

## Fast rates with strongly convex loss functions

Strongly convex loss functions → Same predictor, faster rates

(Still Setting #2, but  $l_t$  only assumed to be convex, not linear)

Exp-concavity: The chosen loss functions  $l_t$  are  $\eta$ -exp-concave if  
 $\forall t, e^{-\eta l_t}$  is concave on  $X$

- Remarks:
- $x \mapsto x^\alpha$  is concave and increasing for  $0 < \alpha \leq 1$ ; thus  $\eta$ -exp-concave functions are also  $\eta'$ -exp-concave for all  $0 < \eta' \leq \eta$
  - $\ln$  is convex and decreasing: exp-concave functions are in particular convex functions

Example: Investment in the stock market, without transaction costs

$$l_t(p) = -\ln\left(\sum_{j=1}^J p_j x_{jt}\right)$$

1-exp-concave { where  $x_{jt}$  denotes the multiplicative evolution of stock  $j$ .  
The  $-\ln$  comes because we consider a cumulative loss (not a reward)

- The square loss  $l_t(p) = \left(\sum_j p_j x_{jt} - y_t\right)^2$   
If  $f_{jt}, y_t \in [0, B]$   $\forall j, t$ , then the  $l_t$  are  $\frac{1}{2B^2}$ -exp concave; it suffices to prove that  $\forall y \in [0, B], \Psi_y: x \in [0, B] \mapsto e^{-\eta(x-y)^2}$  is concave for  $\eta = \frac{1}{2B^2}$

But  $\Psi'_y(x) = -2\eta(x-y)e^{-\eta(x-y)^2}$

$$\Psi''_y(x) = (4\eta^2(x-y)^2 - 2\eta)e^{-\eta(x-y)^2} \leq 0$$

as soon as  $2\eta(x-y)^2 \leq 1$ , which is guaranteed by  $\eta \leq \frac{1}{2B^2}$ .

Same predictor:

with  $l_t^* \stackrel{\text{def}}{=} l_t(s_t)$ ; ie  $p_j^* \propto e^{-\eta \sum_{s_t} l_{js}}$

Theorem: If  $\ell_t$ ,  $\ell_t$  is  $\eta$ -exp-concave, then the exponentially weighted average predictor used with this  $\eta$  is such that

(for all strategies  
to pick the  $\ell_t$ )

$$\sum_{t=1}^T \ell_t(p_t) - \min_{k=1 \dots N} \sum_{t=1}^T \ell_t(s_k) \leq \frac{\ln N}{\eta}$$

Ex: For all sequences of bounded forecasts and observations

$y_t$   $f_j \in [0, B]$  and  $\ell_t$ ,  $y_t \in [0, B]$ ,

$$\frac{1}{T} \sum_{t=1}^T \left( \sum_j p_j f_j - y_t \right)^2 \leq \frac{2B^2 \ln N}{T}$$

$$+ \min_{k=1 \dots N} \frac{1}{T} \sum_{t=1}^T (f_k - y_t)^2$$

In particular, taking  $\sqrt{\cdot}$  and using  $\sqrt{ab} \leq \sqrt{a} + \sqrt{b}$ :

For sequences, RMSE of meta-predictor  $\leq$  RMSE of best predictor +  $B \sqrt{\frac{2 \ln N}{T}}$ .

↳ Hence the name "prediction of individual sequences" or "robust aggregation of forecasts".

Proof:  $\ell_t(p_t) = -\frac{1}{\eta} \ln \underbrace{e^{-\eta \ell_t(p_t)}}_{\geq \sum_j p_j e^{-\eta f_j}} \leq -\frac{1}{\eta} \ln \sum_j p_j e^{-\eta f_j}$

Same telescoping argument:

$$\sum_{t=1}^T \ell_t(p_t) \leq -\frac{1}{\eta} \ln \frac{\sum_{j=1}^N e^{-\eta \sum_{t=1}^T f_j}}{N}$$

and same way of concluding as before.

