

Sequential optimization for general convex sets.

We so far focused our attention on the simplex

$$\mathcal{X} = \{(\mathbf{p}_1, \dots, \mathbf{p}_N) : \forall j, p_j \geq 0 \text{ and } \sum_k p_k = 1\} \subset \mathbb{R}^N$$

Let's now deal with more general convex sets $\mathcal{C} \subset \mathbb{R}^N$, which we assume to be CLOSED.

Protocol: At each round $t=1, 2, \dots$

1. The statistician picks $x_t \in \mathcal{C}$ while the opponent picks simultaneously a convex and differentiable function $l_t : \mathcal{C} \rightarrow \mathbb{R}$
2. x_t and l_t are revealed

Aim: Control $R_T = \sum_{t=1}^T l_t(x_t) - \inf_{x \in \mathcal{C}} \sum_{t=1}^T l_t(x)$ (well, the inf is a min here)

Algorithm: Online gradient descent with fixed learning rate $\eta > 0$



(called OGDX(η))

- Play any $x_1 \in \mathcal{C}$
- For $t=2, 3, \dots,$

$$x_t = \text{Proj}_{\mathcal{C}}(x_{t-1} - \eta \nabla l_{t-1}(x_{t-1}))$$

where $\text{Proj}_{\mathcal{C}}$ is the Euclidean projection onto \mathcal{C} (well defined because \mathcal{C} is closed).

$\|\cdot\|$ is the Euclidean norm.

Theorem: Assume that $\max_{t \leq T} \sup_{x \in \mathcal{C}} \|\nabla l_t(x)\| \leq G$ [bounded gradients]

and that \mathcal{C} is compact, i.e.,

$$\sup_{x, x' \in \mathcal{C}} \|x - x'\| \leq D$$
 [bounded diameter]

Then OGDX(η) is such that

$$\sum_{t=1}^T l_t(x_t) - \min_{x \in \mathcal{C}} \sum_{t=1}^T l_t(x) \leq \frac{D^2}{2\eta} + \frac{\eta G^2 T}{2}$$

↑
inf achieved by continuity + compactness

In particular, for $\eta = \frac{D}{G\sqrt{T}}$, the bound equals $D\sqrt{T}$.

Remark:

In the case with EWA, we were considering a linear function:

$$\nabla p \in X, \quad l_t(p) = \sum_{j=1}^N p_j l_{jt}$$

↳ OG would also be applicable in this setting, though it leads to a suboptimal bound:

diameter of X :

$$D = 2$$

$$\forall p, q \in X, \quad \sum_j (p_j - q_j)^2$$

$$\leq \sum_j |p_j - q_j| \leq \sum_j (p_j + q_j) = 2$$

equality achieved if p, q Dirac masses at $i \neq k$

bound on the gradients: if, e.g., $l_{jt} \in [-M, M] \quad \forall j, t$,

$$\text{then } \nabla l_t(p) = \begin{pmatrix} l_{1t} \\ l_{2t} \\ \vdots \\ l_{Nt} \end{pmatrix}$$

$$\text{and } \|\nabla l_t(p)\|^2 \leq N M^2 = G^2$$

We get the bounds:

- for OG: $D\sqrt{T} = M\sqrt{2NT}$

- for EWA: $(M - (-M))^2 \sqrt{\frac{T}{2} \ln N} = M\sqrt{8T \ln N}$

The dependency in N is suboptimal for OG.

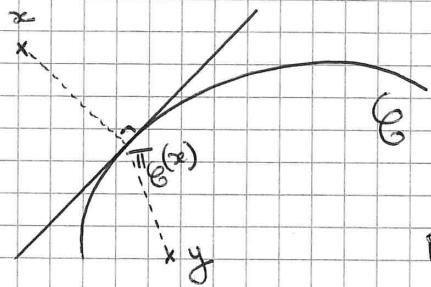
Proof (of the theorem): By convexity (« inégalité des pentes », in French):

$$\forall x \in \mathcal{E}, \quad \forall t, \quad l_t(x_t) - l_t(x) \leq \nabla l_t(x_t) \cdot (x_t - x)$$

Summing over t ,

$$\begin{aligned} R_T &= \sup_{x \in \mathcal{E}} \left\{ \sum_{t=1}^T l_t(x_t) - \sum_{t=1}^T l_t(x) \right\} \\ &\leq \sup_{x \in \mathcal{E}} \left\{ \sum_{t=1}^T \nabla l_t(x_t) \cdot (x_t - x) \right\} \end{aligned}$$

Now, we will need some facts about projections on CLOSED convex sets:



Def: $\Pi_G(x)$ is the (unique) $y^* \in G$ with $\|x - y^*\| \leq \|x - y\| \quad \forall y \in G$

Property 1: if $x \notin G$, then $\forall y \in G, (y - \Pi_G(x)) \cdot (x - \Pi_G(x)) \leq 0$

Property 2: $\forall y \in G, \|y - \Pi_G(x)\| \leq \|y - x\|$

Proof: (Property 1): $t \in [0,1], y_t = \Pi_G(x) + t(y - \Pi_G(x)) \in G$ by convexity
 $\|x - y_t\|^2 = \|x - \Pi_G(x)\|^2 + t^2 \|y - \Pi_G(x)\|^2 - 2t(y - \Pi_G(x)) \cdot (x - \Pi_G(x)) \geq \|x - \Pi_G(x)\|^2$
 Simplify the $\|x - \Pi_G(x)\|^2$, divide by t and $t \rightarrow 0$.

