

Online Learning

Pierre Gaillard

January 16, 2025

INRIA

References

These monographs are available online.

- Cesa-Bianchi and Lugosi 2006
- Shalev-Shwartz et al. 2012
- Orabona 2019
- Hazan et al. 2016
- Lattimore and Szepesvári 2019

Classical Machine Learning

In classical supervised machine learning, the learner

1. observes training data with labels,
2. builds a program to minimize the training error
3. controls the error of new data if they are similar to the training data



→ Learning method → Prediction on test data

Sequential Learning

In some applications, the environment may evolve over time and the data may be available sequentially.

Spam detection: can be seen as a game between spammer and spam filters. Each trying to fool the other one. The data is possibly adversarial.

Necessity to take a robust approach by learning as ones goes along from experiences as more aspects of the problem are observed.

This is the goal of sequential learning (or sequential learning).

Sequential learning

In sequential learning, we do not have any training data.

Data are **acquired and treated on the fly**.

Feedbacks are received and algorithms updated step by step.



This field has received a lot of attention recently because of the possible applications coming from internet:

- ads to display,
- repeated auctions,
- spam detection,
- experts/algorithm aggregation

Setting of an online learning problem/online convex optimization

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta$ \rightarrow full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$ \rightarrow bandit feedback.

Goal. Minimize the cumulative loss:

$$\widehat{L}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t).$$

A simple stochastic model:

- K arms (actions: here price signals)
- Each arm k is associated an **unknown** probability distribution with mean μ_k



Setting: sequentially pick an arm k_t and get reward $X_{k_t,t}$ with mean μ_{k_t}

Goal: maximize the expected cumulative reward

$$\mathbb{E} \left[\sum_{t=1}^T X_{k_t,t} \right]$$

Exploration vs Exploitation trade-off.

Bandit applications

Maximize one's gains in casino? Hopeless ...



Historical motivation (Thomson, 1933): clinical trials, for each patient t in a clinical study

- choose a treatment k_t
- observe response to the treatment $X_{k_t,t}$

Goal: maximize the number of patient healed (or find the best treatment)

Successful because of many applications coming from Internet: recommender systems, online advertisements,...

Setting of an online learning problem – Multi-armed bandits

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t = k_t \in \Theta := \{1, \dots, K\}$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$ (by sampling the arms);
- the player suffers loss $\ell_t(\theta_t) = 1 - X_{k_t, t}$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta \rightarrow$ full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t) = X_{k_t, t} \rightarrow$ bandit feedback.

The goal of the player is to minimize his cumulative loss:

$$\widehat{L}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t).$$

Example 2: 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.

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 $\theta_t := p_t$, we recover the general setting. The inputs correspond to the expert advice $x_t(k)$ that are often revealed before the learner makes his decision p_t .

Example 2: 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.

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.

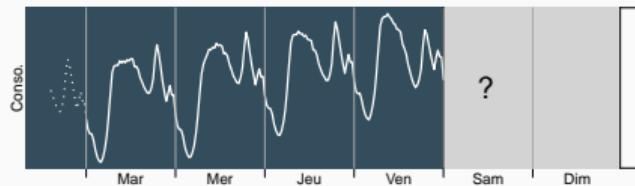
Player's performance is then measured via a loss function $\ell_t(p_t) = \ell(\hat{y}_t, y_t)$ which measures the distance between the prediction \hat{y}_t and the output y_t :

- | | |
|---|--|
| - squared loss $\ell(\hat{y}_t, y_t) = (\hat{y}_t - y_t)^2$ | $\ell(\hat{y}_t, y_t) = \hat{y}_t - y_t / y_t $ |
| - absolute loss $\ell(\hat{y}_t, y_t) = \hat{y}_t - y_t $ | - pinball loss. |
| - absolute percentage of error | |

All these loss functions are convex, which will play an important role in the analysis.

Example: Prediction with expert advice for electricity forecasting

Short term prediction (one day ahead) of the French electricity consumption

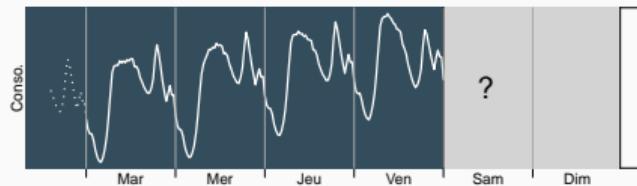


Important because electricity is hard to store.



Example: Prediction with expert advice for electricity forecasting

Short term prediction (one day ahead) of the French electricity consumption

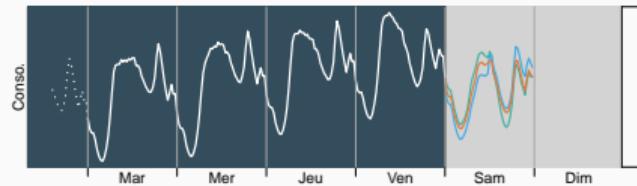


Many experts (statisticians or data scientists)
design prediction models:

Simultaneously, the French electricity market is evolving (electric cars,...)

Example: Prediction with expert advice for electricity forecasting

Short term prediction (one day ahead) of the French electricity consumption

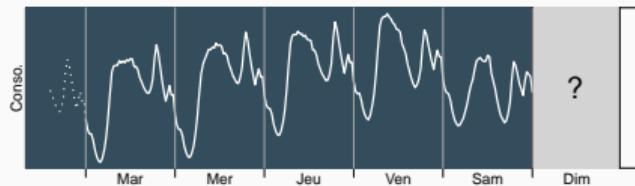


Many experts (statisticians or data scientists)
design prediction models:

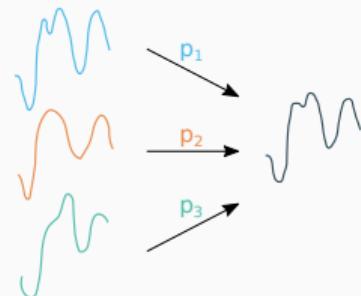
Simultaneously, the French electricity market is evolving (electric cars,...)

Example: Prediction with expert advice for electricity forecasting

Short term prediction (one day ahead) of the French electricity consumption



Combine the predictions using adaptive methods:



Each day,

1. Assign a weight to each expert based on past performance

$$\theta_t = \text{weight vector}$$

2. Predict the weighted average $\hat{y}_t = \langle \theta_t, x_t \rangle$ and suffer loss

$$\ell_t(\theta_t) = (y_t - \hat{y}_t)^2$$

How to measure the performance? The regret

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.

→ 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 (Regret)

The regret of the player with respect to a fixed parameter $\theta^* \in \Theta$ after T time steps is

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

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

Regret decomposition

We have some approximation-estimation decomposition:

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

We will focus on the regret in these lectures.

The goal of the player is to ensure a sublinear regret $\text{Reg}_T = o(T)$ as $T \rightarrow \infty$ and this for any possible sequence of losses ℓ_1, \dots, ℓ_T .

→ the average performance of the player will approach on the long term the one of the best decision.

Adversarial / Stochastic setting

The losses ℓ_t are unknown to the player beforehand and may be:

- **Adversarial setting** (lessons 1, 2, and 3): No stochastic 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_1, \dots, ℓ_T and θ^* . Players's goal is to control the quantity:

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

- **Stochastic setting** (lessons 4, 5, and 6): the losses are generated by some stochastic process (e.g., i.i.d.). The regret bounds hold then in expectation or with high probability.

Why a different loss at every round t ?

This may be caused by many phenomena, e.g. by

- 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 $\ell_t(x) = (x - y_t)^2$.
- noise: the environment is stochastic and the variation over time t models some noise effect.
- 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: what about best θ_t^* at every round?

Regret

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

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

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

where the infimum is inside the sum.

Exercise: Show that the environment can ensure Reg_T^* to be linear in T by choosing properly the loss functions ℓ_t .

Setting of an online learning problem/online convex optimization

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta$ \rightarrow full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$ \rightarrow bandit feedback.

Goal. Minimize the regret

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

Online Linear Optimization

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

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

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 .

How to choose the weights

At round t the player needs to choose a weight vector $p_t \in \Delta_K$.

How to choose the weights? The player should

- give more weight to actions that performed well in the past.
- not give all the weight to the current best action, otherwise it would not work (see Exercise next).

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)

The exponentially weighted average forecaster

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

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

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

Exercise:

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 $\text{Reg}_T \geq \Omega(T)$.

Solution

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

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

Exercise:

1. Show that $p_t(k) > 0$ implies that $k \in \arg \min_j \sum_{s=1}^{t-1} g_s(j)$

Solution

Assume that there exists $k \in [K]$ such that $p_t(k) > 0$ and $k \notin \arg \min_j \sum_{s=1}^{t-1} g_s(j)$. Then, there exists $k' \neq k$ such that $\sum_{s=1}^{t-1} g_s(k') < \sum_{s=1}^{t-1} g_s(k)$. Therefore,

$$\begin{aligned} \sum_{s=1}^{t-1} \ell_s(p_t) &= \sum_{s=1}^{t-1} \sum_{j=1}^K p_s(j) g_s(j) = \sum_{s=1}^{t-1} \sum_{j \neq k} p_s(j) g_s(j) + p_s(k) \sum_{s=1}^{t-1} g_s(k) \\ &> \sum_{s=1}^{t-1} \sum_{j \neq k} p_s(j) g_s(j) + p_s(k') \sum_{s=1}^{t-1} g_s(k) = \sum_{s=1}^{t-1} \ell_s(q_t), \end{aligned}$$

where $q_t(j) = p_t(j)$ if $j \notin \{k, k'\}$ and $q_t(k) = 0$ and $q_t(k') = p_t(k') + q_t(k')$. This yields a contradiction.

Solution

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

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

Exercise:

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 $g_1, \dots, g_T \in [0, 1]^K$ such that $\text{Reg}_T \geq \Omega(T)$.

Solution

It suffices to choose $g_t(k) = 1$ if $p_t(k) > 0$ and $g_t(k) = 0$ otherwise. The cumulative loss of FTL is T while there exists an action with cumulative loss smaller than T/K .

Regret guarantee for EWA

Theorem 1 (Regret bound for EWA)

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

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

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 $\text{Reg}_T \leq 2\sqrt{T \log K}$.

This regret bound is optimal up to constant factor (see [1]).

Exercise: Generalize the above theorem when the losses $g_1, \dots, g_T \in [-B, B]^K$ for some $B > 0$.

[1] Cesa-Bianchi and Lugosi 2006.

Proof (Step 1 - Reformulation of the regret for linear losses)

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

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

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

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

Proof (Step 2 – Upper-bound of W_T)

We denote $W_t(j) = e^{-\eta \sum_{s=1}^t g_t(j)}$ 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 $p_t = (p_t(1), \dots, p_t(K))$.

Proof (Step 2 - Upper-bound of W_T)

Now, using $1 + x \leq e^x$, we get:

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

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

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

Proof (Step 3 – Lower-bound of W_T)

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:

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

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

Regret guarantee for EWA

Theorem 1 (Regret bound for EWA)

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

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

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 $\text{Reg}_T \leq 2\sqrt{T \log K}$.

This regret bound is optimal up to constant factor (see [1]).

Exercise: Generalize the above theorem when the losses $g_1, \dots, g_T \in [-B, B]^K$ for some $B > 0$.

[1] Cesa-Bianchi and Lugosi 2006.

The previous algorithms 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**.

Anytime algorithm: the doubling trick

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 regret of EWA-doubling is then upper-bounded as:

$$\text{Reg}_T \leq 7\sqrt{T \log K}.$$

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 $T \geq 1$.

Proof

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

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

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:

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

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

Reminder of the setting of prediction with expert advice

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.

The goal is to minimize the regret with respect to the best expert

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

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

Reminder of the setting of prediction with expert advice

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.

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 :

- squared loss $\ell(\hat{y}_t, y_t) = (\hat{y}_t - y_t)^2$
- absolute loss $\ell(\hat{y}_t, y_t) = |\hat{y}_t - y_t|$
- absolute percentage of error
- $\ell(\hat{y}_t, y_t) = |\hat{y}_t - y_t| / |y_t|$
- pinball loss.

All these loss functions are convex, how can we apply our analysis for linear losses?

