# An Intuitive Proof of MWU

#### April 10, 2023

- Recall the setting.
  - There are n experts and T days.
  - At each day  $t = 1, 2, \ldots, T$ , the following happen in order:
    - 1. We choose a distribution  $p^t = (p_1^t, \ldots, p_n^t)$  of experts
    - 2. Then, the **adversary** reveals the **loss vector**  $\ell^t = (\ell_1^t, \ldots, \ell_n^t)$
    - 3. Our (expected) loss is  $\ell_A^t = \langle p^t, \ell^t \rangle = \sum_i p_i^t \ell_i^t$ .
- Goal: our total loss  $L_A = \sum_{t=1}^T \ell_A^t$  is as good as the total loss of the best expert  $L_{\star} = \min_i \sum_{t=1}^T \ell_i^t$ .
- We want to analyze the MWU algorithm:

Algorithm 1 The Multiplicative Weights Update (MWU) algorithm  $\overline{MWU(\epsilon)}$ :

- $w_i^1 \leftarrow 1 \quad \forall i$
- For t = 1, ..., T

- Follows expert i with probability  $p_i^t = \frac{w_i^t}{\sum_j w_j^t}$ - After  $\ell^t$  is revealed,  $w_i^{t+1} \leftarrow w_i^t \cdot \exp(-\epsilon \ell_i^t) \quad \forall i$ 

**Theorem 0.1.** For  $0 < \epsilon \leq 1/2$ , MWU( $\epsilon$ ) guarantees that

$$L_A \le L_\star + \epsilon T + \frac{\ln n}{\epsilon}$$

**Theorem 0.2.** Suppose  $\ell^t \in [0,1]^n$  for all t. For  $0 < \epsilon \le 1/2$ , MWU( $\epsilon$ ) guarantees that

$$L_A \le L_\star + \epsilon L_A + \frac{\ln n}{\epsilon}$$

SO

$$L_A \le (1+2\epsilon)(L_\star + \frac{\ln n}{\epsilon})$$

• Most proofs in the literature<sup>1</sup> involved a lot of low-level calculation and the intuition is lost.

<sup>&</sup>lt;sup>1</sup>https://lucatrevisan.github.io/40391/lecture11.pdf

http://www.theoryofcomputing.org/articles/v008a006/v008a006.pdf

- The proof here:
  - All low-level calculation is modularized into the soft-min function.
  - Using soft-min, high-level explanation is very intuitive
- The proof is inspired by the paper of Quanrud<sup>2</sup> and the blog of Zuzic<sup>3</sup> (they work in a more specific context of solving LPs).

## 1 Basic Calculus: Approximating Smooth Functions

- Let  $f : \mathbb{R} \to \mathbb{R}$  be a **smooth** function.
- Suppose you know f(x). Can we approximate  $f(x + \delta)$ ? Use Taylor expansion

$$f(x+\delta) = f(x) + f'(x)\delta + \frac{1}{2}f''(x)\delta^{2} + \dots$$

and, when  $\delta$  is tiny, we have

$$f(x+\delta) \approx f(x) + f'(x)\delta$$

- What if  $f : \mathbb{R}^n \to \mathbb{R}$ ?
  - The first derivative  $f'(\cdot) \in \mathbb{R}^n$  is called the **gradient** of f.
  - The second derivative  $f''(\cdot) \in \mathbb{R}^{n \times n}$  is called the **Hessian** of f.
- Suppose you know  $f(\vec{x})$ . Can we approximate  $f(\vec{x} + \vec{\delta})$ ? Taylor expansion:

$$f(\vec{x} + \vec{\delta}) = f(\vec{x}) + \left\langle f'(\vec{x}), \vec{\delta} \right\rangle + \frac{1}{2} \left\langle \vec{\delta}, f''(\vec{x})\vec{\delta} \right\rangle + \dots$$

and, when  $\|\vec{\delta}\|$  is tiny, we have

$$f(\vec{x}+\delta) \approx f(\vec{x}) + \left\langle f'(\vec{x}), \vec{\delta} \right\rangle$$