(Property 2) $\|y - x\|^2 = \|y - \Pi_G(x) + \Pi_G(x) - x\|^2 = \|y - \Pi_G(x)\|^2 + \|x - \Pi_G(x)\|^2 \leftarrow$ both are ≥ 0
 $+ 2(y - \Pi_G(x)) \cdot (\Pi_G(x) - x) \leftarrow$ (by Prop 1)

In particular, for $t \geq 1$:

$$\begin{aligned} \forall x \in \mathcal{E}, \quad \|x_{t+1} - x\|^2 &= \|\Pi_G(x_t - \eta \nabla l_t(x_t)) - x\|^2 \\ &\stackrel{\text{Property 1}}{\leq} \|x_t - \eta \nabla l_t(x_t) - x\|^2 \\ &= \|x_t - x\|^2 + \eta^2 \|\nabla l_t(x_t)\|^2 - 2\eta \nabla l_t(x_t) \cdot (x_t - x) \end{aligned}$$

That is, $\forall t, \forall x \in \mathcal{E}$,

$$\begin{aligned} &\nabla l_t(x_t) \cdot (x_t - x) \\ &\leq \frac{1}{2\eta} \left(\|x_t - x\|^2 - \|x_{t+1} - x\|^2 \right) + \frac{\eta}{2} \underbrace{\|\nabla l_t(x_t)\|^2}_{\leq G^2} \\ &\quad \text{by the assumption of bounded gradients} \end{aligned}$$

Summing over t :

$$\begin{aligned} &\sum_{t=1}^T \nabla l_t(x_t) \cdot (x_t - x) \\ &\leq \frac{1}{2\eta} \left(\underbrace{\|x_1 - x\|^2}_{\leq D^2} - \underbrace{\|x_{T+1} - x\|^2}_{\geq 0} \right) + \frac{\eta}{2} T G^2 \end{aligned}$$

Adaptation to time T (and other parameters) / for CGD and EWA

The algorithms discussed so far:

EWA with fixed learning rate $\eta > 0$
 CGD with fixed step size $\eta > 0$

required the knowledge of T and of

the range $[m, M]$ of the losses for EWA
 the bound G on the gradients of the losses for CGD

We now study adaptation to these parameters.

There exists a general trick, called the doubling trick which we will detail in 4 pages from now for EWA.

Another popular way of being adaptive is to tune η according to the past observations:

Adaptive CGD:

- Play an arbitrary $x_1 \in \mathcal{E}$, pick $\eta_1 > 0$
- For $t \geq 2$, play $x_t = \text{Proj}_{\mathcal{E}}(x_{t-1} - \eta_{t-1} \nabla l_{t-1}(x_{t-1}))$
 pick $\eta_t > 0$

Theorem: If Adaptive CGD picks a non-increasing sequence $(\eta_t)_t$ with $\eta_t > 0 \forall t$,
 then: $R_T \leq \frac{D^2}{2\eta_1} + \frac{1}{2} \sum_{t=1}^T \eta_t \|\nabla l_t(x_t)\|^2$ with the same notation as in the Theorem for CGD(η)

Corollary: Picking $\eta_t = \begin{cases} 1 & \text{if } V_t = \sum_{s=1}^t \|\nabla l_s(x_s)\|^2 < D^2 \\ D/\sqrt{V_t} & \text{if } V_t \geq D^2 \end{cases}$

Adaptive CGD achieves the bound $\max \left\{ D^2, \frac{3}{2} D \sqrt{\sum_{t=1}^T \|\nabla l_t(x_t)\|^2} \right\}$

without any prior knowledge of neither T nor G

(of course, \mathcal{E} and D are known).

$$= O(DG\sqrt{T})$$

Proof of the theorem and of the corollary:

Th: We extract from the proof of OGD(η) that

$$\forall t \in \mathbb{C}, \quad \nabla l_t(x_t) \cdot (x_t - x) \leq \frac{1}{2\eta_t} (\|x_t - x\|^2 - \|x_{t+1} - x\|^2) + \frac{\eta_t}{2} \|\nabla l_t(x_t)\|^2$$

