

An Optimal Algorithm for Online Non-Convex Learning

LIN YANG, The Chinese University of Hong Kong, Hong Kong

LEI DENG, The Chinese University of Hong Kong, Hong Kong

MOHAMMAD H. HAJIESMAILI, Johns Hopkins University, MD, USA

CHENG TAN, The Chinese University of Hong Kong, Hong Kong

WING SHING WONG, The Chinese University of Hong Kong, Hong Kong

In many online learning paradigms, convexity plays a central role in the derivation and analysis of online learning algorithms. The results, however, fail to be extended to the non-convex settings, while non-convexity is necessitated by a large number of recent applications. The Online Non-Convex Learning (ONCL) problem generalizes the classic Online Convex Optimization (OCO) framework by relaxing the convexity assumption on the cost function (to a Lipschitz continuous function) and the decision set. The state-of-the-art result for the ONCL demonstrates that the classic online exponential weighting algorithm attains a sublinear regret of $O(\sqrt{T \log T})$. The regret lower bound for the OCO, however, is $\Omega(\sqrt{T})$, and to the best of our knowledge, there is no result in the context of the ONCL problem achieving the same bound. This paper proposes the Online Recursive Weighting (ORW) algorithm with regret of $O(\sqrt{T})$, matching the tight regret lower bound for the OCO problem, and fills the regret gap between the state-of-the-art results in the online convex and non-convex optimization problems.

CCS Concepts: • **Theory of computation** → **Online learning algorithms**; **Regret bounds**; • **Computing methodologies** → **Machine learning algorithms**;

Additional Key Words and Phrases: Online non-convex learning, online convex optimization, Lipschitz expert, regret, online recursive weighting

ACM Reference Format:

Lin Yang, Lei Deng, Mohammad H. Hajiesmaili, Cheng Tan, and Wing Shing Wong. 2018. An Optimal Algorithm for Online Non-Convex Learning. *Proc. ACM Meas. Anal. Comput. Syst.* 2, 2, Article 25 (June 2018), 25 pages. <https://doi.org/10.1145/3224420>

1 INTRODUCTION

1.1 Background and Motivation

The Online Convex Optimization (OCO) framework has widely influenced the online learning community since the seminal work by Zinkevich [51]. The OCO is modeled as a repeated game composed of T iterations. At iteration t , the player chooses a point \mathbf{x}_t from a bounded convex decision set $\mathcal{K} \subset \mathbb{R}^n$; after the choice is committed, a bounded convex cost function $f_t : \mathcal{K} \mapsto \mathbb{R}$ is revealed to the player. The goal of the player is to minimize the *regret*, which is defined as the

Authors' addresses: Lin Yang, The Chinese University of Hong Kong, Shatin, NT, Hong Kong, 999077, yl015@ie.cuhk.edu.hk; Lei Deng, The Chinese University of Hong Kong, Shatin, NT, Hong Kong, 999077, ldeng@ie.cuhk.edu.hk; Mohammad H. Hajiesmaili, Johns Hopkins University, 3400 N Charles St. Baltimore, MD, USA, hajiesmaili@jhu.edu; Cheng Tan, The Chinese University of Hong Kong, Shatin, NT, Hong Kong, 999077, tancheng1987love@163.com; Wing Shing Wong, The Chinese University of Hong Kong, Shatin, NT, Hong Kong, 999077, wswong@ie.cuhk.edu.hk.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

2476-1249/2018/6-ART25 \$15.00

<https://doi.org/10.1145/3224420>

difference between the online cumulative cost and the cumulative cost using an optimal offline choice in hindsight. This model can be applied to many real-world problems, such as online routing [7, 45], spam email filtering [17, 42], online metric learning [24], ad selection and content ranking in search engines [9, 35, 47], etc.

A promising result related to the OCO model is that the regret of the state-of-the-art efficient algorithms [18] is sublinear, implying that the average cost difference per iteration converges to zero as the number of iterations goes to infinity. It is well known that the lower bound of the regret for the OCO problem is $O(\sqrt{T})$ [18] and researchers have proposed a large number of online algorithms whose regret attains this lower bound, including the Online Gradient Decent method [51], the Stochastic Gradient Decent method [20, 39, 40, 44], the Online Newton Step, and many regularization-related methods [16, 26] (see the recent survey paper [18], and the references therein).

One of the most natural extensions for the OCO problem is to relax the convexity assumption on the decision set \mathcal{K} and the cost function f_t . This extension brings out the Online Non-Convex Learning (the ONCL) problem, which is necessitated by tons of state-of-the-art applications. For example, in the portfolio selection problem [10, 25], the decision maker (e.g., the trader) chooses a distribution of her wealth allocation over n assets \mathbf{x}_t , at each round. By the end of each round, the adversary chooses the market returns of the assets with positive values. In some specific settings [5, 32, 46], the online portfolio selection problem is a non-convex one due to the non-convex diversification constraints and non-convex transaction costs, and thus the traditional OCO framework fails in modeling such case. In addition, there are extensive machine learning research focusing on non-convex loss functions in large margin classifiers [12, 37, 49]. In [12, 37], non-convex online Support Vector Machine (SVM) models has been studied which adopts a non-convex loss function, called Ramp Loss, to suppress the influence of outliers. In [49], a special non-convex penalty, called the smoothly clipped absolute deviation penalty, is imposed on the hinge loss function in the SVM. Such a new SVM is applied to identify important genes for cancer classification [49].

1.2 Related Results

The ONCL problem is not a new problem and there are plenty prior works on it. Among them, [12] and [15] propose respective heuristic online training algorithms, but neither of them are rigorously shown to satisfy any regret bound. In [19], Hazan and Kale tackle the ONCL problem with submodular cost functions, and propose an online algorithm that attains the regret of $O(\sqrt{T})$. In [50], an online bandit learning problem with non-convex losses is investigated. The cost function is again a special non-convex function, defined as the composition of a non-increasing scalar function with a linear function of small variation. An online algorithm is developed with $\tilde{O}(\text{poly}(d)T^{2/3})$ regret, where $\text{poly}(d)$ stands for a polynomial with respect to the dimension of the decision set d . The most related works to ours are [31] and [36], where by applying the exponential weighting method [6, 14] the regret of $O(\sqrt{T \log T})$ is attained.

In addition, the ONCL problem has been broadly investigated under a similar problem, called the Lipschitz Expert problem [27, 29], which generalizes the traditional Expert problem [14] to metric spaces. For such a problem, a regret of $O(\sqrt{T \log T})$ can be achieved [29]. However, to the best of our knowledge, for the general ONCL problem or the Lipschitz Expert problem, no online algorithm can achieve the regret of $O(\sqrt{T})$ as the well-known lower bound for the OCO problem.

1.3 Our Contributions and Adopted Techniques

This paper tackles the ONCL problem, with non-convex L -Lipschitz cost functions and non-convex decision set. We propose a novel online algorithm, called the Online Recursive Weighting (the

ORW) algorithm, and prove that the regret of the ORW is upper bounded by $O(\sqrt{T})$. The obtained regret bound matches the well-known lower bound of the regret for the OCO [18], hence, the ORW algorithm is asymptotically optimal for the general ONCL problem.

The general idea of the ORW algorithm is to divide the decision set into multiple subsets according to a *grid layered* structure. Any subset at the upper layer is divided into multiple smaller subsets at the lower layers. Our algorithm recursively selects the subset from the topmost layer to the bottommost layer until a decision point is identified. At each layer, the ORW algorithm leverages the classic Exponential Weighting or Hedge algorithm [6, 14] to select a subset at the lower layer. As the core technical contribution, the ORW properly partitions the subsets and sets the subset-selecting probabilities, thereby, it asymptotically achieves $O(\sqrt{T})$ regret.

Then, the ORW algorithm is extended to an adaptive version (the AORW algorithm), which increases the granularity parameter gradually as time goes on. The AORW guarantees $O(\sqrt{T})$ regret, reduces the computational complexity of the ORW, and works properly when the duration of time horizon is unknown to the online player.

1.4 Basic Notations and Organizations

In the rest of this paper, we use calligraphy font to denote sets, e.g., \mathcal{K} . In equations, we use Fraktur font to denote algorithms, e.g., \mathfrak{A} . We use bold math font to denote vectors whose entries are represented by the corresponding normal math font, e.g., $\mathbf{i} = (i_1, i_2, \dots, i_n)$. We further gives the following definition for L -Lipschitz functions.

DEFINITION 1. A function $f : \mathcal{K} \rightarrow \mathbb{R}$ (where $\mathcal{K} \subset \mathbb{R}^n$) is L -Lipschitz if

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{K}. \quad (1)$$

The rest of the paper is organized as follows. Section 2 introduces details on the problem formulation as well as the notation used in this paper. In Section 3, we introduce the proposed ORW algorithm and its regret is analyzed in Section 4. In Sections 5, we extend our algorithm to an adaptive version called AORW. Section 6 reviews the related literature and potential extensions. Finally, Section 7 concludes the paper.

2 PROBLEM SETTING OF ONLINE NON-CONVEX LEARNING

We describe our online non-convex learning problem as a structured repeated game. At each iteration t , the player chooses a decision $\mathbf{x}_t \in \mathcal{K}$, where $\mathcal{K} \subseteq \mathbb{R}^n$ is a bounded decision set whose diameter is D , i.e., $\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \|\mathbf{x} - \mathbf{y}\|_2 = D$. After the player commits to a decision point at iteration t , the adversary chooses a cost function $f_t(\mathbf{x}) \in \mathcal{F}$, where $\mathcal{F} : \mathcal{K} \mapsto \mathbb{R}^+$ is a set of cost functions which are assumed to be non-negative and L -Lipschitz (as defined in Equation (1)).

Note that in the original OCO model, the cost function and decision set are assumed to be convex. In this work, we take into account the general class of cost functions and decision sets which are not necessarily convex.

In our model, we consider the full feedback case, in which the cost function $f_t(\mathbf{x})$ is revealed to the player only after a choice is made at iteration t . At each iteration, the player needs to make online decisions without knowing the current and future cost functions. Fed with the historical cost functions $\mathbf{H}_t = (f_1, f_2, \dots, f_{t-1})$, the decision of an online algorithm \mathfrak{A} at iteration t is denoted as $\mathbf{x}_t = \mathfrak{A}(\mathbf{H}_t)$. The summary of main notations related to problem setting is given in Table 1.

