## Algorithm for a distribution-free bound.

We will study it in the following, more general, setting, called the adversarial setting:

For $t = 1, 2, \dots$ :

- The opponent picks $(g_{1,t}, \dots g_{k,t}) \in [0,1]^k$ while, simultaneously,   *(and thanks to an auxiliary randomization $U_t$)*

- The decision-maker picks $I_t \in \{1 \dots K\}$, possibly at random according to a distribution $p_t = (p_{1,t} \dots p_{k,t})$   $\downarrow$

- The payoff of the decision-maker is $g_{I_t, t}$, and this is the only feedback she gets ( she does not observe $g_{j,t}$ for $j \neq I_t$ )

Aim:   Control the regret   $R_T = \max_{k=1 \dots K} \sum_{t=1}^{T} g_{k,t} - \sum_{t=1}^{T} g_{I_t, t}$

    $\hookrightarrow$ in high probability (preferred but difficult)

    $\hookrightarrow$ in expectation:   $\mathbb{E}[R_T]$   *Reminder: even $g_{k,t}$ depends on $I_1 \dots I_{t-1}$*

Example:   A stochastic opponent picks once for all, $\nu_1 \dots \nu_k \in \mathcal{P}([0,1])$ and draws $g_{k,t} \sim \nu_k$ for all $t$, independently.

Then (as we already discussed):

$$\mathbb{E}[R_T] \geqslant T \max_{k=1 \dots K} \mathbb{E}(\nu_k) - \mathbb{E}\left[\sum_{t=1}^{T} g_{I_t, t}\right] = \bar{R}_T$$

Fundamental observations:   Assume with no loss of generality that the opponent is deterministic.

Then $g_{k,t}$ is $\mathcal{F}_{t-1} \overset{= \sigma(U_1 \dots U_{t-1})}{=} \sigma(I_1 \dots I_{t-1})$ -measurable, and so is $p_t$.

Thus, in $\mathbb{E}[g_{I_t, t} \mid \mathcal{F}_{t-1}]$, the only randomness to integrate over, is the random draw of $I_t$ according to $p_t$ thanks to $U_t$:

$$\mathbb{E}[g_{I_t, t} \mid \mathcal{F}_{t-1}] = \mathbb{E}\left[\sum_{k=1}^{k} g_{k,t} \mathbb{1}_{\{I_t = k\}} \mid \mathcal{F}_{t-1}\right]$$

$$= \sum_{k=1}^{k} g_{k,t} \underbrace{\mathbb{E}[\mathbb{1}_{\{I_t = k\}} \mid \mathcal{F}_{t-1}]}_{= p_{k,t} \text{ by the model}} = \sum_{k=1}^{k} p_{k,t} g_{k,t}$$

This leads to the following important observation:     if $P_{kt} > 0$, then

$$\hat{\ell}_{kt} = \frac{(1-g_{I_t,t})}{P_{kt}} \mathbb{1}_{\{I_t = k\}} \qquad \text{is a (conditionally) unbiased estimator of the pseudo-losses } \ell_{kt} = 1 - g_{kt}$$

indeed, similar calculations show that

$$E\left[\hat{\ell}_{kt} \mid \mathcal{F}_{t-1}\right] = E\left[\frac{1-g_{kt}}{P_{kt}} \mathbb{1}_{\{I_t = k\}} \mid \mathcal{F}_{t-1}\right] = \frac{1-g_{kt}}{P_{kt}} \underbrace{E\left[\mathbb{1}_{\{I_t = k\}} \mid \mathcal{F}_{t-1}\right]}_{= P_{kt}}$$

$$\underset{P_{kt} > 0}{\overset{as}{\longleftarrow}} = 1 - g_{kt}$$

↓ exponentially weighted averages

**Algorithm:**     EWA  (with fixed learning rate)          Parameter: $\eta > 0$

Set $P_1 = (1/k, \dots, 1/k)$,  draw $I_1 \sim P_1$,  get $g_{I_1,1}$,  compute the $\hat{\ell}_{k1}$ for $k = 1..K$

For $t = 2, 3, \dots$

— set $P_{it} = \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_{is}\right) \Big/ \sum_{k=1}^{K} \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_{ks}\right)$

— draw $I_t \sim P_t$,  get $g_{I_t,t}$,  compute the $\hat{\ell}_{kt}$ for $k = 1, \dots K$

**Theorem:**     For all opponents picking payoffs $g_{kt} \in [q,1]$,  for all $T \geq 1$,

$$\overline{R}_T = \max_{k=1..K} E\left[\sum_{t=1}^{T} g_{kt}\right] - E\left[\sum_{t=1}^{T} g_{I_t,t}\right] \leq \frac{\ln K}{\eta} + \eta T \frac{K}{2}$$

In particular,  for $\eta = \sqrt{\frac{2k \ln K}{T}}$,  we get     $\overline{R}_T \leq \sqrt{2 \, T k \ln K}$

**Issue 1:**  The horizon $T$ is unknown in advance,  and even worse, we may have $T \to +\infty$   We will need to study how to better tune $\eta$.

