p_t(j)} \middle| \mathcal{F}_{t-1} \right) \\ &\leq \mathbb{E} \left[\left(\frac{g_t(j) \mathbb{1}_{j=k_t}}{p_t(j)} \right)^2 \middle| \mathcal{F}_{t-1} \right] \leq \mathbb{E} \left[\frac{\mathbb{1}_{j=k_t}}{p_t(j)^2} \middle| \mathcal{F}_{t-1} \right] = \sum_{k=1}^K \frac{p_t(k) \mathbb{1}_{j=k}}{p_t(j)^2} = \frac{1}{p_t(j)}. \end{aligned}$$

Substituting into the previous inequality and using $1 + x \leq e^x$, it yields

$$\mathbb{E} \left[\exp \left(\beta (g_t(j) - \hat{g}_t(j)) \right) \middle| \mathcal{F}_{t-1} \right] \leq \left(1 + \frac{\beta^2}{p_t(j)} \right) e^{-\beta^2/p_t(j)} \leq 1.$$

The proof is concluded by induction

$$\begin{aligned} \mathbb{E} \left[\exp \left(\beta \sum_{t=1}^T (g_t(j) - \hat{g}_t(j)) \right) \right] &= \mathbb{E} \left[\underbrace{\mathbb{E} \left[\exp \left(\beta (g_T(j) - \hat{g}_T(j)) \right) \middle| \mathcal{F}_{T-1} \right]}_{\leq 1} \exp \left(\beta \sum_{t=1}^{T-1} (g_t(j) - \hat{g}_t(j)) \right) \right] \\ &\leq \mathbb{E} \left[\exp \left(\beta \sum_{t=1}^{T-1} (g_t(j) - \hat{g}_t(j)) \right) \right] \leq \dots \leq 1. \end{aligned}$$

□

The issue with the estimators $\hat{g}_t(j) \in (0, +\infty)$ defined in Equation (14) is that they might be unbounded if the weights $p_t(j)$ are close to zero. The condition $\eta \hat{g}_t(j) \leq 1$ which appeared in the proof of EWA cannot hold for any $\eta > 0$. Remark that this was not a problem for EXP3 with the preceding choice (9) because $-\eta g_t(j) \leq 1$ (see the proof of EWA for details).

The next algorithm called EXP3.P, is close to EXP3 but ensures the weights do not vanish to zero by adding an exploration parameter $\gamma > 0$.

EXP3.P

Parameters: $\eta > 0, \beta \in (0, 1), \gamma \in (0, 1)$

Initialize: $p_1 = \left(\frac{1}{K}, \dots, \frac{1}{K} \right)$

For $t = 1, \dots, T$

- draw $k_t \sim p_t$; receive gain $g_t(k_t) = 1 - \ell_t(k_t)$ and observe $g_t(k_t) \in [0, 1]$;
- update for all $k \in \{1, \dots, K\}$

$$p_{t+1}(k) = (1 - \gamma) \frac{e^{\eta \sum_{s=1}^t \hat{g}_s(k)}}{\sum_{j=1}^K e^{\eta \sum_{s=1}^t \hat{g}_s(j)}} + \frac{\gamma}{K},$$

where $\hat{g}_s(k) = \frac{g_s(k) \mathbb{1}_{k=k_s} + \beta}{p_s(k)}$.

The weights $p_t(k)$ of EXP3.P are necessary larger than γ/K and thus $|\eta g_t(j)| \leq 1$ as soon as $\eta(1+\beta)K/\gamma \leq 1$. We get the following high-probability bound on the regret.

Theorem 10. *For well-chosen parameters $\gamma \in (0, 1)$, $\beta \in (0, 1)$ and $\eta > 0$ satisfying $\eta(1 + \beta)K/\gamma \leq 1$, for any $\delta > 0$, the EXP3.P algorithm achieves*

$$R_T \leq 6\sqrt{TK \log K} + \sqrt{\frac{TK}{\log K}} \log(1/\delta).$$

with probability at least $1 - \delta$.

Remark that the above bound leads to a bound on the expected regret, with the choice $\delta = 1/T$ it yields

$$\mathbb{E}[R_T] \leq 6\sqrt{TK \log K} + \sqrt{\frac{TK}{\log K}} \log(T) + 1$$

The logarithmic dependency on T can even be removed using $\mathbb{E}[X] = \int_0^\infty \mathbb{P}(X \geq \varepsilon) d\varepsilon$.

Proof of Theorem 10. Defining the weights that would assign EXP3,

$$q_t(j) \stackrel{\text{def}}{=} \frac{e^{\eta \sum_{s=1}^{t-1} \hat{g}_s(j)}}{\sum_{k=1}^K e^{\eta \sum_{s=1}^{t-1} \hat{g}_s(k)}},$$

we get from Inequality (12) applied with $\hat{g}_t(j)$,

$$\max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j) \leq \sum_{t=1}^T q_t \cdot \hat{g}_t + \eta \sum_{t=1}^T q_t \cdot \hat{g}_t^2 + \frac{\log K}{\eta}.$$

where we used $\eta \hat{g}_t(j) \leq 1$ because $\eta(1 + \beta)K/\gamma \leq 1$. Now, we use that $p_t \stackrel{\text{def}}{=} (1 - \gamma)q_t + \gamma/K$, which entails $q_t = (p_t - \gamma/K)/(1 - \gamma) \leq p_t/(1 - \gamma)$. Substituting into the above inequality

$$(1 - \gamma) \max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j) \leq \sum_{t=1}^T p_t \cdot \hat{g}_t + \eta \sum_{t=1}^T p_t \cdot \hat{g}_t^2 + \frac{\log K}{\eta}. \quad (16)$$

But by definition of \hat{g}_t ,

$$p_t \cdot \hat{g}_t = \sum_{j=1}^K p_t(j) \hat{g}_t(j) = \sum_{j=1}^K (g_t(j) \mathbf{1}_{j = k_t} + \beta) = g_t(k_t) + K\beta.$$

and since $p_t(j) \hat{g}_t(j) \leq (1 + \beta)$,

$$\sum_{t=1}^T p_t \cdot \hat{g}_t^2 \leq (1 + \beta) \sum_{j=1}^K \sum_{t=1}^T \hat{g}_t(j) \leq K(1 + \beta) \max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j) \leq \frac{\gamma}{\eta} \max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j).$$

Therefore, substituting into Inequality (16) gives

$$(1 - \gamma) \max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j) \leq \sum_{t=1}^T g_t(k_t) + K\beta T + \gamma \max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j) + \frac{\log K}{\eta},$$

where we used $(1 + \beta)K \leq \gamma/\eta$. Reorganizing, we get

$$(1 - 2\gamma) \max_{j \in [K]} \sum_{t=1}^T \hat{g}_t(j) \leq \sum_{t=1}^T g_t(k_t) + K\beta T + \frac{\log K}{\eta}.$$

Using Lemma 2 together with a union bound (to have it for all $j \in [K]$), we have with probability $1 - \delta$

$$(1 - 2\gamma) \left(\max_{j \in [K]} \sum_{t=1}^T g_t(j) - \frac{\log(K/\delta)}{\beta} \right) \leq \sum_{t=1}^T g_t(k_t) + K\beta T + \frac{\log K}{\eta},$$

and thus reorganizing and choosing $\gamma \stackrel{\text{def}}{=} 2\eta K \geq \eta(1 + \beta)K$,

$$\max_{j \in [K]} \sum_{t=1}^T g_t(j) - \sum_{t=1}^T g_t(k_t) \leq K\beta T + \frac{\log K}{\eta} + \frac{\log(K/\delta)}{\beta} + 4\eta K T.$$

The proof is concluded by optimizing $\eta \stackrel{\text{def}}{=} (1/2)\sqrt{(\log K)/KT}$ and $\beta \stackrel{\text{def}}{=} \sqrt{(\log K)/(KT)}$. \square

4.2 Adversarial bandits with experts

We turn back to the loss version of the game. We now consider prediction with expert advice in the bandit framework. The setting is the same as the one described in Figure 1, but at the beginning of each round $t \geq 1$, some experts $i = 1, \dots, N$ propose recommendations $h_t(i) \in [K]$. These recommendations may be random and may depend on past actions k_s , $s \leq t - 1$ and past observations $\ell_s(k_s)$. The loss of each expert is given by the loss of the chosen decision $\ell_t(h_t(i))$ but only $\ell_t(k_t)$ is observed by the learner. The goal of the learner is then to be competitive with the best expert on a long run. To do so, it minimizes the pseudo-regret

$$R_T^{\text{exp}} \stackrel{\text{def}}{=} \max_{i=1, \dots, N} \mathbb{E} \left[\sum_{t=1}^T \ell_t(k_t) - \sum_{t=1}^T \ell_t(h_t(i)) \right]$$

with respect to the experts. In order to bound the pseudo-regret, one could consider experts as the set of arms and use EXP3. This would give a bound of order $\sqrt{TN \log N}$. However it does not take into account the information on the reward of all experts that choose the same action $h_t(i) = k_t$.

EXP4

Parameter: $\eta > 0$

Initialize: $q_1 = (\frac{1}{N}, \dots, \frac{1}{N})$.

For each round $t = 1, \dots, n$

1. Get expert advice $h_t(1), \dots, h_t(N) \in [K]$
2. Draw an expert i_t with probability distribution $q_t \in \Delta_N$
3. Choose decision $k_t = h_t(i_t)$
4. Compute the estimated loss for each decision

$$g_t(k) = \frac{\ell_t(k)}{p_t(k)} \mathbf{1}_{k = k_t},$$

where $p_t \stackrel{\text{def}}{=} \sum_{i=1}^N q_t(i) \delta_{\ell_t(i)} \in \Delta_K$.