Prediction with expert advice with 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

$$\text{Reg}_T^{\text{expert}} \stackrel{\text{def}}{=} \sum_{t=1}^T \ell(\hat{y}_t, y_t) - \min_{1 \leq k \leq K} \sum_{t=1}^T \ell(x_t(k), y_t) \leq 2\sqrt{T \log K}$$

where $\hat{y}_t = p_t \cdot x_t$ and where $\eta > 0$ 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

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

The result is then obtained by Theorem 1.

□

Setting of an online learning problem/online convex optimization

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta$ \rightarrow full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$ \rightarrow bandit feedback.

Goal. Minimize the regret

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

Previously: we have seen an algorithm when $\ell_t(\theta) = \langle \theta, g_t \rangle$ and $\Theta = \Delta_K$. How to generalize for convex ℓ_t ?

From linear to convex losses

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

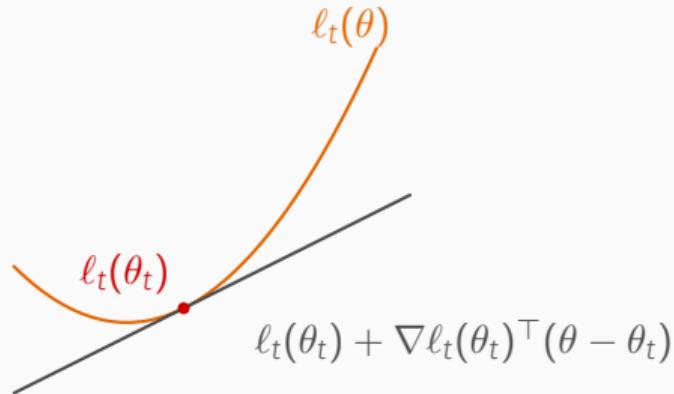
Assumptions and notations: Actions are denoted by p_t (instead of θ_t). The losses are assumed to be convex and Lipschitz

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

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 $g_t = \nabla \ell_t(p_t) \in [-G_\infty, G_\infty]^K$ as loss vectors.

The gradient trick



For $g_t = \nabla \ell_t(\theta_t)$, the linear loss $\tilde{\ell}_t(\theta) = g_t^\top \theta$ satisfies for any $\theta \in \Theta$

$$\ell_t(\theta_t) - \ell_t(\theta) \leq g_t^\top (\theta_t - \theta) \leq \tilde{\ell}_t(\theta_t) - \tilde{\ell}_t(\theta).$$

To prevent infinite regret, need finite $|\tilde{\ell}_t(\theta)|$ and hence bounds on the dual norms of the domain and gradients

$$|\tilde{\ell}_t(\theta)| \leq \|g_t\|_p \|\theta\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1.$$

The Exponentiated Gradient forecaster (EG)

Parameter: $\eta > 0$

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

For $t = 1, \dots, T$

- select p_t ; incur loss $\ell_t(p_t)$ and observe $\ell_t : \Theta \rightarrow [0, 1]$;
- compute the gradient $g_t = \nabla \ell_t(p_t) \in [-G_\infty, G_\infty]^K$
- update for all $k \in \{1, \dots, K\}$

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

Theorem 3

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

$$\text{Reg}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(p_t) - \min_{p \in \Theta} \sum_{t=1}^T \ell_t(p) \leq \eta G_\infty^2 T + \frac{\log K}{\eta}. \quad (4)$$

Therefore, for the choice $\eta = \frac{1}{G_\infty} \sqrt{\frac{\log K}{T}}$, EWA satisfies the regret bound
 $\text{Reg}_T \leq 2G_\infty \sqrt{T \log K}$.

Proof

1. Apply the regret bound of EWA with g_t (see Theorem 1 of last class):

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

Remark that the theorem also holds for loss vectors $g_t \in [-G_\infty, G_\infty]^K$ as soon as $\eta \leq 1/G_\infty$.

Upper-bounding $g_t(j)^2 \leq \|\nabla \ell_t(p_t)\|_\infty^2 \leq G_\infty^2$, substituting $g_t = \nabla \ell_t(p_t)$, this yields for all $p \in \Delta_K$

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

2. Gradient inequality: by convexity of the losses

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

which yields

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

3. Optimize η : $\eta = \frac{1}{G_\infty} \sqrt{\frac{\log K}{T}}$. □

Example: Prediction with expert advice (continued)

Setting: 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 3 with $\ell_t : p \mapsto \ell(p \cdot x_t, y_t)$.

For instance, with the absolute loss, $G_\infty = 1$ and EG satisfies:

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

Hence, on the long run we perform as good as the best convex combination of the experts which may outperform the best expert.

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

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

$$\theta_{t+1} = \text{Proj}_{\Theta} \left(\theta_t - \eta \nabla \ell_t(\theta_t) \right).$$

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

Regret bound for OGD

Online Gradient Descent

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta}(\theta_t - \eta \nabla \ell_t(\theta_t))$$

Theorem 4 (Regret of OGD)

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 $\text{Reg}_T \leq DG\sqrt{T}$.

Comparison of EG and OGD

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 3 and 12).

We saw in Theorem 3 that EG has a regret bound $\text{Reg}_T \leq 2G_\infty \sqrt{T \log K}$. 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 $R_t \leq G_\infty \sqrt{2KT}$. Thus

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

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

Online Gradient Descent

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta}(\theta_t - \eta \nabla \ell_t(\theta_t))$$

Theorem 4 (Regret of OGD)

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 $\text{Reg}_T \leq DG\sqrt{T}$.

Proof (Step 1)

Recall the update of OGD:

$$\text{OGD} : \quad \theta_{t+1} \leftarrow \text{Proj}_{\Theta} \left(\underbrace{\theta_t - \eta \nabla \ell_t(\theta_t)}_{z_t} \right)$$

1. Upper-bound the regret with gradient inequality: by convexity

$$\text{Reg}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \ell_t(\theta^*) \stackrel{\text{Convexity}}{\leqslant} \sum_{t=1}^T \langle \nabla \ell_t(\theta_t), \theta_t - \theta^* \rangle$$

Proof (Step 2)

2. Get a telescoping sum:

$$\begin{aligned}
 \|\theta_{t+1} - \theta^*\|^2 &\stackrel{\text{Projection}}{\leq} \|z_t - \theta^*\|^2 \\
 &= \|\theta_t - \eta \nabla \ell_t(\theta_t) - \theta^*\|^2 \\
 &= \|\theta_t - \theta^*\|^2 + \eta^2 \|\nabla \ell_t(\theta_t)\|^2 - 2\eta \langle \nabla \ell_t(\theta_t), \theta_t - \theta^* \rangle
 \end{aligned}$$

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta} \left(\underbrace{\theta_t - \eta \nabla \ell_t(\theta_t)}_{z_t} \right)$$

Thus,

$$\langle \nabla \ell_t(\theta_t), \theta_t - \theta^* \rangle \leq \frac{1}{2\eta} \left(\|\theta_t - \theta^*\|^2 - \|\theta_{t+1} - \theta^*\|^2 \right) + \frac{\eta}{2} \|\nabla \ell_t(\theta_t)\|^2$$

Summing over $t = 1, \dots, T$ and it telescopes

$$\begin{aligned}
 \text{Reg}_T &\leq \frac{1}{2\eta} \left(\|\theta_1 - \theta^*\|^2 - \cancel{\|\theta_{T+1} - \theta^*\|^2} \right) + \frac{\eta}{2} G^2 T \\
 &\leq \frac{D^2}{2\eta} + \frac{\eta G^2 T}{2}
 \end{aligned}$$

Exercise: 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 $\eta_t = D/(G\sqrt{t})$

Exercise: 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 $\sum_{t=2}^T \|\theta_t^* - \theta_{t-1}^*\| \leq X$

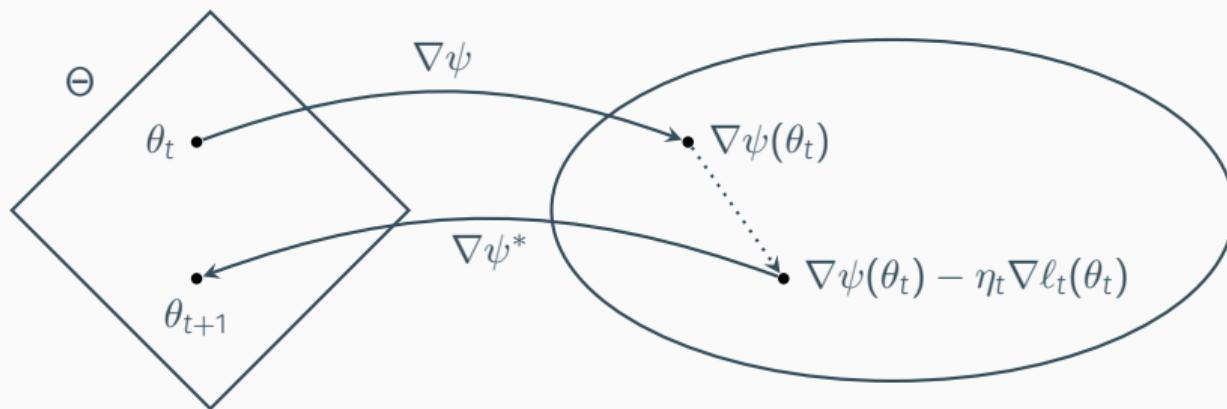
$$\sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(\theta_t^*) \leq \dots$$

Online Mirror Descent (OMD)

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.

Gradient Descent updates are performed into a dual space defined by a function $\psi : \Theta \rightarrow \mathbb{R}$ called mirror map. The goal is to transform the decision space into one more suited to Euclidean geometry.



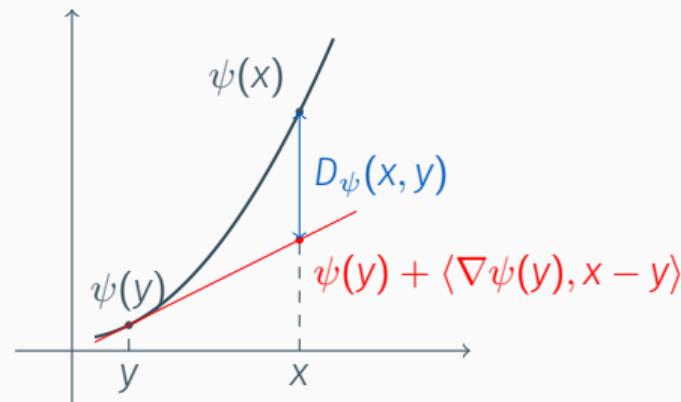
Bregman divergence

Definition (Bregman divergence)

For any continuously differentiable convex function ψ , the Bregman divergence with respect to ψ is defined as

$$D_\psi(x, y) \leq \psi(x) - \psi(y) - \nabla\psi(y) \cdot (x - y) \quad \forall x, y \in \Theta.$$

It is the difference between the value of the regularization function at x and the value of its first order Taylor approximation.



Online Mirror Descent (OMD) – Lazy and agile versions

Online Mirror Descent (OMD)

Parameters: $\eta > 0$, mirror map ψ

Initialize: $\theta_1 \in \Theta$, $z_1 = x_1$

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

$$\nabla \psi(z_{t+1}) = \nabla \psi(\theta_t) - \eta \nabla \ell_t(\theta_t) \quad \leftarrow \text{agile version}$$

$$\nabla \psi(z_{t+1}) = \nabla \psi(z_t) - \eta \nabla \ell_t(\theta_t) \quad \leftarrow \text{lazy version}$$

- project according to the Bregman divergence

$$\theta_{t+1} \in \arg \min_{\theta \in \Theta} D_\psi(\theta, z_{t+1}).$$

Main difference: the **lazy version** performs all gradient steps in the dual space before computing the projection over Θ to form predictions θ_t , while the **agile version** alternates projections and gradient steps in the dual space.

Example: OMD (agile version) with Euclidean regularization = OGD

$$\text{OGD} : \quad \theta_{t+1} \leftarrow \text{Proj}_{\Theta} (\theta_t - \eta \nabla \ell_t(\theta_t))$$

$$\begin{aligned}\text{OMD} : \quad & \nabla \psi(z_{t+1}) = \nabla \psi(\theta_t) - \eta \nabla \ell_t(\theta_t) \\ & \theta_{t+1} \in \arg \min_{\theta \in \Theta} D_\psi(\theta, z_{t+1})\end{aligned}$$

If $\Theta \subset \mathbb{R}^d$, we can choose $\psi(x) = \frac{1}{2} \|x\|^2$.

