

Lecture 6: Online prediction with expert advice

Related to:

- bandits
- stochastic optimization
- classical statistics

References:

"Predict, Learning and Games", M. Cesa-Bianchi and G. Lugosi
 Chapter 2

miscellaneous desc. → Survey: "Convex optim.: Algorithms and complexity" S. Boyd

Survey: "Regret analysis of stoc. and non stoc. multiarmed bandit pb" Beberck and Cesa-Bianchi

- I Setting
- II Convex case
- III Non convex case

I Setting

1.1 The problem

A statistician/learner wants to predict round after round the values of a sequence $y_1, y_2, y_3, \dots \in Y$
 His/her predictions are denoted by

$$\hat{a}_1, \hat{a}_2, \dots \in \mathcal{D} \begin{cases} \text{not necessarily } \mathcal{C} \\ \text{convex subset of a vector space} \end{cases}$$

Sequential "aggregat°": at each round $t \geq 1$, the statistician has access to K base predictions (the "expert advice")

$$a_{1,t}, \dots, a_{K,t} \in \mathcal{D}$$

that can be combine ("aggregate") to choose his/her own predict^o $\hat{a}_\epsilon \in \mathcal{D}$. (typically, \hat{a}_ϵ is a weighted average of the $a_{i,\epsilon}$)

Evaluat^o of the predict^o quality

loss funct^o $l: \mathcal{D} \times \mathcal{Y} \rightarrow \mathbb{R}$

$(a, y) \mapsto l(a, y)$

measure how a "is close" to y

Goal of the statistician: on the long run the goal is to output predict^o that are on average almost as good as those of the best expert in hindsight, i.e., to minimize the "regret".

$$\text{Reg}_T = \underbrace{\sum_{\epsilon=1}^T l(\hat{a}_\epsilon, y_\epsilon)}_{\text{cumulative loss of the statistician}} - \underbrace{\min_{1 \leq i \leq K} \sum_{\epsilon=1}^T l(a_{i,\epsilon}, y_\epsilon)}_{\text{cum. loss of the best expert in hindsight}}$$

cumulative loss of the statistician

cum. loss of the best expert in hindsight

Assumpt^o: ? almost none!

We will prove guarantees of the form $\text{Reg}_T \leq \dots$ for all sequences $(y_\epsilon)_{\epsilon \geq 1}$ in \mathcal{Y} and $(a_\epsilon)_{\epsilon \geq 1}$ in \mathcal{D}^K

$$a_\epsilon = (a_{i,\epsilon})_{1 \leq i \leq K}$$

provided: l is convex in its first argument and bounded.

In particular: $y_\epsilon, \epsilon \geq 1$ not necessarily i.i.d.

All results stated below will be valid in the following scenarios:

- Meta-statistical setting: the sequence $(y_\epsilon)_{\epsilon \geq 1}$ is the realization of a stochastic process and the $a_{i,\epsilon}$ are

sequential estimations coming from k different statistical methods.

→ The goal is to learn as fast as the best statistical method in our toolbox.

• deterministic setting

We make no stochastic assumption on $(g_t)_{t \geq 1}$

$a_{i,t}$ = predict^o output by engineering / physical models.

Example: $g_t =$ max amplitude of the ozone concentration at ground level (t : index of day)

$a_{i,t}$ = predict^o made by chemical-physical models

• adversarial setting:

The g_t and the $a_{i,t}$ can be chosen by a malicious adversary who wants to defeat the statistician

Example: computer security, spam detect^o

Formal recap' of the setting:

At every round $t \geq 1$:

- ① The expert advice $a_{1,t}, \dots, a_{k,t} \in \mathcal{D}$ are revealed to the statistician
- ② The statistician makes her own predict^o $\hat{a}_t \in \mathcal{D}$, using the $a_{i,t}$ but also the past data $(g_s, a_s)_{1 \leq s \leq t-1}$
- ③ The statistician observes $g_t \in \mathcal{Y}$ and suffers the loss $\ell(\hat{a}_t, g_t)$