5. Compute the estimated loss of the experts component-wise $g_t(h_t(i))$
6. Update the probability distribution over the experts component-wise

$$q_{t+1}(i) = \frac{\exp\left(-\eta \sum_{s=1}^t g_s(h_s(i))\right)}{\sum_{j=1}^N \exp\left(\eta \sum_{s=1}^t g_s(h_s(j))\right)}, \quad \forall 1 \leq i \leq N.$$

Theorem 11. *EXP4 with $\eta = \sqrt{\log N/(KT)}$ satisfies $R_T^{\text{exp}} \leq 2\sqrt{TK \log N}$.*

Similarly to the variant EXP3.P, we can define a variant EXP4.P to bound the regret with high probability (and thus the expected regret). Furthermore, the above algorithm (and theorem) can be extended to the case where expert advice are distributions $h_t(i) \in \Delta_K$. The algorithm is the same by sampling k_t according to $h_t(i_t)$ and assigning to expert i the loss $\sum_{k=1}^K h_t(i)(k)g_t(k)$.

Proof. We can apply the analysis of EXP to a learner using distribution q_t over N actions (here experts) with (full-information) losses $g_t(h_t(i))$ for $i \in \{1, \dots, N\}$. We get from Inequality (*)

$$\mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N q_t(i) \cdot g_t(h_t(i)) - \min_{1 \leq j \leq N} \sum_{t=1}^T g_t(h_t(j)) \right] \leq \eta \sum_{t=1}^T \sum_{i=1}^N \mathbb{E} \left[q_t(i) g_t(h_t(i))^2 \right] + \frac{\log N}{\eta}. \quad (17)$$

Remark that $k_t = h_t(i)$ with probability $q_t(i)$ so that, k_t follows the distribution $p_t \stackrel{\text{def}}{=} \sum_{i=1}^N q_t(i) \delta_{h_t(i)}$ knowing the past information $\mathcal{F}_{t-1} \stackrel{\text{def}}{=} \sigma(h_1, i_1, k_1, \dots, i_{t-1}, k_{t-1}, h_t)$. Now, similarly to the proof of EXP3, we compute the expectations. We have for all $k \in [K] \stackrel{\text{def}}{=} \{1, \dots, K\}$

$$\mathbb{E} \left[g_t(k) \middle| \mathcal{F}_{t-1} \right] = \mathbb{E} \left[\frac{\ell_t(k)}{p_t(k)} \mathbf{1}_{k = k_t} \middle| \mathcal{F}_{t-1} \right] = \sum_{j=1}^K p_t(j) \frac{\ell_t(k)}{p_t(k)} \mathbf{1}_{k = j} = \ell_t(k),$$

and thus for all $i \in \{1, \dots, N\}$

$$\mathbb{E} \left[g_t(h_t(i)) \middle| \mathcal{F}_{t-1} \right] = \ell_t(h_t(i)), \quad (18)$$

and

$$\begin{aligned} \mathbb{E} \left[\sum_{i=1}^N q_t(i) \cdot g_t(h_t(i)) \middle| \mathcal{F}_{t-1} \right] &= \sum_{i=1}^N q_t(i) \mathbb{E} \left[g_t(h_t(i)) \middle| \mathcal{F}_{t-1} \right] = \sum_{i=1}^N q_t(i) \ell_t(h_t(i)) \\ &= \mathbb{E} \left[\ell_t(h_t(i_t)) \middle| \mathcal{F}_{t-1} \right] = \mathbb{E} \left[\ell_t(k_t) \middle| \mathcal{F}_{t-1} \right]. \end{aligned} \quad (19)$$

Furthermore,

$$\mathbb{E} \left[g_t(h_t(i))^2 \middle| \mathcal{F}_{t-1} \right] = \sum_{k=1}^K p_t(k) \left(\frac{\ell_t(h_t(i))}{p_t(h_t(i))} \right)^2 \mathbf{1}_{k = h_t(i)} = \frac{\ell_t(h_t(i))^2}{p_t(h_t(i))} \leq \frac{1}{p_t(h_t(i))},$$

and

$$\sum_{i=1}^N q_t(i) \mathbb{E} \left[g_t(h_t(i))^2 \middle| \mathcal{F}_{t-1} \right] \leq \sum_{i=1}^N \frac{q_t(i)}{p_t(h_t(i))} = \mathbb{E} \left[\frac{1}{p_t(h_t(i_t))} \middle| \mathcal{F}_{t-1} \right] = \mathbb{E} \left[\frac{1}{p_t(k_t)} \middle| \mathcal{F}_{t-1} \right] = \sum_{k=1}^K \frac{p_t(k)}{p_t(k)} = K. \quad (20)$$

Substituting (18), (19), and (20) into Inequality (17) and lower-bounding the expected regret with the pseudo-regret, we get

$$\begin{aligned} \bar{R}_T^{\text{exp}} &\stackrel{\text{def}}{=} \max_{1 \leq i \leq N} \mathbb{E} \left[\sum_{t=1}^T \ell_t(k_t) - \ell_t(h_t(i)) \right] \\ &\stackrel{(18),(19)}{=} \max_{1 \leq i \leq N} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N q_t(i) g_t(h_t(i)) - g_t(h_t(i)) \right] \\ &\stackrel{\text{Jensen}}{\leq} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N q_t(i) g_t(h_t(i)) - \min_{1 \leq i \leq N} g_t(h_t(i)) \right] \\ &\stackrel{(17)}{\leq} \eta \sum_{t=1}^T \sum_{i=1}^N \mathbb{E} \left[q_t(i) g_t(h_t(i))^2 \right] + \frac{\log N}{\eta} \\ &\stackrel{(20)}{\leq} \eta K T + \frac{\log N}{\eta}. \end{aligned}$$

Optimizing η concludes the proof. \square

4.3 Adversarial Bandits with side information

A natural extension of the previous setting is by adding side (or contextual) information: this is called contextual bandits. It arises in most applications such as recommendation systems or online advertisement. The side information can then be the cookies of a new user to which we need to recommend a product.

Assume that for each time step $t \geq 1$, before doing its prediction k_t the learner observes a context x_t in a finite set \mathcal{X} of contexts. The learner must then learn the best mapping $g : \mathcal{X} \rightarrow [K]$ and is evaluated with the contextual pseudo-regret:

$$R_T^{\mathcal{X}} \stackrel{\text{def}}{=} \max_{g: \mathcal{X} \rightarrow [K]} \mathbb{E} \left[\sum_{t=1}^T \ell_t(k_t) - \sum_{t=1}^T \ell_t(g(x_t)) \right].$$

Similarly, to the stochastic setting, if \mathcal{X} is finite, a simple algorithm consists in running a different copy $EXP3(c)$ of EXP3 for each context $c \in \mathcal{X}$. We denote by \mathcal{X} -EXP3 this algorithm. At each time step $t \geq 1$, the learner uses $EXP3(x_t)$ to make the prediction. The following theorem follows from Theorem 9.

Theorem 12. *The contextual pseudo-regret of \mathcal{X} -EXP3 is upper-bounded as:*

$$R_T^{\mathcal{X}} \leq 2\sqrt{T|\mathcal{X}|K \log K}.$$

Proof. Applying the proof of the pseudo-regret bound of EXP3 for each instance $x \in \mathcal{X}$:

$$\max_{j \in [K]} \mathbb{E} \left[\sum_{t=1}^T (\ell_t(k_t) - \ell_t(j)) \mathbb{1}_{x_t = x} \right] \leq 2\sqrt{K(\log K)T_x},$$

where $T_x = \sum_{t=1}^T \mathbb{1}_{x_t = x}$. Summing over $x \in \mathcal{X}$,

$$\sum_{x \in \mathcal{X}} \max_{j \in [K]} \mathbb{E} \left[\sum_{t=1}^T (\ell_t(k_t) - \ell_t(j)) \mathbb{1}_{x_t = x} \right] \leq 2 \sum_{x \in \mathcal{X}} \sqrt{K(\log K)T_x} \stackrel{\text{Jensen}}{\leq} 2\sqrt{|\mathcal{X}|K(\log K)T}$$

where the last inequality is by using the concavity of the square root together with $\sum_{s \in \mathcal{X}} T_s = T$. The proof is concluded by remarking that the left-hand side is the contextual pseudo-regret. \square

Similarly to the classical lower-bound $O(\sqrt{TK})$, a lower-bound of order $\sqrt{|\mathcal{X}|KT}$ holds under the assumption that a significant proportion of the contexts are used at least a constant fraction of the T rounds. The above bound is nice but the dependency $|\mathcal{X}|$ might be annoying if \mathcal{X} is large.

Exercise 4.1. Generalize the above algorithm and upper-bound when the context-space is continuous and the loss functions are β -Hölder in the contexts.

