

# Recursive Exponential Weighting for Online non-convex optimization

Speaker: Yi-Shan Wu

Institute of Information Science

Academia Sinica

Taiwan

Feb, 02, 2018

# Online Convex Optimization

For  $t = 1$  to  $T$

- player chooses a point  $x_t$  from  $\mathcal{K} \subset \mathbb{R}^n$  (bounded convex set)
- cost function  $f_t : \mathcal{K} \mapsto \mathbb{R}$  is revealed (bounded convex function)

Regret:

$$R_T = \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t=1}^T f_t(x_t) - \min_{x \in \mathcal{K}} \sum_{t=1}^T f_t(x) \right\}$$

- OCO lower bound:  $O(\sqrt{T})$  obtained by OGD, SGD,  $\dots$
- If further assumes that  $f_t$  is strictly convex, OCO lower bound:  $O(\ln T)$

A natural question: Is it possible to obtain  $O(\sqrt{T})$  regret bound if the convexity assumption on the cost function is relaxed?

# Recursive Exponential Weighting Algorithm

**Assumption:** cost function  $f_t \in \mathcal{F}$ ,  $f_t(x) : \mathcal{K} \mapsto [0, B]$ , and it is  $L$ -Lipschitz

## REW Algorithm:

### ① Set Discretization from $\mathcal{K}$ to $\mathcal{I}$

Idea: Group highly correlated decisions into one set.

$\implies$  reduce the original problem to expert problem with finite experts

### ② Set partition according to layers

Idea: Divide and conquer, recursively digging into subsets that are chosen.

## Analysis:

$$R_{ImC} = \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t=1}^T c_t(I_t) - \min_{i \in \mathcal{I}} \sum_{t=1}^T c_t(i) \right\}$$

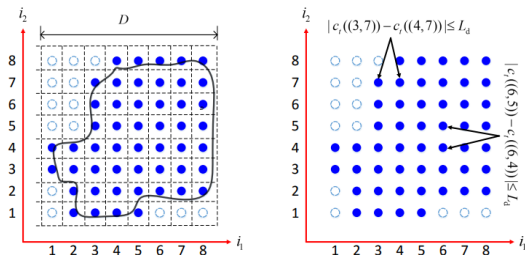
$$R_{ImD} = \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \min_{i \in \mathcal{I}} \sum_{t=1}^T c_t(i) - \min_{x \in \mathcal{K}} \sum_{t=1}^T f_t(x) \right\}$$

# Set Discretization

- Dimension of  $\mathcal{K}$ :  $n$
- edge length of sub-cube:  $D/2^m$
- index of sub-cube  $D_{\mathbf{i}} : \mathbf{i} = (i_1, i_2, \dots, i_n), \quad 1 \leq i_j \leq 2^m, \forall j \in [n]$
- $\mathcal{I} = \{\mathbf{i} : D_{\mathbf{i}} \cap \mathcal{K} \neq \emptyset\} \Rightarrow |\mathcal{I}|$  experts
- cost of  $\mathbf{i} \in \mathcal{I}$  at time  $t : c_t(\mathbf{i})$

$$|c_t(\mathbf{p}) - c_t(\mathbf{q})| \leq L_d \|\mathbf{p} - \mathbf{q}\|_1, \quad \mathbf{p}, \mathbf{q} \in \mathcal{I}$$

where  $L_d = L\sqrt{n}2D/2^m$



# Set Partition for $\mathcal{I}$

$$\mathcal{I}_\ell(\mathbf{i}) = \{\mathbf{p} \in \mathcal{I} : 1 + (i_j - 1)2^{m-\ell} \leq p_j \leq 2^{m-\ell} + (i_j - 1)2^{m-\ell}, j = 1, 2, \dots, n\}.$$