]

What about the more challenging comparison to convex combinations?

$$\inf_{p \in X} \sum_{t=1}^T \ell_t(p)$$

?

Exp-concavity & comparison to convex losses

Algorithm:

(continuous EWA predictor)

$$p_t = \frac{\int_X p e^{-\eta \sum_{s=1}^{t-1} l_s(p)} d\mu(p)}{\int_X e^{-\eta \sum_{s=1}^{t-1} l_s(p)} d\mu(p)}$$

not.  $\int_X p d\mu_t(p)$

where  $\mu$  is the uniform (Lebesgue) measure on  $X$ .

implicitly defining a measure  
here

Remarks:

- $p_1 = (\frac{1}{N}, \dots, \frac{1}{N})$
- how to better grasp  $\mu$ ?  $X$  is included in an affine hyperspace, in which it has a positive Lebesgue measure  $\rightarrow \mu$  is a conditional measure of that
- one can show that if  $X_1, \dots, X_{N-1}$  iid  $\sim U_{[0,1]}$  then  $0 \leq X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(N-1)} \leq 1$  (order statistics) define  $N$  segments

$$(X_{(1)}, X_{(2)} - X_{(1)}, \dots, 1 - X_{(N-1)})$$

whose joint law is  $\mu$ .

- how to compute  $p_t$ ?

↳ grid arguments (with  $S$  intervals) have a complexity of  $O(S^{N+1})$

↳ Monte-Carlo methods? Simulate  $P_1, \dots, P_m$  iid  $\sim \mu$

$$\text{and use } \mathbb{E}_\mu g(p) \approx \frac{1}{m} \sum_{j=1}^m g(P_j)$$

in any case,  
a nightmare  
in practice

(I had to  
implement it once,

see Stoltz & Lugosi 2005)

↳ NCNC methods, see Vempala 1996.

Theorem: For all sequences  $l_t$  of  $\eta$ -exp-concave functions, the continuous EWA predictor satisfies:

$$\sum_{t=1}^T l_t(p_t) - \inf_{p \in \mathcal{X}} \sum_{t=1}^T l_t(p) \leq \frac{1}{\eta} (1 + (N-1) \ln(T))$$

PROOF:

$$\begin{aligned} l_t(p_t) &= -\frac{1}{\eta} \ln e^{-\eta l_t(p_t)} \\ &= -\frac{1}{\eta} \ln \underbrace{e^{-\eta l_t(\int_X p d\mu_t(p))}}_{\geq \int_X e^{-\eta l_t(p)} d\mu_t(p)} \quad \text{by exp-concavity} \end{aligned}$$

thus

$$\begin{aligned} l_t(p_t) &\leq -\frac{1}{\eta} \ln \int_X e^{-\eta l_t(p)} d\mu_t(p) \\ &= -\frac{1}{\eta} \ln \frac{\int_X e^{-\eta \sum_{s=1}^t l_s(p)} d\mu_t(p)}{\int_X e^{-\eta \sum_{s=1}^{t-1} l_s(p)} d\mu_t(p)} \end{aligned}$$

Telescoping argument:

$$\sum_{t=1}^T l_t(p_t) \leq \cancel{-\frac{1}{\eta} \ln \frac{1}{\eta}} - \frac{1}{\eta} \ln \int_X e^{-\eta \sum_{t=1}^T l_t(p)} d\mu_t(p)$$

Fix  $\delta > 0$ , let  $p_\delta^*$  be s.t.  $\inf_p \sum_{t=1}^T l_t(p) \leq \delta + \sum_{t=1}^T l_t(p_\delta^*)$

why is the inf not a min?  $l_t$  is convex thus continuous on the interior of  $\mathcal{X}$  but can be discontinuous on its border

(1D-example:  $x \in [0, 1] \mapsto x^2 + \frac{1}{x+1}$ )

$$\Delta_{\delta, \varepsilon}^* = \left\{ (1-\varepsilon)p_\delta^* + \varepsilon r, \quad r \in \mathcal{X} \right\} \quad \text{for } \varepsilon \in [0, 1]$$