Competing against the best context set

In some cases, one may want to combine bandit algorithms. For example, we could have in hand different context set \mathcal{X} . For each of these sets \mathcal{X} , we can bound the pseudo-regret $R_T^{\mathcal{X}}$ using Theorem 12 with \mathcal{X} -EXP3 of Section 4.3, but we would like to find the best set \mathcal{X} . To do so, we may want to combine with EXP4 different instances of \mathcal{X} -EXP3, each using its own context set \mathcal{X} . We can then combine the bounds of Theorem 11 and Theorem 12 to ensure we are competing with the best possible context set \mathcal{X} . In this case, each instance of \mathcal{X} -EXP3 does not observe their own choice of action but the action chosen by EXP4 which follows a different distribution. The bound of Theorem 11 is valid but the regrets of the experts cannot be bounded using Theorem 12. It is however possible to use a variant of EXP4 to combine bandit algorithms

by adding an exploration parameter. We then lose however in the rate of the regret bound which is then of order

$$\max_{\mathcal{X}} R_T^{\mathcal{X}} \leq \mathcal{O}\left(T^{2/3} (\max_{\mathcal{X}} |\mathcal{X}| K \log K)^{1/3} \sqrt{\log M}\right)$$

where M is the number of context sets \mathcal{X} . We refer to Section 4.2.1 of Bubeck et al. [2012] for more details on this application.

5 Stochastic multi-armed bandits

During the last few chapters, we have reviewed the framework for comprehensive information. We designed algorithms to minimize regret for the different decision spaces Θ and loss assumptions f_t . Most of the algorithms were based on variations in the exponentially weighted average forecaster or online gradient descent. We also found some links with game theory, including the Blackwell approach, two-player zero-sum games, and calibration.

In this chapter, we consider the bandit setting, when the player only observes the performance of $f_t(\theta_t)$ but not $f_t(\theta)$ for $\theta \neq \theta_t$. We will start by providing fundamental results for stochastic bandits with finitely many actions, also called K -armed bandits which basically corresponds to $\Theta = \{1, \dots, K\}$ and i.i.d. loss functions f_t . For a thorough introduction to stochastic bandits we refer the interested student to the monographs Bubeck et al. [2012] or Lattimore and Szepesvári [2020].

5.1 Setting: stochastic bandit with finitely many actions

We state here the setting of stochastic bandits with finitely many actions (also called multi-armed bandit) and fix notations that we will use.

A multi-armed bandit problem is a sequential decision problem defined by a finite set of actions $\Theta \stackrel{\text{def}}{=} \{1, \dots, K\}$ also called *arms*. We assume that there are K unknown sequences $X_{i,1}, X_{i,2}, \dots$ of rewards in $[0, 1]$ associated with each arm $i = 1, \dots, K$. At each round, the player makes a decision by pulling an arm k_t in Θ and observes the corresponding reward¹ $X_{k_t,t}$. The objective of the player is to minimize his cumulative regret:

$$R_T \stackrel{\text{def}}{=} \max_{k=1, \dots, K} \sum_{t=1}^T X_{k,t} - \sum_{t=1}^T X_{k_t,t}.$$

In stochastic bandits, we generally assume the sequences to be i.i.d. Each arm $k = 1, \dots, K$ is associated an unknown probability distribution ν_k over $[0, 1]$ and $X_{k,t} \sim \nu_k$. We also denote

$$\mu_k \stackrel{\text{def}}{=} \mathbb{E}[X_{k,t}], \quad \text{and} \quad \mu^* \in \arg \max_{k=1, \dots, K} \{\mu_k\}.$$

The player aims at finding the arm with the highest mean reward μ_k as quickly as possible. The setting is summarized in Setting 1. Note that we retrieve the setting of online optimization (Setting 1) with the notation $X_{k,t} \stackrel{\text{def}}{=} 1 - f_t(k)$ with i.d.d. loss functions.

Unknown parameters: K probability distributions ν_1, \dots, ν_K on $[0, 1]$

At each time step $t = 1, \dots, T$

- the player chooses an action $k_t \in \Theta = \{1, \dots, K\}$;
- given k_t , the environment draws the reward $X_{k_t,t} \sim \nu_{k_t}$;
- the player only observes the feedback $X_{k_t,t}$.

Setting 1: Setting of stochastic bandit with finitely many actions

Multi-armed bandits have several concrete historical applications in a variety of fields, including ad placement, clinical trials, source routing or game AI. The name bandit refers to the “slot machine” in casinos, and the bandit problem corresponds to a player that inserts coins into different machines and tries to maximize

¹In the bandit community, it is more common to consider rewards rather than losses.

his payoff. The finite arms bandit settings we consider are simple enough to be analyzed and the algorithms can often be generalized to more realistic settings including for example contextual bandits.

Remark. Assume that all arms $\nu_k \sim \mathcal{B}(1/2)$ for $k = 1, \dots, K$. Then, $\mathbb{E}[X_{k,t}] = 1/2$ and $\mathbb{E}[X_{k_t,t}] = 1/2$. But because of fluctuations of random walks, the expected magnitude of the maximum rewards is of order

$$\mathbb{E} \left[\max_{k=1, \dots, K} \sum_{n=1}^T X_{k,n} \right] \approx \sqrt{T \log K}.$$

Therefore, in this case though all arms are optimal, the expected regret is of order $\sqrt{T \log K}$. We will thus consider a more quantity in the stochastic framework called the pseudo-regret which corresponds to competing with the best action in expectation, rather than the optimal action on the sequence of realized rewards.

Definition 4 (Pseudo-regret). *The pseudo-regret is defined as*

$$\bar{R}_T \stackrel{\text{def}}{=} T\mu^* - \mathbb{E} \left[\sum_{t=1}^T \mu_{k_t} \right],$$

where we recall $\mu_k = \mathbb{E}[X_{k,t}]$.

Remark that the pseudo-regret is upper-bounded by the expected regret $\bar{R}_T \leq \mathbb{E}[R_T]$. It is thus harder to design algorithms for the true regret but we will focus here on the pseudo-regret.

Useful notation In the following, we will denote by $\hat{\mu}_k(s)$ the empirical mean of rewards obtained after pulling arm k s times. Let us also denote for all arms $k = 1, \dots, K$ by

$$\Delta_k \stackrel{\text{def}}{=} \mu^* - \mu_k \quad \text{and} \quad N_k(t) \stackrel{\text{def}}{=} \sum_{s=1}^t \mathbb{1}_{k_s=k},$$

respectively the suboptimal gap of arm k and the number of times the arm k was pulled by the player before time t . Then, the pseudo-regret can be rewritten

$$\bar{R}_T = \left(\sum_{k=1}^K \mathbb{E}[N_k(t)] \right) \mu^* - \mathbb{E} \left(\sum_{k=1}^K N_k(t) \mu_k \right) = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(t)]. \quad (21)$$

We recall Hoeffding's inequality that will be used in the proofs. We will often use Azuma-Hoeffding's inequality which is a generalization of Hoeffding's inequality to martingals.

Proposition 3 (Hoeffding's Inequality). *If X_1, \dots, X_n are independent random variables almost surely in $[a, b]$ then for all $\delta \in (0, 1)$ we have*

$$\mathbb{P} \left\{ \sum_{t=1}^n X_k - \mathbb{E} \left[\sum_{t=1}^n X_k \right] \geq (b-a) \sqrt{\frac{n}{2} \log \frac{1}{\delta}} \right\} \leq \delta,$$

or equivalently for all $\varepsilon > 0$

$$\mathbb{P} \left\{ \sum_{t=1}^n X_k - \mathbb{E} \left[\sum_{t=1}^n X_k \right] \geq \varepsilon \right\} \leq \exp \left(- \frac{2\varepsilon^2}{n(b-a)^2} \right).$$

Parameter: $m \geq 1$.

1. Exploration

- For rounds $t = 1, \dots, mK$ explore by drawing each arm m times.
- Compute for each arm k its empirical mean of rewards obtained by pulling arm k m times

$$\hat{\mu}_k(m) = \frac{1}{m} \sum_{s=1}^{Km} \mathbb{1}_{k_t=k} X_{k,t}.$$

2. Exploitation: keep playing the best arm $\arg \max_k \hat{\mu}_k(m)$ for the remaining rounds $t = mK + 1, \dots, T$.

Algorithm 1: Explore-Then-Commit (ETC)

5.2 Explore-Then-Commit (ETC)

Contrary to the full information we examined earlier, the player only observes the rewards of the chosen actions. He must therefore make a trade-off between exploration and exploitation. The first bandit algorithm that we consider is Explore Then Commit (ETC). It consists in first performing an exploration phase of mK length in which each arm is pulled $m \geq 1$ times. Then it exploits by pulling the arm with the best empirical reward for the remaining rounds. It is summarized in Algorithm 1.