Then

$$\nabla \psi(x) = x \quad \text{and} \quad D_\psi(x, y) = \frac{1}{2} \|x - y\|^2.$$

Therefore, the update of OMD becomes $z_{t+1} = \theta_t - \eta \nabla \ell_t(\theta_t)$ and $\theta_{t+1} = \text{Proj}_{\Theta}(z_{t+1})$.

We recover the online gradient descent algorithm.

OMD (agile or lazy) with negative entropy = EG

$$\text{EG : } \begin{aligned} g_t &= \nabla \ell_t(\theta_t) \\ \theta_{t+1}(k) &= \frac{\theta_t(k)e^{-\eta g_t(k)}}{\sum_{j=1}^K \theta_t(j)e^{-\eta g_t(j)}} \end{aligned}$$

$$\text{OMD : } \begin{aligned} \nabla \psi(z_{t+1}) &= \nabla \psi(\theta_t) - \eta \nabla \ell_t(\theta_t) \\ \theta_{t+1} &\in \arg \min_{\theta \in \Theta} D_\psi(\theta, z_{t+1}) \end{aligned}$$

If $\Theta = \Delta_K$. We consider the negative entropy $\psi(x) = \langle x, \log x \rangle$ so that $\nabla \psi(x) = 1 + \log x$.

The update of OMD is then

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

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

$$z_{t+1}(i) = \theta_t(i) e^{-\eta g_t(i)}.$$

The projection to the simplex is a simple renormalization (exercise next), we thus recover EG.

Exercice: Bregmann projection over the simplex is the renormalization

Exercise

Let $z \in \mathbb{R}_+^d$, $\psi(x) = \langle x, \log x \rangle$ and $\theta_* = \arg \min_{\theta \in \Delta_d} D_\psi(\theta, z)$. Show that

1. For any $x \in \Delta_d$, $D_\psi(x, z) = \langle x, \log \frac{x}{z} \rangle + \|z\|_1 - 1 \geq -\log(\|z\|_1) + \|z\|_1 - 1$
2. For all $k \in [d]$, $\theta_* = z / \|z\|_1$.

Exercice: Bregmann projection over the simplex is the renormalization

Exercise

Let $z \in \mathbb{R}_+^d$, $\psi(x) = \langle x, \log x \rangle$ and $\theta_* = \arg \min_{\theta \in \Delta_d} D_\psi(\theta, z)$. Show that

1. For any $x \in \Delta_d$, $D_\psi(x, z) = \langle x, \log \frac{x}{z} \rangle + \|z\|_1 - 1 \geq -\log(\|z\|_1) + \|z\|_1 - 1$
2. For all $k \in [d]$, $\theta_* = z/\|z\|_1$.

Solution:

1.

$$\begin{aligned} D_\psi(x, z) &= \psi(x) - \psi(z) - \langle \nabla \psi(z), x - z \rangle = \langle x, \log x \rangle - \langle z, \log z \rangle - \langle 1 + \log z, x - z \rangle \\ &= \langle x, \log \frac{x}{z} \rangle + \|z\|_1 - 1 \stackrel{\text{Jensen}}{\geq} -\log \left(\left\langle x, \frac{z}{x} \right\rangle \right) + 1 - \|z\|_1 = -\log(\|z\|_1) + \|z\|_1 - 1 \end{aligned}$$

Exercice: Bregmann projection over the simplex is the renormalization

Exercise

Let $z \in \mathbb{R}_+^d$, $\psi(x) = \langle x, \log x \rangle$ and $\theta_* = \arg \min_{\theta \in \Delta_d} D_\psi(\theta, z)$. Show that

1. For any $x \in \Delta_d$, $D_\psi(x, z) = \langle x, \log \frac{x}{z} \rangle + \|z\|_1 - 1 \geq -\log(\|z\|_1) + \|z\|_1 - 1$
2. For all $k \in [d]$, $\theta_* = z/\|z\|_1$.

Solution:

1.

$$\begin{aligned} D_\psi(x, z) &= \psi(x) - \psi(z) - \langle \nabla \psi(z), x - z \rangle = \langle x, \log x \rangle - \langle z, \log z \rangle - \langle 1 + \log z, x - z \rangle \\ &= \langle x, \log \frac{x}{z} \rangle + \|z\|_1 - 1 \stackrel{\text{Jensen}}{\geq} -\log \left(\left\langle x, \frac{z}{x} \right\rangle \right) + 1 - \|z\|_1 = -\log(\|z\|_1) + \|z\|_1 - 1 \end{aligned}$$

2.

$$D_\psi\left(\frac{z}{\|z\|_1}, z\right) = \left\langle \frac{z}{\|z\|_1}, -\log(\|z\|_1) \right\rangle + \|z\|_1 - 1 = -\log(\|z\|_1) + \|z\|_1 - 1$$

Equivalent Formulation of OMD (agile version)

OMD (agile version) is equivalent to the following update

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_t) \right\}$$

Indeed,

$$\begin{aligned} & \arg \min_{\theta \in \Theta} \left\{ \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_t) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ \eta \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \psi(\theta) - \psi(\theta_t) - \langle \nabla \psi(\theta_t), \theta - \theta_t \rangle \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ \langle \eta \nabla \ell_t(\theta_t) - \nabla \psi(\theta_t), \theta \rangle + \psi(x) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ \langle -\nabla \psi(z_{t+1}), \theta \rangle + \psi(x) \right\} \quad \text{since } \nabla \psi(z_{t+1}) = \nabla \psi(\theta_t) - \eta \nabla \ell_t(\theta_t) \\ &= \arg \min_{\theta \in \Theta} \left\{ -\psi(z_{t+1}) - \langle \nabla \psi(z_{t+1}), \theta - z_{t+1} \rangle + \psi(x) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ D_\psi(\theta, z_{t+1}) \right\} \end{aligned}$$

Equivalent formulation of OMD (lazy version) – linearized FTRL

OMD (lazy version) is the linearized version of Follow The Regularized Leader (FTRL)

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_1) \right\}$$

Indeed, for the lazy version $\nabla \psi(z_{t+1}) = \nabla \psi(z_t) - \eta \nabla \ell_t(\theta_t) = \nabla \psi(z_1) - \sum_{s=1}^t \eta \nabla \ell_s(\theta_s)$, and $z_1 = x_1$. Thus

$$\begin{aligned} & \arg \min_{\theta \in \Theta} \left\{ \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_1) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ \eta \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \psi(\theta) - \cancel{\psi(\theta_1)} - \langle \nabla \psi(\theta_1), \theta - \theta_1 \rangle \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ \left\langle \sum_{s=1}^t \eta \nabla \ell_s(\theta_s) - \nabla \psi(\theta_1), \theta \right\rangle + \psi(x) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ \langle -\nabla \psi(z_{t+1}), \theta \rangle + \psi(x) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ D_\psi(\theta, z_{t+1}) \right\} \end{aligned}$$

Regret of OMD (agile or lazy version)

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \begin{array}{ll} \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_t) & \text{(agile)} \\ \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_1) & \text{(lazy)} \end{array} \right\}$$

Theorem 5 (Regret of OMD)

Let $T \geq 1$. Let Θ be a compact and convex, $\theta_1 \in \Theta$, and $\psi : \Theta \rightarrow \mathbb{R}$ be μ -strongly convex mirror map with respect to some norm $\|\cdot\|$. Assume

$$D_\psi(\theta, \theta_1)^{1/2} \leq D \quad \text{and} \quad \|\nabla \ell_t(\theta_t)\|_* \leq G_*.$$

Then, for any convex $\ell_1, \dots, \ell_T : \Theta \rightarrow \mathbb{R}$, the regret of OMD (agile or lazy version) with $\eta = DG_*^{-1} \sqrt{2\mu/T}$ is upper bounded as

$$\text{Reg}_T \leq DG_* \sqrt{\frac{2T}{\mu}}.$$

We exactly the regrets of OGD and EG since $\psi : x \mapsto \frac{1}{2}\|x\|_2^2$ and $\psi : x \mapsto \langle x, \log x \rangle$ are resp. 1-strongly convex with $\|\cdot\|_2$ and $\|\cdot\|_1$ norms.

Proof (agile version)

Let $t \geq 1$ and $\theta \in \Theta$. Denote $\Phi_t(\theta) = \langle \nabla \ell_t(\theta_t), \theta \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_t)$. From the optimality condition on θ_{t+1} , $\langle \nabla \Phi_t(\theta_{t+1}), \theta_{t+1} - \theta \rangle \leq 0$, which entails $\langle \nabla \ell_t(\theta_t), \theta_{t+1} - \theta \rangle \leq \frac{1}{\eta} \langle \nabla \psi(\theta_{t+1}) - \nabla \psi(\theta_t), \theta - \theta_{t+1} \rangle$ and yields

$$\begin{aligned}\ell_t(\theta_t) - \ell_t(\theta) &\leq \langle \nabla \ell_t(\theta_t), \theta_t - \theta \rangle = \langle \nabla \ell_t(\theta_t), \theta_t - \theta_{t+1} \rangle + \langle \nabla \ell_t(\theta_t), \theta_{t+1} - \theta \rangle \\ &\leq \langle \nabla \ell_t(\theta_t), \theta_t - \theta_{t+1} \rangle + \frac{1}{\eta} \langle \nabla \psi(\theta_{t+1}) - \nabla \psi(\theta_t), \theta - \theta_{t+1} \rangle \\ &= \langle \nabla \ell_t(\theta_t), \theta_t - \theta_{t+1} \rangle + \frac{1}{\eta} (D_\psi(\theta, \theta_t) - D_\psi(\theta, \theta_{t+1}) - D_\psi(\theta_{t+1}, \theta_t)) \\ &\leq \langle \nabla \ell_t(\theta_t), \theta_t - \theta_{t+1} \rangle + \frac{1}{\eta} (D_\psi(\theta, \theta_t) - D_\psi(\theta, \theta_{t+1}) - \frac{\mu}{2} \|\theta_t - \theta_{t+1}\|^2),\end{aligned}$$

where the last inequality because ψ is μ -strongly convex. Moreover,

$$\langle \nabla \ell_t(\theta_t), \theta_t - \theta_{t+1} \rangle \stackrel{\text{Cauchy-Schwarz}}{\leq} \|\theta_t - \theta_{t+1}\| \|\nabla \ell_t(\theta_t)\|_* \stackrel{\text{Young's}}{\leq} \frac{\mu}{2\eta} \|\theta_t - \theta_{t+1}\|^2 + \frac{\eta}{2\mu} \|\nabla \ell_t(\theta_t)\|_*^2.$$

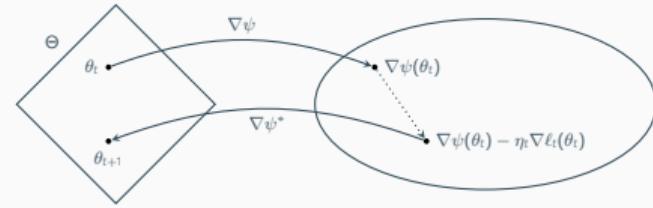
Combining the last two upper-bounds and summing over $t = 1, \dots, T$ entails the stated regret upper bound.

Proof (lazy version)

Left as exercise

Summary on OMD (lazy or agile)

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \begin{array}{ll} \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_t) & \text{(agile)} \\ \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \frac{1}{\eta} D_\psi(\theta, \theta_1) & \text{(lazy)} \end{array} \right.$$



OMD provide efficient algorithms with generic $O(\sqrt{T})$ regret bounds but suffer from some drawbacks

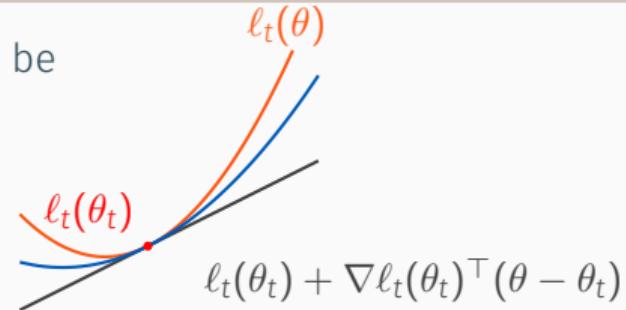
- $O(\sqrt{T})$ is sub-optimal for some losses with curvature (e.g., squared, logistic, or log loss) for which logarithmic regret may be achieved → ONS
- needs bounded gradients which is not the case in some settings (e.g., portfolio selection).
- the best mirror map ψ should be specified by the user → AdaGrad.