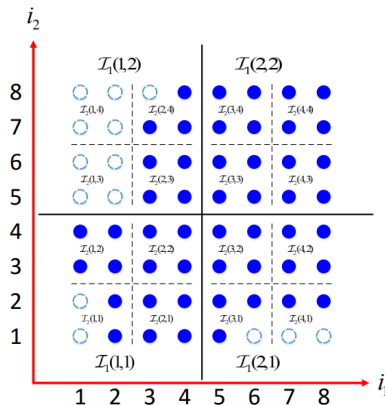


Figure: Set partition when  $n = 2$  and  $m = 3$

# Set Partition for $\mathcal{I}$

- $\mathcal{U}_\ell(\mathbf{i})$  : the layer- $\ell$  subset containing lower-layer subset  $\mathcal{I}_{\ell+1}(\mathbf{i})$

$$\mathcal{U}_\ell(\mathbf{i}) = \{\mathbf{j} : \mathcal{I}_{\ell+1}(\mathbf{i}) \subset \mathcal{I}_\ell(\mathbf{j}) : \mathcal{I}_{\ell+1}(\mathbf{i}) \neq \emptyset\}.$$

**Ex:**  $\mathcal{U}_1(3, 2) = \{(2, 1)\}, \quad (\mathcal{I}_2(3, 2) \subset \mathcal{I}_1(2, 1))$

- $\mathcal{M}_\ell(\mathbf{k})$  : the layer- $\ell$  subsets that contain point  $\mathbf{k}$

**Ex:**  $\mathcal{M}_2(5, 3) = \{(3, 2)\}$

- $\mathcal{D}_\ell(\mathbf{i})$  : the index set of non-empty layer- $(\ell + 1)$  subsets within  $\mathcal{I}_\ell(\mathbf{i})$

$$\mathcal{D}_\ell(\mathbf{i}) = \{\mathbf{j} : \mathcal{I}_{\ell+1}(\mathbf{j}) \subset \mathcal{I}_\ell(\mathbf{i}) : \mathcal{I}_{\ell+1}(\mathbf{j}) \neq \emptyset\}.$$

**Ex:**  $\mathcal{D}_2(3, 2) = \{(5, 3), (6, 3), (5, 4), (6, 4)\}$

$$\bar{C}_{\ell+1,t}(\mathbf{i}) = \sum_{\tau=1}^t \bar{c}_{\ell+1,\tau}(\mathbf{i})$$

```

2: for  $t = 1$  to  $T$  do
3:    $\mathbf{s} = \mathbf{1}$ 
4:   //Recursively select the subsets in all layers
5:   for  $l = 0$  to  $m - 1$  do
6:     Select a non-empty subset  $\mathcal{I}_{l+1}(\mathbf{i}) \subset \mathcal{I}_l(\mathbf{s})$  with probability

```

$$p_t(\mathcal{I}_{l+1}(\mathbf{i})) = \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}))}{\sum_{\mathbf{i} \in \mathcal{D}_l(\mathbf{s})} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}))}$$

```

7:     if subset  $\mathcal{I}_{l+1}(\mathbf{i})$  is selected then
8:       Set  $\mathbf{s} = \mathbf{i}$ 
9:     end if
10:  end for

```

Figure: Selection

# REW-cumulative expected normalized cost

$$\begin{aligned} & \Pr[\mathbf{I}_t = \mathbf{k} | \mathbf{I}_t \in \mathcal{I}_{\ell+1}(\mathbf{i})] \\ &= \prod_{i=\ell+2}^m \Pr[\mathbf{I}_t \in \mathcal{M}_i(\mathbf{k}) | \mathbf{I}_t \in \mathcal{M}_{i-1}(\mathbf{k})] \end{aligned}$$

```
11: //Get the expected normalized cost for all subsets in all layers
12: for  $l = 0$  to  $m - 1$  do
13:   for each non-empty subset  $\mathcal{I}_{l+1}(\mathbf{i})$  in layer- $(l + 1)$  do
14:     Calculate  $\Pr[\mathbf{I}_t = \mathbf{k} | \mathbf{I}_t \in \mathcal{I}_{l+1}(\mathbf{i})]$  based on Equation (3.3)
15:     Calculate
```

$$\begin{aligned} & \bar{c}_{l+1,t}(\mathbf{i}) \\ &= \sum_{\mathbf{k} \in \mathcal{I}_{l+1}(\mathbf{i})} \frac{c_t(\mathbf{k}) - \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k})}{n2^{m-l}L_d} \cdot \Pr[\mathbf{I}_t = \mathbf{k} | \mathbf{I}_t \in \mathcal{I}_{l+1}(\mathbf{i})] \end{aligned}$$

```
16:   end for
17: end for
18: end for
```