Theorem 13 (Thm 6.1, ?). *If $1 \leq m \leq T/K$ then*

$$\bar{R}_T \leq m \sum_{k=1}^K \Delta_k + (T - mK) \sum_{k=1}^K \Delta_k \exp(-m\Delta_k^2).$$

Proof. Assume without loss of generality that the first arm is optimal, i.e., $\mu_1 = \mu^*$ and $\Delta_1 = 0$. From (21), we have

$$\bar{R}_T = \sum_{k=1}^K \Delta_k \mathbb{E}[N_k(t)].$$

Let $k \geq 2$ be a suboptimal arm. Then, the arm k is selected m times during the exploration phase, and $T - mK$ times during the exploitation if k is selected, which implies $\hat{\mu}_k(m) \geq \hat{\mu}_1(m)$. Therefore,

$$\mathbb{E}[N_k(t)] \leq m + (T - mK) \mathbb{P}(\hat{\mu}_k(m) \geq \hat{\mu}_1(m))$$

Now, we can use Hoeffding's inequality to control the probability in the right-hand side. Indeed $\hat{\mu}_k(m)$ and μ_1 are respectively the empirical averages of m i.i.d. random variables in $[0, 1]$ of mean μ_k and $\mu_1 = \mu^*$. Therefore,

$$\begin{aligned} \mathbb{P}(\hat{\mu}_k(m) - \hat{\mu}_1(m) \geq 0) &= \mathbb{P}(\hat{\mu}_k(m) - \hat{\mu}_1(m) - \mu_k + \mu_1 \geq -\mu_k + \mu_1) \\ &= \mathbb{P}(\hat{\mu}_k(m) - \hat{\mu}_1(m) - \mu_k + \mu_1 \geq \Delta_k) \\ &= \mathbb{P}(m\hat{\mu}_k(m) - m\hat{\mu}_1(m) - m\mu_k + m\mu_1 \geq m\Delta_k) \\ &\leq \exp(-m\Delta_k^2). \end{aligned}$$

This implies

$$\bar{R}_T \leq m \sum_{k=1}^K \Delta_k + (T - mK) \sum_{k=1}^K \Delta_k \exp(-m\Delta_k^2).$$

□

The bound in Theorem 13 illustrates the trade-off between exploration and exploitation. If m is large, the exploration is too long and the first term $m \sum_{k=1}^K \Delta_k$ yields a large regret. On the other hand, for small m , there is a large probability to choose a suboptimal arm during the exploitation and the other term might lead to a large regret. The question is which value of m should we choose?

To have an idea, we will consider the case $K = 2$, in which case the bound is

$$\bar{R}_T \leq m\Delta_2 + T\Delta_2 \exp(-m\Delta_2^2).$$

Corollary 3. *If $K = 2$ and $m = \max\{1, \lceil \log(T\Delta_2^2)/\Delta_2^2 \rceil\}$, then*

$$\bar{R}_T \leq \Delta_2 + \frac{1 + \log(T\Delta_2^2)}{\Delta_2}.$$

The above bound is of order $O((\log T)/\Delta_2)$. Such bounds are called distribution-dependent because they heavily depend on the distributions ν_k via Δ_k . If $\Delta_2 \rightarrow 0$, it explodes. However, we also have from (21) that $\bar{R}_T \leq \Delta_2 T$. Therefore, in the worst case for any value of Δ_2 , Corollary 3 yields to the worst-case bound

$$\bar{R}_T \leq \min \left\{ T\Delta_2, \Delta_2 + \frac{1 + \log(T\Delta_2^2)}{\Delta_2} \right\} \lesssim \sqrt{T \log T}. \quad (22)$$

The above bound is close to be optimal. Yet, the issue is that the parameter m depends on Δ_2 and T . If the dependence on T can be dealt with a doubling-trick it is harder to optimize it in Δ_2 . Furthermore, when there are more than two arms, one might want to explore differently the arms. The upper-confidence-bound algorithm that we will see next does not have these issues.

Exercise 5.1. Show that it is possible to achieve the worst-case bound on the pseudo-regret of order $O(T^{2/3})$ by optimizing m independently of Δ (only with T).

Exercise 5.2. Generalize the results of Theorem 13 and 3 when the rewards are not-bounded but σ^2 -sub-Gaussian, i.e., for all $\lambda > 0$

$$\mathbb{E} \left[\exp(\lambda(X_{k,t} - \mathbb{E}[X_{k,t}])) \right] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right).$$

5.3 Upper-Confidence-Bound (UCB)

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC. It does not rely on an initial exploration phase but explores on the fly as rewards are observed. It explores and exploits sequentially throughout the experience. Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

To perform exploration and face uncertainty, the UCB algorithm is based on the *optimism principle*.

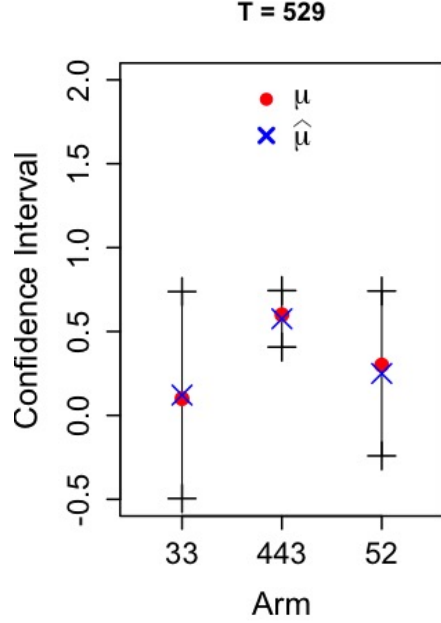
For each arm k , it builds a confidence interval on its expected reward based on past observation

$$I_k(t) = [LCB_k(k), UCB_k(t)]$$

where LCB is the Lower-Confidence-Bound and UCB is the Upper-Confidence-Bound. Then it is *optimistic* acting as if the best possible rewards are the real rewards and chooses the next arm accordingly

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} UCB_k(t).$$

In other words, it pulls the arm with the higher upper-confidence-bound. An example of how UCB works with three arms of means $\mu_1 = 0.1$, $\mu_2 = 0.6$ and $\mu_3 = 0.3$ is plotted in the Figure on the right. The best arm is pulled more often (see x-axis for number of times arms are selected) and his confidence interval is smaller.



The only question is how to design the upper-confidence-bounds. This is based on Hoeffding's inequality. Since the rewards are i.i.d. the distribution of $\hat{\mu}_k(s)$ is equal to the distribution of

$$\frac{1}{s} \sum_{s'=1}^s X_{k,s'},$$

with mean μ_k . Therefore, from Hoeffding's inequality, we have for all arms $k \in \{1, \dots, K\}$, for all $s \geq 1$ and all $\delta \in (0, 1)$

$$\mathbb{P} \left\{ \mu_k \geq \hat{\mu}_k(s) + \sqrt{\frac{\log \frac{1}{\delta}}{2s}} \right\} \leq \delta. \quad (23)$$

where $\hat{\mu}_k(t)$ is the empirical reward of arm k after pulling it t times. Therefore, it is reasonable to choose the upper-confidence bound

$$UCB_t(k) = \begin{cases} \infty & \text{if } N_k(t-1) = 0 \\ \hat{\mu}_k(N_k(t-1)) + \sqrt{\frac{2 \log t}{N_k(t-1)}} & \text{otherwise} \end{cases}$$

The UCB algorithm is described in Algorithm 2.

Theorem 14. *If the distributions ν_k have supports all included in $[0, 1]$ then for all k such that $\Delta_k > 0$*

$$\mathbb{E}[N_k(T)] \leq \frac{8 \log T}{\Delta_k^2} + 2.$$

In particular, this implies that the pseudo-regret of UCB is upper-bounded as

$$\bar{R}_T \leq 2K + \sum_{k: \Delta_k > 0} \frac{8 \log T}{\Delta_k}.$$

Remark. Let us make some remarks about the about upper-bound on the pseudo-regret.

Initialization For rounds $t = 1, \dots, K$ pull arm $k_t = t$

For $t = K + 1, \dots, T$, choose

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} \left\{ \hat{\mu}_k(N_k(t-1)) + \sqrt{\frac{2 \log t}{N_k(t-1)}} \right\},$$

and get reward $X_{k_t, t}$.

Algorithm 2: Upper-Confidence-Bound (UCB)

- UCB has a regret bound of order

$$\bar{R}_T \lesssim \frac{K \log T}{\Delta},$$

where $\Delta = \min_{i: \Delta_i > 0} \Delta_i$. Once again, using that the regret incurred from pulling arm k cannot be larger than $T\Delta_k$, this distribution-dependent upper-bound can be transformed into a distribution-free bound of order $\bar{R}_T \lesssim \sqrt{TK \log T}$. We leave this proof as an exercise.

- This bound is close to optimal since the lower bound is of order $O(\sqrt{KT})$. There exists modification of UCB to get rid of the extra logarithmic term. For instance, the MOSS algorithm (Minimax Optimal Strategy in the Stochastic Case) achieves

$$\bar{R}_T \lesssim \min \left\{ \sqrt{TK}, \frac{K}{\Delta} \log \frac{T\Delta^2}{K} \right\},$$

however it depends on the smallest gap Δ only and not on all gaps Δ_i .

- The assumption that the rewards are independent between arms can be relaxed.
- The assumption that the rewards are in $[0, 1]$ can be relaxed to a sub-Gaussian assumption.
- While a bound on the pseudo-regret is interesting, one would actually want a bound with high probability on

$$\hat{R}_T \stackrel{\text{def}}{=} T\mu^* - \sum_{t=1}^T \mu_{k_t, t}.$$

Using Hoeffding's inequality to control \hat{R}_T with $\bar{R}_T = \mathbb{E}[\hat{R}_T]$ would yield an additional term of order \sqrt{T} due to fluctuations which would dominate $O(K \log T / \Delta)$. Obtaining a bound of order $O(K \log T / \Delta)$ on \hat{R}_T is a challenging problem and not achieved by UCB. Some strategies using the knowledge of T can satisfy it.