→ for such  $p = (1-\varepsilon)p_\delta^* + \varepsilon r$ , we have by exp-concavity:

$$\forall t, \quad e^{-\eta l_t(p)} \geq (1-\varepsilon) e^{-\eta l_t(p_\delta^*)} + \underbrace{\varepsilon e^{-\eta l_t(r)}}_{\geq 0}$$

$$\rightarrow \text{also, } \mu(\Delta_{\delta, \varepsilon}^*) = \varepsilon^{N-1}$$

f. properties of the Lebesgue measure (and ambient dimension  $N-1$ )

Thus

$$\int_{\mathcal{X}} e^{-\eta \sum_{t=1}^T l_t(p)} d\mu(p) \geq (1-\varepsilon)^T \varepsilon^{N-1} e^{-\eta \sum_{t=1}^T l_t(p_s^*)}$$

Substituting in the bound:

$$\sum_{t=1}^T l_t(p_t) \leq -\frac{1}{\eta} \ln((1-\varepsilon)^T \varepsilon^{N-1}) + \underbrace{\sum_{t=1}^T l_t(p_s^*)}_{\leq \delta + \inf_{p \in \mathcal{X}} \sum_{t=1}^T l_t(p)}$$

Let  $\delta \downarrow 0$ , we have proved:

$$\sum_{t=1}^T l_t(p_t) - \inf_{p \in \mathcal{X}} \sum_{t=1}^T l_t(p) \leq \inf_{\mathbb{E}[l] \in \mathcal{C}} -\frac{1}{\eta} \ln((1-\varepsilon)^T \varepsilon^{N-1})$$

Choosing, e.g.,  $\varepsilon = 1/T$ :

$$\begin{aligned} \frac{1}{\eta} (1-\varepsilon)^T &= (1 - \frac{1}{T})^T = (\frac{T-1}{T})^T = (1 - \frac{1}{T})^T \\ &= \exp(T \ln(1 - \frac{1}{T})) \leq e^{-\frac{1}{T}} \end{aligned}$$

hence the stated bound.  $\square$

A better bound?

$$B(\varepsilon) = T \ln(1-\varepsilon) + (N-1) \ln \varepsilon$$

$$B'(\varepsilon) = \frac{-T}{1-\varepsilon} + \frac{N-1}{\varepsilon}$$

$$B''(\varepsilon) < 0 \quad \text{thus the } \varepsilon^* \text{ s.t. } B'(\varepsilon^*) = 0$$

is the global minimum of  $B$

$$\text{We have } B'(\varepsilon^*) = 0 \Leftrightarrow \frac{\varepsilon^*}{N-1} = \frac{1}{T} - \frac{\varepsilon^*}{T} \Leftrightarrow \varepsilon^* = \frac{1/T}{1/T + 1/(N-1)} = \frac{N-1}{T+N-1}$$

We thus get the sharper final bound

$$\begin{aligned} \frac{1}{\eta} \left( T \ln \left( \frac{T+N-1}{T} \right) + (N-1) \ln \left( \frac{T+N-1}{N-1} \right) \right) &= \frac{T+N-1}{\eta} H \left( \frac{N-1}{T+N-1} \right) \\ &= \frac{1}{\eta} \left( (N-1) \underbrace{\ln \left( 1 + \frac{N-1}{T} \right)}_{\leq (N-1)/T} + (N-1) \ln \left( 1 + \frac{T}{N-1} \right) \right) \end{aligned}$$

where  $H(x) = -x \ln x + (1-x) \ln(1-x)$   
is the binary entropy

\* However \* both bounds are  $\sim \frac{N-1}{\eta} \ln T$  as  $T \rightarrow +\infty$ .

↳ And the second bound seems less readable to me...