Table 1. Summary of notations related to problem setting

Notation	Description
t	Index of iteration
T	The number of iterations, $T \geq 1$
\mathcal{T}	Set $\mathcal{T} = \{1, 2, \dots, T\}$
n	The dimension of the decision space
\mathbf{x}_t	Chosen decision point at t
\mathcal{K}	The decision set available for the online player
D	Diameter of the decision set \mathcal{K}
$f_t(\mathbf{x})$	$f_t(\mathbf{x}) : \mathcal{K} \mapsto \mathbb{R}^+$. Cost function at iteration t , known for $t - 1$, unknown for $\tau \geq t$
L	Lipschitz constant for the cost functions
\mathcal{F}	The set of cost functions which are available for adversary
H_t	$H_t = (f_1, f_2, \dots, f_{t-1})$. The historical cost functions available for the online algorithm at iteration t
\mathfrak{A}	An online algorithm which maps from the historical cost functions H_t to a decision point \mathbf{x}_t for all $t = 1, 2, \dots$

We use the pseudo-regret¹ as a common performance metric to evaluate our proposed online algorithm for the ONCL problem, formally defined as

$$\text{regret}_T(\mathfrak{A}) \stackrel{\text{def}}{=} \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t=1}^T \mathbb{E}[f_t(\mathbf{x}_t)] - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}) \right\}, \quad (2)$$

which is the cumulative difference between the expected cost of the online algorithm and the cost of the best fixed offline decision up to T . The expectation in (2) is taken over the randomness of the decision at each iteration t if the online algorithm \mathfrak{A} is randomized.

Our goal is to design an online algorithm for the ONCL problem and try to minimize the regret. Ideally, it is desired to have a sublinear regret with respect to T , i.e., $\text{regret}_T(\mathfrak{A}) = o(T)$. The sublinear regret implies that time-average performance of the online algorithm is as good as the best fixed strategy as time goes to infinity.

In Section 3, we propose an online algorithm for the ONCL problem, called the Online Recursive Weighting (ORW). In Section 4, we analyze the performance of the ORW and demonstrate that its regret is $O(\sqrt{T})$ asymptotically, which attains the lower bound of the regret for the basic OCO problem [18].

3 ONLINE RECURSIVE WEIGHTING ALGORITHM

The Online Recursive Weighting (ORW) algorithm is based on the idea of dividing the entire decision set into several *sampling subsets* so that each subset contains highly correlated elements. An essential tuning parameter for the algorithm is the granularity of the sampling subsets, which determines the number of final sample points (details in Sections 3.1 and 3.2).

To define the granularity of the sampling subsets and the decision policy of the ORW algorithm, we construct a *layered grid structure*. The topmost layer consists of a single grid and hence a single subset including the entire decision set \mathcal{K} . The sampling subsets are the subsets at the bottommost layer. In Section 3.1, we explain the details of constructing the layered grid structure.

¹We call it regret in short in the rest of this paper.

Table 2. Summary of key notations related to the ORW algorithm

Notation	Description
\mathcal{D}	Cover cube of the decision set \mathcal{K}
l	Index of layer
\mathbf{i}_l	Index for a sub-cube/subset at l -th layer, $\mathbf{i}_l = (i_{l,1}, i_{l,2}, \dots, i_{l,n})$
$\mathcal{D}_l(\mathbf{i}_l)$	Layer l sub-cube indexed by \mathbf{i}_l
$\mathcal{K}_l(\mathbf{i}_l)$	Layer l subset indexed by \mathbf{i}_l
\mathcal{I}_l	Index set of nonempty subsets at layer l
$\mathbf{v}_s(\mathbf{i}_l)$	Index of layer s subset that contains $\mathcal{K}_l(\mathbf{i}_l)$
$\mathcal{E}_s(\mathcal{K}_l(\mathbf{i}_l))$	Index set of nonempty layer s subsets within the subset $\mathcal{K}_l(\mathbf{i}_l)$, $\mathbf{i}_l \in \mathcal{I}_l$
m	The index of the layer where the subsets are sampled
$\mathbf{r}_m(\mathbf{i}_m)$	The sample point for the subset $\mathcal{K}_m(\mathbf{i}_m)$
$c_{m,t}(\mathbf{i}_m)$	The cost on the sample point of the subset $\mathcal{K}_m(\mathbf{i}_m)$ at iteration t , i.e., $c_{m,t}(\mathbf{i}_m) \stackrel{\text{def}}{=} f_t(\mathbf{r}_m(\mathbf{i}_m))$
L_m	The maximum cost difference between two points in the same subset at layer m
$\mathcal{I}_{l,t}$	The index of the subset chosen at layer l at iteration t
$\bar{c}_{l,t}(\mathbf{i}_l)$	The <i>expected normalized cost</i> conditioning on that subset $\mathcal{K}_l(\mathbf{i}_l)$ is chosen at layer l at iteration t . For short, we also call it the expected cost of choosing $\mathcal{K}_l(\mathbf{i}_l)$ at layer l at iteration t
$\bar{C}_{l,t}(\mathbf{i}_l)$	The cumulative sum for $\bar{c}_{l,\tau}(\mathbf{i}_l)$ up to iteration t , i.e., $\bar{C}_{l,t}(\mathbf{i}_l) = \sum_{\tau=1}^t \bar{c}_{l,\tau}(\mathbf{i}_l)$. For short, we also call it the cumulative expected cost of choosing $\mathcal{K}_l(\mathbf{i}_l)$ at layer l up to iteration t
\mathfrak{R}	The notation used for the ORW algorithm in the equations (Algorithm 1)

In the next step (Section 3.2), a single element is randomly selected from each sampling subset called *sample point*, and by grouping all these points, we construct the set of sample points. At each iteration, the algorithm selects a *decision point* randomly from the set of sample points according to a tree decision structure. Beginning with the topmost layer, the ORW algorithm probabilistically selects a subset by running Hedge update², and subsequently, the ORW algorithm continues among the subsets of one layer below that lies inside the selected subset of the previous layer.

Designing the tree decision structure to select among sample points is the core technical contribution of the ORW algorithm and is proposed in Section 3.3. The summary of key notations related to the algorithm design is given in Table 2.

3.1 A Layered Grid Structure

Recall that the diameter of bounded decision set \mathcal{K} is D , hence, it is possible to find a bounded cube of length D , denoted by \mathcal{D} , that can cover \mathcal{K} entirely, i.e., $\mathcal{K} \subset \mathcal{D}$. At layer $l \in \mathbb{N}^+$, we partition \mathcal{D} into smaller identical sub-cubes with edge length of size $D/2^l$. The decision set is n -dimensional, hence the total number of sub-cubes is equal to 2^{nl} . For simplicity, each sub-cube is indexed by a distinct n -dimensional vector $\mathbf{i}_l = (i_{l,1}, i_{l,2}, \dots, i_{l,n})$, where $1 \leq i_{l,j} \leq 2^l, j = 1, 2, \dots, n$, and the sub-cube indexed by \mathbf{i}_l is denoted by $\mathcal{D}_l(\mathbf{i}_l)$. One can refer to Fig. 1 for an illustrative example with $n = 2$ and $l = 1$.

²Hedge algorithm maintains weights for each expert/decision point, which are updated according to the observed losses at each step. For details, readers can refer to [14].

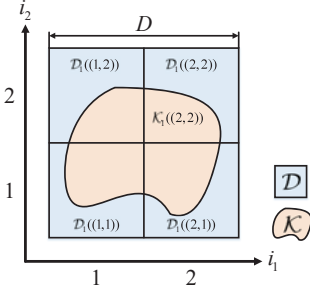


Fig. 1. Layer 1 grid construction for a general decision set \mathcal{K} .

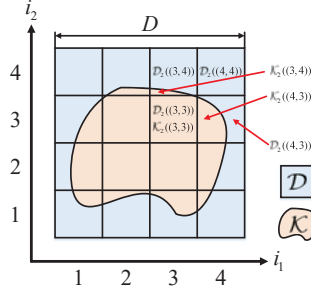


Fig. 2. Index set for overlapped sub-cubes when $n = 2$ and $l = 2$.

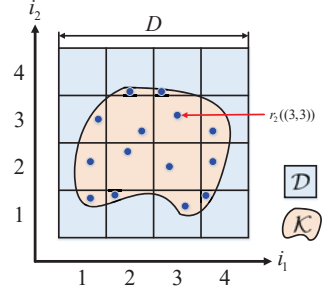


Fig. 3. Sampling for the case when $n = 2$ and $m = 2$.

The notation $\mathcal{K}_l(i_l)$ is used to represent the intersection of the decision set \mathcal{K} and corresponding sub-cube $\mathcal{D}_l(i_l)$, i.e.,

$$\mathcal{K}_l(i_l) \stackrel{\text{def}}{=} \mathcal{D}_l(i_l) \cap \mathcal{K}.$$

By convention, we regard the decision set \mathcal{K} as the only layer 0 subset, denoted by $\mathcal{K}_0(1) = \mathcal{K}$. By the above partitioning structure, in total 2^{nl} subsets at layer l are constructed, whose union is the decision set \mathcal{K} . Note that there could be some empty subsets, and we denote the index set for nonempty layer l subsets by \mathcal{I}_l , i.e.,

$$\mathcal{I}_l \stackrel{\text{def}}{=} \{i_l : \mathcal{K}_l(i_l) \neq \emptyset\}.$$

Any subset indexed by $i_l \in \mathcal{I}_l$, i.e., $\mathcal{K}_l(i_l)$, consists of a group of neighboring lower-layer subsets. Specifically, for any $s > l$, we have

$$\mathcal{K}_l(i_l) = \bigcup_{\substack{i_s : i_{s,j} \geq 1 + (i_{l,j} - 1)2^{s-l}, \\ i_{s,j} \leq i_{l,j}2^{s-l}, \\ j = 1, 2, \dots, n.}} \mathcal{K}_s(i_s).$$

To ease the presentation, we define the following notation:

▷ We use $\mathcal{E}_s(\mathcal{K}_l(i_l))$, $s \geq l$, to denote the index set of nonempty layer s subsets within $\mathcal{K}_l(i_l)$, i.e.,

$$\mathcal{E}_s(\mathcal{K}_l(i_l)) \stackrel{\text{def}}{=} \{i_s \in \mathcal{I}_s : 1 + (i_{l,j} - 1)2^{s-l} \leq i_{s,j} \leq i_{l,j}2^{s-l}, j = 1, 2, \dots, n\}. \quad (3)$$

Indeed, we have $|\mathcal{E}_s(\mathcal{K}_l(i_l))| \leq 2^{(s-l)n}$.

▷ Assume $i_l \in \mathcal{I}_l$ and $s \leq l$. We use $\mathbf{v}_s(i_l)$ to denote an index for some layer s subset, and the elements of $\mathbf{v}_s(i_l)$ satisfy

$$v_{s,j} = \left\lceil \frac{i_{l,j}}{2^{l-s}} \right\rceil, j = 1, 2, \dots, n.$$

Due to the above definition of $\mathbf{v}_s(i_l)$, $i_l \in \mathcal{I}_l$, together with the grid construction method we adopt, it follows that $\mathcal{K}_l(i_l) \subset \mathcal{K}_s(\mathbf{v}_s(i_l))$.