NB: The above online protocol enables us to make out the dependencies between all quantities, eg, $\hat{a}_t = \hat{a}_t(y_{1:t-1}, a_{1:t-1})$

of 'strategy' is a sequence of funct^o $(\hat{a}_t(\cdot))_{t \geq 1}$
 where $\hat{a}_t : (Y \times \mathcal{D})^{t-1} \times \mathcal{D}^k \rightarrow \mathcal{D}$

Goal: minimizing the regret

More precisely, the goal of this lecture is to build a strategy $(\hat{a}_t(\cdot))_{t \geq 1}$

st

$$\sup_{\substack{(y_t)_{t \geq 1} \text{ in } Y \\ (a_t)_{t \geq 1} \text{ in } \mathcal{D}^k}} \left\{ \frac{1}{T} \sum_{t=1}^T \ell(\hat{a}_t, y_t) - \min_{1 \leq i \leq k} \frac{1}{T} \sum_{t=1}^T \ell(a_i, y_t) \right\} \leq o(1)$$

as $T \rightarrow +\infty$

This implies that the average loss of the statistician is almost as small as that of the best expert.

This is true for all sequences $(y_t)_{t \geq 1}$ and $(a_t)_{t \geq 1}$
 \rightarrow This theory is sometimes called 'predict^o of individual sequences', the theorems are valid for every seq.

Examples of \mathcal{D} , Y and ℓ

- Square loss for bounded predict^o and observat^o
 $\mathcal{D} = Y = [-B, B]$; $\ell(a, y) = (a - y)^2$

- Absolute loss for binary classification
 $\mathcal{D} = [0, 1]$, $Y = \{0, 1\}$, $\ell(a, y) = |a - y|$

- Linear loss on the simplex
 $\mathcal{D} = \Delta(k) := \{v \in \mathbb{R}_+^k : \sum_{i=1}^k v_i = 1\}$
 probab. vect. on k elements

$Y = [0, 1]^k$
 $\ell(a, y) = a \cdot y = \sum_{i=1}^k a_i \cdot y_i$

1-2. Why is convex aggregat^o useful?

Imagine the simple problem:

$$\mathcal{D} = \{0, 1\}, \quad y = \{0, 1\},$$

$$\ell(a, y) = |a - y| = \mathbb{1}_{a \neq y} \quad (\text{because } a, y \in \{0, 1\})$$

At each round $t \geq 1$,

$$\begin{cases} a_{1,t} = 0 \\ a_{2,t} = 1 \end{cases}$$

Important comment: whatever $\hat{a}_t \in \{0, 1\}$ there exists $y_t \in \{0, 1\}$ s.t. $\ell(\hat{a}_t, y_t) = 1$ (just take $y_t = 1 - \hat{a}_t$)
so that $\sum_{i=1}^T \ell(\hat{a}_t, y_t) = T$

$$\text{but } \min_{1 \leq i \leq 2} \sum_{t=1}^T \ell(a_{i,t}, y_t) \leq \frac{1}{2} \sum_{i=1}^2 \sum_{t=1}^T \ell(a_{i,t}, y_t) = \frac{T}{2}$$

As a consequence,

$$\begin{aligned} \frac{\text{Reg}_T}{T} &= \frac{1}{T} \sum_{t=1}^T \ell(\hat{a}_t, y_t) - \min_{1 \leq i \leq 2} \frac{1}{T} \sum_{t=1}^T \ell(a_{i,t}, y_t) \\ &\geq \frac{1}{2}, \quad \text{not } \leq o(1)! \end{aligned}$$

The problem is that \mathcal{D} is not convex.

If $\mathcal{D} = [0, 1]$ instead, we could choose $\hat{a}_t \in [0, 1]$ (eg $\frac{1}{2}$)

and probably achieve $\text{Reg}_T \leq \sqrt{T}$ so that

$$\frac{\text{Reg}_T}{T} \lesssim \frac{1}{\sqrt{T}} = o(1) \quad \begin{matrix} \uparrow \\ \text{constants are} \\ \text{omitted} \end{matrix}$$

Take-home msg: convexity is crucial!