### 2 Overall Plan

- Let's forget the MWU algorithm you saw.
  - Suppose that we want to derive a very natural algorithm such that  $L_A \leq L_{\star} + \epsilon T + \frac{\ln n}{\epsilon}$ .
  - The algorithm we derived will be MWU exactly.
- Notation:
  - $\begin{aligned} & L^t = \sum_{t'=1}^t \ell^{t'} \in \mathbb{R}^n \text{ encode total loss of all experts after day } t \text{ (after } \ell^t \text{ is revealed)}. \\ & L^t_\star = \min(L^t) := \min_i(L^t_i) \text{ be the total loss of the best expert after day } t. \\ & L^t_A = \sum_{t'=1}^t \left\langle p^{t'}, \ell^{t'} \right\rangle \text{ is our total loss after day } t. \end{aligned}$
- Wishful hope: whenever  $L_A^t$  increases,  $L_{\star}^t$  increases by at least the same amount.

<sup>&</sup>lt;sup>2</sup>https://epubs.siam.org/doi/abs/10.1137/1.9781611976014.11

<sup>&</sup>lt;sup>3</sup>https://zuza.github.io/Intuitive-Multiplicative-Weights/

- So after T days, we have  $L_A^T \leq L_{\star}^T$ .
- No regret at all. But this is clearly too good to hope for.
- Suppose, magically, there exists a function  $smin(L^t)$  where
  - $\operatorname{smin}(L^t) \approx \min(L^t) = L^t_\star,$
  - but  $\operatorname{smin}(L^t)$  is smooth.
- Actual plan: whenever  $L_A^t$  increases,  $smin(L^t)$  increases by at least the same amount (modulo some small error).
- For each day,
  - The change in smin $(L^t)$  (which is a proxy for  $L^t_{\star}$ ) is

$$\approx \langle \operatorname{smin}'(L^{t-1}), \ell^t \rangle$$

because  $smin(\cdot)$  is smooth.

- By Taylor, we have

$$\operatorname{smin}(L^t) = \operatorname{smin}(L^{t-1} + \ell^t) \approx \operatorname{smin}(L^{t-1}) + \left\langle \operatorname{smin}'(L^{t-1}), \ell^t \right\rangle,$$

- \* The approximation holds when  $\ell^t$  is tiny compared to  $L^{t-1}$ .
- \* We will not assume this assumption in formal proof (but there will be small error term).
- the change in  $L_A^t$  is exactly  $\langle p^t, \ell^t \rangle$  by definition:

$$L_A^t = L_A^{t-1} + \left< p^t, \ell^t \right>$$

- So... how should we set our distribution  $p^t$  on each day t? Clearly, we should set

$$p^t \leftarrow \min'(L^{t-1})$$

so that  $\left< \min'(L^{t-1}), \ell^t \right> = \left< p^t, \ell^t \right>$ .

• So every time  $L_A^t$  increases,  $\operatorname{smin}(L^t)$  increases by almost the same amount. As  $L_{\star}^t \approx \operatorname{smin}(L^t)$ , we should have

 $L_A^T \leq L_{\star}^t + (\text{error from Taylor expansion}) + (\text{error from smin})$ 

## 3 Soft-Max/Min: Smooth version of Max/Min

- Let  $x = (x_1, \ldots, x_n)$  be a vector.
- The soft-max of x is  $\operatorname{smax}(x) = \ln(\sum_i \exp(x_i))$ .
- Why does it behave like max?
  - If  $x_{i_0} \gg x_i$  for all i,  $\exp(x_{i_0}) \approx \sum_i \exp(x_i)$ . So  $\ln(\sum_i \exp(x_i)) \approx x_{i_0}$ .
  - So smax $(x) \approx \max(x)$ .
- The soft-min of x is smin(x) = -smax(-x).