Example 1. Fig. 1 illustrates a simple example of grid construction when $n = 2$ and $l = 1$. At layer 1, the cover cube \mathcal{D} is partitioned into four smaller sub-cubes, each of whom intersects the decision set \mathcal{K} , forming four nonempty subsets, respectively.

$$\mathcal{I}_1 = \{(1, 1), (1, 2), (2, 1), (2, 2)\}.$$

In Fig. 2, the subsets at layer 1 are further divided and 14 fixed nonempty layer 2 subsets are constructed. Particularly, the subset $\mathcal{K}_1((2, 2))$ contains three nonempty layer 2 subsets, i.e., $\mathcal{K}_2((3, 3))$, $\mathcal{K}_2((3, 4))$, and $\mathcal{K}_2((4, 3))$. Thus, we have

$$\mathcal{E}_2(\mathcal{K}_1((2, 2))) = \{(3, 3), (3, 4), (4, 3)\}.$$

According to the definition of $\mathbf{v}_s(\mathbf{i}_l)$, we have

$$\mathbf{v}_1((3, 3)) = \mathbf{v}_1((3, 4)) = \mathbf{v}_1((4, 3)) = (2, 2),$$

and $\mathcal{K}_2((3, 3)) \subset \mathcal{K}_1((2, 2))$, $\mathcal{K}_2((3, 4)) \subset \mathcal{K}_1((2, 2))$, and $\mathcal{K}_2((4, 3)) \subset \mathcal{K}_1((2, 2))$.

3.2 Sampling

Assuming layer 0 contains the original decision set \mathcal{K} , we fix the number of layers to be $m + 1$, where layer m contains at most 2^m nonempty subsets. In the sampling procedure, a *sample point*, $\mathbf{r}_m(\mathbf{i}_m)$, is selected randomly as a representative, from each nonempty subset $\mathcal{K}_m(\mathbf{i}_m)$ at layer m . Fig. 3 depicts a simple example for $n = 2$ and $m = 2$, where the sample points are colored in blue. Due to the bijective correspondence between the sample points and the subsets at layer m , we interchangeably use choosing a sample point from a subset and choosing the corresponding layer m subset.

In Section 3.3, we define a recursive, probabilistic algorithm to select subsets as we go down the layered structure. Correspondingly, this defines a probabilistic approach to choose a *decision point* from the sample points as the final decision of the ORW algorithm.

Once the decision point is selected at each iteration, the cost of the decision, $f_t(\mathbf{r}_m(\mathbf{i}_m))$, along with the function f_t will be revealed.

For notation convenience, we denote the cost of the sample point from the layer m subset, $\mathcal{K}_m(\mathbf{i}_m)$, $\mathbf{i}_m \in \mathcal{I}_m$, at iteration t by $c_{m,t}(\mathbf{i}_m)$, i.e.,

$$c_{m,t}(\mathbf{i}_m) \stackrel{\text{def}}{=} f_t(\mathbf{r}_m(\mathbf{i}_m)).$$

Note that by this sample point construction step, the proposed algorithm reduces the original problem in principle to an Expert problem with $|\mathcal{I}_m|$ experts. As compared to the classic setting of the Expert problem in [8], the difference is that the cost function in our problem is a non-convex L -Lipschitz continuous one. Finally, the following lemma characterizes the cost difference between two sample points.

LEMMA 3.1. Assume $\mathbf{p}, \mathbf{q} \in \mathcal{I}_m$, we have

$$|c_{m,t}(\mathbf{p}) - c_{m,t}(\mathbf{q})| \leq 2L_m \|\mathbf{p} - \mathbf{q}\|_1,$$

where $L_m = \sqrt{n}DL/(2^m)$ and $\|\mathbf{p} - \mathbf{q}\|_1 = \sum_{j=1}^n |p_j - q_j|$.

PROOF. The maximum distance between two points within a subset at layer m is $\sqrt{n}D/(2^m)$. Considering the Lipschitz continuous condition, the maximum cost difference within the same subset at layer m is $L_m = \sqrt{n}DL/(2^m)$. We require that, any two subsets of the same layer, \mathbf{p} and \mathbf{q} , are said to be neighboring subsets if their coordinates satisfy $p_j - q_j \leq 1$ for $j = 1, 2, \dots, n$. Then, the maximum cost difference of two points in the union of any two neighboring layer m subsets is $2\sqrt{n}DL/(2^m)$. Thus, the maximum distance for any two points in the union of two layer m subsets, \mathbf{p} and \mathbf{q} , is $2\sqrt{n}DL/(2^m)\|\mathbf{p} - \mathbf{q}\|_1$. This yields the result in Lemma 3.1. \square

3.3 Recursive Choosing Policy

In this subsection, we propose a novel recursive decision structure to determine the index point for the bottom-layer subset over the index set \mathcal{I}_m .

At iteration t , $t = 1, 2, \dots, T$, the final choice of the ORW algorithm is obtained by recursively choosing a nonempty subset from the topmost layer ($l = 0$) to the bottommost layer ($l = m$). In this way, the action of the ORW algorithm at each iteration t , consists of a sequence of indexes of subsets, i.e., $(I_{0,t}, I_{1,t}, I_{2,t}, \dots, I_{m,t})$. By default, $I_{0,t}$ is set to be 1, referring to the only decision set \mathcal{K} . Additionally, $I_{l,t}$ should satisfy that $I_{l,t} \in \mathcal{E}_l(\mathcal{K}_{l-1}(I_{l-1,t}))$, i.e., $\mathcal{K}_l(I_{l,t}) \subset \mathcal{K}_{l-1}(I_{l-1,t})$ for $l = 1, 2, \dots, m$.

At each layer, a stochastic policy is adopted to select a lower-layer subset. Suppose that, at iteration t and at layer l , $0 \leq l < m$, the ORW algorithm has chosen subset $i_l \in \mathcal{I}_l$, i.e., $I_{l,t} = i_l$. Then, the algorithm proceeds to choose a subset at the next layer (layer $(l + 1)$) in $\mathcal{K}_l(i_l)$ randomly according to a conditional probability.

In the ORW algorithm, and for $i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))$, the conditional probability, $\Pr_{l,l+1}^t = \Pr[I_{l+1,t} = i_{l+1} | I_{l,t} = i_l]$, is based on the cumulative expected cost of the *previous* iterations.

For iteration τ , $0 \leq \tau \leq t-1$, the *expected normalized cost* of subset $\mathcal{K}_l(i_l)$, denoted by $\bar{c}_{l+1,\tau}(i_{l+1})$, is defined as follows

$$\begin{aligned} \bar{c}_{l,\tau}(i_l) &\stackrel{\text{def}}{=} \mathbb{E} \left[\frac{c_{m,\tau}(I_{m,\tau}) - \min_{i_m \in Q} c_{m,\tau}(i_m)}{2^{m-l+1} L_m} | I_{l,\tau} = i_l \right] \\ &= \sum_{j_m \in \mathcal{E}_m(\mathcal{K}_l(i_l))} \frac{c_{m,\tau}(j_m) - \min_{i_m \in Q} c_{m,\tau}(i_m)}{2^{m-l+1} L_m} \cdot \Pr[I_{m,\tau} = j_m | I_{l,\tau} = i_l], \end{aligned} \quad (4)$$

where $Q = \mathcal{E}_m(\mathcal{K}_{l-1}(\mathbf{v}_{l-1}(i_l)))$. Note that $\bar{c}_{l+1,\tau}(i_{l+1})$ can be computed at the end of iteration τ , once the cost function is revealed. In Equation (4), the conditional selection probability for subset $j_m \in \mathcal{I}_m$, i.e., $\Pr[I_{m,\tau} = j_m | I_{l,\tau} = i_l]$ can be obtained by multiplying the selection probabilities of the corresponding subsets at each layer, i.e.,

$$\Pr[I_{m,\tau} = j_m | I_{l,\tau} = i_l] = \prod_{k=l}^{m-1} \Pr[I_{k+1,\tau} = \mathbf{v}_{k+1}(\mathcal{K}_m(j_m)) | I_{k,\tau} = \mathbf{v}_k(\mathcal{K}_m(j_m))]. \quad (5)$$

For a subset $\mathcal{K}_l(i_l)$, the *cumulative expected cost*, denoted by $\bar{C}_{l,t}(i_l)$, is defined as the sum of the *expected normalized cost* from iteration 1 to iteration t , i.e.,

$$\bar{C}_{l,t}(i_l) \stackrel{\text{def}}{=} \sum_{\tau=1}^t \bar{c}_{l,\tau}(i_l). \quad (6)$$

By convention, let the initial value of the cumulative expected cost be zero, i.e., $\bar{C}_{l,0}(i_l) = 0$ for all $i_l \in \mathcal{I}_l$ and for any layer $l = 1, 2, \dots, m$.

At iteration t , the selection probability, $\Pr[I_{l+1,t} = i_{l+1} | I_{l,t} = i_l]$, is proportional to the exponent of the expected cumulative cost of choosing subset i_{l+1} at layer $(l + 1)$ up to iteration $t - 1$, i.e.,

$$\Pr[I_{l+1,t} = i_{l+1} | I_{l,t} = i_l] = \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(i_{l+1}))}{\sum_{i'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(i'_{l+1}))}, \quad (7)$$

where η_t is a positive and decreasing parameter which implies that the changing on $\bar{C}(\cdot)$ has decreasing influence on the choosing probability and thus the decision of the online algorithm gets

Algorithm 1 The Online Recursive Weighting (ORW) Algorithm

Input: index set $\mathcal{I}_m, T, \{\eta_t = \sqrt{\frac{n}{t}}\}$
Output: $\mathbf{r}_m(\mathcal{I}_{m,1}), \mathbf{r}_m(\mathcal{I}_{m,2}), \dots$

- 1: //↑ Output sample points of layer m subsets
- 2: //↓ Initialize the cumulative expected cost for all subsets at all layers
- 3: Set $\bar{C}_{l,0}(\mathbf{i}_l) = 0, \forall \mathbf{i}_l \in \mathcal{I}_l, l = 1, 2, \dots, m$
- 4: **for** $t = 1, 2, \dots, T$ **do**
- 5: $\mathbf{I}_{0,t} = 1$
- 6: //↓ Recursively select subsets at all layers
- 7: **for** $l = 0, 1, \dots, m-1$ **do**
- 8: Randomly select an index $\mathbf{I}_{l+1,t} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{I}_{l,t}))$ according to the probability distribution specified in Equation (7)
- 9: **end for**
- 10: Choose the subset indexed by $\mathbf{I}_{m,t}$ at iteration t
- 11: Choose the sample point for subset $\mathcal{K}_m(\mathbf{I}_{m,t})$, i.e., $\mathbf{r}_m(\mathbf{I}_{m,t})$, as the final decision
- 12: //↓ Get the normalized expected cost for all subsets at all layers based on the revealed cost function $f_t(\mathbf{x})$ at iteration t
- 13: **for** $l = 1, 2, \dots, m$ **do**
- 14: **for** $\mathbf{i}_l \in \mathcal{I}_l$ **do**
- 15: **for** $j_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))$ **do**
- 16: Calculate $\Pr[\mathbf{I}_{m,t} = j_m | \mathbf{I}_{l,t} = \mathbf{i}_l]$ according to Equation (5)
- 17: **end for**
- 18: Calculate $\bar{c}_{l,t}(\mathbf{i}_l)$ according to Equation (4)
- 19: Update $\bar{C}_{l,t}(\mathbf{i}_l) = \bar{C}_{l,t-1}(\mathbf{i}_l) + \bar{c}_{l,t}(\mathbf{i}_l)$
- 20: **end for**
- 21: **end for**
- 22: **end for**

more unalterable as time goes on. Note that the denominator in Equation (7) is a normalizer such that right hand side of Equation (7) is a probability mass function.