II The convex case

2-1 - The exponentially weighted average forecaster ("EWA")

Popular online algorithm / strategy dating back to the 90's
EWA algorithm

- parameter: $\eta > 0$
- At each round $t \geq 1$,
 - * collect the expert advice $a_{i,t}$ $1 \leq i \leq k$
 - * compute the ^{weight} probability vector $p_t \in \Delta(k)$ as

$$p_{i,t} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell(a_{i,s}, g_s)\right)}{\sum_{j=1}^k \exp\left(-\eta \sum_{s=1}^{t-1} \ell(a_{j,s}, g_s)\right)} \quad 1 \leq i \leq k$$

- * compute the average prediction

$$\hat{a}_t = \sum_{i=1}^k p_{i,t} a_{i,t} \in \mathcal{D}$$

↑ because $a_{i,t} \in \mathcal{D}$ and \mathcal{D} is convex

- * receive the observation $g_t \in \mathcal{G}$

Theorem: (Cesa-Bianchi 1999)

Let $T \geq 1$ and $k \geq 2$.

Assume that $\ell: \mathcal{D} \times \mathcal{G} \rightarrow [B_1, B_2]$ is convex in its first argument, i.e. $a \mapsto \ell(a, g)$ is convex for all $g \in \mathcal{G}$.

Then, the EWA algorithm tuned with $\eta > 0$ satisfies, for all sequences $(g_t)_{t \geq 1}$ in \mathcal{G} and $(a_t)_{t \geq 1}$ in \mathcal{D}^k ,

$$\sum_{t=1}^T \ell(\hat{a}_t, g_t) + \min_{1 \leq i \leq K} \sum_{t=1}^T \ell(a_{i,t}, g_t) \leq \frac{\ln K}{\eta} + \frac{\eta T (B_2 - B_1)^2}{8}$$

In particular, the choice of $\eta = \frac{1}{B_2 - B_1} \sqrt{\frac{8 \ln K}{T}}$

yields $\text{Reg}_T \leq (B_2 - B_1) \sqrt{\frac{T \ln K}{2}}$

Consequence: for $\eta = \frac{1}{B_2 - B_1} \sqrt{\frac{8 \ln K}{T}}$, $\frac{\text{Reg}_T}{T} \leq (B_2 - B_1) \sqrt{\frac{\ln K}{2T}} = o(1)$ as $T \rightarrow \infty$

Proof: we set

$$l_{i,t} = \ell(a_{i,t}, g_t)$$

$$L_{i,t} = \sum_{s=1}^t l_{i,s} \quad (\text{cumsum of expert } i)$$

$$W_t = \frac{1}{K} \sum_{i=1}^K e^{-\eta L_{i,t-1}} \quad L_{i,0} = 0 \text{ convention}$$

The idea here is to control the quantity $\ln W_{T+1} - \ln W_1$ in two different ways.

• Upper bound:

$$\ln W_{T+1} - \ln W_1 = \sum_{t=1}^T \ln \left(\frac{W_{t+1}}{W_t} \right) \quad (\text{telescopic sum})$$

$$= \sum_{t=1}^T \ln \left(\frac{\frac{1}{K} \sum_{i=1}^K e^{-\eta L_{i,t}}}{\frac{1}{K} \sum_{i=1}^K e^{-\eta L_{i,t-1}}} \right)$$

$$= \sum_{t=1}^T \ln \left(\frac{\sum_{i=1}^K e^{-\eta L_{i,t-1}} e^{-\eta l_{i,t}}}{\sum_{i=1}^K e^{-\eta L_{i,t-1}}} \right)$$

$P_{i,t}$

Next we provide an upper bound on $\ln \left(\sum_{i=1}^k p_{i,t} e^{-\eta l_{i,t}} \right)$ using Hoeffding's lemma.