## Convex loss functions

### Convex Loss Functions

$$l_t : X \rightarrow [m, M]$$

convex

#### 1. Comparison to the best individual expert

Regret  $R_T = \sum_{t=1}^T l_t(p_t) - \min_{i=1\dots N} \sum_{t=1}^T l_t(s_i)$  where  $s_i$   
Dirac mass at  $i$

is bounded by convexity by

$$\hat{R}_T \leq \sum_{t=1}^T \sum_{j=1}^N p_{tj} l_{jt} - \min_{i=1\dots N} \sum_{t=1}^T l_{it}$$

where we denoted  
 $l_{jt} = l_t(s_j)$

linear upper bound, which we know how to  
control with EWA

#### 2. Comparison to the best convex vector

An inefficient algorithm to do so: the continuous EWA predictor

(↳ we'll see a more efficient way (algorithm with  
a better computational complexity), a better regret  
bound)

1) Show that the continuous EWA strategy is such that:

Against all sequences  $l_t : X \rightarrow [m, M]$  of convex functions,

$$\sum_{t=1}^T l_t(p_t) - \inf_{q \in X} \sum_{t=1}^T l_t(q) \leq \frac{\eta}{8} (M-m)^2 T + \inf_{E \in \mathcal{E}(q)} -\frac{1}{\eta} \ln \left( \varepsilon^{N-i} e^{-\frac{\eta}{2} E(M-m)} \right)$$

2) When  $T, M$  and  $m$  are known, which  $\eta$  is a good choice  
to minimize the obtained upper bound, and can you  
provide an easy to read final upper bound?



## Convex functions & comparison to the best convex vector

↳ efficient forecaster

$\mathcal{X} = \{(p_1, \dots, p_N) : \sum_j p_j = 1, \forall k, p_k \geq 0\}$  is the simplex of convex weight vectors.

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

1. The statistician and the opponent pick simultaneously  $p_t \in \mathcal{X}$  and a convex loss function  $l_t: \mathcal{X} \rightarrow [m, M]$

2.  $l_t$  and  $p_t$  are revealed

↳ Regret:  $R_T = \sum_{t=1}^T l_t(p_t) - \inf_{p \in \mathcal{X}} \sum_{t=1}^T l_t(p)$  to be controlled in a uniform way

Fact: Convex functions are subdifferentiable on the interior of their domain of definition:

Let  $f: \mathcal{S} \rightarrow \mathbb{R}$  be convex, where  $\mathcal{S} \subseteq \mathbb{R}^N$  is convex:

$$\text{for } x \in \mathcal{S}, \exists \nabla f_x \in \mathbb{R}^N \mid \text{for } y \in \mathcal{S}, f(x) - f(y) \leq \nabla f_x \cdot (x-y)$$

$\mathcal{E}_f(x) = \{\text{the set of all possible such } \nabla f_x\}$  is called the subgradient of  $f$  at  $x$

If  $f$  is differentiable at  $x$ , then  $\mathcal{E}_f(x) = \{\nabla f(x)\}$ .

Application: If  $p_t \in \mathcal{X}$  (ie  $p_{jt} \geq 0 \forall j$ ) then  $\exists \nabla l_t(p_t) \in \mathbb{R}^N \mid \forall p \in \mathcal{X}, l_t(p_t) - l_t(p) \leq \nabla l_t(p_t) \cdot (p_t - p)$

Example: Meta-statistical framework:  $l_t(p) = \left( \sum_j p_j f_{jt} - y_t \right)^2$

$l_t$  is differentiable,  $\nabla l_t(p) = \left( 2 \left( \sum_j p_j f_{jt} - y_t \right) f_{jt} \right)_j$

↳ Gradients in  $[-G, G]^N$  with  $G = 2B^2$  if  $f_{jt}, y_t \in [0, B]$

Strategy: Exponentiated Gradients (EG) with learning rate  $\eta > 0$

$$p_t = (\gamma_1, \dots, \gamma_N) \quad \text{and} \quad \text{denoting } \tilde{l}_{js} = j\text{-th component of } \partial l_s(p_t);$$

$$p_{jt} = \exp(-\eta \sum_{s=1}^{t-1} \tilde{l}_{js}) / \sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} \tilde{l}_{ks})$$

↑ interpreted as pseudo-loss

Strategy easy to implement as soon as subgradients can be easily computed.