- Note  $\min(x) = -\max(-x)$ .
- We will consider the scaled version:  $smin_{\epsilon}(x) = -\frac{1}{\epsilon}smax(-\epsilon x)$ .

**Proposition 3.1.** We prove this at the end.

- $\min(x) \frac{\ln n}{\epsilon} \le \min_{\epsilon}(x) \le \min(x)$
- $\operatorname{smin}_{\epsilon}'(x)_i = \frac{\exp(-\epsilon x_i)}{\sum_j \exp(-\epsilon x_j)}$
- When  $\epsilon \leq 1/2$  and  $\|\delta\|_{\infty} \leq 1$ ,  $\min_{\epsilon}(x+\delta) \geq \min_{\epsilon}(x) + \langle \min'_{\epsilon}(x), \delta \rangle \epsilon \frac{\sum_{i} \exp(-\epsilon x_{i})\delta_{i}^{2}}{\sum_{j} \exp(-\epsilon x_{j})}$

## 4 Deriving MWU exactly

- We use  $\operatorname{smin}_{\epsilon}(L^t)$  to approximate  $L^t_{\star}$ .
- Recall that we want to set  $p^t$  as the gradient of soft-min:

$$p^t \leftarrow \min_{\epsilon}'(L^{t-1})$$

 $\mathbf{SO}$ 

$$p_i^t = \frac{\exp(-\epsilon L_i^{t-1})}{\sum_j \exp(-\epsilon L_j^{t-1})}$$

- But what is this?
  - Recall the MWU algorithm: the weight of expert *i* is  $w_i^{t+1} = w_i^t \cdot \exp(-\epsilon \ell_i^t)$ .
  - So  $w_i^t = \exp(-\epsilon L_i^{t-1})$ .
  - Therefore,

$$p_i^t = \frac{w_i^t}{\sum_j w_j^t}.$$

- We just derived the MWU algorithm!
- Let's analyze it.
- We analyze how much  $L_A^t$  and  $smin_{\epsilon}(L^t)$  change each day.
  - $L_A^t$  increase exactly  $\langle p^t, \ell^t \rangle$  by definition:

$$L_A^t - L_A^{t-1} = \left\langle p^t, \ell^t \right\rangle$$

- How about smin<sub> $\epsilon$ </sub>( $L^t$ )? Via Taylor's expansion, we have

$$\begin{aligned} \min_{\epsilon}(L^{t}) - \min_{\epsilon}(L^{t-1}) &= \min_{\epsilon}(L^{t-1} + \ell^{t}) - \min_{\epsilon}(L^{t-1}) \\ &\geq \left\langle \min_{\epsilon}'(L^{t-1}), \ell^{t} \right\rangle - \epsilon \frac{\sum_{i} \exp(-\epsilon L_{i}^{t-1})(\ell_{i}^{t})^{2}}{\sum_{j} \exp(-\epsilon L_{j}^{t-1})} \quad \text{as } \epsilon \leq 1/2, \|\ell^{t}\| \leq 1. \\ &\geq \left\langle p^{t}, \ell^{t} \right\rangle - \epsilon \left\langle p^{t}, (\ell^{t})^{2} \right\rangle \end{aligned}$$

where  $(\ell^t)^2 := (\ell^t_i)^2$  for all i.

- That is,  $\min_{\epsilon}(L^t)$  increases at least as much as  $L_A^t$  increases (except some error term  $\epsilon \langle p^t, (\ell^t)^2 \rangle$ ).
- Summing over all t, the sum telescopes to

$$\operatorname{smin}_{\epsilon}(L^{T}) - \operatorname{smin}_{\epsilon}(L^{0}) \ge \sum_{t=1}^{T} \langle p^{t}, \ell^{t} \rangle - \epsilon \langle p^{t}, (\ell^{t})^{2} \rangle$$
$$= L_{A}^{T} - \epsilon \sum_{t=1}^{T} \langle p^{t}, (\ell^{t})^{2} \rangle$$