Hoeffding's Lemma

If $X \in [a, b]$ a.s., then $\forall \lambda \in \mathbb{R}$, $\ln \mathbb{E} \left[e^{\lambda(X - \mathbb{E}[X])} \right] \leq \frac{\lambda^2 (b-a)^2}{8}$

We define a r.v. I_t s.t. $\mathbb{P}(I_t = i) = p_{i,t}$ and $X = l_{I_t, t}$ which is a r.v. with values in $[B_1, B_2]$ ($l: \mathcal{I} \times \mathcal{I} \rightarrow [B_1, B_2]$)

We use Hoeffding's lemma with X and $\lambda = -\eta$

$$\ln \mathbb{E} \left[e^{-\eta(X - \mathbb{E}[X])} \right] \leq \frac{\eta^2 (B_2 - B_1)^2}{8}$$

$$\text{but } \mathbb{E}[X] = \mathbb{E} \left[l_{I_t, t} \right] = \sum_{i=1}^k p_{i,t} l_{i,t} = p_t \cdot l_t$$

and $\ln \mathbb{E} \left[e^{-\eta(X - \mathbb{E}[X])} \right]$ $l_t = (l_{i,t})_{i=1, \dots, k} \in \mathbb{R}^k$

$$= \ln \left(\sum_{i=1}^k p_{i,t} e^{-\eta(l_{i,t} - p_t \cdot l_t)} \right)$$

$$= \eta p_t \cdot l_t + \ln \left(\sum_{i=1}^k p_{i,t} e^{-\eta l_{i,t}} \right)$$

Therefore,

$$\ln \left(\sum_{i=1}^k p_{i,t} e^{-\eta l_{i,t}} \right) \leq \eta p_t \cdot l_t + \frac{\eta^2 (B_2 - B_1)^2}{8}$$

Finally, summing over $t = 1, \dots, T$

$$\ln W_{T+1} - \ln W_1 \leq \eta \sum_{t=1}^T p_t \cdot l_t + \frac{T \eta^2 (B_2 - B_1)^2}{8}$$

• Lower bound

$$\begin{aligned} \ln W_{T+1} - \ln W_1 &= \ln W_{T+1} - 0 = \ln \left(\sum_{i=1}^K e^{-\eta L_{i,T}} \right) - \ln K \\ &\geq \ln \left(\max_{1 \leq i \leq K} e^{-\eta L_{i,T}} \right) - \ln K \\ &= -\eta \min_{1 \leq i \leq K} L_{i,T} - \ln K \end{aligned}$$

• Collecting the upper and lower bounds.

$$\begin{aligned} -\eta \min_{1 \leq i \leq K} L_{i,T} - \ln K &\leq \ln W_{T+1} - \ln W_1 \\ &\leq -\eta \sum_{t=1}^T p_t \cdot \ell_t + \frac{T\eta^2 (B_2 - B_1)^2}{8} \end{aligned}$$

Rearranging bounds and dividing by η , we get:

$$(1) \quad \left\{ \begin{aligned} \sum_{t=1}^T p_t \cdot \ell_t - \min_{1 \leq i \leq K} \sum_{t=1}^T \ell_{i,t} &\leq \frac{\ln K}{\eta} + \frac{T\eta (B_2 - B_1)^2}{8} \\ &P(a_{i,t}, g_t) \end{aligned} \right.$$

will be used later

To conclude, note that

$$\begin{aligned} \ell(\vec{a}_t, g_t) &\stackrel{\text{def of } \vec{a}_t}{=} \ell \left(\sum_{i=1}^K p_{i,t} a_{i,t}, g_t \right) \\ &\leq \sum_{i=1}^K p_{i,t} \ell(a_{i,t}, g_t) \\ &= p_t \cdot \ell_t \quad \text{since } \ell(\cdot, g_t) \text{ is convex} \end{aligned}$$

□

Remarks:

- Important: how to tune the parameter η in practice? The values of T or $B_2 - B_1$ may be unknown in practice, so that $\eta = \frac{1}{B_2 - B_1} \sqrt{\frac{8mk}{T}}$ may be unfeasible

take η_ϵ with $B_2, B_1 \sim$ empirical min and max (proof similar)

\rightarrow if T is unknown, taking $\eta_\epsilon = \frac{1}{B_2 - B_1} \sqrt{\frac{8mk}{\epsilon}}$ works as

well i.e. it yields a similar regret bound

\rightarrow if $B_2 - B_1$ is unknown, we can tune η_ϵ as a funct^o of the past data

\rightarrow similar regret bound, without needing to know T or $B_2 - B_1$. [bounds, e.g., $\text{Reg}_T \leq \ln k$ instead of $\text{Reg}_T \leq \sqrt{T \ln k}$]

- For special loss functions ℓ (e.g., strongly convex), it is possible to prove better regret

2-2 Applicat^o to $\nabla \ell_\epsilon$; links with gradient descent

Assume in this subsection that $\mathcal{D} = \Delta(K)$

$\ell: \Delta(K) \times \mathcal{Y} \rightarrow \mathbb{R}$ is differentiable in its 1st arg.

We write $\nabla \ell(a, y)$ for the gradient of $a \mapsto \ell(a, y)$ at point a .

$$a_{i,\epsilon} = (0, 0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^k \\ = e_i$$

So, if $p_\epsilon \in \Delta(K)$ then $\ell\left(\sum_{i=1}^k p_{i,\epsilon} a_{i,\epsilon}, y_\epsilon\right) = \ell(p_\epsilon, y_\epsilon)$

(stated otherwise, a pt $p_\epsilon \in \Delta(K)$ can be seen as a convex combinat^o of the corners e_i of $\Delta(K)$)

Algorithm ("EG": Exponential gradient)

• parameter $\eta > 0$

• $\forall t \geq 1,$

$$p_{i,t} := \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \nabla_i \ell(p_s, g_s)\right)}{\sum_{i=1}^K \exp\left(-\eta \sum_{s=1}^{t-1} \nabla_i \ell(p_s, g_s)\right)}, \quad 1 \leq i \leq K$$

where ∇_i is the i -th coordinate of the gradient.

Theorem (Cesa-Bianchi 1999)

Let $T \geq 1$ and $K \geq 2$

Assume $\ell: \Delta(K) \times \mathcal{Y} \rightarrow \mathbb{R}$ is differentiable in $\Delta(K)$, and $\|\nabla \ell\|_{\infty} \leq G$, then the EG algorithm tuned with $\eta > 0$ satisfies, $\forall (g_t)_{t=1}^T$ in \mathcal{Y} .

$$\begin{aligned} \text{Reg}_T^{\text{conv}} &= \sum_{t=1}^T \ell(p_t, g_t) - \inf_{w \in \Delta(K)} \sum_{t=1}^T \ell(w, g_t) \\ &\leq \frac{\ln K}{\eta} + \frac{T \eta G^2}{2} \end{aligned}$$

In particular choosing $\eta = \frac{1}{G} \sqrt{\frac{2 \ln K}{T}}$ leads to

$$\text{Reg}_T^{\text{conv}} \leq G \sqrt{2 T \ln K}$$

NB: here the goal was more ambitious: instead of targeting $\min_{1 \leq i \leq K} \sum_{t=1}^T \ell(e_i, g_t)$ we compete against the best convex combinat^o of the e_i , $1 \leq i \leq K$, ie we target $\inf_{w \in \Delta(K)} \sum_{t=1}^T \ell(w, g_t)$

which is always \leq (and sometimes \ll) than

$$\min_{1 \leq i \leq K} \sum_{t=1}^T \ell(e_i, g_t)$$

(take w in $\{e_i: 1 \leq i \leq K\}$)