Theorem: Assume that the  $l_t$  are picked such that

$$\forall t \forall p \in \mathcal{X} \quad (\partial l_t(p))_j \in [-G, G]$$

Then the regret of EG tuned with  $\eta > 0$  is controlled as

$$R_T = \sum_{t=1}^T l_t(p_t) - \inf_{p \in \mathcal{X}} \sum_{t=1}^T l_t(p) \leq \frac{\ln N}{\eta} + \eta \frac{G^2}{2}$$

In particular, the choice  $\eta = \frac{1}{G} \sqrt{\frac{2 \ln N}{T}}$  leads to  $R_T \leq G \sqrt{2T \ln N}$ .

Proof: Fix a  $p$ . Since by construction (cf. exponential weights)  $p_t \in \mathcal{X}$ , we have by convexity

$$l_t(p_t) - l_t(p) \leq \partial l_t(p_t) \cdot (p_t - p) = \underbrace{\sum_{j=1}^N p_{jt} \tilde{l}_{jt}}_{\text{not.}} - \underbrace{\sum_{j=1}^N p_j \tilde{l}_{jt}}$$

→ Reduction to linear losses: let's pretend that the opponent picks vectors of losses

$$(\tilde{l}_{1t}, \dots, \tilde{l}_{Nt}) \in [-G, G]$$

By the above inequality,

$$\begin{aligned} R_T &\leq \sum_{t=1}^T \sum_{j=1}^N p_{jt} \tilde{l}_{jt} - \inf_{p \in \mathcal{X}} \sum_{t=1}^T \sum_{j=1}^N p_j \tilde{l}_{jt} \\ &= \sum_{t=1}^T \sum_{j=1}^N p_{jt} \tilde{l}_{jt} - \min_{\tilde{p}_t} \sum_{t=1}^N \tilde{l}_{jt} \leq \frac{\ln N}{\eta} + \eta \frac{(2G)^2}{8} T \end{aligned}$$

This concludes the proof.

Bound proved last week in the linear case

Application: Oracle inequalities  $\rightarrow$  From individual sequences to stochastic sequences

Question: Let  $(Y_1, \dots, Y_T)$  be a sequence of iid random variables taking values in an arbitrary set  $\mathcal{Y}$ ; and  $Q: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a loss function, convex in its first argument, with bounded subgradients:

with either  
 $Q \geq 0$   
or  $\forall \theta, Q(\theta, y) \in \mathbb{L}^1$   
for all expectations below  
to be defined

$$G = \sup_{y \in \mathcal{Y}} \sup_{\theta \in \mathcal{X}} \| \partial Q(\theta, y) \|_\infty < +\infty$$

Aim: Construct  $\hat{\theta}_T = \hat{\theta}_T(Y_1, \dots, Y_T)$  such that

$$\mathbb{E}[Q(\hat{\theta}_T, Y)] \leq \inf_{\theta \in \mathcal{X}} \mathbb{E}[Q(\theta, Y)] + \varepsilon_T$$

where  $\hat{Y}$  is independent of the  $Y_t$  with the same distribution and  $\varepsilon_T \rightarrow 0$ .

(The expectation  $\mathbb{E}$  is w.r.t.  $Y_1, \dots, Y_T$  and  $\hat{Y}$ .)

Typical machine learning method:

Empirical risk minimization, possibly in a regularized way:

$$\hat{\theta}_T \in \arg \min_{\theta \in \mathcal{X}} \left\{ \frac{1}{T} \sum_{t=1}^T Q(\theta, Y_t) + d \text{reg}(\theta) \right\}$$

with  $\text{reg}(\theta) = \|\theta\|_2$  or  $\|\theta\|_1$  or...  
and  $d$  to be tuned (called the regularization factor)

Our method: (as in our proof of Sion's lemma)

1) Pretend data is sequential, while of course it is batch

!

$$\tilde{\theta}_1 = (1, \dots, 1) \quad \text{and by induction, with } \eta = \frac{1}{G} \sqrt{\frac{2 \ln T}{T}}$$

Strongly depends on the order picked

$$\tilde{\theta}_{t+1} = \exp(-\eta \sum_{s=1}^{t+1} (\partial Q(\tilde{\theta}_s, Y_s))_j) / \sum_{k=1}^n \exp(-\eta \sum_{s=1}^{t+1} (\partial Q(\tilde{\theta}_s, Y_s))_k)$$

2) Consider an average:

$$\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \tilde{\theta}_t$$

(note: we do not use  $\bar{\cdot}$ !)