Note that the ORW algorithm updates the cumulative expected cost at the end of each iteration. Once the conditional probabilities are defined and the cost function at iteration t is revealed, we then calculate the expected normalized cost and further update the cumulative expected cost for each subset, which is

$$\bar{C}_{l,t}(\mathbf{i}_l) = \bar{C}_{l,t-1}(\mathbf{i}_l) + \bar{c}_{l,t}(\mathbf{i}_l). \quad (8)$$

Then, the updated cumulative expected cost will be used to calculate the choosing probability at the next iteration as shown in Equation (7). The summary of the proposed ORW algorithm is listed as Algorithm 1.

Remark 1. (Technical differences with the traditional weighting methods) For the expert learning problem, the intuitive idea to attain a sublinear regret is to allocate more preference to the expert of smaller cumulative cost in a stochastic manner [34], which can be realized by the Exponential Weighting algorithm (or Hedge algorithm). The Hedge algorithm observes the costs on each point, and updates the weight of a point based on its own cost only. Such a point-by-point weighting method fails in utilizing the cost correlation among neighboring decision nodes and attains only sub-optimal performance within continuum decision space which could consist of infinitely many decision points. Than the Hedge algorithm, the ORW algorithm makes better use of the correlation of neighboring points by grouping highly correlated decision points as a high-level decision. As

expressed in Equation (4), a common evaluation is conducted for each decision group, and a particular amount of common preference will be allocated to the group based on its cumulative performance. Furthermore, by constructing a tree decision structure, the ORW algorithm chooses among at most 2^n “experts” or decision points at each layer, avoiding the performance degradation due to the large number of decision points that the Hedge algorithm might have been faced with.

4 REGRET ANALYSIS FOR THE ONLINE RECURSIVE WEIGHTING ALGORITHM

Since the ONCL problem is a more general model than the OCO, the lower bound for the OCO is still valid for the non-convex case. In the OCO framework, a well-known lower bound of the regret is $\Omega(\sqrt{T})$, and the reader can refer to [18] for the sketch of the proof. In this section, we analyze the regret of the ORW algorithm (represented by \mathfrak{R} in equations) and demonstrate it matches the above lower bound. The main technical results are summarized in the following theorem.

THEOREM 4.1. *With $\eta_t = \sqrt{\frac{n}{t}}$, the ORW algorithm guarantees that*

$$\text{regret}_T(\mathfrak{R}) \leq 2(\ln 2 + 1)nDL\sqrt{T} + 2\ln 2 \cdot nDL + \frac{\sqrt{n}}{2^m}DLT.$$

If m is further set to $\lceil \log_2 \sqrt{nT} \rceil$, it follows that

$$\text{regret}_T(\mathfrak{R}) < (4n + 1)DL\sqrt{T} + 2nDL.$$

Remark 2. (Dimensional dependency) Theorem 4.1 implies that the ORW algorithm attains a regret of $O(n\sqrt{T})$, with a mild polynomial dependency on the dimension. Note that the regret of the Hedge algorithm on a continuum [31] has dimension dependency and is shown to be $O(\sqrt{nT \log T})$ (see Corollary 2 in [31]), which has a milder dependency on the dimension as compared to our algorithm. This milder dependency is due to an additional assumption on the cost function in [31] to be always bounded by a constant no matter how large the dimension is. In our setting, we do not have such assumption on cost functions, and the regret is $O(n\sqrt{T})$.

Remark 3. [31] shows that when the decision set is uniformly fat, the Hedge algorithm can attain a regret of $O(\sqrt{T \log T})$. An additional advantage of our ORW algorithm is that the only assumption on decision set is that it is bounded.

To carry out the analysis, we split the regret analysis of the ORW algorithm into two parts:

▷ The first part is the regret due to the “imperfect choice” among sample points, i.e.,

$$\text{regret}_{\text{ImC}, T}(\mathfrak{R}) \stackrel{\text{def}}{=} \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t \in \mathcal{T}} \mathbb{E} [c_{m,t}(I_{m,t})] - \min_{i_m \in I_m} \sum_{t \in \mathcal{T}} c_{m,t}(i_m) \right\},$$

where the first term is the cumulative cost incurred by the online algorithm (whose choice at iteration t is denoted by $I_{m,t}$), and the second term is the minimum cumulative cost among sampled points.

▷ The second part of the regret is introduced by “imperfect discretization”, which is expressed as

$$\text{regret}_{\text{ImD}, T}(\mathfrak{R}) \stackrel{\text{def}}{=} \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \min_{i_m \in I_m} \sum_{t \in \mathcal{T}} c_{m,t}(i_m) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t \in \mathcal{T}} f_t(\mathbf{x}) \right\},$$

where the second term is the minimum cumulative cost over the decision set \mathcal{K} . Since the supreme of sum is less than or equal to the sum of supreme, we have the following lemma.

LEMMA 4.2. $\text{regret}_T(\mathfrak{R}) \leq \text{regret}_{\text{ImC}, T}(\mathfrak{R}) + \text{regret}_{\text{ImD}, T}(\mathfrak{R})$.

In the following lemma, we derive a bound for the first-part regret of the ORW algorithm, which is related to the subproblem of choosing a point over the set \mathcal{I}_m . The regret analysis for the imperfect discretization is given in Lemma 4.4.

LEMMA 4.3. *With $\eta_t = \sqrt{\frac{n}{t}}$, the ORW algorithm guarantees that*

$$\text{regret}_{\text{ImC}, T}(\mathfrak{R}) \leq 2(\ln 2 + 1)nDL\sqrt{T} + 2\ln 2 \cdot nDL.$$

PROOF. The ORW algorithm has a recursive decision-making structure to determine the final decision. Correspondingly, the regret due to “imperfect choice” can be further splitted into multiple pieces which are introduced at each layer.

Suppose a nonempty subset $\mathcal{K}_l(i_l)$ is chosen at layer $l \in \{0, 1, 2, \dots, m-1\}$. In the next step, the ORW algorithm will further choose a subset whose index lies in $\mathcal{E}_{l+1}(\mathcal{K}_l(i_l))$. Among the subsets of $\mathcal{K}_l(i_l)$, there exists a local optimal subset (for example, i''_{l+1}), whose performance is known in hindsight and potentially a regret loss due to imperfect choice at layer l will be incurred by the online algorithm. Equation (9) bounds such cost difference (or regret loss) between $\sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l,t} = i_l]$ and $\sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l+1,t} = i''_{l+1}]$ at the l -th layer, where i''_{l+1} is any element in $\mathcal{E}_{l+1}(\mathcal{K}_l(i_l))$.

In Equation (9), equality (E1) is obtained simply by the law of total probability. Equality (E2) is due to the definition of $\bar{c}_{l+1,t}(\mathbf{i}_{l+1})$ in (4). Equality (E3) is based on the conditional probability in (7). Inequality (E4) is due to the following inequality,

$$\begin{aligned} & \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1})) \cdot \bar{c}_{l+1,t}(\mathbf{i}_{l+1})}{\sum_{i'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} \\ & \leq -\frac{1}{\eta_t} \ln \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1})) \cdot \exp(-\eta_t \bar{c}_{l+1,t}(\mathbf{i}_{l+1}))}{\sum_{i'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} + \frac{\eta_t}{2}. \end{aligned} \quad (10)$$

The proof of Equation (10) is given in Appendix A.1. Equality (E5) is due to the fact that $\bar{C}_{l+1,t}(\mathbf{i}_{l+1}) = \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1}) + \bar{c}_{l+1,t}(\mathbf{i}_{l+1})$. Equality (E6) simplifies the expression for the log-sum-exp function by defining

$$\Phi_t(\alpha) \stackrel{\text{def}}{=} -\frac{1}{\alpha} \ln \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\alpha \bar{C}_{l+1,t}(\mathbf{i}_{l+1})),$$

where $\alpha > 0$. Equality (E7) rearranges the first set of terms and Inequality (E8) uses the following bound result

$$\sum_{t=1}^{T-1} [\Phi_t(\eta_t) - \Phi_t(\eta_{t+1})] \leq \sum_{t=1}^{T-1} n \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \leq \ln 2 \cdot \sqrt{nT}, \quad (11)$$

which is proved in Appendix A.2. Inequality (E9) entails the following two results:

$$-\Phi_0(\eta_1) = \frac{1}{\eta_1} \ln |\mathcal{E}_{l+1}(\mathcal{K}_l(i_l))| = \frac{1}{\sqrt{n}} \ln |\mathcal{E}_{l+1}(\mathcal{K}_l(i_l))| \leq \frac{1}{\sqrt{n}} \ln 2^n \leq \ln 2 \cdot \sqrt{n},$$

and

$$\sum_{t=1}^T \frac{\eta_t}{2} < \sum_{t=1}^T \frac{\sqrt{n}}{(\sqrt{t-1} + \sqrt{t})} = \sum_{t=1}^T \sqrt{n} \cdot (\sqrt{t} - \sqrt{t-1}) = \sqrt{nT}. \quad (12)$$