**Issue 2:**  We don't deal yet with $E[R_T]$ …     → we might do so in the exam!

Before we actually prove this Theorem, let us state and prove some lemmas in what is called the full-information setting.

Two fundamental results in the adversarial, full-information setting

Setting:          for $t = 1, 2, \ldots$

- the opponent picks losses $(\ell_{1t} \ldots \ell_{Kt}) \in [0,1]^K$

  while, simultaneously,

- the decision-maker picks $I_t \in \{1..K\}$, possibly at random according to a distribution $p_t = (p_{1t} \ldots p_{Kt})$ and thanks to an auxiliary randomization $U_t$

- the decision-maker incurs a loss $\ell_{I_t, t}$

- she gets to observe the full vector $(\ell_{1t} \ldots \ell_{Kt})$

Aim:          Control the regret          $R_T = \sum_{t=1}^{T} \ell_{I_t, t} - \min_{k=1..K} \sum_{t=1}^{T} \ell_{kt}$

Algorithm:          EWA (exponentially weighted averages)          — parameter: $\eta > 0$

for $t \geq 1$,
for $k \in \{1..K\}$,          $p_{kt} = \exp\left(-\eta \sum_{s=1}^{t-1} \ell_{ks}\right) \Big/ \sum_{i=1}^{k} \exp\left(-\eta \sum_{s=1}^{t-1} \ell_{is}\right)$

with the convention that an empty sum is null;          so that          $p_1 = (1/K, \ldots, 1/K)$.

Theorem 1:          For all opponents picking losses $\ell_{kt} \in [m, M]$, for all $T \geq 1$,

$$\tilde{R}_T = \sum_{t=1}^{T} \sum_{i=1}^{k} p_{it} \ell_{it} - \min_{k=1..K} \sum_{t=1}^{T} \ell_{kt} \leq \frac{\ln K}{\eta} + \eta T \frac{(M-m)^2}{8}$$

In particular,          $\tilde{R}_T \leq (M-m)\sqrt{\frac{T}{2} \ln K}$          for          $\eta = \frac{1}{M-m} \sqrt{\frac{8 \ln K}{T}}$

Practical use of the bound:

The Hoeffding-Azuma inequality ensures that for all $T$, with probability at least $1-\delta$, where $\delta > 0$:

$$R_T \leq \tilde{R}_T + (M-m)\sqrt{\frac{T}{2} \ln \frac{1}{\delta}}$$

**Proof:** for a given $t$,

$$\sum_{i=1}^{k} p_{it}\, \ell_{it} \;\leq\; -\frac{1}{\eta} \ln\left( \sum_{i=1}^{k} p_{it}\, e^{-\eta \ell_{it}} \right) + \frac{\eta}{8}(M-m)^2$$

**Indeed:** Hoeffding's lemma ensures that if $Z$ is a $[m,M]$-valued random variable, then

$$\ln \mathbb{E}[e^{-\eta Z}] \leq -\eta\, \mathbb{E}[Z] + \frac{\eta^2}{8}(M-m)^2$$

or put differently,

$$\mathbb{E}[Z] \leq -\frac{1}{\eta}\ln \mathbb{E}[e^{-\eta Z}] + \frac{\eta}{8}(M-m)^2$$

Thus,

$$\sum_{t=1}^{T}\sum_{i=1}^{k} p_{it}\,\ell_{it} \;\leq\; \eta\frac{T}{8}(M-m)^2 - \frac{1}{\eta}\sum_{t=1}^{T}\ln\underbrace{\sum_{i=1}^{k} p_{it}\, e^{-\eta \ell_{it}}}_{=\; \dfrac{\sum\limits_{k=1}^{k}\exp\left(-\eta\sum\limits_{s=1}^{t}\ell_{ks}\right)}{\sum\limits_{k=1}^{k}\exp\left(-\eta\sum\limits_{s=1}^{t-1}\ell_{ks}\right)}} \quad \begin{array}{l}\text{by def.}\\ \text{of}\\ p_t\end{array}$$

$$\underset{\substack{\text{telescoping}\\ \text{sum}}}{\leq}\; \eta\, T(M-m)^2/8 - \frac{1}{\eta}\ln\frac{\sum\limits_{k=1}^{k}\exp\left(-\eta\sum\limits_{t=1}^{T}\ell_{kt}\right)}{K}$$