Curvature

When losses have curvature, logarithmic regret may be achieved.

- Convexity:

$$\ell_t(\theta) \geq \ell_t(\theta_t) + \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle$$



- Strong convexity:

$$\ell_t(\theta) \geq \ell_t(\theta_t) + \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{\gamma}{2} \|\theta - \theta_t\|^2$$

- Exp-concavity: $\theta \mapsto \exp(-\eta \ell_t(\theta))$ is concave, which implies (exercise)

$$\ell_t(\theta) \geq \ell_t(\theta_t) + \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{\gamma}{2} \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle^2$$

for $\gamma = \frac{1}{2} \min \left\{ \eta, \frac{1}{4G_* D} \right\}$.

Exp-concavity

Definition (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.

Properties:

- Exp-concavity \Rightarrow convexity because $-\log$ is convex and decreasing.
- Strong convexity + bounded domain and gradients \Rightarrow exp-concavity.
- η -exp-concavity \Rightarrow η' -exp-concavity for $0 \leq \eta' \leq \eta$.

Many losses are exp-concave:

- strongly convex losses on bounded domain
 - logistic loss
 - relative entropy
- squared loss is $\frac{1}{2Y^2}$ -expconcave on $[0, Y]$

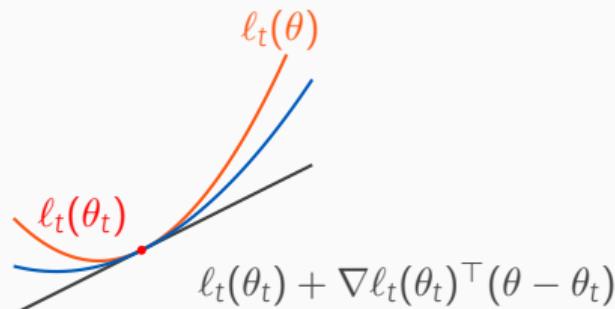
Exercise: Prove the above the above facts.

Exercise: Exp-concavity implies a quadratic lower-bound

Show that if the loss ℓ_t is η -exp-concave with gradients bounded by G and a domain diameter bounded by D , then it can be lower-bounded by a quadratic approximation: for all $\theta, \theta_t \in \Theta$,

$$\ell_t(\theta) \geq \ell_t(\theta_t) + \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle + \frac{\gamma}{2} \langle \nabla \ell_t(\theta_t), \theta - \theta_t \rangle^2$$

for $\gamma = \frac{1}{2} \min \left\{ \eta, \frac{1}{4G_* D} \right\}$.



Constant regret for exp-concave loss functions

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

$$\text{Reg}_T^{\text{expert}} \stackrel{\text{def}}{=} \sum_{t=1}^T \ell(\hat{y}_t, y_t) - \min_{1 \leq k \leq K} \sum_{t=1}^T \ell(x_t(k), y_t) \leq \frac{\log K}{\eta},$$

for all $T \geq 1$.

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

Proof (Step 1)

We define $W_t(i) = e^{-\eta \sum_{s=1}^t g_s(i)}$ and $W_t = \sum_{i=1}^N W_t(i)$. We have

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

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

$$W_T \leq K \exp \left(-\eta \sum_{t=1}^T \ell(\hat{y}_t, y_t) \right). \quad (5)$$

Proof (Step 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 (5) and taking the log concludes the proof.

Continuous EWA

Can we obtain a regret with respect to the best combination of experts

$$\min_p \sum_{t=1}^T \ell_t(p)$$

instead of the regret with respect to the best fixed expert?

Continuous EWA

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

where μ is the uniform (Lebesgue) measure on $\Theta = \Delta_K$.

Regret bound for continuous EWA

Theorem 6 (Regret of continuous EWA)

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

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

Nice theoretical result but hard to implement because of the integral.

In practice, p_t can be computed by using $(1/T)$ -discretization grid of Θ (bad complexity of order $T^K!$) or by using Monte-Carlo methods to approximate the integral.

Proof (Step 1 – Upper-bound of W_T)

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

$$\begin{aligned}
 W_T &= \int_{\Theta} e^{-\eta \sum_{t=1}^T \ell_t(p)} d\mu(p) \\
 &= W_{T-1} \int_{\Theta} \frac{W_{T-1}(p)}{W_{T-1}} e^{-\eta \ell_T(p)} d\mu(p) \\
 &= W_{T-1} \int_{\Theta} e^{-\eta \ell_T(p)} d\hat{\mu}_{T-1}(p) && \leftarrow p_T = \int_{\Theta} p d\hat{\mu}_{T-1}(p) \\
 &\leq W_{T-1} \exp(-\eta \ell_T(p_T)) && \leftarrow \eta\text{-exp-concavity} \\
 &\leq \exp\left(-\eta \sum_{t=1}^T \ell_t(p_t)\right), && \leftarrow \text{induction}
 \end{aligned} \tag{6}$$

Proof (Step 2 – Lower bound of W_T)

For simplicity, we assume that ℓ_t are continuous. Therefore the infimum is a minimum and let $p^* \in \arg \min_{p \in \Theta} \sum_{t=1}^T \ell_t(p)$ and define

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

By exp-concavity 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 $p \in \Theta_\varepsilon$

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

Integrating both parts over Θ_ε and using $\mu(\Theta_\varepsilon) = \varepsilon^{K-1} \mu(\Theta)$ (exercise) we get

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

Proof (Step 3 – Conclusion)

Combining with (6), using $W_0 = \mu(\Theta)$, taking the log and reorganizing the terms yields

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

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

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

Logarithmic regret for OGD under strong-convexity

Online Gradient Descent:

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta} (\theta_t - \eta_t \nabla \ell_t(\theta_t))$$

Theorem 7 (Regret of OGD under strong-convexity)

Let $D, G, \gamma > 0$. Assume that $\max_{\theta, \theta' \in \Theta} \|\theta - \theta'\| \leq D$ and. Then for any sequence ℓ_1, \dots, ℓ_T of γ -strongly convex differentiable loss functions such that $\max_{\theta \in \Theta} \|\nabla \ell_t(\theta)\| \leq G$, the regret of OGD with $\eta_t = 1/(\gamma t)$ satisfies

$$\text{Reg}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \leq \frac{G^2}{2\gamma} (1 + \log T).$$

Proof

1. Upper-bound the regret with strong convexity:

$$\text{Reg}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t) - \ell_t(\theta^*) \stackrel{\text{Strong Convexity}}{\leqslant} \sum_{t=1}^T \langle \nabla \ell_t(\theta_t), \theta_t - \theta^* \rangle - \frac{\gamma}{2} \|\theta_t - \theta^*\|^2$$

2. Upper-bound the gradient term as for OGD analysis

$$\langle \nabla \ell_t(\theta_t), \theta_t - \theta^* \rangle \leqslant \frac{1}{2\eta_t} \left(\|\theta_t - \theta^*\|^2 - \|\theta_{t+1} - \theta^*\|^2 \right) + \frac{\eta_t}{2} \|\nabla \ell_t(\theta_t)\|^2$$

3. Substitute in the previous inequality and conclude

$$\begin{aligned} \text{Reg}_T &\leqslant \sum_{t=1}^T \frac{1}{2\eta_t} \left(\|\theta_t - \theta^*\|^2 - \|\theta_{t+1} - \theta^*\|^2 \right) + \frac{\eta_t G^2}{2} - \frac{\gamma}{2} \|\theta_t - \theta^*\|^2 \\ &= \frac{1}{2} \sum_{t=1}^T \left(\frac{1}{\eta_t} - \cancel{\frac{1}{\eta_{t-1}}} - \gamma \right) \|\theta_t - \theta^*\|^2 + \frac{G^2}{2} \sum_{t=1}^T \frac{1}{\gamma t} \\ &\leqslant \frac{G^2}{\gamma} (1 + \log T) \end{aligned}$$

Online Newton Step (ONS)

Idea: Include the quadratic upper-bound of the loss into FTRL:

$$\theta_t = \arg \min_{\theta \in \Theta} \left\{ \sum_{s=1}^{t-1} \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \frac{\gamma}{2} \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle^2 + \frac{\lambda}{2} \|\theta\|^2 \right\}.$$

Theorem 8 (Regret of ONS)

Let $T \geq 1$. Let $\Theta \subset \mathbb{R}^d$ be a centered convex set with diameter $D > 0$. Then, for any G -Lipschitz, η -exp-concave losses $\ell_1, \dots, \ell_T : \Theta \rightarrow \mathbb{R}$, the regret of ONS with $\lambda > 0$ and $\gamma \leq \frac{1}{2} \min \left\{ \eta, \frac{1}{4GD} \right\}$ satisfies

$$\text{Reg}_T \leq \frac{\lambda}{2} D^2 + \frac{d}{2\gamma} \log \left(1 + \frac{\gamma T G^2}{d\lambda} \right).$$

If the losses are η -exp-concave \rightarrow logarithmic regret.

Complexity: Ignoring the projection step (may be difficult), it requires the inversion of a $d \times d$ matrices and could be done in $O(d^3)$ operations per round.

Proof (Step 1)

For any $t \geq 1$, define the matrix $A_t = \gamma \sum_{s=1}^t \nabla \ell_t(\theta_t) \nabla \ell_t(\theta_t)^\top + \lambda I$ and set for any $\theta \in \Theta$

$$\psi_t(\theta) = \frac{\gamma}{2} \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle^2 + \frac{\lambda}{2} \|\theta\|^2 \quad \text{and} \quad \Phi_t(\theta) = \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \psi_t(\theta).$$

Then, ψ_t and Φ_t are 1-strongly convex with respect to the norm $\|\cdot\|_{A_t}$.

The ONS update is then

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \sum_{s=1}^t \langle \nabla \ell_s(\theta_s), \theta - \theta_s \rangle + \psi_t(\theta) \right\} = \arg \min_{\theta \in \Theta} \Phi_t(\theta).$$

On the one hand, by Cauchy-Schwarz inequality

$$\sum_{t=1}^T \langle \nabla \ell_t(\theta_t), \theta_t - \theta_{t+1} \rangle \leq \frac{1}{2} \sum_{t=1}^T \|\nabla \ell_t(\theta_t)\|_{A_t^{-1}}^2 + \frac{1}{2} \sum_{t=1}^T \|\theta_t - \theta_{t+1}\|_{A_t}^2. \quad (7)$$

Proof (Step 2)

On the other hand, fixing $\theta \in \Theta$, by definition of $\Phi_t(x)$

$$\langle \nabla \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle = \Phi_t(\theta_{t+1}) - \Phi_{t-1}(\theta_{t+1}) - \frac{\gamma}{2} \langle \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle^2 \quad (8)$$

Thus

$$\begin{aligned} \sum_{t=1}^T \langle \nabla \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle &= -\Phi_T(\theta) + \psi_T(\theta) + \sum_{t=1}^T \langle \nabla \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle \quad \leftarrow \text{by Definition of } \Phi_T(\theta) \\ &= -\Phi_T(\theta) + \psi_T(\theta) + \sum_{t=1}^T \Phi_t(\theta_{t+1}) - \Phi_{t-1}(\theta_{t+1}) - \frac{\gamma}{2} \langle \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle^2 \quad \leftarrow \text{by (8)} \\ &= \cancel{\Phi_T(\theta_{T+1}) - \Phi_T(\theta) - \Phi_0(\theta_1)} + \psi_T(\theta) + \sum_{t=1}^T \Phi_{t-1}(\theta_t) - \Phi_{t-1}(\theta_{t+1}) - \frac{\gamma}{2} \langle \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle^2 \\ &\leq \psi_T(\theta) - \frac{1}{2} \sum_{t=1}^T \|\theta_{t+1} - \theta_t\|_{A_{t-1}}^2 - \frac{\gamma}{2} \langle \ell_t(\theta_t), \theta_{t+1} - \theta_t \rangle^2 = \psi_T(\theta) - \frac{1}{2} \sum_{t=1}^T \|\theta_{t+1} - \theta_t\|_{A_t}^2. \end{aligned} \quad (9)$$

where the last inequality is because $\Phi_T(\theta_{T+1}) \leq \Phi_T(\theta)$ and by strong convexity of Φ_{t-1} and by the optimality condition of $\theta_t = \arg \min_{\theta \in \Theta} \Phi_{t-1}(\theta)$ we have $\Phi_{t-1}(\theta_t) - \Phi_{t-1}(\theta) \leq -\frac{1}{2} \|\theta_t - \theta\|_{A_{t-1}}^2$.