$$\begin{aligned}
& \sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l,t} = \mathbf{i}_l] \\
\stackrel{(E1)}{=} & \sum_{t \in \mathcal{T}} \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \Pr[\mathbf{I}_{l+1,t} = \mathbf{i}_{l+1} | \mathbf{I}_{l,t} = \mathbf{i}_l] \cdot \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l+1,t} = \mathbf{i}_{l+1}] \\
\stackrel{(E2)}{=} & \sum_{t \in \mathcal{T}} \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \Pr[\mathbf{I}_{l+1,t} = \mathbf{i}_{l+1} | \mathbf{I}_{l,t} = \mathbf{i}_l] \cdot \left[\bar{c}_{l+1,t}(\mathbf{i}_{l+1}) \cdot 2^{m-l} L_m + \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \right] \\
= & \sum_{t \in \mathcal{T}} \left\{ \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \Pr[\mathbf{I}_{l+1,t} = \mathbf{i}_{l+1} | \mathbf{I}_{l,t} = \mathbf{i}_l] \cdot \bar{c}_{l+1,t}(\mathbf{i}_{l+1}) \cdot 2^{m-l} L_m + \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \Pr[\mathbf{I}_{l+1,t} = \mathbf{i}_{l+1} | \mathbf{I}_{l,t} = \mathbf{i}_l] \cdot \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \right\} \\
= & \sum_{t \in \mathcal{T}} \left\{ \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \Pr[\mathbf{I}_{l+1,t} = \mathbf{i}_{l+1} | \mathbf{I}_{l,t} = \mathbf{i}_l] \cdot \bar{c}_{l+1,t}(\mathbf{i}_{l+1}) \cdot 2^{m-l} L_m + \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \right\} \\
\stackrel{(E3)}{=} & \sum_{t \in \mathcal{T}} \left[\sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1}))}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} \cdot \bar{c}_{l+1,t}(\mathbf{i}_{l+1}) \right] \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E4)}{\leq} & \sum_{t \in \mathcal{T}} \left\{ \underbrace{\left[-\frac{1}{\eta_t} \ln \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1})) \cdot \exp(-\eta_t \bar{c}_{l+1,t}(\mathbf{i}_{l+1}))}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} + \frac{\eta_t}{2} \right]}_{\text{See (10), which is proved in Appendix A.1}} \cdot 2^{m-l} L_m + \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \right\} \\
\stackrel{(E5)}{=} & \sum_{t \in \mathcal{T}} \left[-\frac{1}{\eta_t} \ln \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t}(\mathbf{i}_{l+1}))}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} + \frac{\eta_t}{2} \right] \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E6)}{=} & \sum_{t \in \mathcal{T}} \left[\Phi_t(\eta_t) - \Phi_{t-1}(\eta_t) + \frac{\eta_t}{2} \right] \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E7)}{=} & \left\{ \Phi_T(\eta_T) + \underbrace{\sum_{t=1}^{T-1} (\Phi_t(\eta_t) - \Phi_t(\eta_{t+1})) - \Phi_0(\eta_1)}_{\text{See (11) which is proved in Appendix A.2}} + \sum_{t=1}^T \frac{\eta_t}{2} \right\} \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E8)}{\leq} & \left\{ \underbrace{\Phi_T(\eta_T) + \ln 2 \cdot \sqrt{nT}}_{\text{See (11) which is proved in Appendix A.2}} - \Phi_0(\eta_1) + \sum_{t=1}^T \frac{\eta_t}{2} \right\} \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E9)}{\leq} & \left\{ \underbrace{\Phi_T(\eta_T) + \ln 2 \cdot \sqrt{nT} + \ln 2 \cdot \sqrt{n} + \sqrt{nT}}_{\text{See (13) which is proved in Appendix A.3}} \right\} \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E10)}{\leq} & \left\{ \sum_{t \in \mathcal{T}} \mathbb{E} \left[\frac{c_{m,t}(\mathbf{I}_{m,t}) - \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m)}{2^{m-l} L_m} | \mathbf{I}_{l+1,t} = \mathbf{i}''_{l+1} \right] + (\ln 2 + 1) \sqrt{nT} + \ln 2 \cdot \sqrt{n} \right\} \cdot 2^{m-l} L_m + \sum_{t \in \mathcal{T}} \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_l(\mathbf{i}_l))} c_{m,t}(\mathbf{i}_m) \\
\stackrel{(E11)}{=} & \sum_{t \in \mathcal{T}} \mathbb{E} [c_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l+1,t} = \mathbf{i}''_{l+1}] + \left((\ln 2 + 1) \sqrt{nT} + \ln 2 \cdot \sqrt{n} \right) \cdot 2^{m-l} L_m \\
\stackrel{(E12)}{=} & \sum_{t \in \mathcal{T}} \mathbb{E} [c_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l+1,t} = \mathbf{i}''_{l+1}] + \underbrace{\left((\ln 2 + 1) nDL\sqrt{T} + \ln 2 \cdot nDL \right)}_{\text{Denoted as } \phi} \cdot \frac{1}{2^l}, \quad \forall \mathbf{i}''_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))
\end{aligned}$$

(9)

Inequality (E10) is due to the following inequality

$$\Phi_T(\eta_T) \leq \sum_{t \in \mathcal{T}} \mathbb{E} \left[\frac{c_{m,t}(\mathbf{I}_{m,t}) - \min_{i_m \in \mathcal{E}_m(\mathcal{K}_l(i_l))} c_{m,t}(i_m)}{2^{m-l} L_m} | I_{l+1,t} = i_{l+1}'' \right], \quad (13)$$

which is proved in Appendix A.3. Equality (E11) combines the first term and the last term. Equality (E12) follows from the fact that $L_m = \sqrt{n}DL/2^m$.

For any cost functions (f_1, f_2, \dots, f_T) , let i_m^* be the index of the sample point with the smallest cumulative cost, i.e.,

$$i_m^* \stackrel{\text{def}}{=} \arg \min_{i_m \in \mathcal{I}_m} \sum_{t \in \mathcal{T}} c_{m,t}(i_m).$$

By repeatedly applying the result in Equation (9), we have

$$\begin{aligned} & \sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t})] - \min_{i_m \in \mathcal{I}_m} \sum_{t \in \mathcal{T}} c_{m,t}(i_m) \\ &= \sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | I_{0,t} = 1] - \sum_{t \in \mathcal{T}} c_{m,t}(i_m^*) \\ &\leq \sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | I_{1,t} = \mathbf{v}_1(i_m^*)] + \phi \cdot \frac{1}{2^0} - \sum_{t \in \mathcal{T}} c_{m,t}(i_m^*) \\ &\leq \sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | I_{2,t} = \mathbf{v}_2(i_m^*)] + \phi \cdot \frac{1}{2^1} + \phi \cdot \frac{1}{2^0} - \sum_{t \in \mathcal{T}} c_{m,t}(i_m^*) \\ &\vdots \\ &\leq \sum_{t \in \mathcal{T}} \mathbb{E}[c_{m,t}(\mathbf{I}_{m,t}) | I_{m,t} = \mathbf{v}_m(i_m^*)] + \phi \cdot \left(\frac{1}{2^{m-1}} + \dots + \frac{1}{2^0} \right) - \sum_{t \in \mathcal{T}} c_{m,t}(i_m^*) \\ &= \phi \cdot \left(\frac{1}{2^{m-1}} + \dots + \frac{1}{2^0} \right) = 2\phi \cdot \left[1 - \frac{1}{2^m} \right] \leq 2\phi. \end{aligned} \quad (14)$$

Since (14) holds for any cost functions (f_1, f_2, \dots, f_T) , we have

$$\text{regret}_{\text{lmC},T}(\mathfrak{R}) \leq 2\phi = 2(\ln 2 + 1)nDL\sqrt{T} + 2 \ln 2 \cdot nDL.$$

This completes the proof. \square

Remark 4. At each layer, the ORW algorithm chooses among at most 2^n subsets. Moreover, due to the Lipschitz condition, the expected normalized cost difference among subsets to be chosen at each layer decreases exponentially as the algorithm goes down. These two facts result in the exponentially decreasing of the regret loss due to “imperfect choice” at each layer. This is also implied by Equation (9) in the proof of Lemma 4.3. Thus, the total regret loss of “imperfect choice”, $\text{regret}_{\text{lmC},T}$ is always upper bounded by $O(\sqrt{T})$ in timescale, no matter how many layers there are.

Moreover, the regret loss due to imperfect discretization can be reduced with larger m . Now, we state the the following lemma.

LEMMA 4.4. *The ORW algorithm guarantees that*

$$\text{regret}_{\text{lmD},T}(\mathfrak{R}) \leq \frac{\sqrt{n}}{2^m} DLT.$$

PROOF. For any cost functions (f_1, f_2, \dots, f_T) , let

$$\mathbf{x}^* \stackrel{\text{def}}{=} \arg \min_{\mathbf{x} \in \mathcal{K}} \sum_{t \in \mathcal{T}} f_t(\mathbf{x}).$$

Suppose that layer m subset containing \mathbf{x}^* is indexed by \mathbf{i}_m^* . The sample point for subset $\mathcal{K}_m(\mathbf{i}_m^*)$ is $\mathbf{r}_m(\mathbf{i}_m^*) \in \mathcal{K}$. Clearly, $\mathbf{r}_m(\mathbf{i}_m^*) = (r_1, r_2, \dots, r_n)$ and $\mathbf{x}^* = (x_1, x_2, \dots, x_n)$ are in the same sub-cube whose edge length is $\frac{D}{2^m}$. Thus, we have

$$\begin{aligned} \|\mathbf{r}_m(\mathbf{i}_m^*) - \mathbf{x}^*\|_2 &= \sqrt{(r_1 - x_1)^2 + (r_2 - x_2)^2 + \dots + (r_n - x_n)^2} \\ &\leq \sqrt{\left(\frac{D}{2^m}\right)^2 + \left(\frac{D}{2^m}\right)^2 + \dots + \left(\frac{D}{2^m}\right)^2} = \frac{D\sqrt{n}}{2^m}. \end{aligned}$$

Then, we have

$$\begin{aligned} &\min_{\mathbf{i}_m \in \mathcal{I}_m} \sum_{t \in \mathcal{T}} c_{m,t}(\mathbf{i}_m) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t \in \mathcal{T}} f_t(\mathbf{x}) \\ &\leq \sum_{t \in \mathcal{T}} c_{m,t}(\mathbf{i}_m^*) - \sum_{t \in \mathcal{T}} f_t(\mathbf{x}^*) \\ &= \sum_{t \in \mathcal{T}} f_t(\mathbf{r}_m(\mathbf{i}_m^*)) - f_t(\mathbf{x}^*) \leq TL \|\mathbf{r}_m(\mathbf{i}_m^*) - \mathbf{x}^*\|_2 \leq \frac{\sqrt{n}}{2^m} DLT, \end{aligned} \tag{15}$$

where the first inequality is due to the Lipschitz condition. Since (15) holds for any cost functions (f_1, f_2, \dots, f_T) , we conclude that $\text{regret}_{\text{ImD},T}(\mathcal{R}) \leq \frac{\sqrt{n}}{2^m} DLT$. This completes the proof. \square

Putting together the results in lemmas 4.2, 4.3, and 4.4, the results in Theorem 4.1 is proved.

5 ADAPTIVE SAMPLING

The ORW algorithm described in Section 3, specifies a granularity parameter, m , and samples the costs at the bottommost layer (layer m). With m larger than $\log_2 \sqrt{nT}$, the ORW algorithm guarantees the regret to be $O(\sqrt{T})$. In many cases, however, the duration interval is long and unknown to the online player, and setting a large m to guarantee $O(\sqrt{T})$ regret might be costly to compute. In order to handle cases with unknown T and reduce this complexity, in this section, we devise an adaptive version of the ORW algorithm (called the AORW algorithm), which increases the granularity parameter gradually. The AORW algorithm is denoted by $\mathcal{R}^{\text{adpt}}$ in equations.