• So

$$\begin{split} L_A^T &\leq \min_{\epsilon} (L^T) - \min_{\epsilon} (L^0) + \epsilon \sum_{t=1}^T \left\langle p^t, (\ell^t)^2 \right\rangle \\ &\leq \min(L^T) + \frac{\ln n}{\epsilon} + \epsilon \sum_{t=1}^T \left\langle p^t, (\ell^t)^2 \right\rangle \\ &\leq L_\star^T + \frac{\ln n}{\epsilon} + \epsilon T \end{split}$$

because  $\langle p^t, (\ell^t)^2 \rangle \leq 1$ . This proves Theorem 0.1.

• If  $\ell^t_i \in [0, 1]$ , then we have  $(\ell^t_i)^2 \leq \ell^t_i$ . So

$$\sum_{t=1}^{T} \left\langle p^{t}, (\ell^{t})^{2} \right\rangle \leq \sum_{t=1}^{T} \left\langle p^{t}, \ell^{t} \right\rangle = L_{A}^{T}$$

and so we get

$$L_A^T \le L_\star^T + \frac{\ln n}{\epsilon} + \epsilon L_A^T,$$

which proves Theorem 0.2

### 5 Recap: How to Derive MWU

- How to design an algorithm for choosing experts with no regret property?
- Wishful hope: whenever  $L_A^t$  increases,  $L_{\star}^t = \min(L^t)$  increases by at least the same amount.
  - But min is not smooth.
- Actual plan: whenever  $L_A^t$  increases,  $\min(L^t) \approx \min(L^t)$  increases by at least the same amount (modulo some small error).
  - Since smin is smooth, we now can approximate the change of  $\operatorname{smin}(L^t)$  per day, which is just  $\langle \operatorname{smin}'_{\epsilon}(L^{t-1}), \ell^t \rangle$
  - Since the change of  $L_A^t = \langle p^t, \ell^t \rangle$ , we set  $p^t = \operatorname{smin}'_{\epsilon}(L^{t-1})$ . So

$$p_i^t = \frac{\exp(-\epsilon L_i^{t-1})}{\sum_j \exp(-\epsilon L_j^{t-1})}.$$

- This is the same thing as saying that  $p^t \sim w^t$  where  $w^t = \exp(-\epsilon L_i^{t-1})$ .
- To have  $w^t = \exp(-\epsilon L_i^{t-1})$ , we just need to multiplicatively update weights

$$w_i^{t+1} \leftarrow w_i^t \cdot \exp(-\epsilon \ell_i^t)$$

• This is the exactly what happens in MWU. Done.

## 6 Low-Level Calculation about Soft-max

• Let  $x \in \mathbb{R}^n$ . Define  $sx(x) = \sum_j exp(x_j)$  for "sum-exp".

#### 6.1 Upper/Lower Bounds

- $\operatorname{smin}_{\epsilon}(x) = \min(x)$  when x concentrates on only 1 entry.
- $\operatorname{smin}_{\epsilon}(x) = \min(x) \frac{\ln n}{\epsilon}$  when x has uniform value on all entries.
- For other x, we have  $\min(x) \frac{\ln n}{\epsilon} \leq \min(x) \leq \min(x)$  but I omit the proof.

#### 6.2 Gradient

• We just compute

$$(\operatorname{smin}_{\epsilon}'(x))_{i} = \frac{\partial}{\partial x_{i}} \frac{-1}{\epsilon} (\ln \operatorname{sx}(-\epsilon x))$$
$$= \frac{-1}{\epsilon} \cdot \frac{1}{\operatorname{sx}(-\epsilon x)} \cdot \frac{\partial}{\partial x_{i}} \operatorname{sx}(-\epsilon x)$$
$$= \frac{-1}{\epsilon} \cdot \frac{1}{\operatorname{sx}(-\epsilon x)} \cdot \exp(-\epsilon x_{i}) \cdot (-\epsilon)$$
$$= \frac{1}{\operatorname{sx}(-\epsilon x)} \cdot \exp(-\epsilon x_{i})$$