Proof (Step 3)

Combining (7) and (9) into the quadratic upper-bound yields

$$\sum_{t=1}^T \ell_t(\theta_t) - \ell_t(\theta) \stackrel{\text{Exp-concavity}}{\leqslant} \frac{\lambda}{2} \|\theta\|^2 - \psi_T(\theta) + \sum_{t=1}^T \langle \nabla \ell_t(\theta_t), \theta_t - \theta \rangle \leqslant \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{2} \sum_{t=1}^T \|\nabla \ell_t(\theta_t)\|_{A_t^{-1}}^2. \quad (10)$$

Lemma 1 (Lemma 11.11 of Cesa-Bianchi and Lugosi 2006)

For any full rank matrix A and any vector $x \in \mathbb{R}^d$, then

$$x^\top (A + xx^\top)^{-1} x = 1 - \det(A)/\det(A + xx^\top)$$

Thus

$$\sum_{t=1}^T \left(1 - \frac{\det(A_{t-1})}{\det(A_t)}\right)^{1-u} \stackrel{1-u \leq -\log u}{\leqslant} \sum_{t=1}^T \log \frac{\det(A_t)}{\det(A_{t-1})} = \log \frac{\det(A_{T+1})}{\det(A_0)} \leq d \log \left(1 + \frac{\gamma T G^2}{d\lambda}\right),$$

where the last line used that $\|\nabla \ell_t(\theta_t)\| \leq G$. Substituting into (10) concludes the proof.

Setting of an online learning problem/online convex optimization

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta$ \rightarrow full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$ \rightarrow bandit feedback.

Goal: Minimize the regret

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

Previous results for full-information feedback

The Exponentially Weighted Average (EWA) forecaster

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

achieves a cumulative regret $\text{Reg}_T \lesssim \sqrt{T \log K}$ when the set of actions is the K -dimensional simplex and for linear losses $\ell_t(p) = p^\top g_t$ with $g_t \in [-1, 1]^K$.

In particular, we saw the intermediate regret-bound if $-\eta g_t(k) \leq 1$

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

Note that the loss vectors g_t may depend on past information $p_1, g_1, \dots, g_{t-1}, p_t$.

Adversarial multi-armed bandit and pseudo-regret

Setting: $\Theta = \{1, \dots, K\}$. At round t , the player chooses an action $k_t \in \{1, \dots, K\}$ and suffers and observes the loss $\ell_t(k_t) \in [0, 1]$ only.

Regret with respect to action $k \in [K]$ by

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

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

Definition (Pseudo-regret)

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

Oblivious vs adaptive adversary

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

The expectation is 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: $\bar{\text{Reg}}_T = \mathbb{E}[\text{Reg}_T]$.
- **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 $\max_{k \in [K]} \mathbb{E}[\text{Reg}_T(k)] \leq \mathbb{E}[\max_{k \in [K]} \text{Reg}_T(k)]$ and thus $\bar{\text{Reg}}_T \leq \mathbb{E}[\text{Reg}_T]$.

How to use EWA for bandits?

The Exponentially Weighted Average (EWA) forecaster

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

(EWA)

Question: Can we use directly $p_t(k)$ as defined by EWA with $g_t = (\ell_t(1), \dots, \ell_t(K))$ and sample $k_t \sim p_t$ as we did for random EWA?

Yes No

How to use EWA for bandits?

The Exponentially Weighted Average (EWA) forecaster

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

Question: Can we use directly $p_t(k)$ as defined by EWA with $g_t = (\ell_t(1), \dots, \ell_t(K))$ and sample $k_t \sim p_t$ as we did for random EWA?

Answer: No, since the player does not observe $\ell_t(k)$ for $k \neq k_t$ and cannot compute p_t .

How to use EWA for bandits?

The Exponentially Weighted Average (EWA) forecaster

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

Question: What about setting using $\ell_t(k)$ if we observe it and 0 otherwise:

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

Yes No

How to use EWA for bandits?

The Exponentially Weighted Average (EWA) forecaster

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

Question: What about setting using $\ell_t(k)$ if we observe it and 0 otherwise:

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

Answer: No, because this estimate would be 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.

How to use EWA for bandits?

The Exponentially Weighted Average (EWA) forecaster

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

Therefore, we choose

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

which leads to the algorithm EXP3 detailed below.

Exponential Weights for bandits

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 $k \in \{1, \dots, K\}$

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

$$p_{t+1}(k) = \frac{e^{-\eta \sum_{s=1}^t g_s(k)}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^t g_s(j)}}, \quad \text{where } g_s(k) = \frac{\ell_s(k)}{p_s(k)} \mathbb{1}\{k = k_s\} \quad (\text{EXP3})$$

Theorem 9

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

$$\bar{\text{Reg}}_T \stackrel{\text{def}}{=} \max_{k \in [K]} \mathbb{E} \left[\sum_{t=1}^T \ell_t(k_t) - \sum_{t=1}^T \ell_t(k) \right] \leq 2\sqrt{KT \log K}.$$

Proof

Applying EWA to the estimated losses $g_t(j)$ that are completely observed and taking the expectation:

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

The rest of the proof consists in computing the expectations:

$$\mathbb{E}[p_t \cdot g_t] = \mathbb{E}[\ell_t(k_t)], \quad \mathbb{E}[g_t(j)] = \mathbb{E}[\ell_t(j)] \quad \text{and} \quad \mathbb{E}[p_t \cdot g_t^2] \leq K \quad (11)$$

Proof

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.

1) Proof that $\mathbb{E}[g_t(j)] = \mathbb{E}[\ell_t(j)]$

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

2) Proof that $\mathbb{E}[p_t \cdot g_t] = \mathbb{E}[\ell_t(k_t)]$

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

Proof

Therefore, using

$$\mathbb{E}[p_t \cdot g_t] = \mathbb{E}[\ell_t(k_t)] \quad \text{and} \quad \mathbb{E}[g_t(j)] = \mathbb{E}[\ell_t(j)] \quad (12)$$

we have

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

Proof

3) Proof that $\mathbb{E}[p_t \cdot g_t^2] \leq K$

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

4) Conclusion. Substituting into Inequality $(*)$ yields

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

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

Limit of the result

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

$$\mathbb{E} \left[\min_j \sum_{t=1}^T g_t(j) \right] \not\leq \mathbb{E} \left[\min_j \sum_{t=1}^T \ell_t(j) \right]. \quad (13)$$

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

EXP3.P

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

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

For $t = 1, \dots, T$

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

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

$$\text{where } g_s(k) = \frac{r_s(k) \mathbb{1}\{k=k_s\} + \beta}{p_s(k)}.$$

The weights $p_t(k)$ of EXP3.P are necessarily larger than γ/K and thus $|\eta g_t(j)| \leq 1$ as soon as $\eta(1 + \beta)K/\gamma \leq 1$.

Regret bound for Exp3.P

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

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

with probability at least $1 - \delta$.

With the choice $\delta = 1/T$ it yields

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

Setting of adversarial bandits with experts

Setting

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

- N experts propose recommendations $h_t(i) \in [K]$ for $i \in [N]$
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player chooses an action $k_t \in [K]$
- the player suffers loss $\ell_t(k_t)$
- the player observes the loss of the chosen action only: $\ell_t(k_t)$

Goal: compete with the best expert, i.e., minimize

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

with respect to the experts.

By using EXP3 on the set of experts instead of the set of actions, we would get

$$\bar{\text{Reg}}_T \leq \sqrt{T N \log N}.$$

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

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

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 $\text{Reg}_T^{\text{exp}} \leq 2\sqrt{TK \log N}$.

Proof left as exercise.

Beyond finite set of actions?

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta$ \rightarrow full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$ \rightarrow bandit feedback.

The goal of the player is to minimize his cumulative loss:

$$\widehat{L}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t).$$

This lecture: we saw variants of EXP3 when Θ is finite.

What if the losses ℓ_t are convex but Θ is any bounded convex set in \mathbb{R}^d ?

In the full information setting (when gradient can be observed), we saw OGD algorithm:

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta} (\theta_t - \eta \nabla \ell_t(\theta_t))$$

Theorem 12 (Regret of OGD)

Let $D, G, \eta > 0$. Assume that Θ has diameter bounded by D and the convex losses have sub-Gradients bounded by G in ℓ_2 -norm, the regret of OGD satisfies

$$\sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \leq DG\sqrt{T}.$$

How to adapt this algorithm to the bandit setting? That is, when only $\ell_t(\theta_t)$ are observed and not $\nabla \ell_t(\theta_t)$?

Point-wise gradient estimators

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta} (\theta_t - \eta \nabla \ell_t(\theta_t))$$

Similarly to EXP3, the idea is to replace the gradient in OGD with **unbiased estimators**. That is try to find an observable random variable \hat{g}_t that satisfies

$$\mathbb{E}[\hat{g}_t] \approx \nabla \ell_t(\theta_t)$$

Example: one-dimensional gradient estimate

$$\ell'(x) = \lim_{\delta \rightarrow 0} \frac{\ell(x + \delta) - \ell(x - \delta)}{2\delta}.$$

Thus we can define

$$\hat{g}(x) = \begin{cases} \frac{\ell(x+\delta)}{\delta} & \text{with proba } \frac{1}{2} \\ -\frac{\ell(x-\delta)}{\delta} & \text{with proba } \frac{1}{2} \end{cases} \quad \text{which yields} \quad \mathbb{E}[\hat{g}(x)] = \frac{\ell(x + \delta) - \ell(x - \delta)}{2\delta}.$$

Thus **in expectation**, for small δ , $\hat{g}(x)$ approximates $\ell'(x)$.

Point-wise gradient estimators: multi-dimensional case

We show here how the one-dimensional pointwise gradient estimator can be extended to the multi-dimensional case.

We define $\hat{\ell}_t$ to be a smoothed version of the loss:

$$\hat{\ell}_t(\theta) = \mathbb{E}_v [\ell_t(\theta + \delta v)]$$

where $v \sim \text{Unif}(\mathbb{B})$. If δ is small, $\hat{\ell}_t$ is a good approximation of ℓ_t .

Lemma 2

Let $\hat{\ell}_t(\theta) = \mathbb{E}[\ell_t(\theta + \delta v)]$ where $v \sim \text{Unif}(\mathbb{B})$ be a smoothed version of the loss, then

$$\mathbb{E}_u \left[\frac{d}{\delta} \ell_t(\theta + \delta u) u \right] = \nabla \hat{\ell}_t(\theta).$$

Proof.

Left as exercise. See Lem. 6.7, Hazan et al. 2016. □

OGD without Gradients

Similarly to EXP3, the idea is to replace the gradient in OGD with **unbiased estimators**.