Proof. Without loss of generality let us assume that the first arm is optimal, i.e., $\mu_1 = \mu^*$ and $\Delta_1 = 0$. We show below that if $k_t = k$, then at least one of the following three inequalities must be satisfied

$$\mu^* > \hat{\mu}_1(N_1(t-1)) + \sqrt{\frac{2 \log t}{N_1(t-1)}} \quad \leftarrow \mu^* \text{ larger than UCB} \quad (\text{i})$$

$$\mu_k < \hat{\mu}_k(N_k(t-1)) - \sqrt{\frac{2 \log t}{N_k(t-1)}} \quad \leftarrow \mu_k \text{ smaller than LCB} \quad (\text{ii})$$

$$N_k(t-1) \leq \frac{8 \log t}{\Delta_k^2} \quad \leftarrow k \text{ not played enough yet} \quad (\text{iii})$$

Indeed, otherwise assume that the three inequalities are all false than

$$\begin{aligned}
\widehat{\mu}_1(N_1(t-1)) + \sqrt{\frac{2 \log t}{N_1(t-1)}} &\geq \mu^* && \leftarrow \text{not (i)} \\
&\geq \mu_k + \Delta_k && \leftarrow \text{Def of } \Delta_k \\
&> \mu_k + 2\sqrt{\frac{2 \log t}{N_k(t-1)}} && \leftarrow \text{not (iii)} \\
&\geq \widehat{\mu}_k(N_k(t-1)) + \sqrt{\frac{2 \log t}{N_k(t-1)}} && \leftarrow \text{not (ii)}.
\end{aligned}$$

This contradicts the fact that $k_t = k$ (see Algorithm 2). Therefore, denoting $u = \lfloor \frac{8 \log T}{\Delta_k^2} \rfloor$, we have

$$\begin{aligned}
\mathbb{E}[N_k(T)] &= \sum_{t=1}^T \mathbb{E}[\mathbf{1}_{k_t=k}] = u + \sum_{t=u+1}^T \mathbb{P}\{k_t = k \text{ and (iii) is false}\} \\
&= u + \sum_{t=u+1}^T \left(\mathbb{P}\{(i) \text{ or (ii)}\} \right) \\
&\leq u + \sum_{t=u+1}^T \left(\mathbb{P}\{(i)\} + \mathbb{P}\{(ii)\} \right). \tag{24}
\end{aligned}$$

Therefore, it suffices to control the probabilities of (i) and (ii), which we do now. At round $t \geq 1$,

$$\begin{aligned}
\mathbb{P}\{(i)\} &\leq \mathbb{P}\left\{ \exists s \in \{1, \dots, t-1\}, \text{ such that } \mu^* > \widehat{\mu}_1(s) + \sqrt{\frac{2 \log t}{s}} \right\} \\
&\leq \sum_{s=1}^t \mathbb{P}\left\{ \mu^* > \widehat{\mu}_1(s) + \sqrt{\frac{\log(1/t^{-4})}{2s}} \right\} \\
&\stackrel{(23)}{\leq} \sum_{s=1}^t t^{-4} = t^{-3}.
\end{aligned}$$

By symmetry, the same applies for $\mathbb{P}\{(ii)\} \leq t^{-3}$. Combining into (24), it concludes the proof of the first inequality

$$\mathbb{E}[N_k(T)] \leq \frac{8 \log T}{\Delta_k^2} + 2 \sum_{t=u+1}^T t^{-3} \leq \frac{8 \log T}{\Delta_k^2} + 2.$$

The upper-bound on the pseudo-regret follows from (21). \square

5.4 Other algorithms

Other algorithms exist in the literature. The best known are ε -greedy and Thompson sampling.

5.4.1 ε -greedy

The idea of ε -greedy is very simple: first choose a parameter $\varepsilon \in (0, 1)$, then at each round, select the arm with the highest empirical mean with probability ε (i.e., be greedy), and explore by playing a random arm with probability ε . It works quite well in practice and is used in many application because of its simple implementation (in particular in reinforcement learning). Choosing $\varepsilon \approx K/(\Delta^2 T)$ yields to an upper-bound of order $K \log T/\Delta^2$. However it requires the knowledge of Δ .

5.4.2 Thompson Sampling

Thomson sampling was the first algorithm proposed for bandits by Thomson in 1933. It assumes a uniform prior over the expected rewards $\mu_i \in (0, 1)$, then at each round $t \geq 1$, for each arm $\pi_{k,t}$, it

- computes $\hat{\nu}_{k,t}$ the posterior distribution of the rewards of arm k given the rewards observed so far;
- samples $\theta_{k,t} \sim \hat{\nu}_{k,t}$ independently;
- selects $k_t \in \arg \max_{k \in \{1, \dots, K\}} \theta_{k,t}$.

Thomson sampling has a similar upper-bound of order $O(K \log T / \Delta)$ than the one achieved by UCB. An advantage over UCB is the possibility of incorporating easily prior knowledge on the arms.

5.5 Lower bounds for multi-armed bandit

In this section, we essentially state that the regret bound of UCB

$$\bar{R}_T \lesssim \min \left\{ \sqrt{KT \log T}, \sum_{k: \Delta_k > 0} \frac{\log T}{\Delta_k} \right\}$$

is close to optimal regret for multi-armed bandit.

5.5.1 Distribution-free lower bound

The next theorem shows that the previous results are not improvable (up to log factors).

Theorem 15 (Lower bound). *For any forecaster, there exists distributions ν_1, \dots, ν_K such that*

$$\bar{R}_T \gtrsim \sqrt{KT}.$$

The complete proof can be found in Bubeck et al. [2012]. We only present here the high-level idea of the proof. We design the adversary as follows: it generates i.i.d. Bernoulli rewards such that $\mathbb{E}[X_{k,t}] = \frac{1}{2}$ for all $k \in \{1, \dots, K\}$ except for one arm k^* where $\mathbb{E}[X_{k^*,t}] = \frac{1}{2} + \varepsilon$.

- Fact 1: to distinguish between a Bernoulli of parameter $1/2$ and a Bernoulli of parameter $1/2 + \varepsilon$, one needs $1/\varepsilon^2$ samples. This result can be proved formally by using Pinsker's inequality. The intuition goes as follows. From the Central Limit Theorem (or the distribution of a Binomial random variable), after T_k observations of an arm k , one can estimate its mean with an error of order $1/\sqrt{T_k}$. In other words, to estimate it with an error smaller than ε , one needs $T_k \approx \varepsilon^{-2}$ observations.
- Fact 2: at least one arm is sampled less than T/K times.

Assume that this arm is k^* , then the learner cannot distinguish it with other arms as soon as $T_{k^*} \leq T/K \leq \varepsilon^{-2}$, which corresponds to $\varepsilon \leq \sqrt{K/T}$. Choosing $\varepsilon = \sqrt{K/T}$, the pseudo-regret is then at least $(1 - 1/K)T\varepsilon \approx T\varepsilon \approx \sqrt{KT}$.

5.5.2 Distribution-dependent lower bound

Here, we show that the distribution dependent upper bound is not also optimal in the case of Bernoulli rewards.

A caveat with distribution dependent lower bounds is that for any distribution, there exists an algorithm with no-regret. For instance, consider a distribution ν_1, \dots, ν_K such that ν_1 is optimal (i.e., $\mu_1 = \max_k \mu_k$), then the algorithm that pull always the first arm will have zero regret. Yet such an algorithm will incur linear regret for some other distributions.

Hence, the following theorem states that any algorithm that incur sublinear regret for all distributions, achieves at best a pseudo regret of the same order of the one satisfied by UCB. The proof can be found in Bubeck et al. [2012].

Theorem 16 (Thm 2.2. Bubeck et al. [2012]). *Consider a strategy such that $\mathbb{E}[N_k(T)] = O(T^a)$ for any Bernoulli distributions, all suboptimal arms k and some $a > 0$. Then, for any Bernoulli distributions with means μ_k , we have*

$$\liminf_{T \rightarrow \infty} \frac{\bar{R}_T}{\log T} \geq \sum_{k: \Delta_k > 0} \frac{\mu^*(1 - \mu^*)}{\Delta_k}.$$

Note that the only difference with UCB is the factor $\mu^*(1 - \mu^*)$ which corresponds to the variance of the best arm. In the case of Bernoulli noise, the KL-UCB algorithm can take advantage of the knowledge that the rewards are Bernoulli to close this gap.

6 Contextual bandits

During last chapter, we considered the finite-armed bandit setting and saw several algorithms (ETC, UCB, ε -greedy, Thomson sampling) that achieve sublinear pseudo regret. UCB achieves for instance

$$\bar{R}_T \lesssim \min \left\{ \sum_{k:\Delta_k>0}^K \frac{\log T}{\Delta_k}, \sqrt{TK \log T} \right\},$$

where $\Delta_k \stackrel{\text{def}}{=} \mu^* - \mu_k$ is the suboptimality gap of arm k . The first bound is distribution dependent (it depends on the gaps Δ_k) and is of order $O(\log T)$ while the second bound is distribution free but is of order $O(\sqrt{T})$. In this chapter, we consider the more practical setting of contextual bandits, in which the learner observes a context $c_t \in \mathcal{C}$ before choosing the action k_t .

In most applications, before choosing an action k_t the player observes some context $c_t \in \mathcal{C}$.

For instance, consider a bandit problem in which the player needs to display ads on his website. At each new visitor, the player chooses an add to display and observes if the visitor click on it. The reward is one if there is a click and 0 otherwise. In this case, the player can see the cookie of the visitor before choosing the ad. A first step towards contextual bandits, is to consider continuous sets of actions \mathcal{X} , which may correspond to mapping between context and arms.

6.1 Continuous stochastic bandits

Let first generalize the finite-armed bandit setting to continuous set of arms in Setting 2.

Unknown parameters: $\nu(\theta)$, for each $\theta \in [0, 1]^d$, a probability distribution on $[0, 1]$ with expectation $\mu(\theta) \in [0, 1]$.

At each time step $t = 1, \dots, T$

- the player chooses an action $\theta_t \in \Theta \subseteq [0, 1]^d$;
- given θ_t , the environment draws the reward $Y_t \sim \nu(\theta_t)$ independently from the past;
- the player only observes the feedback Y_t .