Guarantees:

$$\mathcal{E}_T = G \sqrt{\frac{2 \ln N}{T}}$$

Proof: By the theorem on EG (exponentiated gradient) and given our choice of  $\gamma$ :

$$\sum_{t=1}^T Q(\tilde{\theta}_t, y_t) - \inf_{\theta \in \mathcal{X}} \sum_{t=1}^T Q(\theta, y_t) \leq G \sqrt{2T \ln N}$$

Thus: therefore,

$$\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t) \leq \frac{1}{T} \sum_{t=1}^T Q(\theta, y_t) + \mathcal{E}_T$$

with  $\mathcal{E}_T = G \sqrt{\frac{2 \ln N}{T}}$

Now, since  $y_1, \dots, y_T$  and  $y$  are iid:

$$\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\theta, y_t)\right] = \mathbb{E}[Q(\theta, y)]$$

We conclude the proof by showing that  $\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t)\right] \geq \mathbb{E}[Q(\hat{\theta}_T, y)]$

Indeed:  $\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t)\right] = \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y)\right]$

$\tilde{\theta}_t$  only depends on  $y_1, \dots, y_{t-1}$  and

$(y_1, \dots, y_{t-1}, y_t)$  and  $(y_1, \dots, y_{t-1}, y)$  have the same distribution

by convexity and  $\mathbb{E}$  because the second argument is now the same for everyone!

We have shown: thus,

$$\mathbb{E}[Q(\hat{\theta}_T, y)] \leq \mathbb{E}[Q(\theta, y)] + \mathcal{E}_T$$

from which the conclusion follows.

Extension: Oracle inequality for stationary data.

A stationary sequence  $(y_1, y_2, \dots)$  is by definition a sequence of random variables such that

$\forall k \geq 1, \quad \forall t \geq 1, \quad (y_1, \dots, y_k) \text{ and } (y_{t+k}, \dots, y_{t+2k})$   
have the same distribution.

The strategy is much more complicated to state  
and we need first a piece of notation



- \* For a sequence of observations  $z_1, \dots, z_t \in \mathcal{Y}$ , EG outputs the following weights, which we re-define as functions of the information available:

$$\psi_1^{\text{EG}} = (\gamma_1, \dots, \gamma_N)$$

and for  $t \geq 2$ :

$$\psi_j^{\text{EG}}(z_1, \dots, z_{t-1}) = \frac{\exp(-\eta \sum_{s=1}^{t-1} (\Delta Q(\psi_s^{\text{EG}}(z_1, \dots, z_{s-1}), z_s)_j))}{\sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} (\Delta Q(\psi_k^{\text{EG}}(z_1, \dots, z_{s-1}), z_s))_k)}$$

j-th component of

$$\psi_t^{\text{EG}}(z_1, \dots, z_{t-1}) \in \mathbb{R}^N$$

!  $\psi_t^{\text{EG}}$  is an extremely complicated function of its arguments  $z_1, \dots, z_{t-1}$   
(cf. the various calls to  $\psi_s^{\text{EG}}$ , with  $s \leq t-1$ , in its definition).