Summing over t :

$$\begin{aligned} \forall t \in \mathbb{C}, \quad \sum_{t=1}^T l_t(x_t) - \sum_{t=1}^T l_t(x) &\leq \sum_{t=1}^T \nabla l_t(x_t) \cdot (x_t - x) \\ &\leq \underbrace{\sum_{t=1}^T \frac{1}{2\eta_t} (\|x_t - x\|^2 - \|x_{t+1} - x\|^2)}_{\leftarrow \text{We need only to show that this sum is}} + \frac{1}{2} \sum_{t=1}^T \eta_t \|\nabla l_t(x_t)\|^2 \\ &\leq \frac{D^2}{2\eta_T} \end{aligned}$$

Abel's transform:

$$\begin{aligned} &\sum_{t=1}^T \frac{1}{2\eta_t} (\|x_t - x\|^2 - \|x_{t+1} - x\|^2) \\ &= \underbrace{\frac{1}{2\eta_1} \|x_1 - x\|^2}_{\leq D^2} + \sum_{t=2}^T \left(\underbrace{\frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}}}_{\geq 0 \text{ because } (\eta_t) \text{ is non-increasing}} \right) \underbrace{\|\nabla l_t(x_t)\|^2}_{\leq D^2} - \underbrace{\frac{1}{2\eta_T} \|x_{T+1} - x\|^2}_{\leq 0} \\ &\leq \frac{D^2}{2\eta_T} \text{ after telescoping} \end{aligned}$$

Cor: Recall that we picked:

⚠ Pay attention to the indexing: η_t is used to pick x_{t+1} and may thus depend on the past till and including round t .

$$\forall t \geq 1, \quad \eta_t = \begin{cases} 1 & \text{if } V_t = \sum_{s=1}^t \|\nabla l_s(x_s)\|^2 < D^2 \\ \frac{D}{\sqrt{V_t}} & \text{if } V_t \geq D^2 \end{cases}$$

(Some practice actually explains why this is a natural choice \rightarrow we replace the theoretically optimal but

impossible in practice choice $\eta_t = D/V_T$ by its “sequential estimation”;
and use a threshold to avoid too large values of η_t for small t .)

The regret bound reads

$$\frac{D^2}{2\eta_T} + \frac{1}{2} \sum_{\substack{t: \\ V_t < D^2}} \|v\ell_t(x_t)\|^2$$

$$+ \frac{1}{2} \sum_{\substack{t: \\ V_t > D^2}} \frac{D}{\sqrt{V_t}} \|v\ell_t(x_t)\|^2$$

We denote G the largest $t \leq T$ such that $V_t < D^2$. Then, if $G < T$, we have $V_G \geq D^2$.

The above bound can be rewritten as:

$$\frac{D^2}{2\eta_T} + \underbrace{\frac{1}{2} \sum_{\substack{t: \\ V_t < D^2}} \|v\ell_t(x_t)\|^2}_{= \frac{V_G}{2} < D^2/2} + \frac{1}{2} \sum_{t=G+1}^T \frac{D}{\sqrt{V_t}} \|v\ell_t(x_t)\|^2$$

possibly valid sum
(if $G = T$)

$$= \underbrace{\frac{1}{2} \sum_{t=G+1}^T \frac{D}{\sqrt{V_t}} (V_t - V_{t-1})}_{\leq D \sum_{t=G+1}^T (\sqrt{V_t} - \sqrt{V_{t-1}})} + \underbrace{\frac{D}{\sqrt{V_T}} (\sqrt{V_T} + \sqrt{V_{G+1}})(\sqrt{V_G} - \sqrt{V_{G+1}})}_{\leq 2\sqrt{V_G}}$$

$$\leq D \sum_{t=G+1}^T (\sqrt{V_t} - \sqrt{V_{t-1}})$$

$$= D(\sqrt{V_T} - \sqrt{V_G})$$

We conclude by distinguishing two cases:

* $G = T$: $\eta_T = 1$ and the bound is $\leq \frac{D^2}{2} + \frac{D^2}{2} = D^2$

* $G < T$: $\eta_T = D/\sqrt{V_T}$ and the bound is

$$\frac{D^2}{2D\sqrt{V_T}} + \frac{V_G}{2} + D\sqrt{V_T} - D\sqrt{V_G} \leq \frac{3D\sqrt{V_T}}{2}$$

$\leq \frac{D\sqrt{V_G}}{2}$ as $V_G < D^2$

↳ A readable and simple final bound is thus the one we indicated:

$$\max \left\{ D^2, \frac{3}{2} D \sqrt{V_T} \right\} = \max \left\{ D^2, \frac{3}{2} D \sqrt{\sum_{t=1}^T \|v\ell_t(x_t)\|^2} \right\}$$

Calibration of η for the EWA strategy.

Reminder: by considering subgradients and pseudo-bases, we can reduce the case of convex losses to the one of linear losses. It thus suffices to explain how to sequentially calibrate η in the latter case.