### 6.3 Smooth

• When  $\epsilon \leq 1/2$  and  $\|\delta\|_{\infty} \leq 1$ , we have

$$\operatorname{smin}_{\epsilon}(x+\delta) \ge \operatorname{smin}_{\epsilon}(x) + \left\langle \operatorname{smin}'_{\epsilon}(x), \delta \right\rangle - \epsilon \frac{\sum_{i} \exp(-\epsilon x_{i}) \delta_{i}^{2}}{\operatorname{sx}(-\epsilon x)}$$

### • This is because

$$-\operatorname{smin}_{\epsilon}(x+\delta) = \frac{1}{\epsilon} \ln\left(\sum_{i} \exp(-\epsilon x_{i} - \epsilon \delta_{i})\right)$$

$$= \frac{1}{\epsilon} \ln\left(\sum_{i} \exp(-\epsilon x_{i}) \exp(-\epsilon \delta_{i})\right)$$

$$\leq \frac{1}{\epsilon} \ln\left(\sum_{i} \exp(-\epsilon x_{i})(1 - \epsilon \delta_{i} + (\epsilon \delta_{i})^{2}\right) \qquad e^{a} \leq 1 + a + a^{2} \forall |a| \leq 1/2$$

$$= \frac{1}{\epsilon} \ln\left(\sum_{i} \exp(-\epsilon x_{i}) - \epsilon \sum_{i} \exp(-\epsilon x_{i})\delta_{i} + \epsilon^{2} \sum_{i} \exp(-\epsilon x_{i})\delta_{i}^{2}\right)$$

$$= \frac{1}{\epsilon} \ln\left(\operatorname{sx}(-\epsilon x)(1 - \epsilon \frac{\sum_{i} \exp(-\epsilon x_{i})\delta_{i}}{\operatorname{sx}(-\epsilon x)} + \epsilon^{2} \frac{\sum_{i} \exp(-\epsilon x_{i})\delta_{i}^{2}}{\operatorname{sx}(-\epsilon x)}\right)$$

$$= -\operatorname{smin}_{\epsilon}(x) + \frac{1}{\epsilon} \ln\left(1 - \epsilon \langle \operatorname{smin}_{\epsilon}'(x), \delta \rangle + \epsilon^{2} \frac{\sum_{i} \exp(-\epsilon x_{i})\delta_{i}^{2}}{\operatorname{sx}(-\epsilon x)}\right)$$

$$= -\operatorname{smin}_{\epsilon}(x) + \frac{1}{\epsilon} \left(-\epsilon \langle \operatorname{smin}_{\epsilon}'(x), \delta \rangle + \epsilon^{2} \frac{\sum_{i} \exp(-\epsilon x_{i})\delta_{i}^{2}}{\operatorname{sx}(-\epsilon x)}\right)$$

$$\ln(1 + a) \leq \epsilon$$

$$= -\operatorname{smax}_{\epsilon}(x) - \langle \operatorname{smax}_{\epsilon}'(x), \delta \rangle + \epsilon \frac{\sum_{i} \exp(-\epsilon x_{i})\delta_{i}^{2}}{\operatorname{sx}(-\epsilon x)}$$

Question 6.1. Can you prove this?

$$\operatorname{smin}_{\epsilon}(x+\delta) \ge \operatorname{smin}_{\epsilon}(x) + \left\langle \operatorname{smin}_{\epsilon}'(x), \delta \right\rangle - \delta^{\top} \operatorname{smax}_{\epsilon}''(x)\delta.$$

I don't know how to prove it.

- This is more beautiful, cleaner, and stronger than the bound we proved above.
- This is because, for any y

$$y^{\top} \operatorname{smin}_{\epsilon}''(x) y > \epsilon \frac{\sum_{i} \exp(-\epsilon x_{i}) y_{i}^{2}}{\sum_{j} \exp(-\epsilon x_{j})}.$$