5.1 Online Recursive Weighting Algorithm with Adaptive Sampling

The AORW algorithm adopts an increasing granularity parameter m_t . We denote the time interval that satisfies $m_t \geq l$ as \mathcal{T}_l or $[t_l, T]$. With increasing m_t , the index and the index set for the sampling subsets are updated to the new notation, \mathbf{i}_{m_t} and \mathcal{I}_{m_t} , respectively. $\mathbf{r}_{m_t}(\mathbf{i}_{m_t})$ is used to denote the sample point for the subset of index \mathbf{i}_{m_t} , and the cost on the sampled point of the subset $\mathbf{i}_{m_t} \in \mathcal{I}_{m_t}$ at iteration t is denoted as

$$c'_{m_t,t}(\mathbf{i}_{m_t}) \stackrel{\text{def}}{=} f_t(\mathbf{r}_{m_t}(\mathbf{i}_{m_t})).$$

The AORW algorithm maintains the *cumulative expected cost*, $\tilde{C}'_{l,t}(\mathbf{i}_l)$, for each subset at layers $\{1, 2, \dots, m_t\}$. The cumulative expected cost is initialized to be zero, i.e., $\tilde{C}'_{l,0}(\mathbf{i}_l) = 0$. During run time, the cumulative expected cost is used to determine the choosing probability of subsets at layers $\{1, 2, \dots, m_t\}$. The selection probability for the subset $\mathcal{K}_{l+1}(\mathbf{i}_{l+1})$ conditioning on that $\mathcal{K}_l(\mathbf{i}_l)$ has

$$\begin{aligned}
\bar{c}'_{l,t}(\mathbf{i}_l) &\stackrel{\text{def}}{=} \mathbb{E} \left[\frac{c_{m_t,t}(\mathbf{I}_{m_t,t}) - \min_{\mathbf{i}_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_{l-1}(\mathbf{v}_{l-1}(\mathbf{i}_l)))} c_{m_t,t}(\mathbf{i}_{m_t})}{2^{m_t-l+1} L_{m_t}} \mid \mathbf{I}_{l,t} = \mathbf{i}_l \right] \\
&= \sum_{\mathbf{j}_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(\mathbf{i}_l))} \frac{c_{m_t,t}(\mathbf{j}_{m_t}) - \min_{\mathbf{i}_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_{l-1}(\mathbf{v}_{l-1}(\mathbf{i}_l)))} c_{m_t,t}(\mathbf{i}_{m_t})}{2^{m_t-l+1} L_{m_t}} \cdot \Pr [\mathbf{I}_{m_t,t} = \mathbf{j}_{m_t} \mid \mathbf{I}_{l,t} = \mathbf{i}_l].
\end{aligned} \tag{17}$$

been selected is

$$\Pr [\mathbf{I}_{l+1,t} = \mathbf{i}_{l+1} \mid \mathbf{I}_{l,t} = \mathbf{i}_l] = \frac{\exp \left(-\eta_t \bar{C}'_{l+1,t-1}(\mathbf{i}_{l+1}) \right)}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp \left(-\eta_t \bar{C}'_{l+1,t-1}(\mathbf{i}'_{l+1}) \right)}. \tag{16}$$

Since there is a subset at each layer chosen in a recursive manner from the top layer to the bottom layer, the action sequence of the AORW algorithm at iteration t can be correspondingly denoted by a vector, $(\mathbf{I}_{0,t}, \mathbf{I}_{1,t}, \mathbf{I}_{0,t}, \dots, \mathbf{I}_{m_t,t})$, where $\mathbf{I}_{0,t}$ is set to be 1 by default.

After the costs on the sampled points of the bottom-layer subsets are revealed, the AORW algorithm calculates the *expected normalized cost* for each subset at layers $\{1, 2, \dots, m_t\}$, which is defined in Equation (17).

Then, the calculated *expected normalized cost* is used to update the cumulative expected cost, i.e.,

$$\bar{C}'_{l,t}(\mathbf{i}_l) = \bar{C}'_{l,t-1}(\mathbf{i}_l) + \bar{c}'_{l,t}(\mathbf{i}_l). \tag{18}$$

For the case when $m_{t+1} = m_t + 1$, the AORW algorithm will add the cumulative expected cost for those subsets at the new layer to maintenance, with previous values over $[1, t_{m_{t+1}} - 1]$ set to be 0, i.e.,

$$\bar{c}'_{m_{t+1},\tau}(\mathbf{i}_{m_{t+1}}) = 0, \text{ for } \tau = 0, 1, 2, \dots, t_{m_{t+1}} - 1.$$

The summary of the AORW algorithm is listed in Algorithm 2.

5.2 Regret Analysis

By adaptive sampling, the algorithm parameter does not depend on the length of time interval. Meanwhile, an interesting result of the regret analysis is that adaptive sampling does not degrade the regret bound. We state the main result in the following theorem.

THEOREM 5.1. *With $\eta_t = \sqrt{\frac{n}{t}}$ and m_t being set to $\lceil \log_2 \sqrt{nt} \rceil$, the AORW algorithm achieves the following regret bound*

$$\text{regret}_T(\mathcal{R}^{\text{adpt}}) \leq (6n + 2\sqrt{n}) DL\sqrt{T}.$$

Let $\mathbf{i}_{m_t}^*$, $t \in \mathcal{T}$ indicate the bottom subset that contains the optimal decision point \mathbf{x}^* at iteration t , i.e.,

$$\mathbf{i}_{m_t}^* : \mathbf{x}^* \in \mathcal{K}_{m_t}(\mathbf{i}_{m_t}^*), \quad t \in \mathcal{T}.$$

Similar to the regret analysis in Section 4, we define the regret due to the imperfect choice among sample points as

$$\text{regret}_{\text{ImC},T}(\mathcal{R}^{\text{adpt}}) \stackrel{\text{def}}{=} \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t \in \mathcal{T}} \mathbb{E} [c'_{m_t,t}(\mathbf{I}_{m_t,t})] - \sum_{t \in \mathcal{T}} c'_{m_t,t}(\mathbf{i}_{m_t}^*) \right\},$$

and that due to the imperfect sampling as

$$\text{regret}_{\text{ImD},T}(\mathfrak{R}^{\text{adpt}}) \stackrel{\text{def}}{=} \sup_{f_1, \dots, f_T \in \mathcal{F}} \left\{ \sum_{t \in \mathcal{T}} c'_{m_t, t}(\mathbf{i}_{m_t}^*) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t \in \mathcal{T}} f_t(\mathbf{x}) \right\}.$$

Similar to Lemma 3.1, we have

$$\text{regret}_T(\mathfrak{R}^{\text{adpt}}) \leq \text{regret}_{\text{ImC},T}(\mathfrak{R}^{\text{adpt}}) + \text{regret}_{\text{ImD},T}(\mathfrak{R}^{\text{adpt}}).$$

In the following two lemmas, we bound the regret loss due to imperfect choice and sampling, respectively.

LEMMA 5.2. *With $\eta_t = \sqrt{\frac{n}{t}}$ and the $m_t = \lceil \log_2 \sqrt{t} \rceil$, the AORW algorithm guarantees that*

$$\text{regret}_{\text{ImC},T}(\mathfrak{R}^{\text{adpt}}) \leq 6DLn\sqrt{T}.$$

The proof of Lemma 5.2 is analogous to that of Lemma 4.3. The details are given in Appendix A.4.

LEMMA 5.3. *With $m_t = \lceil \log_2 \sqrt{t} \rceil$, the AORW algorithm guarantees that*

$$\text{regret}_{\text{ImD},T}(\mathfrak{R}^{\text{adpt}}) \leq 2DL\sqrt{nT}.$$

PROOF. For any cost functions (f_1, f_2, \dots, f_T) , let

$$\mathbf{x}^* \stackrel{\text{def}}{=} \arg \min_{\mathbf{x} \in \mathcal{K}} \sum_{t \in \mathcal{T}} f_t(\mathbf{x}).$$

Suppose that the bottom subset containing point \mathbf{x}^* is indexed by $\mathbf{i}_{m_t}^*$, whose sample point is $\mathbf{r}_{m_t}(\mathbf{i}_{m_t}^*) \in \mathcal{K}$. Clearly, $\mathbf{r}_{m_t}(\mathbf{i}_{m_t}^*) = (r_{m_t,1}, r_{m_t,2}, \dots, r_{m_t,n})$ and $\mathbf{x}^* = (x_1, x_2, \dots, x_n)$ are in the same sub-cube whose edge length is $\frac{D}{2^{m_t}}$. Thus, we have

$$\begin{aligned} \|\mathbf{r}_{m_t}(\mathbf{i}_{m_t}^*) - \mathbf{x}^*\|_2 &= \sqrt{(r_{m_t,1} - x_1)^2 + (r_{m_t,2} - x_2)^2 + \dots + (r_{m_t,n} - x_n)^2} \\ &\leq \sqrt{\left(\frac{D}{2^{m_t}}\right)^2 + \left(\frac{D}{2^{m_t}}\right)^2 + \dots + \left(\frac{D}{2^{m_t}}\right)^2} = \frac{D\sqrt{n}}{2^{m_t}} \\ &\leq D\sqrt{\frac{n}{t}}. \end{aligned}$$

Then we have

$$\begin{aligned} \sum_{t \in \mathcal{T}} c'_{m_t, t}(\mathbf{i}_{m_t}^*) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t \in \mathcal{T}} f_t(\mathbf{x}) &= \sum_{t \in \mathcal{T}} c'_{m_t, t}(\mathbf{i}_{m_t}^*) - \sum_{t \in \mathcal{T}} f_t(\mathbf{x}^*) \\ &= \sum_{t \in \mathcal{T}} f_t(\mathbf{r}_{m_t}(\mathbf{i}_{m_t}^*)) - f_t(\mathbf{x}^*) \\ &\stackrel{(E1)}{\leq} L \sum_{t \in \mathcal{T}} \|\mathbf{r}_{m_t}(\mathbf{i}_{m_t}^*) - \mathbf{x}^*\|_2 \leq 2DL\sqrt{nT}, \end{aligned} \tag{19}$$

where inequality (E1) is due to the Lipschitz condition. Since (19) holds for any cost functions (f_1, f_2, \dots, f_T) , we conclude that $\text{regret}_{\text{ImD},T}(\mathfrak{R}) \leq 2DL\sqrt{nT}$. This completes the proof. \square

The results in lemmas 5.2 and 5.3 immediately prove the results in Theorem 5.1.