Setting: At each round $t = 1, 2, \dots, T$,

- The statistician and the opponent simultaneously pick $p_t \in \mathcal{X}$ and $\ell_t = (\ell_{1t}, \dots, \ell_{Nt}) \in [m, M]^N$
- p_t and ℓ_t are publicly revealed

Aim: Control the regret $R_T = \sum_{t=1}^T \sum_{j=1}^N p_{jt} \ell_{jt} - \min_{k=1, \dots, N} \sum_{t=1}^T \ell_{kt}$

Algorithm: EWA with fixed learning rate $\eta > 0$

$$\forall t \geq 1, \quad p_t(\eta) = \exp(-\eta \sum_{k=1}^{t-1} \ell_{ks}) / \sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} \ell_{ks})$$

Performance bound: $R_T \leq \frac{\ln N}{\eta} + \eta \frac{(M-m)^2}{8} T$

For the choice $\eta^* = \frac{1}{M-m} \sqrt{\frac{8 \ln N}{T}}$ we get $R_T \leq \sqrt{\frac{1}{2} \ln N}$

but there are two issues with this choice:

1. T, m and M are not always known in advance
2. This η^* has often a poor performance in practice

Solutions for 1. : “doubling trick” and $p_t(\eta_t)$
 where η_t varies over time

2. : $p_t(\eta_t)$ again, but with different choices for η_t

Calibration of η for the EWA strategy: much harder!

The cheap way: the doubling trick.

Suppose m and M are known (we'll see how to deal with m and M unknown later on). Then: for each $r = 1, 2, \dots$

RESTART EWA with $\eta_r = \frac{1}{M-m} \sqrt{\frac{8 \ln N}{2^r}}$ for rounds $t = 2^r+1, \dots, 2^{r+1}$

(we deal with rounds $t=1$ and $t=2$ by picking uniform weights)

Fix $T \geq 2$: Denote by r_T the smallest $r \geq 1$ such that $T \leq 2^{r+1}$

(in particular, $2^{r_T} < T$)

We have

$$\begin{aligned}
 R_T &= \sum_{t=1}^T \sum_{j=1}^N p_j t l_{jt} - \min_k \sum_{t=1}^T l_{kt} \\
 &\leq 2(M-m) + \sum_{r=1}^{r_T-1} \left(\sum_{t=2^r+1}^{2^{r+1}} \sum_j p_j t l_{jt} - \min_k \sum_{t=2^r+1}^{2^{r+1}} l_{kt} \right) \\
 &\quad + \sum_{t=2^{r_T}+1}^T \sum_j p_j t l_{jt} - \min_k \sum_{t=2^{r_T}+1}^T l_{kt} \\
 &\quad \text{↑ performance bound of EWA with } \eta_{r_T} : \\
 &\quad \text{is } \leq \frac{\ln N}{\eta_{r_T}} + \eta_{r_T} (M-m)^2 \frac{(T-2^{r_T})}{8} \\
 &\quad \leq \frac{\ln N}{\eta_{r_T}} + \eta_{r_T} (M-m)^2 \frac{2^{r_T}}{8} \\
 &\quad = (M-m) \sqrt{\frac{2^{r_T}}{2} \ln N}
 \end{aligned}$$

max of a sum \leq sum of the max.

We have proved:

$$\begin{aligned}
 R_T &\leq 2(M-m) + \left(\sum_{r=1}^{r_T} \sqrt{2^r} \right) \sqrt{\frac{\ln N}{2}} \times (M-m) \\
 &= \sqrt{2} \times \frac{(\sqrt{2})^{r_T} - 1}{\sqrt{2} - 1}
 \end{aligned}$$

Since $2^{r_T} < T$, we have $(\sqrt{2})^{r_T} < \sqrt{T}$

The final bound reads :

$$R_T \leq \frac{(M-m)}{(\sqrt{2}-1)} \sqrt{T \ln N} + 2(M-m)$$

measures the price
for the adaptivity in T if
one uses the doubling trick

the main difference w.r.t. bound
when T is known is an additional
 $\sqrt{2}/(\sqrt{2}-1) \approx 3.42$

NOTE: An extension of the above argument is possible when m and M are also unknown : we take fresh starts whenever "t doubles" or whenever the current "estimates" of m and M are exceeded (and in the fresh start, their value is doubled). But it's messy to write : I don't even dare giving it to you as an exercise... because I would need to write the solution !

Here's a more elegant way to cope with the problem, as we did for OGJ :

The smarter way: post-data adaptive choice of η

Let $p_t(\eta) = (p_{1t}(\eta), \dots, p_{Nt}(\eta))$ denote the vector of weights recommended by ENA with fixed $\eta > 0$ at round t : $p_{jt}(\eta) = \exp(-\eta \sum_{s=1}^t l_{js}) / \sum_{i=1}^N \exp(-\eta \sum_{s=1}^t l_{is})$.

We now consider a rule to select $\eta_t > 0$ based on the post-information $l_{j,t-1}, j \in \{1 \dots N\}$, and use the weights $p_t(\eta_t)$ at round t .
 (Note that the choice of η_t is actually irrelevant...)