Proof: uses the previous proof!
 By convexity of $w \mapsto \ell(w, y_\epsilon)$

$$\ell(p_\epsilon, y_\epsilon) - \ell(w, y_\epsilon) \leq \underbrace{\nabla \ell(p_\epsilon, y_\epsilon)}_{\ell'_\epsilon} \cdot (p_\epsilon - w)$$

(We define $\ell'_\epsilon := \nabla \ell(p_\epsilon, y_\epsilon)$)

$$\leq p_\epsilon \cdot \ell'_\epsilon - w \cdot \ell'_\epsilon$$

Therefore,

$$\begin{aligned} & \sum_{\epsilon=1}^T \ell(p_\epsilon, y_\epsilon) - \inf_{w \in \Delta(k)} \sum_{\epsilon=1}^T \ell(w, y_\epsilon) \\ & \leq \max_{w \in \Delta(k)} \left\{ \sum_{\epsilon=1}^T p_\epsilon \cdot \ell'_\epsilon - \sum_{\epsilon=1}^T w \cdot \ell'_\epsilon \right\} \\ & = \sum_{\epsilon=1}^T p_\epsilon \cdot \ell'_\epsilon - \min_{w \in \Delta(k)} \left\{ w \cdot \sum_{\epsilon=1}^T \ell'_\epsilon \right\} \\ & = \sum_{\epsilon=1}^T p_\epsilon \cdot \ell'_\epsilon - \min_{1 \leq i \leq k} \sum_{\epsilon=1}^T \ell'_{i, \epsilon} \end{aligned}$$

by (1) $\leq \frac{rk}{m} + \frac{mT(B_2 - B_1)^2}{8}$ ($w \cdot \sum_{\epsilon=1}^T \ell'_\epsilon$ attains its max on corner of $\Delta(k)$ i.e. on e_i)

because, by definit^o of EG algo, $p_{i, \epsilon} \propto \exp(-m \sum_{s=1}^{\epsilon-1} p_{i, s})$ corresponds to EWA algo used with

$$\ell_{i, \epsilon} = \nabla_i \ell(p_\epsilon, y_\epsilon)$$

$$\|\ell_{i, \epsilon}\| = \|\nabla_i \ell(p_\epsilon, y_\epsilon)\| \leq \|\nabla \ell\|_\infty \leq 6$$

$$\text{so } B_2 - B_1 \leq 26$$

$$\leq \frac{rk}{m} + \frac{mT6^2}{2}$$

□

Take-home msg: Thanks to convexity and by using the gradients of ℓ , we were able to compete against a larger set.

Exercise: Generalize the algo and the thm when

- D is any convex set
- $a_{i,t} \in D$ arbitrary

$$\text{Reg}_T^{\text{conv}} = \sum_{t=1}^T \ell(a_t^*, y_t) - \inf_{w \in \Delta(K)} \sum_{t=1}^T \ell\left(\sum_{i=1}^K w_i a_{i,t}, y_t\right)$$

Links with gradient descent

The EG algorithm can be equivalently written as $\forall t \geq 1$,

(i) $q_t = (q_{i,t})_{1 \leq i \leq K} \in \mathbb{R}_+^K$ st

$$q_{i,t} = p_{i,t-1} e^{-\eta \nabla_x \ell(p_{t-1}, y_{t-1})}$$

(ii) $p_t = \frac{q_t}{\sum_{i=1}^K q_{i,t}} \in \Delta(K)$

Note that:

(i) is a gradient descent applied to the log-likelihood.