The player wants to minimize its pseudo-regret defined as

$$\bar{R}_T \stackrel{\text{def}}{=} T\mu^* - \mathbb{E} \left[\sum_{t=1}^T \mu(\theta_t) \right],$$

where $\mu^* = \sup_{\theta \in \Theta} \mu(\theta)$.

Setting 2: Setting of stochastic bandit with continuous set of actions

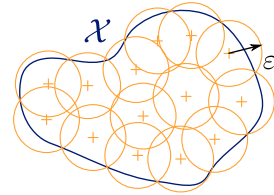
Similarly to what we did in the full-information setting with EWA, if the expectation function μ is β -Hölder: i.e., there exists $c > 0$

$$\forall \theta, \theta' \in \mathcal{X} \quad |\mu(\theta) - \mu(\theta')| \leq c \|\theta - \theta'\|^\beta,$$

then we may discretize the action space Θ and run any discrete bandit algorithm (UCB, ε -greedy, ...).

Theorem 17. *Let $\beta > 0$ and $\varepsilon > 0$. Assume that μ is β -Hölder. If UCB is run on an ε -covering of minimal cardinal of $\Theta \subset [0, 1]^d$, then it satisfies*

$$\bar{R}_T \lesssim T\varepsilon^\beta + \sqrt{\frac{T \log(T)}{\varepsilon^d}}.$$



In particular for $\varepsilon \approx \left(\frac{\log T}{T}\right)^{\frac{1}{2\beta+d}}$, we have $\bar{R}_T \lesssim T \left(\frac{\log T}{T}\right)^{\frac{\beta}{2\beta+d}}$.

Proof. An optimal ε -covering of $[0, 1]^d$ has cardinal of order $K \approx \varepsilon^{-d}$. Let $x^* \in \arg \max_{\theta \in \Theta} \mu(\theta)$ (we assume that it exists) and $\tilde{\theta}^*$ its ε -approximation, then the distribution-free upper-bound of UCB yields

$$\bar{R}_T \lesssim T(\mu(\theta^*) - \mu(\tilde{\theta}^*)) + \sqrt{KT \log T} \approx cT\varepsilon^\beta + \sqrt{\varepsilon^{-d}T \log T}.$$

The second part of the theorem is obtained by optimizing ε . □

To build the discretization, both β and T need to be known in advance. The horizon T can be calibrated online through a “doubling trick” (left as exercise). The parameter β may be tuned through bandit with experts (or bandits where arms are bandit algorithms) that we may see in next lecture (see Exp4 algorithm).

Note that the per-round complexity of such an algorithm is of order $\varepsilon^{-d} \approx T^{\frac{d}{2\beta+d}}$. Quite surprisingly it does not explodes with the dimension d and is always smaller than T . This is due to the fact that the higher the dimension d is, the worse will be the regret bound, and the cruder needs the discretization to be.

6.1.1 Contextual bandits through discretization

No we consider the following contextual bandit setting in which the player has a finite decision set $\Theta = \{1, \dots, K\}$ but observes a context $x_t \in \mathcal{X}$ before choosing his action.

Unknown parameters: $\nu(k, x)$, for each arm $k \in \{1, \dots, K\}$ and context $x \in \mathcal{X}$, a probability distribution on $[0, 1]$ with expectation $\mu(k, x) \in [0, 1]$.

At each time step $t = 1, \dots, T$

- the environment chooses $x_t \in \mathcal{X}$ and reveals it to the player;
- the player chooses an action $k_t \in \{1, \dots, K\}$;
- given k_t , the environment draws the reward $Y_t \sim \nu(k_t, x_t)$ independently from the past;
- the player only observes the feedback Y_t .

The player wants to minimize its pseudo-regret defined as

$$\bar{R}_T \stackrel{\text{def}}{=} \mathbb{E} \left[\sum_{t=1}^T \mu^*(x_t) - \sum_{t=1}^T Y_t \right],$$

where $\mu(k, x) = \mathbb{E}_{Y \sim \nu(k, x)}[Y]$ and $\mu^*(x) = \max_{k=1, \dots, K} \mu(k, x)$.

Setting 3: Setting of contextual stochastic bandit

Finite set of contexts If the set of context is finite $\mathcal{X} \stackrel{\text{def}}{=} \{1, \dots, |\mathcal{X}|\}$ we can denote

$$\bar{R}_T(c) \stackrel{\text{def}}{=} \mathbb{E} \left[\sum_{t=1}^T (\mu^*(x) - \mu(k_t, x)) \mathbb{1}_{x_t=x} \right]$$

the pseudo-regret due to context $x \in \mathcal{X}$. Then applying a separate instance of UCB (or any bandit algorithm) for each context $x \in \mathcal{X}$, we get by using the distribution-free upper-bound of UCB

$$\bar{R}_T(x) \lesssim \sqrt{T_x K \log T_x}, \quad \text{where} \quad T_x \stackrel{\text{def}}{=} \sum_{t=1}^T \mathbb{1}_{x_t=x}.$$

Note that because T_x are not known in advance it is important that the bound of UCB is anytime (i.e., that UCB does not need to know the horizon). The total pseudo-regret of UCB is then obtained by summing over all contexts

$$\bar{R}_T = \sum_{x \in \mathcal{X}} \bar{R}_T(x) \lesssim \sum_{x \in \mathcal{X}} \sqrt{T_x K \log T} \leq \sqrt{|\mathcal{X}| T K \log T},$$

where the last inequality is by Jensen's inequality using the concavity of the square root and $\sum_{x \in \mathcal{X}} T_x = T$.

Continuous set of contexts If the set of context is continuous $\mathcal{X} \subset [0, 1]^d$, one needs again to make assumption on the distributions $\nu(k, x)$ which needs to vary smoothly in x . Doing so, one may discretize the set of context with an ε -covering of \mathcal{X} of size $N \approx \varepsilon^{-d}$ and run an independent instance of UCB in each of the N bins.

Theorem 18. *Let $\beta > 0$ and $\varepsilon > 0$. Assume that $x \mapsto \mu(k, x)$ is β -Hölder for all $k \in \mathcal{X}$. If UCB is independently run in each bin of an optimal ε -covering of \mathcal{X} , then*

$$\bar{R}_T \lesssim T \varepsilon^\beta + \sqrt{\frac{KT \log T}{\varepsilon^d}}.$$

In particular for ε well-optimized, we have $\bar{R}_T \lesssim T \left(\frac{K \log T}{T} \right)^{\frac{\beta}{2\beta+d}}$.

Remark that in all these regret bounds, the suboptimal $\log T$ term can be removed by using MOSS (a minimax optimal variant of UCB).

Better rates using distribution-dependent bound? In the above results, we used the distribution-free regret bound of UCB. Because, if the function $\mu(\cdot, x)$ varies smoothly with x , there should be some context with zero suboptimality gaps. Yet, it is possible to get better rates by assuming the following α -margin assumption. It controls the suboptimality gap with high probability: the contexts x_t are i.i.d. and satisfy for all $\delta \in (0, 1)$

$$\mathbb{P} \left\{ \min_{k: \Delta(k, x_t) > 0} \Delta(k, x_t) < \delta \right\} \leq \square \delta^\alpha \quad (25)$$

where $\Delta(k, x_t) \stackrel{\text{def}}{=} \mu^*(x) - \mu(k, x)$ and \square is some constant. Note that the larger the value of α is the easier is the problem.

Theorem 19 (Theorem 4.1, Perchet and Rigollet [2013]). *Let $\alpha \in (0, 1)$, $\beta > 0$ and $\varepsilon > 0$. Assume that $c \mapsto \mu(k, x)$ is β -Hölder for all $k \in \mathcal{X}$ and that the α -margin assumption (25) holds. Running a bandit algorithm (similar to UCB) independently in each bin of an optimal ε -covering of \mathcal{X} , we get*

$$\bar{R}_T \lesssim T \left(\frac{K \log K}{T} \right)^{\frac{\beta(\alpha+1)}{2\beta+d}},$$

for optimized ε .

The proof (for another algorithm than UCB) may be found in Perchet and Rigollet [2013]. We see that the factor α improves the rate of convergence with respect to the previous rate.

6.1.2 Stochastic Linear bandits

Contextual bandits that we just saw generalizes multi-armed bandits by allowing contexts. However, the dimension of the context space significantly worsen the regret rate from \sqrt{T} to $T^{\frac{d+1}{d+2}}$ for Lipschitz rewards for instance (β -Hölder with $\beta = 1$). In this part, we will see *Stochastic linear bandits*, in which we assume

Unknown parameter: $\mu^* \in \mathbb{R}^d$.

At each time step $t = 1, \dots, T$

- the environment chooses $\Theta_t \subseteq \mathbb{R}^d$ the decision set;
- the player chooses an action $\theta_t \in \Theta_t$;
- given θ_t , the environment draws the reward

$$Y_t = \theta_t \cdot \mu^* + \varepsilon_t$$

where ε_t is i.i.d. 1-subgaussian noise.

- the player only observes the feedback Y_t .

The player wants to minimize its pseudo-regret defined as

$$\bar{R}_T \stackrel{\text{def}}{=} \mathbb{E} \left[\sum_{t=1}^T \max_{\theta \in \Theta_t} \theta \cdot \mu^* - \sum_{t=1}^T Y_t \right].$$

Setting 4: Setting of stochastic linear bandit

the rewards to have a linear structure. This includes rich classes of models and allows better regret of order $O(\sqrt{T})$.