OGD without gradients

For $t = 1, \dots, T$

- Draw $u_t \in \mathbb{S}$ uniformly at random in the unit sphere
- Set $\hat{\theta}_t = \theta_t + \delta u_t$ a random perturbation of the current point θ_t
- Play $\hat{\theta}_t$
- Estimate the gradient in θ_t with

$$\hat{g}_t = \frac{d}{\delta} \ell_t(\hat{\theta}_t) u_t$$

- Update

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta_\delta} (\theta_t - \eta \hat{g}_t)$$

where $\Theta_\delta = \{\theta \in \Theta \text{ s.t. } \theta + \delta u \in \Theta \text{ } \forall u \in \mathbb{S}\}$

Regret of OGD without gradients

OGD without gradients:

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta_\delta} (\theta_t - \eta \hat{g}_t) \quad \text{where } \hat{g}_t = \frac{d}{\eta} \ell_t(\hat{\theta}_t) u_t \text{ and } \hat{\theta}_t = \theta_t + \delta u_t$$

Theorem 13

If the losses are in $[-1, 1]$ and G -Lipschitz, OGD without gradients with parameters $\delta = \min\{D, (1/2)\sqrt{Dd/G}T^{-1/4}\}$ and $\eta = D\delta/(dT^{1/2})$ satisfies the expected regret bound

$$\sum_{t=1}^T \mathbb{E}[\ell_t(\hat{\theta}_t)] - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \leq 2d\sqrt{T} + 2\sqrt{GDd}T^{3/4}.$$

Proof (Step 1)

Denote

$$\theta^* \in \arg \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \quad \text{and} \quad \theta_\delta^* = \text{Proj}_{\Theta_\delta}(\theta^*) .$$

Then,

$$\|\theta^* - \theta_\delta^*\| \leq \delta$$

Thus, if the losses are G -Lipschitz

$$\begin{aligned} \text{Reg}_T &:= \sum_{t=1}^T \mathbb{E}[\ell_t(\hat{\theta}_t)] - \sum_{t=1}^T \ell_t(\theta^*) \leq \sum_{t=1}^T \mathbb{E}[\ell_t(\hat{\theta}_t)] - \sum_{t=1}^T \ell_t(\theta_\delta^*) \\ &\leq \sum_{t=1}^T \mathbb{E}[\ell_t(\theta_t)] - \sum_{t=1}^T \ell_t(\theta_\delta^*) + \delta T G \\ &\leq \sum_{t=1}^T \mathbb{E}[\hat{\ell}_t(\theta_t)] - \sum_{t=1}^T \hat{\ell}_t(\theta_\delta^*) + 3\delta T G \end{aligned} \tag{*}$$

where $\hat{\ell}_t(\theta) = \mathbb{E}_v[\ell_t(\theta + \delta v)]$ with $v \sim \text{Unif}(\mathbb{B})$ are the smoothed versions of the losses.

Proof (Step 2)

Now, recall that the algorithm runs OGD with \hat{g}_t in place of the gradients:

$$\theta_{t+1} \leftarrow \text{Proj}_{\Theta_\delta} (\theta_t - \eta \hat{g}_t)$$

Defining the pseudo-loss $h_t(\theta) = \hat{\ell}_t(\theta) + (\hat{g}_t - \nabla \hat{\ell}_t(\theta_t))^\top \theta$, we can see that

$$\nabla h_t(\theta_t) = \nabla \hat{\ell}_t(\theta_t) + \hat{g}_t - \nabla \hat{\ell}_t(\theta_t) = \hat{g}_t.$$

Therefore, the algorithm actually runs OGD on the losses h_t and thus satisfies the OGD regret bound (see Lecture 2)

$$\sum_{t=1}^T h_t(\theta_t) - \sum_{t=1}^T h_t(\theta_\delta^*) \leq \frac{D^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\hat{g}_t\|^2.$$

Furthermore, by construction of the gradient estimator, we have $\mathbb{E}_{u_t} [\hat{g}_t] = \nabla \hat{\ell}_t(\theta_t)$, which yields

$$\mathbb{E}_{u_t} [h_t(\theta_t)] = \hat{\ell}_t(\theta_t) \quad \text{and} \quad \mathbb{E}_{u_t} [h_t(\theta_\delta^*)] = \hat{\ell}_t(\theta_\delta^*)$$

Thus taking the expectation in the previous regret bound entails

$$\sum_{t=1}^T \mathbb{E}[\hat{\ell}_t(\theta_t)] - \sum_{t=1}^T \hat{\ell}_t(\theta_\delta^*) = \mathbb{E} \left[\sum_{t=1}^T h_t(\theta_t) - \sum_{t=1}^T h_t(\theta_\delta^*) \right] \leq \frac{D^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \mathbb{E}[\|\hat{g}_t\|^2] \tag{**}$$

Proof (Step 3)

Combining the two bounds (*) and (**) that we have proved, we get

$$\text{Reg}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \mathbb{E}[\|\hat{g}_t\|^2] + 3\delta T G$$

Then, since $|\ell_t(\theta)| \leq 1$ for all $\theta \in \Theta$,

$$\|\hat{g}_t\|^2 = \left(\frac{d}{\delta} \ell_t(\hat{\theta}_t) \right)^2 \leq \frac{d^2}{\delta^2}$$

This finally yields the regret

$$\text{Reg}_T \leq \frac{D^2}{2\eta} + \frac{\eta d^2 T}{2\delta^2} + 3\delta T G \leq 2d\sqrt{T} + 2\sqrt{GDd}T^{3/4}$$

for the choices of δ and η .

Convex bandits is still an active research area with many open problems.

The above regret bound of order $O(T^{3/4})$ is suboptimal.

More complicated methods can achieve $O(\sqrt{T})$ regret but with sub-optimal dependence on d and worst computational complexities.

More information can be found in Hazan et al. 2016.

Online learning / adversarial bandit

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$ for all $\theta \in \Theta$ \rightarrow full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$ \rightarrow bandit feedback.

The goal of the player is to minimize his cumulative loss:

$$\widehat{L}_T \stackrel{\text{def}}{=} \sum_{t=1}^T \ell_t(\theta_t).$$

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an arm $k_t \in \Theta$ (compact decision/parameter set, most often $\{1, \dots, K\}$);
- the player observes
 - the **rewards** of every arm: $X_t^k \sim \nu_k$ for all $k \in \Theta \rightarrow$ full-information feedback
 - the **reward** of the chosen arm only: $X_t^{k_t} \sim \nu_{k_t} \rightarrow$ bandit feedback.

The goal of the player is to maximize their cumulative reward.

Regret?

We could use the definition of the regret from adversarial bandits:

Definition (Regret, attempt 1)

$$\text{Reg}_T = \max_k \sum_{t=1}^T X_t^k - \sum_{t=1}^T X_t^{k_t} .$$

Let's see why we don't use that definition.

Notations and assumptions:

- The arm set is $[K] = \{1, \dots, K\}$.
- $\mu^k = \mathbb{E}_{X \sim \nu_k}[X]$, assumed finite for all arms k .
- $\mu^* = \max_{k \in [K]} \mu^k$.

The first notion of regret is inadequate

$$\text{Reg}_T = \max_k \sum_{t=1}^T X_t^k - \sum_{t=1}^T X_t^{k_t}.$$

ν_k Bernoulli(1/2) for all $k \in [K]$. $\mu^k = 1/2$ for all k .

All arms are the same \rightarrow there is no bad choice and **no bad algorithm**.

But:

$$\begin{aligned}\mathbb{E} \text{Reg}_T &= \mathbb{E} \left[\max_{k \in [K]} \sum_{t=1}^T X_t^k \right] - T/2 \\ &= \mathbb{E} \left[\max_{k \in [K]} \sum_{t=1}^T (X_t^k - 1/2) \right] \\ &\approx \sqrt{T \log K}\end{aligned}$$

(See any course/book/wikipedia article on symmetric random walks).

Regret definition

We want a regret notion that does not blow up with stochastic fluctuations.

Definition ((Pseudo)-Regret)

The regret is defined as

$$\text{Reg}_T = \max_k \sum_{t=1}^T \mu^k - \sum_{t=1}^T \mu^{k_t} = T\mu^* - \sum_{t=1}^T \mu^{k_t} .$$

Recall that $\mu^k = \mathbb{E}_{X \sim \nu_k}[X]$.

Most often, we bound the **expected regret** $\mathbb{E}[\text{Reg}_T]$.

Note that the expectation here is over the random rewards and the randomness of the algorithm, if there is any.

Regret decomposition

Suppose that the set of arms is finite: $[K]$.

Define the **gap** of arm $k \in [K]$ by $\Delta_k = \mu^* - \mu^k$.

$$\text{Reg}_T = T\mu^* - \sum_{t=1}^T \mu^{k_t} = \sum_{t=1}^T (\mu^* - \mu^{k_t}) = \sum_{t=1}^T \Delta_{k_t} = \sum_{k=1}^K N_T^k \Delta_k ,$$

where $N_T^k = \sum_{t=1}^T \mathbb{I}\{k_t = k\}$ is the number of pulls of arm k up to time T .

Bounding the regret \Leftrightarrow bounding the number of pulls of bad arms

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

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an arm $k_t \in \Theta$ (compact decision/parameter set, most often $\{1, \dots, K\}$);
- the player observes the reward of the chosen arm only: $X_t^{k_t} \sim \nu_{k_t}$ (independent of other rewards).

The goal of the player is to minimize their expected regret: $\mathbb{E}[\text{Reg}_T] = \sum_{k=1}^K \mathbb{E}[N_T^k] \Delta_k$.

Setting variants:

- **Contextual bandit:** $X_t^{k_t} \sim \nu_{k_t}(x_t)$, for a known context x_t
- **Linear bandit:** $\nu_{k_t} = \mathcal{N}(\theta^\top x_{k_t}, 1)$
- **Structured bandit:** the algorithm knows constraints on $(\mu^k)_{k \in [K]}$, e.g. Lipschitz, linear, monotone...

Goal variants: instead of minimizing the regret, we want to

- **Minimize the simple regret:** return an arm at time T , and minimize its expected gap.
- **Identify the best arm:** return an arm at time T , and minimize the probability that its not one of the best ones.

Relaxed assumptions: rewards not independent, distributions changing with time, etc.

Convergence to the mean

Main idea: we can estimate the mean of the arms with the empirical mean.

Let $(X_s)_{s \in \mathbb{N}}$ be iid random variables with $\mathbb{E}[|X_1|] < \infty$ and expected value $\mathbb{E}[X_1] = \mu$.

Let $\bar{X}_t = \sum_{s=1}^t X_s$.

Theorem 14 (Strong law of large numbers)

$\bar{X}_t \xrightarrow{\text{a.s.}} \mu$, that is $\mathbb{P}(\bar{X}_t \rightarrow \mu) = 1$.

Theorem 15 (Central limit theorem)

If $\mathbb{V}[X] = \sigma^2 < \infty$, then $\sqrt{t}(\bar{X}_t - \mu) \xrightarrow{d} \mathcal{N}(\mu, \sigma^2)$.

Problem: those are asymptotic results.

Main question: if I have 15 samples of arm k , how reliable is my estimate for μ^k ?

Concentration inequalities

Our main tools are **concentration inequalities**: bounds on the probability that the empirical mean (or another statistic) is far from its expected value.

Theorem 16 (Hoeffding's inequality)

If X_1, \dots, X_t are independent random variables almost surely in $[a, b]$ then for all $\delta \in (0, 1)$ we have

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

Equivalently, for all $\varepsilon \geq 0$,

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

Proof

Proof under a sub-Gaussian assumption. Exercise: bounded implies sub-Gaussian.

Assumption: for all s , X_s is σ^2 -sub-Gaussian, which means that for all $\lambda \in \mathbb{R}$,

$$\mathbb{E}[e^{\lambda(X_s - \mu_s)}] \leq e^{\frac{1}{2}\sigma^2\lambda^2}.$$

Proof

Warning: random number of samples

In the analysis of bandit algorithms, we want to bound $\hat{\mu}_t^k - \mu^k$, where
 $\hat{\mu}_t^k = \frac{1}{N_t^k} \sum_{s=1}^t X_s^{k_s} \mathbb{I}\{k_s = k\}$.

k_s is a random variable that depends on all previous rewards.

Issue: $\hat{\mu}_t^k$ is a sum of a **random number** of random variables which are **not independent**.

- $\hat{\mu}_t^k$ is **not** unbiased.
- $\hat{\mu}_t^k$ is **not** a sum of a fixed number of independent random variables.
- **Hoeffding's inequality does not apply.**

How to avoid the difficulty: union bounds, or martingale arguments (see proofs later in the course).

Goal: minimize $\mathbb{E}[\text{Reg}_T] = T\mu^* - \sum_{t=1}^T \mu^{k_t}$.

Since the empirical mean of an arm concentrate around its expected value, can we simply pull the arm with highest empirical mean?