Figure: cumulative cost



# Analysis

$$R_T = R_{ImC} + R_{ImD}$$

$$\begin{aligned} R_{ImC} &= \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t=1}^T c_t(l_t) - \min_{i \in \mathcal{I}} \sum_{t=1}^T c_t(i) \right\} \\ &\leq (4n^2 + \frac{1}{2}n)2^m L_d \sqrt{T} + 2n^2 \cdot 2^m L_d. \end{aligned}$$

$$\begin{aligned} R_{ImD} &= \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \min_{i \in \mathcal{I}} \sum_{t=1}^T c_t(i) - \min_{x \in \mathcal{K}} \sum_{t=1}^T f_t(x) \right\} \\ &\leq \frac{1}{2} L_d T \end{aligned}$$

Note that  $L_d = \frac{2L\sqrt{nD}}{2^m}$

Choose  $m = \log_2 \sqrt{T}$ , then we obtain  $O(n^2 \sqrt{T})$

$$\begin{aligned}
&= \sum_{t \in T} \mathbb{E}[c_t(l_t)] - \sum_{t \in T} c_t(i^*) \\
&= \sum_{t \in T} \mathbb{E}[c_t(l_t) | l_t \in \mathcal{M}_0(i^*)] - \sum_{t \in T} c_t(i^*) \\
&\leq \sum_{t \in T} \mathbb{E}[c_t(l_t) | l_t \in \mathcal{M}_1(i^*)] + (2n + \frac{1}{4})\sqrt{T}n2^mL_d + \frac{\ln |\mathcal{D}_\ell(0)|}{\eta_1} - \sum_{t \in T} c_t(i^*) \\
&\vdots \\
&\leq (4n^2 + \frac{n}{2})2^mL_d\sqrt{T} + 2n^22^mL_d
\end{aligned}$$

# REW-cumulative expected normalized cost

$$\begin{aligned}
& \sum_{t \in \mathcal{T}} \mathbb{E}[c_t(\mathbf{I}_t) | \mathbf{I}_t \in \mathcal{I}_l(\mathbf{i})] \\
&= \sum_{t \in \mathcal{T}} \sum_{j \in \mathcal{D}_l(\mathbf{i})} \Pr[\mathbf{I}_t \in \mathcal{I}_{l+1}(\mathbf{j}) | \mathbf{I}_t \in \mathcal{I}_l(\mathbf{i})] \mathbb{E}[c_t(\mathbf{I}_t) | \mathbf{I}_t \in \mathcal{I}_{l+1}(\mathbf{j})] \\
&\stackrel{(a)}{=} \sum_{t \in \mathcal{T}} \sum_{j \in \mathcal{D}_l(\mathbf{i})} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{j}))}{\sum_{j \in \mathcal{D}_l(\mathbf{i})} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{j}))} \cdot \left[ \bar{c}_{l+1,t}(\mathbf{j}) \cdot n 2^{m-l} L_d + \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k}) \right] \\
&\stackrel{(b)}{\leq} \sum_{t \in \mathcal{T}} \left[ -\frac{1}{\eta_t} \ln \sum_{j \in \mathcal{D}_l(\mathbf{i})} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{j}))}{\sum_{j \in \mathcal{D}_l(\mathbf{i})} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{j}))} \exp(-\eta_t \bar{c}_{l+1,t}(\mathbf{j})) + \frac{\eta_t}{8} \cdot 1^2 \right] \cdot n 2^{m-l} L_d \\
&\quad + \sum_{t \in \mathcal{T}} \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k}) \\
&= \sum_{t \in \mathcal{T}} \left[ -\frac{1}{\eta_t} \ln \sum_{j \in \mathcal{D}_l(\mathbf{i})} \frac{\exp(-\eta_t \bar{C}_{l+1,t}(\mathbf{j}))}{\sum_{j \in \mathcal{D}_l(\mathbf{i})} \exp(-\eta_t \bar{C}_{l+1,t}(\mathbf{j}))} + \frac{\eta_t}{8} \cdot 1^2 \right] \cdot n 2^{m-l} L_d + \sum_{t \in \mathcal{T}} \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k}) \\
&\stackrel{(c)}{=} \sum_{t \in \mathcal{T}} \left[ \Phi_t(\eta_t) - \Phi_{t-1}(\eta_t) + \frac{\eta_t}{8} \cdot 1^2 \right] \cdot n 2^{m-l} L_d + \sum_{t \in \mathcal{T}} \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k}) \\
&\stackrel{(d)}{\leq} \left\{ \Phi_T(\eta_T) + \sum_{t=2}^T \left[ n \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \frac{\eta_t}{8} \right] + \frac{1}{\eta_1} \ln |\mathcal{D}_l(\mathbf{i})| \right\} \cdot n 2^{m-l} L_d + \sum_{t \in \mathcal{T}} \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k}) \\
&\stackrel{(e)}{\leq} \left\{ \sum_{t \in \mathcal{T}} \mathbb{E} \left[ \frac{c_t(\mathbf{I}_t) - \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k})}{n 2^{m-l} L_d} | \mathbf{I}_t \in \mathcal{I}_{l+1}(\mathbf{j}) \right] \right. \\
&\quad \left. + \frac{1}{\eta_T} n + \sum_{t=2}^T \left[ n \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \frac{\eta_t}{8} \right] + \frac{1}{\eta_1} \ln |\mathcal{D}_l(\mathbf{i})| \right\} \cdot n 2^{m-l} L_d + \sum_{t \in \mathcal{T}} \min_{\mathbf{k} \in \mathcal{I}_l(\mathbf{i})} c_t(\mathbf{k}) \\
&\stackrel{(f)}{\leq} \sum_{t \in \mathcal{T}} \mathbb{E}[c_t(\mathbf{I}_t) | \mathbf{I}_t \in \mathcal{I}_{l+1}(\mathbf{j})] + \left( 2n + \frac{1}{4} \right) \sqrt{T} \cdot n 2^{m-l} L_d + \frac{1}{\eta_1} \ln |\mathcal{D}_l(\mathbf{i})| \cdot n 2^{m-l} L_d
\end{aligned}$$