The setting of stochastic linear bandits is described in Setting 4. For simplicity, the noise ε_t is assumed to be i.i.d. and 1-subgaussian noise: i.e., $\mathbb{E}[\varepsilon_t] = 0$ and

$$\forall \lambda > 0, \quad \mathbb{E}[\exp(\lambda \varepsilon_t)] \leq \exp(\lambda^2/2)$$

almost surely. Note that we could consider σ^2 -subgaussian noise, or make it depend on the past $\mathcal{F}_t = \sigma(x_1, \varepsilon_1, \dots, x_t, \varepsilon_t)$ with $\mathbb{E}[\varepsilon_t | \mathcal{F}_t] = 0$.

Particular cases: why is this setting interesting? Different choices of decision sets \mathcal{X}_t lead to different settings of stochastic bandits:

- *Finite-armed bandit*: if $\Theta_t = (e_1, \dots, e_d)$ where e_i are the unit vectors in \mathbb{R}^d and $\mu^* = (\mu_1, \dots, \mu_d)$, we recover the setting of finite-armed bandit.
- *Contextual linear bandit*: we can recover a particular case of Setting 3, if $x_t \in \mathcal{X}$ is a context observed by the player and the reward function μ is of the form

$$\mu(\theta, x) = \psi(\theta, x) \cdot \mu^*, \quad \forall (\theta, x) \in [K] \times \mathcal{X},$$

for some unknown parameter $\mu^* \in \mathbb{R}^d$ and *feature map* $\psi : [K] \times \mathcal{X} \rightarrow \mathbb{R}^d$. For example, assume that you are a website which wants to display ads to visitors. The context x_t can be the cookie of the visitor containing information about what he likes, the actions are the possible ads to be displayed and the reward tells if there is a click. If the possible interests of the visitor are grouped in finite categories (such as traveling), so are the ads (in groups of products, such as flight tickets), the feature maps could contained all the combinations between interests and groups of products. The unknown vector θ^* would be tell which interests and groups of products are positively correlated. Of course the feature map could be created using any methods (such as deep-learning or splines).

- *Combinatorial bandit*: if $\Theta_t \subseteq \{0, 1\}^d$ yields to combinatorial bandit problems. For instance, the decision set corresponds to possible paths in a graph, the vector μ^* assigns to each edge a reward corresponding to its cost and the goal is to find the smallest path with smallest cost.

Algorithm: LinUCB As we saw earlier with UCB, the “optimism principle” is a good option for bandit problems to explore. The LinUCB algorithm is based on the same principle:

1. Build confidence set that contain μ^* : $\mu^* \in C_t$ with high probability
2. Build confidence upper-bound on the rewards: for all $\theta \in \Theta_t$

$$\text{UCB}_t(\theta) = \max_{\mu \in C_t} \theta \cdot \mu \quad (26)$$

3. Be optimistic: act as if the best possible rewards were the true rewards

$$\theta_t \in \arg \max_{\theta \in \Theta_t} \text{UCB}_t(\theta). \quad (27)$$

Therefore the only remaining question is how to build the confidence set $C_t \subseteq \mathbb{R}^d$? They should contain μ^* with high probability but be as small as possible. Given the observed rewards the key is thus to estimate the parameter μ^* . Denoting by I_d the $d \times d$ identity matrix and picking $\lambda > 0$, we can estimate μ^* with *regularized least square*

$$\hat{\mu}_t \stackrel{\text{def}}{=} \arg \min_{\mu \in \mathbb{R}^d} \left\{ \sum_{s=1}^t (Y_s - \theta_s \cdot \mu)^2 + \lambda \|\mu\|^2 \right\} = V_t^{-1} \sum_{s=1}^t \theta_s Y_s,$$

where $V_t \stackrel{\text{def}}{=} \lambda I_d + \sum_{s=1}^t \theta_s \theta_s^\top$. We have the following result whose proof can be found in Lattimore and Szepesvári [2020].

Lemma 3 (Theorem 20.2, Lattimore and Szepesvári [2020]). *Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$, if $\max_{\theta \in \Theta_t} \|\theta\|_2 \leq 1$, for all $t \geq 1$*

$$\|\hat{\mu}_t - \mu^*\|_{V_t} \leq \sqrt{\lambda} \|\mu^*\| + \sqrt{2 \log(1/\delta) + d \log \left(1 + \frac{T}{\lambda} \right)} \stackrel{\text{def}}{=} \beta(\delta),$$

where $\|\mu\|_{V_t}^2 = \mu^\top V_t \mu$.

The above lemma, states that with probability $1 - \delta$, for all $t \geq 1$,

$$\mu^* \in C_t, \quad \text{where } C_t \stackrel{\text{def}}{=} \left\{ \theta \in \mathbb{R}^d : \|\mu - \hat{\mu}_{t-1}\|_{V_{t-1}} \leq \beta(\delta/T) \right\}. \quad (28)$$

Proof of Lemma 3. The proof relies on Laplace’s method on super-martingales which is a standard argument to provide confidence bounds on a self-normalized sum of conditionally centered random vectors. We have

$$\hat{\mu}_t = V_t^{-1} \sum_{s=1}^t \theta_s Y_s = V_t^{-1} \sum_{s=1}^t \theta_s (\theta_s^\top \mu^* + \varepsilon_s) = V_t^{-1} (V_t - \lambda I_d) \mu^* + M_t = \mu^* - \lambda V_t^{-1} \mu^* + V_t^{-1} M_t,$$

where we introduced $M_t = \sum_{s=1}^t \theta_s \varepsilon_s$, which is a martingale with respect to $\mathcal{F}_t = \sigma(\varepsilon_1, \dots, \varepsilon_t)$. Therefore, by triangle inequality

$$\|V_t^{-1/2} (\hat{\mu}_t - \mu^*)\| = \left\| -\lambda V_t^{-1/2} \mu^* + V_t^{-1/2} M_t \right\| \leq \lambda \|V_t^{-1/2} \mu^*\| + \|V_t^{-1/2} M_t\|.$$

On the one hand, given that all eigenvalues of the symmetric matrix V_t are larger than λ , all eigenvalues of $V_t^{-1/2}$ are smaller than $1/\sqrt{\lambda}$ and thus

$$\lambda \|V_t^{-1/2} \mu^*\| \leq \lambda \frac{1}{\sqrt{\lambda}} \|\mu^*\| = \sqrt{\lambda} \|\mu^*\|.$$

We now prove, on the other hand, that with probability at least $1 - \delta$

$$\|V_t^{-1/2} M_t\| \leq \sqrt{2 \log \frac{1}{\delta} + d \log \frac{1}{\lambda} + \log \det(V_t)}.$$

Upper-bounding $\log \det(V_t) \leq d \log(\lambda + t)$ (since all the eigenvalues of V_t are smaller than $\lambda + t$) will then conclude the proof of the Theorem.

Step 1: Introducing super-martingales. For all $\nu \in \mathbb{R}^d$, we consider

$$S_{t,\nu} = \exp\left(\nu^\top M_t - \frac{1}{2}\nu^\top V_t \nu\right)$$

and now show that it is an \mathcal{F}_t -super-martingale. First, note that since the common distribution of the $\varepsilon_1, \dots, \varepsilon_t$ is 1-sub-Gaussian, the for all \mathcal{F}_{t-1} -measurable random variable ν_{t-1}

$$\mathbb{E}\left[e^{\nu_{t-1}^\top \varepsilon_t} \mid \mathcal{F}_{t-1}\right] \leq e^{\frac{\nu_{t-1}^2}{2}}.$$

Now,

$$\mathbb{E}\left[S_{t,\nu} \mid \mathcal{F}_{t-1}\right] = S_{t-1,\nu} \mathbb{E}\left[\exp\left(\nu^\top \theta_t \varepsilon_t - \frac{1}{2}\nu^\top \theta_t \theta_t^\top \nu\right) \mid \mathcal{F}_{t-1}\right] \leq S_{t-1,\nu}.$$

Note that rewriting $S_{t,\nu}$ in its vertex form is, with $m = V_t^{-1} M_t$:

$$S_{t,\nu} = \exp\left(\frac{1}{2}(\nu - m)^\top V_t (\nu - m)\right) \times \exp\left(\frac{1}{2}\|V_t^{-1/2} M_t\|^2\right).$$

Step 2: Laplace's method-integrating $S_{t,\nu}$ over $\nu \in \mathbb{R}^d$. The basic observation behind this method is that (given the vertex form) $S_{t,\nu}$ is maximal at $\nu = m = V_t^{-1} M_t$ and then equals $\exp\left(\frac{1}{2}\|V_t^{-1/2} M_t\|^2\right)$, which is (a transformation of) the quantity to control. Now, because the exp function quickly vanishes, the integral over $\nu \in \mathbb{R}^d$ is close to its maximum. We therefore consider

$$\bar{S}_t = \int_{\mathbb{R}^d} S_{t,\nu} d\nu.$$

We will make repeated uses of the fact that the Gaussian density function

$$\nu \mapsto \frac{1}{\sqrt{\det(2\pi C)}} \exp\left(-\frac{1}{2}(\nu - m)^\top C^{-1}(\nu - m)\right),$$

where $m \in \mathbb{R}^d$ and C is a symmetric positive definite matrix, integrate to 1 over \mathbb{R}^d . This gives us the first rewriting

$$\bar{S}_t = \sqrt{\det(2\pi V_t^{-1})} \exp\left(\frac{1}{2}\|V_t^{-1/2} M_t\|^2\right).$$