- \* We define  $\hat{\theta}_t = \psi_t^{\text{EG}}(y_{T+2-t}, \dots, y_T)$  for  $t \geq 2$   
and  $\hat{\theta}_1 = (\gamma_1, \dots, \gamma_N)$

We recommend:

$$\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \hat{\theta}_t$$

! Computing  $\hat{\theta}_t$  is costly: the computation of each  $\hat{\theta}_t$  requires  $t-1$  steps, namely the computation of the

$$\psi_s^{\text{EG}}(y_{T+2-t}, \dots, y_{T-t}) \quad \text{for } 2 \leq s \leq t$$

$\nearrow$   
no such intermediate quantity appears twice

$$\Rightarrow \text{about } \sum_{t=1}^T (t-1) = O(T^2) \text{ computation steps required.}$$

What does EG guarantee?

$$\forall \theta \in \mathcal{X}, \quad \frac{1}{T} \sum_{t=1}^T Q(\psi_t^{\text{EG}}(y_1, \dots, y_{t-1}), y_t) \leq \frac{1}{T} \sum_{t=1}^T Q(\theta, y_t) + \varepsilon_T$$

still  $\mathcal{O}(\sqrt{2\ln T})$

We take expectations in both sides to conclude:

- stationarity implies in particular that all the  $y_t$  have the same distribution (but they are not independent in general); hence

$$\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\theta, y_t)\right] = \mathbb{E}[Q(\theta, y_T)]$$

- it also implies that  $(\psi_t^{\text{EG}}(y_1, \dots, y_{t-1}), y_t)$  and  $(\psi_t^{\text{EG}}(y_{T+2-t}, \dots, y_T), y_T) = (\hat{\theta}_t, y_T)$  have the same distribution, so that

$$\begin{aligned} & \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\psi_t^{\text{EG}}(y_1, \dots, y_{t-1}), y_t)\right] \\ &= \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T Q(\hat{\theta}_t, y_T)\right] \geq \mathbb{E}[Q(\hat{\theta}_T, y_T)] \end{aligned}$$

↑  
convexity  
of  $\theta \mapsto Q(\theta, y_T)$

- thus we proved, for our choice of  $\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \hat{\theta}_t$

$$\mathbb{E}[Q(\hat{\theta}_T, y_T)] \leq \inf_{\theta \in \mathcal{X}} \mathbb{E}[Q(\theta, y_T)] + \varepsilon_T$$

same  $\varepsilon_T$  as  
for the iid case.

## Sequential optimization for general convex sets

## Sequential optimization for general convex sets.

We so far focused our attention on the simplex

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

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  $R_T = \sum_{t=1}^T l_t(x_t) - \inf_{x \in \mathcal{C}} \sum_{t=1}^T l_t(x)$

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



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

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

where  $\Pi_{\mathcal{C}}$  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.,

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

Then  $OGD(\eta)$  is such that

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

↑  
inf achieved by continuity + compactness

In particular, for  $\eta = \frac{D}{G\sqrt{T}}$ ,

the bound equals  $DG\sqrt{T}$ .

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

$$\nabla \ell(\mathbf{p}), \quad \ell(\mathbf{p}) = \sum_{j=1}^N p_j \ell_j$$

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

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

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

equality achieved if  $p_i, q_i$  Dirac masses at  $i \neq k$

- bound on the gradients: if, e.g.,  $\ell_j \in [-M, M] \forall j$ ,

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

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

We get the bounds:

- for OGD:  $DG\sqrt{T} = M\sqrt{2NT}$

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

The dependency in  $N$  is suboptimal for OGD.

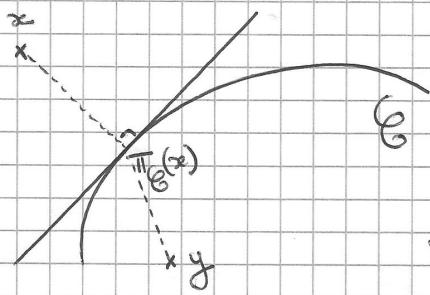
Proof (of the theorem): By convexity ( $\Leftarrow$  inégalité des pentes, in French):

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

Summing over  $t$ ,

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

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



Def:  $\Pi_C(x)$  is the (unique)  $y^* \in C$  with  $\|x-y^*\| \leq \|x-y\| \forall y \in C$