Definition (Follow-The-Leader)

The FTL algorithm first explores each arm once $k_t = t$ for $k \leq K$ and then pulls arm $k_t = \arg \max_{k \in [K]} \hat{\mu}_{t-1}^k$ for all $t \geq K + 1$.

Full information: yes, FTL is optimal.

Bandit: answer is no, FTL does not work. It has linear expected regret in most settings.

FTL still fails

Explore then commit

Need to not only **exploit**, but also **explore**.

Explore-Then-Commit

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

$$\hat{\mu}_{mK}^k = \frac{1}{m} \sum_{s=1}^{Km} \mathbb{I}\{k_s = k\} X_s^k.$$

2. **Exploitation:** keep playing the best arm $\arg \max_k \hat{\mu}_{mK}^k$ for the remaining rounds $t = mK + 1, \dots, T$.

Theorem 17 (Thm 6.1, Lattimore and Szepesvári 2019)

If all distributions are bounded in $[0, 1]$ and $1 \leq m \leq T/K$ then ETC has expected regret

$$\mathbb{E}[\text{Reg}_T] \leq m \sum_{k=1}^K \Delta_k + (T - mK) \sum_{k=1}^K \Delta_k \exp(-m\Delta_k^2).$$

- m too large \Rightarrow too much exploration, linear regret.
- m too small \Rightarrow too little exploration, linear regret.
- What m should we choose?

Proof

Proof

Finding the right trade-off

Two arms bandit: arm 1 is the best arm, arm 2 has gap Δ .

ETC verifies

$$\mathbb{E}[\text{Reg}_T] \leq m\Delta + (T - 2m)\Delta e^{-m\Delta^2}.$$

Theorem 18

If $K = 2$ and $m = \max\{1, \left\lceil \frac{\log(T\Delta^2)}{\Delta^2} \right\rceil\}$, then

$$\mathbb{E}[\text{Reg}_T] \leq \Delta + \frac{1 + \log(T\Delta^2)}{\Delta}.$$

This is a **distribution dependent** bound, meaning that it depends on the gap.

Issue with those bounds: meaningless if Δ is small.

ETC verifies

$$\mathbb{E}[\text{Reg}_T] \leq m\Delta + (T - 2m)\Delta e^{-m\Delta^2}.$$

Theorem 19

If $K = 2$ and $m = \max \left\{ 1, \left\lceil \frac{\log(T\Delta^2)}{\Delta^2} \right\rceil \right\}$, then

$$\mathbb{E}[\text{Reg}_T] \leq \min \left\{ \Delta + \frac{1 + \log(T\Delta^2)}{\Delta}, T\Delta \right\} \lesssim \sqrt{T \log T}.$$

This is close to optimal: we can prove a lower bound of order \sqrt{T} .

Problems:

- m depends on Δ , which is unknown.
- What can we do for $K > 2$?

Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

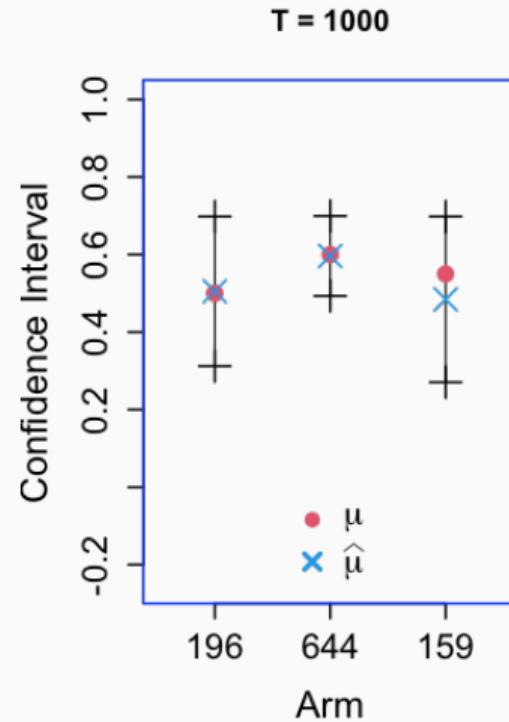
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

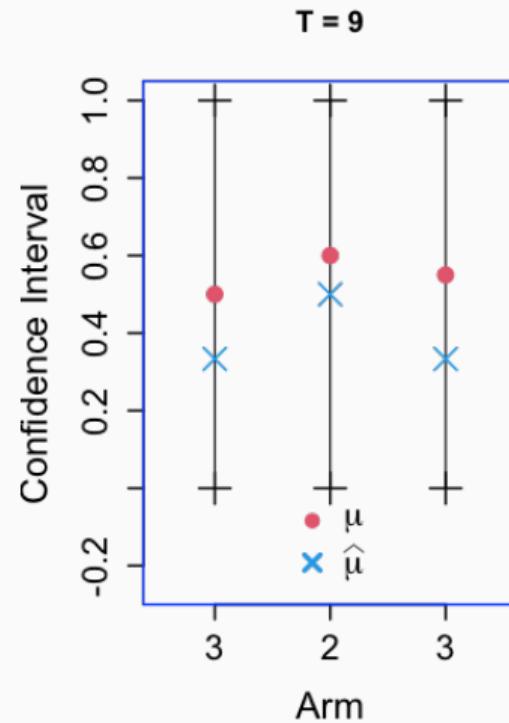
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

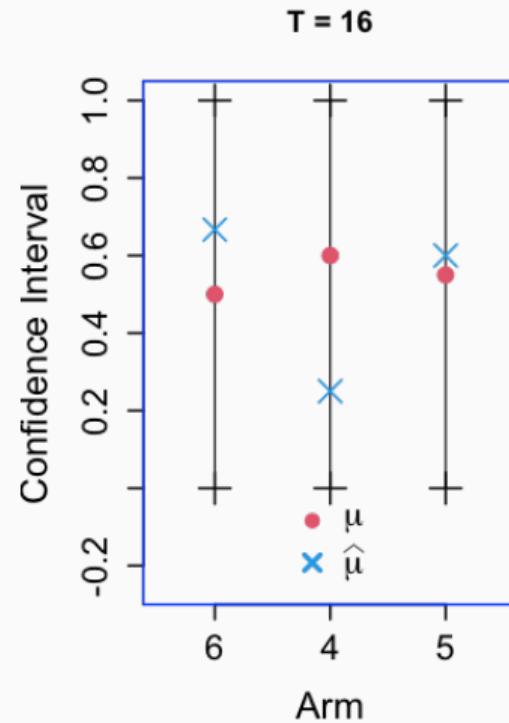
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

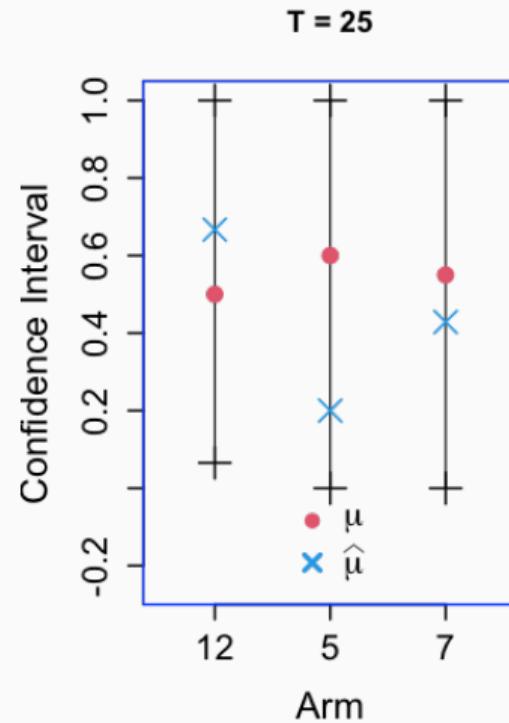
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

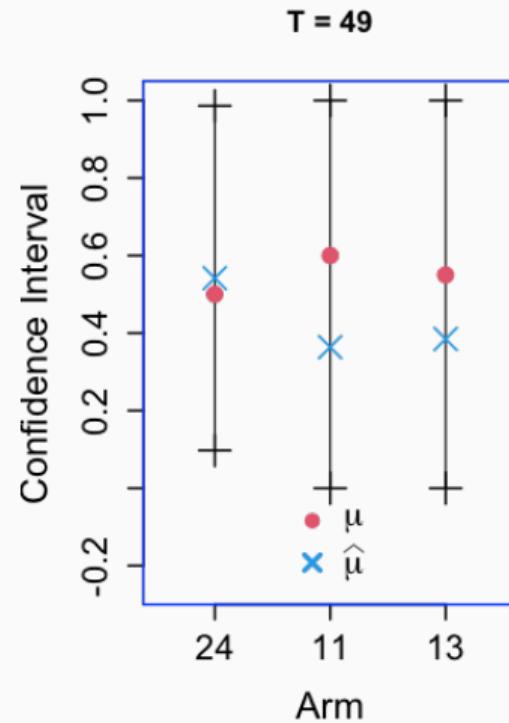
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

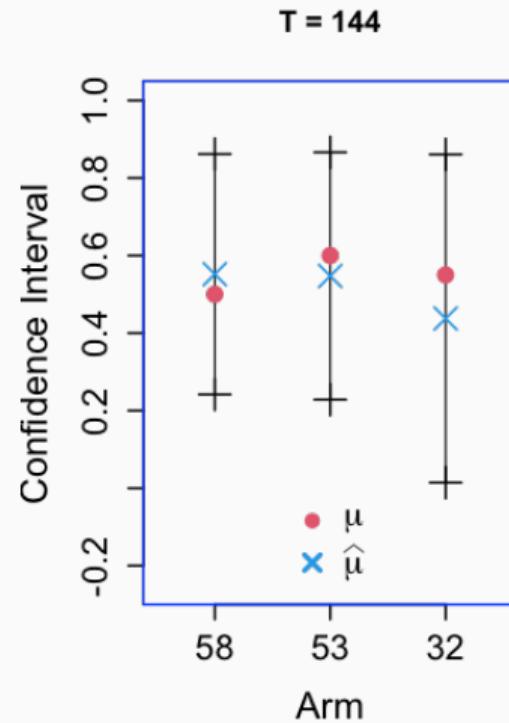
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

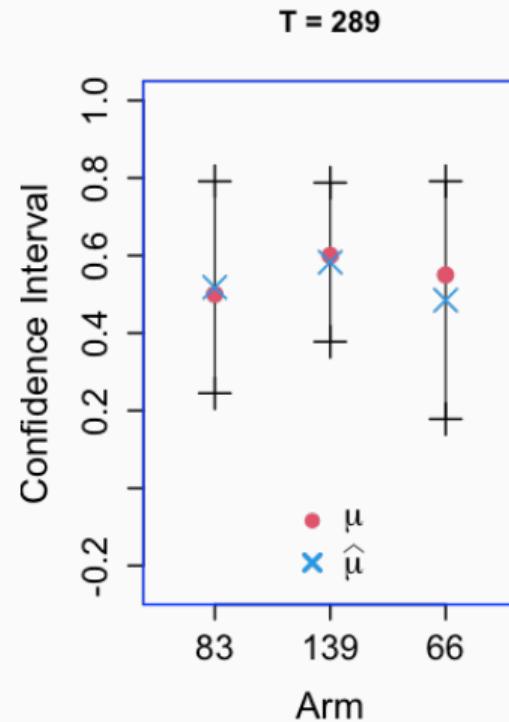
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

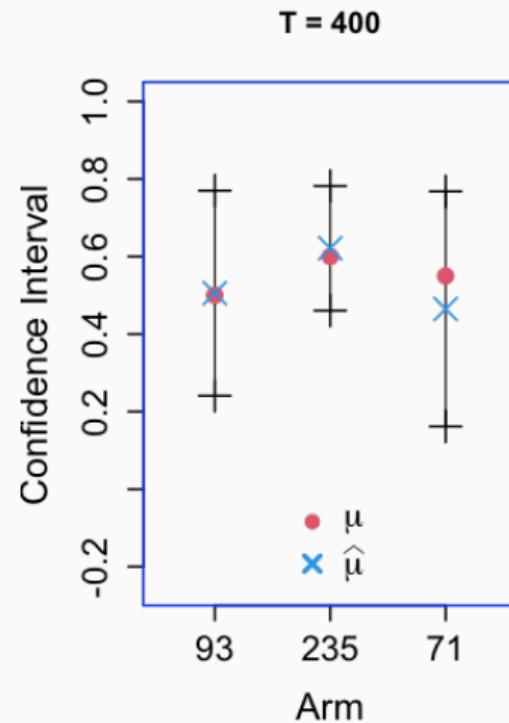
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Being optimistic

The UCB (Upper-Confidence-Bound) algorithm is a very popular bandit algorithm that has several advantages over ETC:

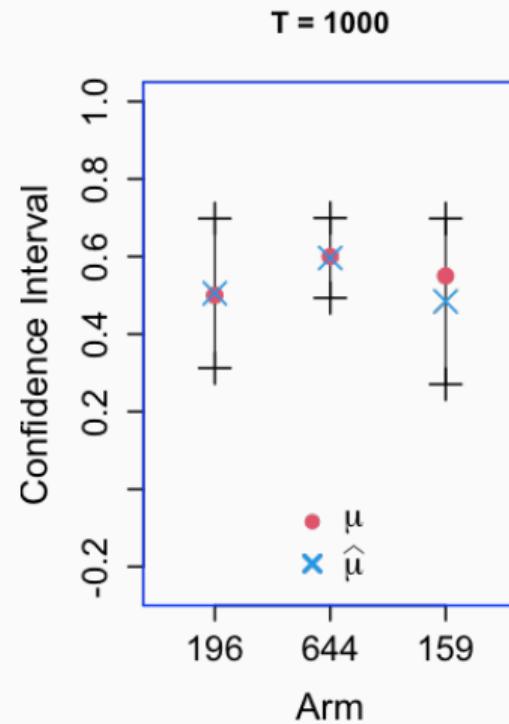
- It does not rely on an initial exploration phase but explores on the fly as rewards are observed.
- Unlike ETC, it does not require knowledge of gaps and behaves well when there are more than two arms.