Lemma: If the selected η_t are non-increasing, then for all $f_{it} \in \mathbb{R}$ (not necessarily bounded),

$$\begin{aligned} R_T &= \sum_{t=1}^T \sum_{j=1}^N p_{jt} f_{jt} - \min_k \sum_{t=1}^T f_{kt} \\ &\leq \frac{\ln N}{\eta_T} + \sum_{t=1}^T s_t \quad \text{where} \\ s_t &= \sum_{j=1}^N p_{jt} f_{jt} + \frac{1}{\eta_t} \ln \left(\sum_{j=1}^N p_{jt} e^{-\eta_t f_{jt}} \right) \end{aligned}$$

Note: s_t was usually bounded via Hoeffding's inequality :

$$\text{if } f_{jt} \in [m_j, M_j] \forall j, \text{ then } s_t \leq \frac{\eta_t (M-m)^2}{8} \rightarrow \text{but we will also consider here other bounds.}$$

Exercise 1:

(easy: see
it as a
warm-up!)

If m and M are known, which η_t could you pick and which $(M-m)\sqrt{T \ln N}$ would you get? (ie, which constant would you get?).

Proof (of the lemma):

By definition of δ_t ,

$$\sum_{j=1}^N p_j e^{-\eta_t l_{jt}} = \delta_t - \frac{1}{\eta_t} \ln \left(\sum_j p_j e^{-\eta_t l_{jt}} \right)$$

The question is thus to control

$$\sum_{t=1}^T -\frac{1}{\eta_t} \ln \left(\sum_j p_j e^{-\eta_t l_{jt}} \right) \quad \uparrow \left[\leq \ln N + \min_k \sum_{t=1}^T \eta_t l_{kt} \right]$$

i.e. can we get some telescoping

An initial transformation is useful:

$$\eta_{t+1} \leq \eta_t \text{ so that } x \mapsto x^{\eta_t/\eta_{t+1}}$$

is convex:

$$\frac{1}{N} \sum_j p_j e^{-\eta_t l_{jt}} = \frac{1}{N} \sum_j \left(p_j e^{-\eta_t \eta_{t+1} l_{jt}} \right)^{\eta_t/\eta_{t+1}}$$

$$\geq \left(\frac{1}{N} \sum_j p_j e^{-\eta_t \eta_{t+1} l_{jt}} \right)^{\eta_t/\eta_{t+1}}$$

$$= \frac{1}{\eta_t/\eta_{t+1}} \left(\sum_j \frac{e^{-\eta_t \sum_{s=1}^{t-1} l_{js}}}{\left(\sum_k e^{-\eta_t \sum_{s=1}^{t-1} l_{ks}} \right)^{\eta_{t+1}/\eta_t}} e^{-\eta_t l_{jt}} \right)^{\eta_t/\eta_{t+1}}$$

$$= \frac{1}{N^{\eta_t/\eta_{t+1}}} \frac{\left(\sum_j \exp \left(-\eta_{t+1} \sum_{s=1}^{t-1} l_{js} \right) \right)^{\eta_t/\eta_{t+1}}}{\sum_k \exp \left(-\eta_t \sum_{s=1}^{t-1} l_{ks} \right)}$$

Therefore,

$$-\frac{1}{\eta_t} \ln \left(\sum_j p_j e^{-\eta_t l_{jt}} \right) \leq -\frac{1}{\eta_t} \ln \left(N^{\eta_t/\eta_{t+1}} \frac{\left(\sum_j e^{-\eta_{t+1} \sum_{s=1}^{t-1} l_{js}} \right)^{\eta_t/\eta_{t+1}}}{\sum_k e^{-\eta_t \sum_{s=1}^{t-1} l_{ks}}} \right)$$

Summing over $t=1, \dots, T$:

$$\sum_{t=1}^T -\frac{1}{\eta_t} \ln \left(\sum_j p_j e^{-\eta_t l_{jt}} \right) \leq \sum_{t=1}^T \left(-\frac{1}{\eta_t} \ln \left(N^{\eta_t/\eta_{t+1}} \right) - \frac{1}{\eta_t} \frac{\eta_t}{\eta_{t+1}} \ln \left(\sum_j e^{-\eta_{t+1} \sum_{s=1}^{t-1} l_{js}} \right) \right)$$

$$+ \frac{1}{\eta_T} \ln \left(\sum_k e^{-\eta_T \sum_{s=1}^{T-1} l_{ks}} \right)$$

$$= \sum_{t=1}^T \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln N - \frac{1}{\eta_{T+1}} \ln \left(\sum_j e^{-\eta_{T+1} \sum_{s=1}^{T-1} l_{js}} \right) + \frac{1}{\eta_1} \ln \left(\sum_k e^{-\eta_1 \times 0} \right)$$

$$\begin{aligned}
 &= \frac{\ln N}{\eta_{TH}} - \frac{1}{\eta_{TH}} \ln \left(\sum_j e^{-\eta_{TH} \sum_{s=1}^T l_{js}} \right) \\
 &\quad \leq \min_{k \leq N} \sum_{t=1}^T l_{kt}
 \end{aligned}$$

we can always assume that $\eta_{TH} = \eta_T$
 (if the rule picked another η_{TH} , just consider for the sake of the proof $\eta_{TH} = \eta_T$)

Comments: We did 2 things with this lemma:

- 1 - have η possibly depend on the past (as long as $\eta_t \rightarrow$ over time)
- 2 - consider sharper bounds than the ones obtained via Hoeffding's lemma: it suffices to focus one's attention on the δ_t .