The proof is concluded by rearranging the inequality and by resorting to the lower bound

$$\ln\left(\sum_{k=1}^{k}\exp\left(-\eta\sum_{t=1}^{T}\ell_{kt}\right)\right) \;\geq\; \ln\left(\max_{k=1..K}\exp\left(-\eta\sum_{t=1}^{T}\ell_{kt}\right)\right)$$

$$= \max_{k=1...K} -\eta\sum_{t=1}^{T}\ell_{kt}$$

$$= -\eta\min_{k=1...K}\sum_{t=1}^{T}\ell_{kt}$$

Our estimated pseudo-losses $\hat{\ell}_{kt} = \dfrac{1-g_{I_t t}}{p_{It}}\,\mathbb{1}_{\{I_t = k\}}$ are $\geq 0$ but can be arbitrarily large (as we don't impose any lower bound on the $p_{kt}$). Thus, to study EWA in the bandit setting, we rather resort to the following result:

**Theorem 2:** For all opponents picking losses $\ell_{kt} \geq 0$, for all $T \geq 1$,

$$\tilde{R}_T = \sum_{t=1}^{T}\sum_{i=1}^{k} p_{it}\,\ell_{it} - \min_{k=1..K}\sum_{t=1}^{T}\ell_{kt}$$

$$\leq \frac{\ln K}{\eta} + \frac{\eta}{2}\sum_{t=1}^{T}\sum_{i=1}^{k} p_{it}\,\ell_{it}^2$$

**Proof:** $e^{-x} \leq 1 - x + \dfrac{x^2}{2}$ for all $x \geq 0$

so that, for all $\eta \geq 0$ and all random variables $Z \geq 0$ a.s.,

$$\mathbb{E}\left[e^{-\eta Z}\right] \leq 1 - \eta\,\mathbb{E}[Z] + \frac{\eta^2}{2}\,\mathbb{E}[Z^2]$$

$\Bigg)$ $\ln(1+u) \leq u$
for all $u > -1$

$$\ln\left(\mathbb{E}\left[e^{-\eta Z}\right]\right) \leq -\eta\,\mathbb{E}[Z] + \frac{\eta^2}{2}\,\mathbb{E}[Z^2]$$

and finally:
$$\mathbb{E}[Z] \leq -\frac{1}{\eta}\ln\mathbb{E}\left[e^{-\eta Z}\right] + \frac{\eta}{2}\,\mathbb{E}[Z^2]$$

In particular, for all $t \geq 1$,
$$\sum_{i=1}^{K} p_{i,t}\,\ell_{i,t} \leq -\frac{1}{\eta}\ln\left(\sum_{i=1}^{K} p_{i,t}\,e^{-\eta\,\ell_{i,t}}\right) + \frac{\eta}{2}\sum_{i=1}^{K} p_{i,t}\,\ell_{i,t}^2$$

from which the bound follows by summation over $t$ and via the same manipulations as the ones performed at the end of the previous proof.

Back to EWA on estimated pseudo-losses $\hat{\ell}_{kt}$ in the bandit setting

Theorem 2 ensures that

$$\sum_{t=1}^{T} \sum_{i=1}^{k} p_{it} \hat{\ell}_{it} - \min_{k=1..K} \sum_{t=1}^{T} \hat{\ell}_{kt} \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \sum_{j=1}^{k} p_{jt} \hat{\ell}_{jt}^2$$

Now, $\sum_{i=1}^{k} p_{it} \hat{\ell}_{it} = \sum_{i=1}^{k} p_{it} \frac{1-g_{I_t t}}{p_{it}} \mathbb{1}_{\{I_t = j\}} = (1 - g_{I_t t}) \sum_{i=1}^{k} \mathbb{1}_{\{I_t = j\}}$

$$= 1 - g_{I_t t}$$

Similarly, $\sum_{j=1}^{k} p_{jt} \hat{\ell}_{jt}^2 = \sum_{j=1}^{k} p_{jt} \frac{1 - g_{I_t t}}{p_{jt}^2} \mathbb{1}_{\{I_t = j\}} \leq \sum_{j=1}^{k} \frac{1}{p_{jt}} \mathbb{1}_{\{I_t = j\}}$