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

$$I_t^k = [L_t^k, U_t^k].$$

It is **optimistic**, acting as if the best possible rewards are the real rewards:

$$k_t \in \arg \max_{k \in \{1, \dots, K\}} U_t^k.$$



Confidence intervals

How to design the upper confidence bounds?

→ concentration inequalities. Here **Hoeffding's inequality**.

Theorem 20 (Hoeffding's inequality)

If X_1, \dots, X_t are independent random variables almost surely in $[a, b]$ with same mean μ then for all $\delta \in (0, 1)$ we have

$$\mathbb{P} \left(\frac{1}{t} \sum_{s=1}^t X_s - \mu \geq \sqrt{\frac{(b-a)^2}{2t} \log \frac{1}{\delta}} \right) \leq \delta.$$

Careful: UCB is adaptive, hence $\hat{\mu}_t$ is not exactly a sum of independent random variables. But we will make it work.

For rewards in $[0, 1]$: $U_t^k = \hat{\mu}_{t-1}^k + \sqrt{\frac{2 \log t}{N_{t-1}^k}}$

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

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

$$k_t \in \arg \max_{k \in [K]} \left\{ \hat{\mu}_{t-1}^k + \sqrt{\frac{2 \log t}{N_{t-1}^k}} \right\},$$

and get reward $X_t^{k_t}$.

Regret Bound

Theorem 21

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

$$\mathbb{E}[N_T^k] \leq \frac{8 \log T}{\Delta_k^2} + 2.$$

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

$$\mathbb{E}[\text{Reg}_T] \leq 2K + \sum_{k: \Delta_k > 0} \frac{8 \log T}{\Delta_k}.$$

Remarks :

- we can also prove $\mathbb{E}[\text{Reg}_T] \lesssim \sqrt{KT \log(T)}$. Close to the optimal $O(\sqrt{KT})$.
- Deals with multiple gaps, without any knowledge of the gaps, unlike ETC.
- Bounded can be replaced by sub-Gaussian.

Proof start

Idea: if the means belong to the confidence intervals and the arms are pulled enough, the algorithm cannot pull a suboptimal arm.

We prove that if $k_t = k \neq *$, then one of these inequalities must be false:

$$\mu^* \leq \hat{\mu}_{t-1}^* + \sqrt{\frac{2 \log t}{N_{t-1}^*}} \quad \leftarrow \mu^* \text{ smaller than UCB} \quad (\text{i})$$

$$\mu^k \geq \hat{\mu}_{t-1}^k - \sqrt{\frac{2 \log t}{N_{t-1}^k}} \quad \leftarrow \mu_k \text{ larger than LCB} \quad (\text{ii})$$

$$N_{t-1}^k \geq \frac{8 \log t}{\Delta_k^2} \quad \leftarrow k \text{ played enough} \quad (\text{iii})$$

Proof 2

$$\mu^* \leq \hat{\mu}_{t-1}^* + \sqrt{\frac{2 \log t}{N_{t-1}^*}} \quad \text{and} \quad \mu^k \geq \hat{\mu}_{t-1}^k - \sqrt{\frac{2 \log t}{N_{t-1}^k}} \quad \text{and} \quad N_{t-1}^k \geq \frac{8 \log t}{\Delta_k^2}$$

Prove that if k is pulled at t , then there is a contradiction.

Proof 3: decomposition wrt events

One of these is false:

$$\mu^* \leq \hat{\mu}_{t-1}^* + \sqrt{\frac{2 \log t}{N_{t-1}^*}} \quad ; \quad \mu^k \geq \hat{\mu}_{t-1}^k - \sqrt{\frac{2 \log t}{N_{t-1}^k}} \quad ; \quad N_{t-1}^k \geq \frac{8 \log t}{\Delta_k^2}$$

$$\text{Then: } \mathbb{E}[N_T^k] \leq u + \sum_{t=u+1}^T \left(\mathbb{P}\{(\text{i}) \text{ is false}\} + \mathbb{P}\{(\text{ii}) \text{ is false}\} \right) \quad \text{for } u = \left\lceil \frac{8 \log T}{\Delta_k^2} \right\rceil.$$

Proof 4: probability of the concentration event

We show: $\mathbb{P}(\mu^k < \hat{\mu}_{t-1}^k - \sqrt{\frac{2 \log t}{N_{t-1}^k}}) \leq t^{-3}$.

Proof summary

For $u = \frac{8 \log T}{\Delta_k^2}$, $\mathbb{E}[N_T^k] \leq u + \sum_{t=u+1}^T \left(\mathbb{P}\{\mu^* > \hat{\mu}_{t-1}^* + \sqrt{\frac{2 \log t}{N_{t-1}^*}}\} + \mathbb{P}\{\mu^k < \hat{\mu}_{t-1}^k - \sqrt{\frac{2 \log t}{N_{t-1}^k}}\} \right)$.

Each of these probabilities is smaller than t^{-3} .

$$\mathbb{E}[\text{Reg}_T] \leq \frac{8 \log T}{\Delta_k^2} + 2 \sum_{t=u+1}^T \frac{1}{t^3} \leq \frac{8 \log T}{\Delta_k^2} + 2.$$

The bound of the regret then comes from $\mathbb{E}[\text{Reg}_T] = \sum_k \mathbb{E}[N_T^k] \Delta_k$.

ε -greedy

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 .

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 .

Other Algorithms: Thompson Sampling

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 $t \geq 1$, it

- computes $\hat{\nu}_{k,t}$ the posterior distribution of the rewards of each 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. Somewhat different proof techniques.

An advantage over UCB is the possibility of incorporating easily prior knowledge on the arms.

UCB proved easier to adapt to structured bandits (it can be hard to sample a posterior conditioned on structural information).

Stochastic Linear Bandits - Motivation

Main motivation: use contexts.

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

$$X_t = \theta_t^\top \mu^* + \varepsilon_t$$

where ε_t is i.i.d. 1-subgaussian noise. ($\forall \lambda > 0$, $\mathbb{E}[\exp(\lambda \varepsilon_t)] \leq \exp(\lambda^2/2)$)

- the player only observes the feedback X_t .

The player wants to minimize its expected regret defined as

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

Examples

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

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

for some unknown parameter $\mu^* \in \mathbb{R}^d$ and feature map $\psi : [K] \times \mathcal{X} \rightarrow \mathbb{R}^d$.

- Combinatorial bandit: $\Theta_t \subseteq \{0, 1\}^d \rightarrow$ combinatorial bandit problems. Example: decision set = 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 - UCB for linear bandits.

- Build confidence region for the parameter: \mathcal{C}_t such that $\mu^* \in \mathcal{C}_t$ with high probability.
- Build confidence bounds for the arm means: $U_t^\theta = \max_{\mu \in \mathcal{C}_t} \theta^\top \mu$.
- Be optimistic: pull $\theta_t = \arg \max_\theta U_t^\theta$.

Main question: how do we get \mathcal{C}_t ?

Confidence region

After time t , the algorithm observed:

$$X_1 = \theta_1^\top \mu^* + \varepsilon_1$$

$$X_2 = \theta_2^\top \mu^* + \varepsilon_2$$

...

$$X_t = \theta_t^\top \mu^* + \varepsilon_t$$

The unknown parameter we want to estimate is μ^* .

Denoting by I_d the $d \times d$ identity matrix and picking $\lambda > 0$, we can estimate μ^* with **regularized least square**

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

where $V_t \stackrel{\text{def}}{=} \lambda I_d + \sum_{s=1}^t \theta_s \theta_s^\top$.

Lemma 3

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

Conclusion: with probability $1 - \delta$, for all $t \geq 1$,

$$\mu^* \in C_t, \quad \text{where} \quad C_t \stackrel{\text{def}}{=} \left\{ \theta \in \mathbb{R}^d : \|\mu - \hat{\mu}_{t-1}\|_{V_{t-1}} \leq \beta(\delta/T) \right\}. \quad (14)$$

Theorem 22

Let $T \geq 1$ and $\mu^* \in \mathbb{R}^d$. Assume that for all $\theta \in \cup_{t=1}^T \Theta_t$, $|\theta^\top \mu^*| \leq 1$, $\|\mu^*\| \leq 1$ and $\|\theta\| \leq 1$, then LinUCB with C_t defined as above satisfies the regret bound

$$\mathbb{E} \text{Reg}_T \leq \square_\lambda d \sqrt{T} \log(T),$$

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

Remark:

- $O(\sqrt{T})$: the exponent does not depend on d .

Proof

With probability $1 - 1/T$, for all $t \geq 1$,

$$\mu^* \in C_t, \quad \text{where} \quad C_t \stackrel{\text{def}}{=} \left\{ \theta \in \mathbb{R}^d : \|\mu - \hat{\mu}_{t-1}\|_{V_{t-1}} \leq \beta(1/T^2) \right\}.$$

Proof 2

Proof 3

LinUCB with C_t defined as above satisfies the regret bound

$$\mathbb{E} \text{Reg}_T \lesssim d\sqrt{T} \log(T),$$

To prove it, we assumed the following lemma:

Lemma 4

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

Under additional assumptions, it is possible to improve the regret bound $O(d\sqrt{T} \log T)$.

- If the set of available actions at time t is fixed and finite; i.e., $\theta_t \in \Theta$ where $|\Theta| = K$. Then, it is possible to achieve

$$\mathbb{E} \text{Reg}_T \leq \square \sqrt{Td \log(TK)},$$

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} .

- 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 $\tilde{O}(\sqrt{dm_0 T})$.

Thank you!

- [1] Nicolo Cesa-Bianchi and Gábor Lugosi. **Prediction, learning, and games**. Cambridge university press, 2006.
- [2] Elad Hazan et al. “**Introduction to online convex optimization**”. In: Foundations and Trends® in Optimization 2.3-4 (2016), pp. 157–325.
- [3] Tor Lattimore and Csaba Szepesvári. “**Bandit algorithms**”. In: preprint (2019).
- [4] Francesco Orabona. “**A modern introduction to online learning**”. In: arXiv preprint arXiv:1912.13213 (2019).
- [5] Shai Shalev-Shwartz et al. “**Online learning and online convex optimization**”. In: Foundations and Trends® in Machine Learning 4.2 (2012), pp. 107–194.
- [6] Martin Zinkevich. “**Online convex programming and generalized infinitesimal gradient ascent**”. In: Proceedings of the 20th International Conference on Machine Learning (ICML-03). 2003, pp. 928–936.