Let's try to get sharper (second-order) bounds on the δ_t (Hoeffding's bound is a zero-order bound: something extremely uniform, that only depends on the scale $[m, M]$ and not on the concentration).

Intuition: $e^x \approx 1 + x + \frac{x^2}{2}$ when x is small

$$\begin{aligned}
 &\text{thus if } \eta_t \text{ is small} \quad -\frac{1}{\eta_t} \ln \left(\sum_j p_{jt} e^{-\eta_t l_{jt}} \right) \approx -\frac{1}{\eta_t} \ln \left(1 - \sum_j \eta_t p_{jt} l_{jt} \right. \\
 &\quad \left. + \frac{1}{2} \sum_j p_{jt} (\eta_t l_{jt})^2 \right) \\
 &\ln(1+u) \approx u \text{ when } u \text{ is small} \\
 &\approx \sum_j p_{jt} l_{jt} - \frac{\eta_t}{2} \sum_j p_{jt} l_{jt}^2
 \end{aligned}$$

$$\text{That is, } \delta_t \approx \frac{\eta_t}{2} \sum_j p_{jt} l_{jt}^2$$

Actually, we could do better:

$$\begin{aligned} \delta_t &= -\frac{1}{\eta_t} \ln \left(\sum_j p_{jt} e^{-\eta_t (l_{jt} - \sum_k p_{kt} l_{kt})} \right) \\ &\stackrel{\eta_t \text{ large}}{\approx} -\frac{1}{\eta_t} \ln \left(1 - \sum_j p_{jt} \left(\eta_t (l_{jt} - \sum_k p_{kt} l_{kt}) \right) + \frac{\eta_t^2}{2} \sum_j p_{jt} \left(l_{jt} - \sum_k p_{kt} l_{kt} \right)^2 \right) \\ &\approx \frac{\eta_t}{2} \underbrace{\sum_j p_{jt} \left(l_{jt} - \sum_k p_{kt} l_{kt} \right)^2}_{\text{not. } v_t \text{ (a variance-like term)}} = \frac{\eta_t}{2} v_t \end{aligned}$$

The regret bound is $\approx \frac{\ln N}{\eta_T} + \sum_{t=1}^T \frac{\eta_t}{2} v_t$

Choosing $\eta_t \approx \frac{\sqrt{\ln N}}{\sqrt{v_t}}$ where $v_{t-1} = \sum_{s=1}^{t-1} v_s$

of $\Omega(\sqrt{N})$
case
we're interested in
the max
size
of
the
subset
in
 \mathcal{B}_t

we have $\sum \eta_t v_t \approx \sqrt{\ln N} \sum_t \frac{1}{\sqrt{v_t}} (v_t - v_{t-1})$

$$= \sqrt{\ln N} \sum_t \frac{1}{\sqrt{v_{t-1}}} \underbrace{(\sqrt{v_t} + \sqrt{v_{t-1}})(\sqrt{v_t} - \sqrt{v_{t-1}})}_{> 0}$$

After telescoping $\sum_t \eta_t v_t \leq C \sqrt{v_T \ln N}$

and the final regret bound is $\leq C' \sqrt{v_T \ln N}$

↳ For a decade (2005-2014) this proof was somewhat wobbly and suboptimal, and led to a non-homogeneous bound.

Recently, in 2014, de Roon, van Erven, Grünwald & Koolen could at last make it 100% precise and homogeneous.

Their key argument to link δ_t and η_t, v_t is Bernstein's inequality.

Lemma [Bernstein's inequality]: X random variable taking values in $[m, M]$
 Then: $\forall \eta > 0$, $\ln E[e^{\eta X}] \leq \eta E[X] + \frac{(e^{\eta(M-m)} - 1 - \eta(M-m)) \text{Var}(X)}{(M-m)^2}$

Proof: $\psi: x \in \mathbb{R} \mapsto (e^x - x - 1)/x^2$ is increasing over \mathbb{R}
 $\eta > 0$, $X \leq M$ and $E[X] \geq m$ so that $\eta(X - E[X]) \leq \eta(M-m)$
 then $\psi(\eta(X - E[X])) \leq \psi(\eta(M-m))$

That is, $e^{\eta(X - E[X])} - \eta(X - E[X]) - 1 \leq \psi^2(X - E[X]^2) \frac{e^{\eta(M-m)} - \eta(M-m) - 1}{\eta^2(M-m)^2}$

Taking expectations: $E[e^{\eta X}] \times e^{-\eta E[X]} \leq 1 + \frac{e^{\eta(M-m)} - \eta(M-m) - 1}{(M-m)^2} \text{Var}(X)$

The proof is concluded by taking logarithms and using $\ln(u) \leq u - 1$.