Thus we have

(*) $$\sum_{t=1}^{T} (1 - g_{I_t t}) - \min_{k=1..K} \sum_{t=1}^{T} \hat{\ell}_{kt} \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \sum_{j=1}^{k} \frac{\mathbb{1}_{\{I_t = j\}}}{p_{jt}}$$

By the tower rule, $\mathbb{E}[\hat{\ell}_{kt}] = \mathbb{E}[\mathbb{E}[\hat{\ell}_{kt} | \mathcal{F}_{t-1}]] = \mathbb{E}[1 - g_{kt}]$

and $\mathbb{E}\left[\frac{\mathbb{1}_{\{I_t = j\}}}{p_{jt}}\right] = \mathbb{E}\left[\frac{1}{p_{jt}} \underbrace{\mathbb{E}[\mathbb{1}_{\{I_t = j\}} | \mathcal{F}_{t-1}]}_{= p_{jt}}\right] = 1$

So that, by taking expectations $\mathbb{E}[\ ]$ in both sides of (*) and using that $\mathbb{E}[\min ...] \leq \min \mathbb{E}[...]$, we get

$$\mathbb{E}\left[\sum_{t=1}^{T} (1 - g_{I_t t}) - \min_{k=1..k} \sum_{t=1}^{T} \hat{\ell}_{kt}\right] \leq \frac{\ln K}{\eta} + \frac{\eta}{2} KT$$

$$\geq \mathbb{E}\left[\sum_{t=1}^{T} (1 - g_{I_t t})\right] - \min_{k=1..k} \mathbb{E}\left[\sum_{t=1}^{T} (1 - g_{kt})\right]$$

$$= \max_{k=1..K} \mathbb{E}\left[\sum_{t=1}^{T} g_{kt}\right] - \mathbb{E}\left[\sum_{t=1}^{T} g_{I_t t}\right]$$

Application of the "EWA forecaster" / Sion's lemma.

<u>Statement:</u> Let $X, Y$ two convex sets, $f: X \times Y \to [0, M]$ a function s.t. $\forall x \in X$, $f(x, \cdot)$ is concave
$\forall y \in Y$, $f(\cdot, y)$ is convex

then (under additional regularity assumptions):

$$\inf_{x \in X} \sup_{y \in Y} f(x, y) = \sup_{y \in Y} \inf_{x \in X} f(x, y).$$

<u>Proof:</u>

1) $\geqslant$ always holds : $\forall x, \forall y, \quad f(x, y) \geqslant \inf_{x' \in X} f(x', y)$

taking $\sup_{y \in Y}$ in both sides: $\forall x, \quad \sup_{y \in Y} f(x, y) \geqslant \sup_{y \in Y} \inf_{x' \in X} f(x', y)$

Conclude by taking the $\inf_{x \in X}$ (the right-hand side is a constant independent of $x$).

2) A (fictitious) statistician and a (fictitious) opponent <u>play</u> as follows :

First, the statistician sets $N \geqslant 2$ and $x^{(1)} \dots x^{(N)}$ in $X$, as well as $T \geqslant 1$

Then, at each round, they simultaneously pick
$$x_t = \sum_{j=1}^{N} p_{jt} \, x^{(j)} \in X \quad \text{and} \quad y_t \in Y$$

with $\eta = \frac{1}{M} \sqrt{\frac{8 \ln N}{T}}$

How? $p_{jt} = \exp\left(-\eta \sum_{s=1}^{t-1} f(x^{(j)}, y_s)\right) \Big/ \sum_{k=1}^{N} \exp\left(-\eta \sum_{s=1}^{t-1} f(x^{(k)}, y_s)\right)$

and (since $p_{jt}$ only depends on the past, the opponent can compute it and pick :)

$y_t$ s.t. (by definition of sup) $\quad f(x_t, y_t) \geqslant \sup_{y \in Y} f(x_t, y) - \frac{1}{T}$

By definition of the exponentially weighted average strategy with learning rate $\eta > 0$,

$$\sum_{t=1}^{T} \underbrace{\sum_{j=1}^{N} P_{jt}\, f(x^{(j)}, y_t)}_{\substack{\text{corresponds} \\ \text{to } g_t \in [0, M]}} \quad - \quad \min_{k=1,\dots N} \sum_{t=1}^{T} f(x^{(k)}, y_t) \qquad \leq \quad \underbrace{\frac{\ln N}{\eta} + \eta \frac{M^2}{8} T}_{= M\sqrt{\frac{T}{2}\ln N}} \qquad \left.\begin{array}{c} \text{given} \\ \text{our} \\ \text{choice} \\ \text{for } \eta \end{array}\right.$$