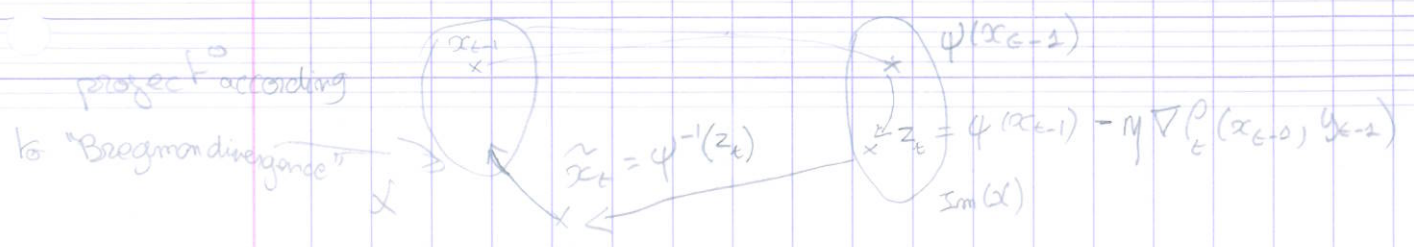
$$\ln(q_{i,t}) = \ln(p_{i,t-1}) - \eta \nabla_x \ell(p_{t-1}, y_{t-1})$$

(ii) is a project^o step

$$q_t \in \underset{p \in \Delta(K)}{\text{argmin}} \text{KL}(p, q_t) \quad \leftarrow \sum_{i=1}^K p_i \ln\left(\frac{p_i}{q_{i,t}}\right)$$

Exercise \uparrow

Therefore, the EG algorithm is a special case of "mirror descent"



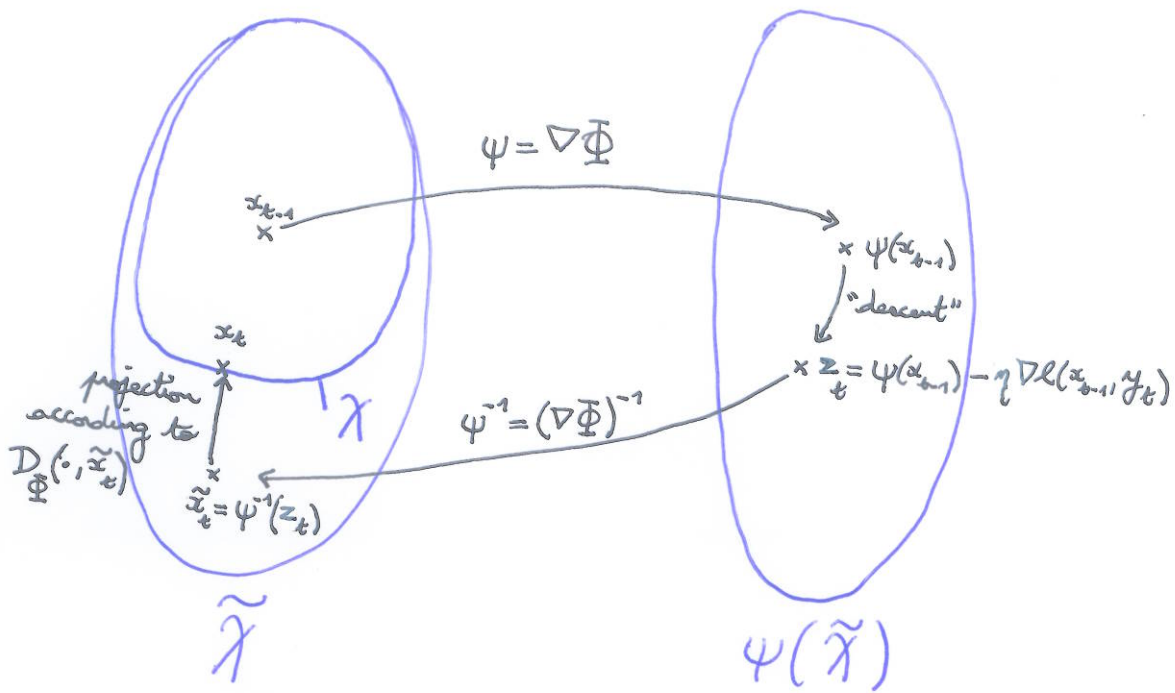


Fig: Goline Mirror Descent

In the above $\mathcal{J} : \mathcal{X} = \Delta(k)$

$$\psi : \mathcal{X} \rightarrow \mathbb{R}_+^k$$

$$(x_i)_{1 \leq i \leq k} \mapsto (\ln(x_i))_{1 \leq i \leq k}$$

• gradient step: $\ln p_{t+1} = \ln p_{t,\epsilon} - \eta \nabla_x \ell(p_t, y_t)$