Algorithm: EWA with at each round $t \geq 1$ $\eta_t = \frac{\ln N}{\sum_{s=1}^{t-1} \delta_s}$ with the convention that $\frac{\ln 0}{0} = +\infty$

Note: $p_1 = (\frac{1}{N}, \dots, \frac{1}{N})$ and $\eta_1 = +\infty$; since the value of η_1 is irrelevant this merely means that $\eta_1 \geq \eta_2$ and

$$\begin{aligned} \delta_1 &= \frac{1}{N} \sum_j l_{j1} + \frac{1}{\eta_1} \ln \left(\sum_j \frac{1}{N} e^{-\eta_1 l_{j1}} \right) \\ &= \frac{1}{N} \sum_j l_{j1} - \min_k l_{k1} \end{aligned}$$

$\delta_j > 0$ except if all losses are equal, in which case no regret is suffered and the indexing of time could well start at 2...

With no loss of generality we will assume $\delta_1 > 0$.

Theorem: The above strategy ensures that:

$$\forall m \leq M, \quad \forall l_{it} \in [m, M], \quad R_T = \sum_{t,i,j} p_{jt} l_{it} - \min_k \sum_{t=1}^T l_{kt}$$

$$\leq 2 \sqrt{\sum_{t=1}^T \sum_j (p_{jt} (l_{jt} - \sum_k p_{kt} l_{kt}))^2 \ln N} + (M-m) \left(2 + \frac{4}{3} \ln N \right)$$

Proof: $\eta_t \downarrow$ thus $R_T \leq \frac{\ln N}{\eta_T} + \sum_{t=1}^T \delta_t$ ($\eta_1 = +\infty$ means that we do here $\eta_T \rightarrow +\infty$)

By Bernstein's inequality: $\delta_t \leq \frac{1}{\eta_t} \frac{e^{\eta_t(M-m)} - e^{\eta_t(1-m)-1}}{(M-m)^2} v_t$ for $t \geq 2$

(with $X_t = -\ell_{jt}$ with probability p_{jt} and $\eta_t > 0$)

$$v_t = \sum_j (\ell_{jt} - \frac{1}{n} \sum_k \ell_{kt})^2$$

Equivalently: $v_t \geq \frac{\eta_t(M-m)}{e^{\eta_t(M-m)} - e^{\eta_t(M-m)-1}} (M-m) \delta_t$ for $t \geq 2$

$$= \frac{2}{\eta_t(M-m)} (M-m) \delta_t = (M-m) \delta_t \times f(\eta_t(M-m))$$

where $f(x) = \frac{2}{x} - \frac{x}{e^x - x - 1}$ for $x \geq 0$

One can show that f is decreasing, so that $f(x) \leq f(0) = \frac{2}{3}$ $\forall x \geq 0$

Therefore, $v_t \geq \frac{2\delta_t}{\eta_t} + \frac{2}{3}(M-m)\delta_t$ for $t \geq 2$, but also valid for $t=1$ as $\eta_1 \rightarrow +\infty$

Then,
$$\left(\sum_{t=1}^T \delta_t \right)^2 = \sum_{t=1}^T \left(\left(\sum_{s=1}^t \delta_s \right)^2 - \left(\sum_{s=1}^{t-1} \delta_s \right)^2 \right)$$

↑
telescoping
 $= (a+b)(a-b)$

$$= \sum_{t=1}^T \delta_t \left(\delta_t + 2 \sum_{s=1}^{t-1} \delta_s \right)$$

$\underbrace{\delta_t}_{\text{direct bound: this is } \leq M-m} + \underbrace{2 \sum_{s=1}^{t-1} \delta_s}_{\text{by def of alg, this}} = \frac{\ln N}{\eta_t}, \text{ even for } t=1$

$$\leq \sum_{t=1}^T \left((M-m) \delta_t + \frac{2\delta_t \ln N}{\eta_t} \right)$$

$$\leq v_T \ln N + \frac{2}{3} (M-m) \delta_T \ln N$$

Summarizing:
$$\left(\sum_{t=1}^T \delta_t \right)^2 \leq V_T \ln N + (M-m) \left(1 + \frac{2}{3} \ln N \right) \sum_{t=1}^T \delta_t$$

↑
reminder:
 $V_T = \sum_{t \leq T} v_t$

This is a 2nd-order inequality: $x^2 \leq b + ax \Rightarrow x \leq a + \sqrt{b}$
for $a, b \geq 0$

[Indeed, α is then smaller than the larger root of the 2nd-order polynomial $u \mapsto u^2 - au - b$, that is,

$$\alpha \leq \frac{a + \sqrt{a^2 + 4b}}{2} \leq a + \sqrt{b} \quad \text{using } \sqrt{t+t'} \leq \sqrt{t} + \sqrt{t'}.]$$

We thus get :

$$\sum_{t=1}^T \delta_t \leq (M-m) \left(1 + \frac{2}{3} \ln N\right) + \sqrt{v_f \ln N}$$

On the other hand: $\frac{\ln N}{\eta_T} \leq \frac{\ln N}{\eta_m} = \sum_{t=1}^T \delta_t$

Final bound: $R_T \leq 2 \sum_{t=1}^T \delta_t$.