From the convexity of $f(\cdot, y_t)$, we finally get:

$$\sum_{t=1}^{T} f\Big(\underbrace{\sum_{j} P_{jt}\, x^{(j)}}_{\substack{= x_t \\ \text{by def.}}}, y_t\Big) \quad - \quad \min_{k=1\cdots N} \sum_{t=1}^{T} f(x^{(k)}, y_t) \qquad \leq \quad M\sqrt{\frac{T}{2}\ln N} \qquad (\ast)$$

where $\bar{x} = \frac{1}{T}\sum_{t=1}^{T} x_t$

3)
$$\inf_{x} \sup_{y} f(x, y) \quad \leq \quad \sup_{y} f(\bar{x}, y) \qquad \left.\begin{array}{c} f(\cdot, y) \\ \text{convex} \\ \forall y \end{array}\right.$$

$$\leq \quad \sup_{y} \frac{1}{T}\sum_{t=1}^{T} f(x_t, y)$$

$$\underset{\substack{\sup \Sigma \\ \leq \Sigma \sup}}{} \quad \leq \quad \frac{1}{T}\sum_{t=1}^{T} \sup_{y} f(x_t, y) \quad \leq \quad \frac{1}{T}\sum_{t=1}^{T} f(x_t, y_t) \quad +\frac{1}{T}$$

$$\underset{\substack{\text{def of} \\ y_t}}{\uparrow}$$

4) $\quad \dfrac{1}{T}\displaystyle\sum_{t=1}^{T} f(x_t, y_t)$

$$\leq \quad \underbrace{M\sqrt{\frac{\ln N}{2T}}}_{= O(1/\sqrt{T})} \quad + \quad \min_{k=1\cdots N} \frac{1}{T}\sum_{t=1}^{T} f(x^{(k)}, y_t) \qquad \left.\begin{array}{c} \text{by} \\ (\ast) \end{array}\right.$$

$$\leq \quad O(1/\sqrt{T}) \quad + \quad \underbrace{\min_{k} f(x^{(k)}, \bar{y})}_{\leq \sup_{y\in\mathcal{Y}} \min_{k} f(x^{(k)}, y)} \qquad \left.\begin{array}{c} \text{by concavity} \\ \text{of} \\ f(x^{(k)}, \cdot), \\ \text{where} \\ \bar{y} = \frac{1}{T}\sum_{t=1}^{T} y_t \end{array}\right.$$

5) In sections (1)-(4)     $T, N$ and $x^{(1)}\dots x^{(N)}$ were fixed, but we can play with them !     We proved

$$\inf_{x} \sup_{y} f(x, y) \quad \leq \quad \frac{1}{T} + O(1/\sqrt{T}) \quad + \quad \sup_{y} \min_{k} f(x^{(k)}, y)$$

Letting $\qquad T \longrightarrow +\infty$ :

$$\inf_x \sup_y f(x,y) \leq \sup_y \min_k f(x^{(k)},y)$$

This holds for all $N$ and all $x^{(1)} \dots x^{(N)}$ in $\mathcal{X}$ :

$$\inf_x \sup_y f(x,y) \leq \inf_{N \geq 1} \inf_{\{x^{(1)} \dots x^{(N)}\} \subset \mathcal{X}} \sup_{y \in \mathcal{Y}} \min_{k=1 \dots N} f(x^{(k)},y)$$

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}$

is clearly $\qquad$ ?

$$\geq \sup_y \inf_x f(x,y) \quad \text{but maybe} \leq \sup_y \inf_x f(x,y) \quad \text{so that}$$
we have
can equality ?

We will now state
and use regularity / topological assumptions.

Assume  • $\mathcal{X}, \mathcal{Y}$ are metric spaces, with distances $d_X$ and $d_Y$

  • $f : \mathcal{X} \times \mathcal{Y} \to [q, M]$ uniformly continuous :

$$\forall \varepsilon > 0, \quad \exists \delta > 0 \mid \quad d_X(x,x') + d_Y(y,y') \leq \delta$$
$$\Rightarrow \left| f(x,y) - f(x',y') \right| \leq \varepsilon$$

  • $\mathcal{X}$ compact : $\qquad \forall \delta > 0, \qquad \exists N$ and $x^{(1)} \dots x^{(N)}$ s.t.

$$\mathcal{X} \subset \bigcup_{j=1}^{N} B(x^{(j)}, \delta)$$

Given $\varepsilon > 0$ and the associated $\delta > 0$ :

$$\forall x,y, \qquad f(x,y) \geq \min_{j=1 \dots N} f(x^{(j)}, y) - \varepsilon \qquad \text{by uniform continuity}$$

Taking $\inf_x$ then $\sup_y$ :

$$\sup_y \inf_x f(x,y) \geq \sup_y \min_j f(x^{(j)},y) - \varepsilon ,$$

$$\geq \inf_N \inf_{\{x^{(j)}\}} \sup_y \min_j f(x^{(j)},y)$$
$$- \varepsilon$$

and we let $\varepsilon \searrow 0$

(Perhaps you can find even better - weaker - assumptions ? If so, let me know !
while still having a smooth and easy-to-read proof ... )

Sequential optimization for general convex sets.