Algorithm 2 Online Recursive Weighting with Adaptive Sampling**Input:** $\{\eta_t = \sqrt{\frac{n}{t}}\}$, $\{m_t = \lceil \log_2 \sqrt{t} \rceil\}$ **Output:** $r_{m_1}(I_{m_1,1}), r_{m_2}(I_{m_2,2}), \dots$

```

1: //↑ Output sample points of layer  $m$  subsets
2: //↓ Initialize the normalized expected cost for all subsets at all layers
3: Set  $\bar{C}_{l,0}(i_l) = 0, \forall i_l \in \mathcal{I}_l, l = 1, 2, \dots, m_1$ 
4: for  $t = 1, 2, \dots, T$  do
5:    $I_{0,t} = 1$ 
6:   //↓ Recursively select the subsets at all layers
7:   for  $l = 0, 1, \dots, m_t - 1$  do
8:     Select an index  $I_{l+1,t} \in \mathcal{E}_{l+1}(I_{l,t})$  according to the probability distribution calculated in
       Equation (16)
9:   end for
10:  Choose the subset indexed by  $I_{m_t,t}$  at iteration  $t$ 
11:  Choose the sample point for subset  $\mathcal{K}_{m_t}(I_{m_t,t})$ , i.e.,  $r_{m_t}(I_{m_t,t})$ , as the final decision
12:  //↓ Get the normalized expected cost for all subsets at all layers based
    on the revealed cost function  $f_t(x)$  at iteration  $t$ 
13:  for  $l = 1, 2, \dots, m_t$  do
14:    for  $i_l \in \mathcal{I}_l$  do
15:      for  $j_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))$  do
16:        Calculate  $\Pr[I_{m_t,t} = j_{m_t} | I_{l,t} = i_l]$  according to Equation (5)
17:      end for
18:      Calculate  $\bar{c}_{l,t}(i_l)$  according to Equation (17)
19:      Update  $\bar{C}'_{l,t}(i_l) = \bar{C}'_{l,t-1}(i_l) + \bar{c}'_{l,t}(i_l)$ 
20:    end for
21:  end for
22:  //↓ Initialize the cumulative expected cost for all subsets at the new
    layer
23:  if  $m_{t+1} > m_t$  then
24:    for  $i_{m_{t+1}} \in \mathcal{I}_{m_{t+1}}$  do
25:      Set  $\bar{C}'_{m_{t+1},t}(i_{m_{t+1}}) = 0$ 
26:    end for
27:  end if
28: end for

```

6 EXTENSIONS AND STATE-OF-THE-ART RESULTS**6.1 Online Optimization in a Decentralized Environment**

Recently, Hosseini [23] and Lee [33] generalized the classic OCO problem to a decentralized optimization framework within a network of agents. In their work, consensus-based gradient-descent algorithms were proposed for distributed online optimization. In their setting, each agent aims to drive its individual average regret, which is the average over time of the regret function evaluated at this agent's estimation for the choice that the whole network should make, to zero. An interesting message of their work is that the $O(\sqrt{T})$ regret can still be attained by leveraging the communication among agents. In addition, [43] addresses decentralized online optimization in non-stationary environments using mirror decent, and in [4], distributed online optimization is studied for strongly convex objective functions over time-varying networks.

In parallel, a promising future work for our approach is to implement the proposed weighting method to solve the distributed ONCL problem. The extension is natural, while techniques to be adopted might be rather different. That is due to the fact that the ORW algorithm has to maintain an estimation for each subset and it might be very costly to exchange such information. Thus, in the opinion of the authors, the main challenge to implement such a weighting method within a decentralized environment is to alleviate the communication overhead among agents.

6.2 Partial Information Feedback

We emphasize that our online non-convex learning problem is based on *full information feedback*. That is, the whole cost function will be revealed at the end of each iteration. It is an interesting and important future direction to consider *partial information feedback* where only the cost value of the player's choice is revealed. The *partial information feedback* extension is motivated by many real-world systems in which the observer is not co-located with the controller and the feedback information is noisy, partial, or incomplete due to limited communication bandwidth. Some examples are online routing in data networks [7], power control in cellular networks [48], and the ad placement problem on a web page [47].

The bandit information feedback setting has been investigated in a huge number of works in the OCO framework such as [1, 2, 11, 13, 21, 22, 41], which is called the Bandit Convex Optimization (BCO) problem. Among those works, [1, 11, 13, 41] reported to attain a sublinear regret for the BCO problem, respectively, but none of them have attained the lower bound of the regret. For the special case of strongly convex and smooth losses, [2] obtained a regret of $\tilde{O}(\sqrt{T})$ in the unconstrained case, and [21] obtained the same rate even in the constrained case. Recently, a new algorithm was reported by Hazan *et al.* to attain a regret of $(\ln T)^{2d} \sqrt{T}$ [22]. This is the first algorithm to attain a $\tilde{O}(\sqrt{T})$ regret for the OCO model with bandit feedback. Different from the above works, [3, 28–30, 38] studied the model where the (expected) payoff function satisfies a Lipschitz condition with respect to the metric. [3, 30, 38] investigate this problem in a few specific metric spaces such as a one-dimensional real interval, while in [28, 29], the action set is from a general metric space. Their model is more general and the results in [29] show that there is an algorithm whose regret on any instance satisfies $R(T) = \tilde{O}(T^{\frac{d+1}{d+2}})$, where d is the dimension of the action set.

Despite of the above results, the optimal online algorithm and tight regret bound for the online non-convex optimization problem with bandit feedback are still open, calling for more investigation from the community.

7 CONCLUSION

In this paper, we investigated the online non-convex learning problem, which removes the convexity assumption of the cost function of the classic online convex optimization problem. This generalization is necessitated by a huge amount of important applications which incur non-convex penalties, such as in portfolio selection and SVM training problems. Without the convexity assumption, it is far more challenging to design an efficient online algorithm with sublinear regret. The classic exponential weighting online algorithm is shown to attain a sublinear regret of $O(\sqrt{T} \log T)$ and the convex optimization methods even fail to converge to the optimum. Our analysis prove that by properly partitioning subsets and using the online recursive weighting method, the regret can be reduced to match the known lower bound for OCO, i.e., $O(\sqrt{T})$.

ACKNOWLEDGEMENT

The authors would like to thank Professor Giulia Fanti for her kind shepherding and the anonymous reviewers of SIGMETRICS 2018 for their insightful comments and suggestions to improve the quality of the paper.

REFERENCES

- [1] J. Abernethy, E. Hazan, and A. Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *COLT*, pages 263–274, 2008.
- [2] A. Agarwal, O. Dekel, and L. Xiao. Optimal algorithms for online convex optimization with multi-point bandit feedback. In *COLT*, pages 28–40. Citeseer, 2010.
- [3] R. Agrawal. The continuum-armed bandit problem. *SIAM journal on control and optimization*, 33(6):1926–1951, 1995.
- [4] M. Akbari, B. Gharesifard, and T. Linder. Distributed online convex optimization on time-varying directed graphs. *IEEE Transactions on Control of Network Systems*, 2015.
- [5] D. Ardia, K. Boudt, P. Carl, K. M. Mullen, and B. Peterson. Differential evolution (deoptim) for non-convex portfolio optimization. 2010.
- [6] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [7] B. Awerbuch and R. Kleinberg. Online linear optimization and adaptive routing. *Journal of Computer and System Sciences*, 74(1):97–114, 2008.
- [8] N. Cesa-Bianchi, Y. Freund, D. Haussler, D. P. Helmbold, R. E. Schapire, and M. K. Warmuth. How to use expert advice. *Journal of the ACM (JACM)*, 44(3):427–485, 1997.
- [9] R. Combes, S. Magureanu, A. Proutiere, and C. Laroche. Learning to rank: Regret lower bounds and efficient algorithms. *ACM SIGMETRICS Performance Evaluation Review*, 43(1):231–244, 2015.
- [10] T. Cover. Universal portfolios. *Mathematical finance*, 1(1):1–29, 1991.
- [11] O. Dekel, R. Eldan, and T. Koren. Bandit smooth convex optimization: Improving the bias-variance tradeoff. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2926–2934, 2015.
- [12] S. Ertekin, L. Bottou, and C. Giles. Non-convex online support vector machines. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(2):368–381, 2011.
- [13] A. Flaxman, A. Kalai, and H. McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms (SODA)*, pages 385–394, 2005.
- [14] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *European conference on computational learning theory*, pages 23–37. Springer, 1995.
- [15] G. Gasso, L. Pappaioannou, M. Spivak, and L. Bottou. Batch and online learning algorithms for nonconvex neyman-pearson classification. *ACM Transactions on Intelligent Systems and Technology*, 2(3):28, 2011.
- [16] A. J. Grove, N. Littlestone, and D. Schuurmans. General convergence results for linear discriminant updates. *Machine Learning*, 43(3):173–210, 2001.
- [17] T. Guzella and W. Caminhas. A review of machine learning approaches to spam filtering. *Expert Systems with Applications*, 36(7):10206–10222, 2009.
- [18] E. Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3–4):157–325, 2016.
- [19] E. Hazan and S. Kale. Online submodular minimization. *Journal of Machine Learning Research*, 13(Oct):2903–2922, 2012.
- [20] E. Hazan and S. Kale. Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization. *Journal of Machine Learning Research*, 15(1):2489–2512, 2014.
- [21] E. Hazan and K. Levy. Bandit convex optimization: Towards tight bounds. In *Advances in Neural Information Processing Systems (NIPS)*, pages 784–792, 2014.
- [22] E. Hazan and Y. Li. An optimal algorithm for bandit convex optimization. *arXiv preprint arXiv:1603.04350*.
- [23] S. Hosseini, A. Chapman, and M. Mesbahi. Online distributed convex optimization on dynamic networks. *IEEE Transactions on Automatic Control*, 61(11):3545–3550, 2016.
- [24] P. Jain, B. Kulis, I. S. Dhillon, and K. Grauman. Online metric learning and fast similarity search. In *Advances in neural information processing systems*, pages 761–768, 2009.
- [25] A. Kalai and S. Vempala. Efficient algorithms for universal portfolios. *Journal of Machine Learning Research*, 3(Nov):423–440, 2002.
- [26] J. Kivinen and M. K. Warmuth. Relative loss bounds for multidimensional regression problems. In *Advances in neural information processing systems (NIPS)*, pages 287–293, 1998.