• $\psi : e^{\cdot}$

More generally:
"online mirror descent"

Potential $\bar{\Phi} : \tilde{\mathcal{X}} \rightarrow \mathbb{R}$ and convex + differentiable on $\tilde{\mathcal{X}}$

$$\psi = \nabla \bar{\Phi}$$

Bregman divergence: the fct

$$D_{\bar{\Phi}} : \tilde{\mathcal{X}} \times \tilde{\mathcal{X}} \rightarrow \mathbb{R}$$

$$(y, x) \mapsto \bar{\Phi}(y) - \bar{\Phi}(x) - \nabla \bar{\Phi}(x) \cdot (y - x)$$

(difference between $\bar{\Phi}(y)$ and its 1st order Taylor approx around x)
 $\rightarrow \geq 0$ since $\bar{\Phi}$ convex

Examples: 1) $\tilde{\mathcal{X}} = \mathbb{R}_+^k$ $\bar{\Phi}(x) = \sum_{i=1}^k x_i (\ln x_i - 1)$

$$D_{\bar{\Phi}}(y, x) = \sum_{i=1}^k y_i \ln \frac{y_i}{x_i} - \sum_{i=1}^k (y_i - x_i) ; \nabla \bar{\Phi}(x) = (\ln x_i)_{1 \leq i \leq k}$$

2) $\mathcal{X} = \mathbb{R}^k$ $\bar{\Phi}(x) = \frac{1}{2} \|x\|_2^2$, $\psi(x) = \nabla \bar{\Phi}(x) = x$

$$D_{\bar{\Phi}}(y, x) = \frac{1}{2} \|y\|_2^2 - \frac{1}{2} \|x\|_2^2 - x \cdot (y - x)$$

$$= \frac{1}{2} \|y - x\|_2^2$$

gradient step here is the classical gradient step

Project^o step: Euclidian project onto \mathcal{X} .

Good news: \exists theorem controlling the regret of mirror descent for many possible potential fct $\bar{\Phi}$.

III The non-convex case

See notes.

D is not assumed to be convex.

How can we choose \hat{a}_t as a convex combination of the $a_{i,t}$?
doesn't exist!

Instead, we can use a randomized algorithm:

$\forall t \geq 1$

• compute $p_t \in \Delta(k)$ as in EWA

• instead of defining $\hat{a}_t = \sum_{i=1}^k p_{i,t} a_{i,t}$ we draw

$I_t \sim p_t$ (conditionally on the past) and predict $\hat{a}_t = a_{I_t, t}$

In expectation, everything is the same!

$$\begin{aligned} \mathbb{E} \left[\ell(a_{I_t, t}, g_t) \mid \underbrace{I_1, \dots, I_{t-1}}_{\text{past}} \right] \\ = \sum_{i=1}^k p_{i,t} \ell(a_{i,t}, g_t) = p_t \cdot \ell_t \end{aligned}$$

Therefore, the inequality (1)

$$\begin{aligned} \sum_{t=1}^T p_t \cdot \ell_t - \min_{1 \leq i \leq k} \sum_{t=1}^T \ell(a_{i,t}, g_t) \\ \leq \frac{\ln k}{\eta} + \frac{T\eta (B_2 - B_1)^2}{8} \end{aligned}$$

implies here that

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \ell(a_{I_t, t}, g_t) \right] - \min_{1 \leq i \leq k} \sum_{t=1}^T \ell(a_{i,t}, g_t) \\ \leq \frac{\ln k}{\eta} + \frac{T\eta (B_2 - B_1)^2}{8} \end{aligned}$$

assuming that $\ell: D \times Y \rightarrow [B_1, B_2]$ (no convexity) and if $a_{i,t}$ and g_t are deterministic.