We so far focused our attention on the simplex
$$\mathcal{X} = \{ (p_1 \dots p_N) : \forall j, \ p_j \geq 0 \quad \text{and} \quad \sum_k p_k = 1 \} \subset \mathbb{R}^N$$

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 $\ell_t : \mathcal{C} \to \mathbb{R}$

     2. $x_t$ and $\ell_t$ are revealed

Aim:       Control $R_T = \sum\limits_{t=1}^{T} \ell_t(x_t) - \inf\limits_{x \in \mathcal{C}} \sum\limits_{t=1}^{T} \ell_t(x)$

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

     — Play any $x_1 \in \mathcal{C}$
     — For $t = 2, 3, \dots,$
$$x_t = \Pi_{\mathcal{C}} \left( x_{t-1} - \eta \, \nabla \ell_{t-1}(x_{t-1}) \right)$$

Called OGD($\eta$)

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

Theorem:     Assume that $\max\limits_{t \leq T} \sup\limits_{x \in \mathcal{C}} \| \nabla \ell_t(x) \| \leq G$     [ bounded gradients ]

and that $\mathcal{C}$ is compact, i.e., $\sup\limits_{x, x' \in \mathcal{C}} \| x - x' \| \leq D$     [ bounded diameter ]

Then OGD($\eta$) is such that $\sum\limits_{t=1}^{T} \ell_t(x_t) - \min\limits_{x \in \mathcal{C}} \sum\limits_{t=1}^{T} \ell_t(x) \leq \dfrac{D^2}{2\eta} + \dfrac{\eta \, G^2 T}{2}$

(inf achieved by continuity + compactness)

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

<u>Remark:</u>    In the case with EWA, we were considering a linear function:

$$\forall p \in \mathcal{X}, \qquad \ell_t(p) = \sum_{j=1}^{N} p_j \cdot \ell_{jt}$$

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

diameter of $\mathcal{X}$ :    $\forall p, q \in \mathcal{X},$    $\sum_j (p_j - q_j)^2$
$D^2 = 2$

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

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

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

then    $\nabla \ell_t(p) = \begin{pmatrix} \ell_{1t} \\ \vdots \\ \ell_{Nt} \end{pmatrix}$

and    $\| \nabla \ell_t(p) \|^2 \leq N M^2 =: G^2$

We get the bounds:
- for OGD :    $DG\sqrt{T} = M\sqrt{2NT}$
- for EWA :    $(M - (-M))^2 \sqrt{\dfrac{T}{2} \ln N} = M\sqrt{8T \ln N}$

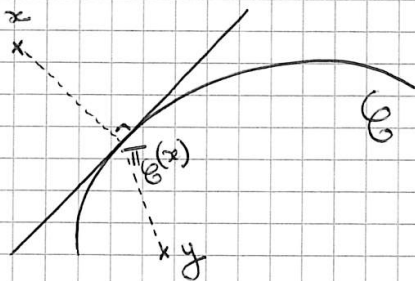The dependency in $N$ is suboptimal for OGD.

<u>Proof</u> (of the theorem): By convexity («inégalité des pentes» in French):

$$\forall x \in \mathcal{C}, \quad \forall t, \qquad \ell_t(x_t) - \ell_t(x) \leq \nabla \ell_t(x_t) \cdot (x_t - x)$$

Summing over $t$,

$$R_T = \sup_{x \in \mathcal{C}} \left\{ \sum_{t=1}^{T} \ell_t(x_t) - \sum_{t=1}^{T} \ell_t(x) \right\}$$

$$\leq \sup_{x \in \mathcal{C}} \left\{ \sum_{t=1}^{T} \nabla \ell_t(x_t) \cdot (x_t - x) \right\}$$

Now, we will need some facts about projections on <u>CLOSED</u> convex sets:



Facts: If $x \notin \mathcal{C}$, then

$$\forall y \in \mathcal{C}, \qquad (y - \Pi_{\mathcal{C}}(x)) \cdot (x - \Pi_{\mathcal{C}}(x)) \leq 0$$

Thus, $\forall y \in \mathcal{C}$, $\qquad \| y - \Pi_{\mathcal{C}}(x) \| \leq \| y - x \|$.

$$\left[ \text{Indeed:} \quad \| y - x \|^2 = \| y - \Pi_{\mathcal{C}}(x) + \Pi_{\mathcal{C}}(x) - x \|^2 \right.$$

$$= \| y - \Pi_{\mathcal{C}}(x) \|^2 + \| x - \Pi_{\mathcal{C}}(x) \|^2 \quad \text{both are } \geq 0$$

$$\left. + 2 (y - \Pi_{\mathcal{C}}(x)) \cdot (\Pi_{\mathcal{C}}(x) - x) \right.$$

In particular, for $t \geq 1$:    def of OGD ↓

$$\forall x \in \mathcal{C}, \qquad \| x_{t+1} - x \|^2 = \| \Pi_{\mathcal{C}}(x_t - \eta \nabla \ell_t(x_t)) - x \|^2$$

$$\underset{\substack{\text{by the}\\\text{facts above}}}{\leq} \| x_t - \eta \nabla \ell_t(x_t) - x \|^2$$

$$= \| x_t - x \|^2 + \eta^2 \| \nabla \ell_t(x_t) \|^2 - 2\eta \, \nabla \ell_t(x_t) \cdot (x_t - x)$$

That is, $\qquad \forall t, \forall x \in \mathcal{C}, \qquad \nabla \ell_t(x_t) \cdot (x_t - x)$

$$\leq \frac{1}{2\eta} \left( \| x_t - x \|^2 - \| x_{t+1} - x \|^2 \right) + \frac{\eta}{2} \underbrace{\| \nabla \ell_t(x_t) \|^2}_{\leq G^2}$$

by the assumption of bounded gradients

Summing over $t$: $\qquad \sum_{t=1}^{T} \nabla \ell_t(x_t) \cdot (x_t - x)$

$$\leq \frac{1}{2\eta} \left( \underbrace{\| x_1 - x \|^2}_{\leq D^2} - \underbrace{\| x_{T+1} - x \|^2}_{\geq 0} \right) + \frac{\eta}{2} T G^2$$

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

The algorithms discussed so far :

　　EWA　with fixed learning rate $\eta > 0$
　　OGD　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 OGD

We now study adaptation to these parameters.

There exists a general trick called the doubling trick which we will detail in 3 pages from now for EWA.

Another popular way of being adaptive is to tune $\eta$ according to the past observations:

　　Adaptive OGD :

　　　　－ Play an arbitrary $x_1 \in \mathcal{E}$,　pick $\eta_1 > 0$

　　　　－ For $t \geqslant 2$,　play　$x_t = \Pi_{\mathcal{E}} \left( x_{t-1} - \eta_{t-1} \nabla \ell_{t-1}(x_{t-1}) \right)$

　　　　　　　　　　pick $\eta_t > 0$

__Theorem__:　If Adaptive OGD picks a non-increasing sequence $(\eta_t)_t$ with $\eta_t > 0 \; \forall t$,

　　then:　$R_T \leqslant \dfrac{D^2}{2\eta_T} + \dfrac{1}{2} \sum_{t=1}^{T} \eta_t \left\| \nabla \ell_t(x_t) \right\|^2$　　with the same notation as in the theorem for OGD($\eta$)

__Corollary__:　Picking　$\eta_t = \begin{cases} 1 & \text{if } V_t = \sum_{s=1}^{t} \left\| \nabla \ell_s(x_s) \right\|^2 < D^2 \\[2mm] D/\sqrt{V_t} & \text{if } V_t \geqslant D^2 \end{cases}$

Adaptive OGD achieves the bound　$\max \left\{ D^2, \; \dfrac{3}{2} D \sqrt{\sum_{t=1}^{T} \left\| \nabla \ell_t(x_t) \right\|^2} \right\}$

　　　　　　　　　　　　　　　　$= O\left( D G \sqrt{T} \right)$

without any prior knowledge of neither $T$ nor $G$

( of $\mathcal{E}$ thus $D$ are known).