Second, by the Fubini-Tonelli theorem and the super-martingale property

$$\mathbb{E}[S_{t,\nu}] \leq \mathbb{E}[S_{0,\nu}] = \exp(-\lambda\|\nu\|^2/2),$$

we also have

$$\mathbb{E}[\bar{S}_t] \leq \int_{\mathbb{R}^d} \exp(-\lambda\|\nu\|^2/2) d\nu = \sqrt{\det(2\pi\lambda^{-1}I_d)}.$$

Combining the two statements, we proved

$$\mathbb{E}\left[\exp\left(\frac{1}{2}\|V_t^{-1/2} M_t\|^2\right)\right] \leq \sqrt{\frac{\det(V_t)}{\lambda^d}}.$$

Step 3: Markov-Chernov bound. For $u > 0$,

$$\begin{aligned} \mathbb{P}\left(\|V_t^{-1/2} M_t\| > u\right) &= \mathbb{P}\left(\frac{\|V_t^{-1/2} M_t\|^2}{2} > \frac{u^2}{2}\right) \\ &\leq \exp\left(-\frac{1}{2}u^2\right) \mathbb{E}\left[\exp\left(\frac{1}{2}\|V_t^{-1/2} M_t\|^2\right)\right] \leq \exp\left(-\frac{u^2}{2} + \frac{1}{2} \log \frac{\det(V_t)}{\lambda^d}\right) = \delta, \end{aligned}$$

for the claimed choice

$$u = \sqrt{2 \log \frac{1}{\delta} + d \log \frac{1}{\lambda} + \log \det(V_t)}.$$

□

Theorem 20 (Corollary 19.2, ?). *Let $T \geq 1$ and $\mu^* \in \mathbb{R}^d$. Assume that for all $\theta \in \cup_{t=1}^T \Theta_t$, $|\theta \cdot \mu^*| \leq 1$, $\|\mu^*\| \leq 1$ and $\|\theta_t\| \leq 1$, then LinUCB with C_t defined in (28) satisfies the pseudo-regret bound*

$$\bar{R}_T \leq \square_\lambda d \sqrt{T} \log(T),$$

where \square_λ is a constant that may depend on λ .

Proof. Let $\delta = 1/T$. By Lemma 3, with probability $1 - 1/T$,

$$\forall t \geq 1, \quad \mu^* \in C_t. \quad (29)$$

Step 1: Small instantaneous regrets under the event (29). Assume that (29) holds. Let

$$\theta_t^* \stackrel{\text{def}}{=} \max_{\theta \in \Theta_t} \theta \cdot \mu^* \quad \text{and} \quad r_t \stackrel{\text{def}}{=} (\theta_t^* - \theta_t) \cdot \mu^*$$

be respectively the optimal decision and the instantaneous regret at round t . We also define

$$\tilde{\mu}_t \in \arg \max_{\mu \in C_t} \{\theta_t \cdot \mu\}.$$

Since $\mu^* \in C_t$, we have

$$\theta_t^* \cdot \mu^* \leq \max_{\mu \in C_t} \{\theta_t^* \cdot \mu\} \stackrel{(26)}{=} UCB_t(\theta_t^*) \stackrel{(27)}{\leq} UCB_t(\theta_t) = \max_{\mu \in C_t} \{\theta_t \cdot \mu\} = \theta_t \cdot \tilde{\mu}_t,$$

which entails because μ^* and $\tilde{\mu}_t$ belong to C_t ,

$$r_t \stackrel{\text{def}}{=} (\theta_t^* - \theta_t) \cdot \mu^* \leq \theta_t \cdot (\tilde{\mu}_t - \mu^*) \stackrel{\text{Cauchy-Schwarz}}{\leq} \|\theta_t\|_{V_{t-1}^{-1}} \|\tilde{\mu}_t - \mu^*\|_{V_{t-1}} \leq 2 \|\theta_t\|_{V_{t-1}^{-1}} \beta(1/T^2).$$

Therefore, summing over $t = 1, \dots, T$ and using $r_t \leq 2$, we have

$$\begin{aligned} R_T &\stackrel{\text{def}}{=} \sum_{t=1}^T r_t \leq \sqrt{T \sum_{t=1}^T r_t^2} \quad \leftarrow \text{Jensen's inequality} \\ &\leq 2 \sqrt{T \sum_{t=1}^T \min \left\{ \|\theta_t\|_{V_{t-1}^{-1}}^2 \beta(1/T^2)^2, 1 \right\}} \\ &\leq 2\beta(1/T^2) \sqrt{T \sum_{t=1}^T \min \left\{ \|\theta_t\|_{V_{t-1}^{-1}}^2, 1 \right\}} \quad \leftarrow \beta_T(1/T^2) \geq 1 \\ &\leq 2\beta(1/T^2) \sqrt{T \sum_{t=1}^T \log \left(1 + \|\theta_t\|_{V_{t-1}^{-1}}^2 \right)} \quad \leftarrow \min\{u, 1\} \leq 2 \log(1+u). \end{aligned}$$

But, we have

$$\begin{aligned}
1 + \|\theta_t\|_{V_{t-1}^{-1}}^2 &= \det(1 + \|\theta_t\|_{V_{t-1}^{-1}}^2) \\
&= \det\left(V_{t-1}^{-1}(V_{t-1} + V_{t-1}^{1/2}\|\theta_t\|_{V_{t-1}^{-1}}^2 V_{t-1}^{1/2})\right) \leftarrow \text{using } \det(I + AB) = \det(I + BA) \\
&= \det\left(V_{t-1}^{-1}(V_{t-1} + \theta_t \theta_t^\top)\right) \leftarrow \|\theta_t\|_{V_{t-1}^{-1}} = V_{t-1}^{-1/2} \theta_t \theta_t^\top V_{t-1}^{-1/2} \\
&= \det(V_{t-1}^{-1} V_t) \leftarrow V_t = V_{t-1} + \theta_t \theta_t^\top \\
&= \frac{\det(V_t)}{\det(V_{t-1})} \leftarrow \det(A^{-1}B) = \frac{\det(B)}{\det(A)}.
\end{aligned}$$

Substituting into the regret bound, the sum telescopes and it entails

$$R_T \leq 2\beta(1/T^2) \sqrt{T \log\left(\frac{\det(V_T)}{\det(V_0)}\right)}.$$

Then, using $V_0 \stackrel{\text{def}}{=} \lambda I_d$ and since $V_T = \lambda I_d + \sum_{t=1}^T \theta_t \theta_t^\top$ with $\|\theta_t\| \leq 1$, all eigenvalues of V_T lie in $[\lambda, \lambda + T]$ which yields

$$\det(V_0) = \lambda^d \quad \text{and} \quad \det(V_T) \leq (\lambda + T)^d.$$

Plugging back into the previous upper-bound and using that $\beta(1/T^2) \leq \square_\lambda \sqrt{d \log T}$

$$R_T \leq 2\sqrt{dT\beta(1/T^2) \log\left(1 + \frac{T}{\lambda}\right)} \leq \square_\lambda d\sqrt{T} \log T.$$

Part 2: without the event (28) We because $r_t \leq 2$, almost surely $R_T \leq 2T$, and we have

$$\begin{aligned}
\bar{R}_T &= \mathbb{E}[R_T] \leq \mathbb{E}\left[R_T \mid \text{Event (28)}\right] \mathbb{P}\{\text{Event (28)}\} + 2T(1 - \mathbb{P}\{\text{Event (28)}\}) \\
&\leq \square d\sqrt{T} \log T + 2.
\end{aligned}$$

This concludes the proof. \square

Better regret with assumptions It is worth pointing out that if we make additional assumptions, it is possible to improve the regret bound $O(d\sqrt{T} \log T)$. A first setting corresponds to the case where the set of available actions at time t is fixed and finite; i.e., the learner needs to choose $\theta_t \in \Theta$ where $|\Theta| = K$. Then, it is possible to achieve

$$R_T \leq \square \sqrt{Td \log(TK)},$$

which improves the previous bound by a factor $\sqrt{d}/\log(K)$ and improves the classical bound of UCB $O(\sqrt{TK \log T})$ by a factor K/\sqrt{d} . These improvements can be significant when $K \gg d \gg 1$. We refer the curious reader to [Lattimore and Szepesvári, 2020, Chapter 22].

Another possible improvement when $d \gg 1$ is to assume that μ^* is m_0 -sparse (i.e., most of its components are zero). Then under assumptions, one can get a regret of order $\tilde{O}(\sqrt{dm_0 T})$.

6.2 Other possible extensions of bandits

Note that there exist many different extensions of stochastic bandits to make it more realistic or with improved regret.

- *Bandit with delays*: For instance, consider the example of the website which wants to display ads. The website does not observe if there is no click, he needs to fix some time after which he consider that the visitor will not click, and if the visitor stays long on the webpage, the website may need to display ads to other visitors before getting the rewards. There is thus delayed feedback the website needs to deal with.

- *Non stationarity*
- *Combinatorial bandits*
- *Dueling bandits*
- ...

We refer the interested student to the monograph Lattimore and Szepesvári [2020] for more information on these settings. Next week, we will deal with adversarial bandits.

References

- Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- Jyrki Kivinen and Manfred K Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *information and computation*, 132(1):1–63, 1997.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Vianney Perchet and Philippe Rigollet. The multi-armed bandit problem with covariates. *The Annals of Statistics*, pages 693–721, 2013.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 928–936, 2003.