Property 1: If  $x \notin C$ , then  $\forall y \in C$ ,  $(y - \Pi_C(x)) \cdot (x - \Pi_C(x)) \leq 0$

Property 2:  $\forall y \in C$ ,  $\|y - \Pi_C(x)\| \leq \|x - x\|$

Proof: (Property 1):  $t \in [0,1]$ ,  $y_t = \Pi_C(x) + t(y - \Pi_C(x)) \in C$  by convexity

$$\begin{aligned} \|x - y_t\|^2 &= \|x - \Pi_C(x) + t(y - \Pi_C(x))\|^2 + t^2 \|y - \Pi_C(x)\|^2 \\ &\quad - 2t(y - \Pi_C(x)) \cdot (x - \Pi_C(x)) \geq \|x - \Pi_C(x)\|^2 \end{aligned}$$

Simplify the  $\|x - \Pi_C(x)\|^2$ , divide by  $t$  and  $t \rightarrow 0$ .

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

$$+ 2(y - \Pi_C(x)) \cdot (\Pi_C(x) - x) \quad \text{(by Prop 1)}$$

In particular, for  $t \geq 1$ :

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

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

$$\begin{aligned} &\nabla l_t(x_t) \cdot (x_t - x) \\ &\leq \frac{1}{2\eta} \left( \|x_t - x\|^2 - \|x_{th} - x\|^2 \right) + \frac{1}{2} \|\nabla l_t(x_t)\|^2 \end{aligned}$$

Summing over  $t$ :

$$\sum_{t=1}^T \nabla l_t(x_t) \cdot (x_t - x) \leq G^2$$

by the assumption of bounded gradients

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

Adaptation to time  $T$  (and other parameters)

A popular way of being adaptive is to tune  $\eta$  according to the past observations.

Adaptive OGD:

- Play an arbitrary  $x_1 \in \mathcal{G}$ , pick  $\eta_1 > 0$
- For  $t \geq 2$ , play  $x_t = \Pi_{\mathcal{G}}(x_{t-1} - \eta_t \nabla f_{t-1}(x_{t-1}))$   
pick  $\eta_t > 0$

Theorem: If Adaptive OGD picks a non-increasing sequence  $(\eta_t)_t$  with  $\eta_t > 0$  then:

$$R_T \leq \frac{D^2}{2\eta_1} + \frac{1}{2} \sum_{t=1}^T \eta_t \|\nabla f_t(x_t)\|^2 \quad \text{with the same notation as in the Theorem for OGD}(\eta).$$

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

Adaptive OGD achieves the bound

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

without any prior knowledge of neither  $T$  nor  $\mathcal{G}$

$$\leq \max \left\{ D^2, \frac{3}{2} D G \sqrt{T} \right\}$$

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

Proof of the theorem and of the corollary:

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

$$\text{True} \ L, \quad \nabla l_t(x_t) \cdot (x_t - x) \leq \frac{1}{2\eta_t} (\|x_t - x\|^2 - \|x_{t+1} - x\|^2) + \frac{\eta_t}{2} \|\nabla l_t(x_t)\|^2$$

Summing over  $t$ :

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

Abel's transform:

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

Cor: Recall that we picked:

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

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

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

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

The regret bound reads

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

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

We denote  $G$  the largest  $t \leq T$  such that  $V_t < D^2$ . Then, if  $G < T$ , we have  $V_G \geq D^2$ .  
(possibly  $t=0$  if  $V_1 \geq D^2$ )

The above bound can be rewritten as:

$$\frac{D^2}{2\eta_T} + \frac{1}{2} \sum_{\substack{t: \\ V_t < D^2}} \|v_t(x_t)\|^2 + \frac{1}{2} \sum_{t=G+1}^T \frac{D}{\sqrt{V_t}} \|v_t(x_t)\|^2$$

$$= \underbrace{\frac{V_G}{2} < \frac{D^2}{2}}_{\text{possibly } V_G \text{ small } (\text{if } G=T)}$$

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

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

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

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

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

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