Proof of the theorem and of the corollary:

Th: We extract from the proof of OGD($\eta$) that

$$\forall x \in \mathcal{C}, \qquad \nabla \ell_t(x_t) \cdot (x_t - x) \leq \frac{1}{2\eta_t}\left( \|x_t - x\|^2 - \|x_{t+1} - x\|^2 \right)$$
$$+ \frac{\eta_t}{2}\|\nabla \ell_t(x_t)\|^2$$

Summing over $t$:

$$\forall x \in \mathcal{C}, \qquad \sum_{t=1}^{T} \ell_t(x_t) - \sum_{t=1}^{T} \ell_t(x) \leq \sum_{t=1}^{T} \nabla \ell_t(x_t) \cdot (x_t - x)$$

$$\leq \underbrace{\sum_{t=1}^{T} \frac{1}{2\eta_t}\left( \|x_t - x\|^2 - \|x_{t+1} - x\|^2 \right)}_{} $$

$$+ \frac{1}{2}\sum_{t=1}^{T} \eta_t \|\nabla \ell_t(x_t)\|^2$$

We need only to show
that this sum is
$\leq D^2/2\eta_T$

Abel's transform:

$$\sum_{t=1}^{T} \frac{1}{2\eta_t}\left( \|x_t - x\|^2 - \|x_{t+1} - x\|^2 \right)$$

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

$$\leq \frac{D^2}{2\eta_T} \qquad \text{after telescoping}$$

Cor: Recall that we picked:

Pay attention to the
indexation: $\eta_t$ is used
to pick $x_{t+1}$ and
may thus depend on the
past till and including
round $t$

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

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

impossible in practice choice $\eta_t = \dfrac{D}{V_T}$ by its "sequential estimation";
and use a threshold to avoid too large values of $\eta_t$ for small $t$.)

The regret bound reads $\qquad \dfrac{D^2}{2\eta_T} + \dfrac{1}{2} \sum\limits_{\substack{t: \\ V_t < D^2}} \| \nabla \ell_t(x_t) \|^2$

$$+ \dfrac{1}{2} \sum\limits_{\substack{t: \\ V_t > D^2}} \dfrac{D}{\sqrt{V_t}} \| \nabla \ell_t(x_t) \|^2$$

We denote $\tau$ the largest $t \le T$ such that $V_t < D^2$. Then, if $\tau < T$, we
have $V_{\tau+1} \geqslant D^2$.

The above bound can be rewritten as:

$$\dfrac{D^2}{2\eta_T} + \underbrace{\dfrac{1}{2} \sum\limits_{t:\, V_t < D^2} \| \nabla \ell_t(x_t) \|^2}_{= \frac{V_\tau}{2} < \frac{D^2}{2}} + \dfrac{1}{2} \sum\limits_{\substack{t = \tau+1 \\ \uparrow}}^{T} \dfrac{D}{\sqrt{V_t}} \| \nabla \ell_t(x_t) \|^2$$

possibly void sum (if $\tau = T$)

$$= \underbrace{\dfrac{1}{2} \sum\limits_{t=\tau+1}^{T} \dfrac{D}{\sqrt{V_t}} \left( V_t - V_{t-1} \right) = \dfrac{1}{2} \sum\limits_{t=\tau+1}^{T} \dfrac{D}{\sqrt{V_t}} \left( \sqrt{V_t} + \sqrt{V_{t-1}} \right)\underbrace{\left( \sqrt{V_t} - \sqrt{V_{t-1}} \right)}_{\leqslant 2\sqrt{V_t}}}$$

$$\leqslant D \sum\limits_{t=\tau+1}^{T} \left( \sqrt{V_t} - \sqrt{V_{t-1}} \right)$$

$$= D \left( \sqrt{V_T} - \sqrt{V_\tau} \right)$$

We conclude by distinguishing two cases:
 * $\tau = T$: $\qquad \eta_T = 1 \quad$ and the bound is $\leqslant \dfrac{D^2}{2} + \dfrac{D^2}{2} = D^2$

 * $\tau < T$: $\qquad \eta_T = \dfrac{D}{\sqrt{V_T}} \quad$ and the bound is

$$\dfrac{D^2}{2\,D/\sqrt{V_T}} + \dfrac{V_\tau}{2} + D\sqrt{V_T} - D\sqrt{V_\tau} \leqslant \dfrac{3D\sqrt{V_T}}{2}$$

$$\underbrace{\qquad\qquad}_{\leqslant \frac{D\sqrt{V_\tau}}{2} \text{ as } V_\tau < D^2}$$

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

$$\max\left\{ D^2, \ \dfrac{3}{2} D\sqrt{V_T} \right\} = \max\left\{ D^2, \ \dfrac{3}{2} D \sqrt{\sum\limits_{t=1}^{T} \| \nabla \ell_t(x_t) \|^2} \right\}$$