Exercise 2. The proof above seems unnatural to me; I do not really see the link to the intuition I give you.

→ If you can find a more natural exposition (still starting from Bernstein's inequality), you will get bonus points at the exam!

Corollary: The algorithm above is such that

$$\forall T, \forall m \leq M, \forall f_T \in [m, M], \sum_{t=1}^T p_t f_t = \min_k \sum_{t=1}^T p_t f_{kt} \leq (M-m) \sqrt{T \ln N} + (M-m) \left(2 + \frac{4}{3} \ln N\right)$$

not known by the algorithm

↳ The algorithm is adaptive to T and to the bounds m and M .

Proof: $v_f = \text{some variance of a random variable } \in [m, M] \leq \frac{(M-m)^2}{4}$, as in Hoeffding's lemma.

Exercise 3. Improvement for small losses.

This exercise shows that the \sqrt{T} order of magnitude for the regret can be improved on easy problems. The results below are of a nature similar to the "fast rates" results in classification/statistical

learning.

(1) Assume that losses l_{jt} are non-negative : $l_{jt} \geq 0$, i.e., $l_{jt} \in [0, M]$.

Show that

$$\sum_{t=1}^T \sum_{j=1}^N p_{jt} \left(l_{jt} - \frac{1}{k} \sum_{k=1}^N l_{kt} \right)^2 \leq M \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt}.$$

(2) Exhibit a regret bound of order $\sqrt{\min_{k=1 \dots N} \sum_{t=1}^T l_{kt} \ln N}$

(3) Specify the value of this bound when $\exists k \mid \sum_{t=1}^T l_{kt} = 0$
(the "perfect expert" case)

↳ Have the bound be
really reader-friendly, please!

Exercise 4. « Impossible tuning »

Consider the Prod forecaster : $p_1 = (\gamma_1, \dots, \gamma_N)$

and for $t \geq 2$,

$$p_{jt} = \frac{\prod_{s=1}^{t-1} (1 - \eta f_{js})}{\sum_{k=1}^N \prod_{s=1}^{t-1} (1 - \eta f_{ks})}$$

It appears as a first-order
approximation to ENA.

(1) Shows that $\sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt} - \sum_{t=1}^T l_{kt} \leq \frac{\ln N}{\eta} + \eta \sum_{t=1}^T l_{kt}^2$

when losses are such that $l_{jt} \leq M$ and when η has been chosen

so that $\eta \leq 1/2M$

Hint: use $\ln(1+u) \geq u - u^2$ $\forall u \geq -\frac{1}{2}$ and $\ln(1+u) \leq u$ $\forall u > -1$

(2) Can you tune η (by doubling trick, by having it depend on time, etc.)
so that:

$$\sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt} \leq \min_{k=1 \dots N} \left\{ \sum_{t=1}^T l_{kt} + C \sqrt{\sum_{t=1}^T l_{kt}^2 \ln N} \right\}$$

(for some universal constant C to be defined by the analysis.) ?

|| IF SO: YOU'LL GET 20/20 AT MY EXAM AND WE WILL WRITE A PAPER
TOGETHER!
↳ a paper by Grünwald / Koolen / van Erven gets such a bound last for a different algorithm

How do we pick the η_t in practice for EWA?

What my PhD students and I do:

$$\hat{L}_t(\eta) = \sum_{s=1}^t \sum_{j=1}^N p_{js}(\eta) l_{js}$$

cumulative loss of
EWA with fixed $\eta > 0$
on rounds 1 to t

Then we select $\eta_t \in \arg\min_{\eta > 0} \hat{L}_{t-1}(\eta)$

Trade off between

- not overfitting (still EWA prediction)
- fitting η according to the true stochasticity / non-stochasticity of the data (the more stochastic, the larger η ; the more adversarial, the smaller η).

Exhibits good empirical performance (much better than when η tuned as prescribed by theory).

Open questions x2 (would be good topics for PhD work):

1. One can show that in general, the regret of the above strategy can be linear. Can you find a calibration strategy that works well both in practice and in theory for EWA?
2. Can you quantify how to set η depending on the degree of stochasticity of the data?

Note: So how different the above calibration strategy is from, e.g.) cross-validation \rightarrow We really exploit here the sequential fashion of our setting!