$$\mathbf{MWC} : \max_{f \in \mathcal{F}} \sum_{r \in R} x_r(f). \quad (1)$$

An equivalent formulation is

$$\begin{aligned} \mathbf{MWC - EQ} : \quad & \max_{p \geq 0} \sum_{f \in \mathcal{F}} p_f \sum_{r \in R} x_r(f) \\ & \text{s.t.} \quad \sum_{f \in \mathcal{F}} p_f = 1, \end{aligned} \quad (2)$$

MWC: maximum weighted configuration

$$\max_{f \in \mathcal{F}} \sum_{r \in R} x_r(f) \approx \frac{1}{\beta} \log \left( \sum_{f \in \mathcal{F}} \exp \left( \beta \sum_{r \in R} x_r(f) \right) \right) \triangleq g_{\beta}(\mathbf{x}), \quad (3)$$

where  $\beta$  is a positive constant and  $\mathbf{x} \triangleq [\sum_{r \in R} x_i(f), f \in \mathcal{F}]$ .

# log-sum-exp approximation

**Theorem 1:** For the log-sum-exp function  $\bar{g}_\beta(\mathbf{x})$ , we have

- its conjugate function<sup>3</sup> is given by

$$g_\beta^*(\mathbf{p}) = \begin{cases} \frac{1}{\beta} \sum_{f \in \mathcal{F}} p_f \log p_f & \text{if } \mathbf{p} \geq 0 \text{ and } \mathbf{1}^T \mathbf{p} = 1 \\ \infty & \text{otherwise.} \end{cases} \quad (7)$$

- it is a convex and closed function; hence, the conjugate of its conjugate  $g_\beta^*(\mathbf{p})$  is itself, i.e.,  $g_\beta(\mathbf{x}) = g_\beta^{**}(\mathbf{x})$ . Specifically,

$$\begin{aligned} g_\beta(\mathbf{x}) &= \max_{\mathbf{p} \geq 0} \sum_{f \in \mathcal{F}} p_f \sum_{r \in R} x_r(f) - \frac{1}{\beta} \sum_{f \in \mathcal{F}} p_f \log p_f \quad (8) \\ \text{s.t. } &\sum_{f \in \mathcal{F}} p_f = 1. \end{aligned}$$

By solving KKT condition,

we have

$$p_f^*(\mathbf{x}) = \frac{\exp(\beta \sum_{r \in R} x_r(f))}{\sum_{f' \in \mathcal{F}} \exp(\beta \sum_{r \in R} x_r(f'))}, \forall f \in \mathcal{F}. \quad (12)$$