- [27] R. Kleinberg and A. Slivkins. Sharp dichotomies for regret minimization in metric spaces. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 827–846. Society for Industrial and Applied Mathematics, 2010.
- [28] R. Kleinberg, A. Slivkins, and E. Upfal. Multi-armed bandits in metric spaces. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 681–690. ACM, 2008.
- [29] R. Kleinberg, A. Slivkins, and E. Upfal. Bandits and experts in metric spaces. *arXiv preprint arXiv:1312.1277*, 2013.
- [30] R. D. Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In *Advances in Neural Information Processing Systems*, pages 697–704, 2005.
- [31] W. Krichene, M. Balandat, C. Tomlin, and A. Bayen. The hedge algorithm on a continuum. In *the 32nd International Conference on Machine Learning (ICML-15)*, pages 824–832, 2015.
- [32] P. Krokmal, J. Palmquist, and S. Uryasev. Portfolio optimization with conditional value-at-risk objective and constraints. *Journal of risk*, 4:43–68, 2002.
- [33] S. Lee, A. Nedich, and M. Raginsky. Stochastic dual averaging for decentralized online optimization on time-varying communication graphs. *IEEE Transactions on Automatic Control*, 2017.
- [34] N. Littlestone and M. Warmuth. The weighted majority algorithm. *Information and computation*, 108(2):212–261, 1994.
- [35] S. Magureanu, A. Proutiere, M. Isaksson, and B. Zhang. Online learning of optimally diverse rankings. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(2):32, 2017.
- [36] O.-A. Maillard and R. Munos. Online learning in adversarial lipschitz environments. *Machine Learning and Knowledge Discovery in Databases*, pages 305–320, 2010.
- [37] L. Mason, P. L. Bartlett, and J. Baxter. Improved generalization through explicit optimization of margins. *Machine Learning*, 38(3):243–255, 2000.
- [38] S. Pandey, D. Agarwal, D. Chakrabarti, and V. Josifovski. Bandits for taxonomies: A model-based approach. In *Proceedings of the 2007 SIAM International Conference on Data Mining*, pages 216–227. SIAM, 2007.
- [39] A. Rakhlin, O. Shamir, and K. Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 449–456, 2012.
- [40] H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [41] A. Saha and A. Tewari. Improved regret guarantees for online smooth convex optimization with bandit feedback. In *AISTATS*, pages 636–642, 2011.
- [42] D. Sculley and G. Wachman. Relaxed online svms for spam filtering. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 415–422, 2007.
- [43] S. Shahrampour and A. Jadbabaie. Distributed online optimization in dynamic environments using mirror descent. *IEEE Transactions on Automatic Control*, 2017.
- [44] O. Shamir and T. Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *International Conference on Machine Learning*, pages 71–79, 2013.
- [45] M. S. Talebi, Z. Zou, R. Combes, A. Proutiere, and M. Johansson. Stochastic online shortest path routing: The value of feedback. *IEEE Transactions on Automatic Control*, 2017.
- [46] S. Uryasev. Conditional value-at-risk: Optimization algorithms and applications. In *Computational Intelligence for Financial Engineering, 2000.(CIFEr) Proceedings of the IEEE/IAFE/INFORMS 2000 Conference on*, pages 49–57. IEEE, 2000.
- [47] F. Wauthier, M. Jordan, and N. Jovic. Efficient ranking from pairwise comparisons. In *International Conference on Machine Learning (ICML)*, pages 109–117, 2013.
- [48] W. Wong and C. Sung. Robust convergence of low-data rate-distributed controllers. *IEEE transactions on automatic control*, 49(1):82–87, 2004.
- [49] H. H. Zhang, J. Ahn, X. Lin, and C. Park. Gene selection using support vector machines with non-convex penalty. *bioinformatics*, 22(1):88–95, 2005.
- [50] L. Zhang, T. Yang, R. Jin, and Z. Zhou. Online bandit learning for a special class of non-convex losses. In *AAAI*, pages 3158–3164, 2015.
- [51] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning (ICML)*, pages 928–936, 2003.

A PROOFS

A.1 Proof of Inequality (10)

The proof relies on the following Hoeffding’s Lemma.

LEMMA A.1. (*Hoeffding's Lemma*) Let Z be any real-valued random variable with expected value $\mathbb{E}[Z] = 0$ and such that $a \leq Z \leq b$ almost surely. Then, for all $\lambda \in \mathbb{R}$,

$$\mathbb{E}[\exp(\lambda Z)] \leq \exp\left(\frac{\lambda^2(b-a)^2}{8}\right).$$

We thus define a random variable Y whose probability mass function is specified as

$$\Pr[Y = \bar{c}_{l+1,t}(\mathbf{i}_{l+1})] = \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1}))}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))}, \quad (20)$$

for $\forall \mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))$.

Since $\bar{c}_{l+1,t}(\mathbf{i}_{l+1}) \in [0, 1]$, $\forall \mathbf{i}_{l+1}$, we have that $0 \leq Y \leq 1$ for sure. However, random variable Y may not have zero mean. We thus denote another random variable $Z = Y - \mathbb{E}[Y]$. Clearly, $\mathbb{E}[Z] = 0$ and $-1 \leq Z \leq 1$ for sure. We further let $\lambda = -\eta_t$. Thus, by Hoeffding's Lemma, we have

$$\mathbb{E}[\exp(-\eta_t Z)] \leq \exp\left(\frac{(-\eta_t)^2(1-(-1))^2}{8}\right) = \exp\left(\frac{\eta_t^2}{2}\right). \quad (21)$$

In addition, we have

$$\mathbb{E}[\exp(-\eta_t Z)] = \mathbb{E}[\exp(-\eta_t(Y - \mathbb{E}[Y]))] = \mathbb{E}[\exp(-\eta_t Y)] \cdot \exp[\eta_t \mathbb{E}[Y]]. \quad (22)$$

Combining (21) and (22), we get

$$\mathbb{E}[\exp(-\eta_t Y)] \cdot \exp(\eta_t \mathbb{E}[Y]) \leq \exp\left(\frac{\eta_t^2}{2}\right). \quad (23)$$

Taking logarithm and rearranging the items in (23), we have

$$\mathbb{E}[Y] \leq -\frac{1}{\eta_t} \ln \mathbb{E}[\exp(-\eta_t Y)] + \frac{\eta_t}{2}. \quad (24)$$

Now applying the probability mass function (20) into (24), we get that

$$\begin{aligned} \mathbb{E}[Y] &= \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1})) \cdot \bar{c}_{l+1,t}(\mathbf{i}_{l+1})}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} \\ &\leq -\frac{1}{\eta_t} \ln \mathbb{E}[\exp(-\eta_t Y)] + \frac{\eta_t}{2}. \\ &= -\frac{1}{\eta_t} \ln \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \frac{\exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}_{l+1})) \cdot \exp(-\eta_t \bar{c}_{l+1,t}(\mathbf{i}_{l+1}))}{\sum_{\mathbf{i}'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_t \bar{C}_{l+1,t-1}(\mathbf{i}'_{l+1}))} + \frac{\eta_t}{2}. \end{aligned}$$

A.2 Proof of Equality (11)

To prove (11), we only need to show that

$$\Phi_t(\eta_t) - \Phi_t(\eta_{t+1}) \leq \ln 2 \cdot n \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right). \quad (25)$$

The second inequality of (11) simply follows from $\eta_t = \sqrt{\frac{n}{t}}$. We next prove (25).

We call a vector $\mathbf{z} = (z_1, z_2, \dots, z_m) \in \mathbb{R}^m$ a *non-increasing non-negative vector* if $z_1 \geq z_2 \geq \dots \geq z_m \geq 0$. We then define a function

$$\Psi(a, \mathbf{z}) \stackrel{\text{def}}{=} -\frac{1}{a} \ln \sum_{j=1}^m \exp(-az_j), \quad a > 0,$$

where $\mathbf{z} = (z_1, z_2, \dots, z_m)$ is a non-increasing non-negative vector in \mathbb{R}^m . We further define a function $\tilde{h}(\bullet)$ mapping from a non-increasing non-negative vector to another non-increasing non-negative vector, where $\tilde{h}(\mathbf{z})$ changes all largest elements of \mathbf{z} into the second-largest element of \mathbf{z} . Namely, if the first k elements \mathbf{x} are the largest, i.e., $z_1 = z_2 = \dots = z_k > z_{k+1} \geq \dots \geq z_m$, then

$$\tilde{h}(\mathbf{z}) = \underbrace{(z_{k+1}, z_{k+1}, \dots, z_{k+1})}_{\text{In total } k}, z_{k+1}, z_{k+2}, \dots, z_m).$$

Clearly, if we apply function $\tilde{h}(\bullet)$ to a non-increasing non-negative vector \mathbf{x} for m times, we get a constant vector $\mathbf{z}_{\min} = (z_m, z_m, \dots, z_m)$.

We have the following lemma.

LEMMA A.2. For $a > b > 0$, any non-increasing non-negative vector \mathbf{z} , we have

$$\Psi(a, \mathbf{z}) - \Psi(b, \mathbf{z}) \leq \Psi(a, \tilde{h}(\mathbf{z})) - \Psi(b, \tilde{h}(\mathbf{z})). \quad (26)$$

PROOF. (i) If \mathbf{z} is a constant vector, i.e., $z_1 = z_2 = \dots = z_m$, then $\tilde{h}(\mathbf{z}) = \mathbf{z}$ and thus (26) holds as an equality.

(ii) If \mathbf{z} is not a constant vector, we assume that the largest element of \mathbf{z} is z and there are in total $k < m$ largest elements in \mathbf{z} , i.e., $z = z_1 = z_2 = \dots = z_k > z_{k+1} \geq z_{k+2} \geq \dots \geq z_m$.

Then

$$\begin{aligned} & \Psi(a, \mathbf{z}) - \Psi(b, \mathbf{z}) \\ &= -\frac{1}{a} \ln \left(k \exp(-az) + \sum_{j=k+1}^m \exp(-az_j) \right) + \frac{1}{b} \ln \left(k \exp(-bz) + \sum_{j=k+1}^m \exp(-bz_j) \right) \stackrel{\text{def}}{=} \Upsilon(z), \end{aligned} \quad (27)$$

where the last equality defines a function with respect to z satisfying $z > z_{k+1}$. Taking the derivative on $\Upsilon(z)$ with respect to z , we get

$$\frac{\partial \Upsilon(z)}{\partial z} = \frac{k \exp(-az)}{k \exp(-az) + \sum_{j=k+1}^m \exp(-az_j)} - \frac{k \exp(-bz)}{k \exp(-bz) + \sum_{j=k+1}^m \exp(-bz_j)}.$$

Define

$$\Theta(\sigma) \stackrel{\text{def}}{=} \frac{k \exp(-\sigma z)}{k \exp(-\sigma z) + \sum_{j=k+1}^m \exp(-\sigma z_j)}.$$

Then we have that

$$\frac{\partial \Upsilon(z)}{\partial z} = \Theta(a) - \Theta(b). \quad (28)$$

Taking derivative on $\Theta(\sigma)$ with respect to σ , we get

$$\frac{\partial \Theta(\sigma)}{\partial \sigma} = \frac{k \exp(-\sigma z) \sum_{j=k+1}^m (z_j - z) \exp(-\sigma z_j)}{\left(k \exp(-\sigma z) + \sum_{j=k+1}^m \exp(-\sigma z_j) \right)^2}.$$

Since $z > z_j$, for $j = k+1, k+2, \dots, m$, we have $\frac{\partial \Theta(a)}{\partial a} < 0$. Because $a > b$, so we have

$$\frac{\partial \Upsilon(z)}{\partial z} = \Theta(a) - \Theta(b) < 0, \quad \forall z > z_{i+1}.$$

Thus, $\Upsilon(z) < \Upsilon(z_{i+1})$. By noting that $\Upsilon(z_{i+1}) = \Psi(a, \tilde{h}(z)) - \Psi(b, \tilde{h}(z))$, we thus have

$$\Psi(a, z) - \Psi(b, z) < \Psi(a, \tilde{h}(z)) - \Psi(b, \tilde{h}(z)).$$

Part (i) and part (ii) complete the proof. \square

Based on Lemma A.2, we can immediately obtain the following result.

LEMMA A.3. *For $a > b > 0$ and any non-increasing non-negative vector $\mathbf{z} \in \mathbb{R}^m$, we have*

$$\Psi(a, \mathbf{z}) - \Psi(b, \mathbf{z}) \leq \left(\frac{1}{b} - \frac{1}{a} \right) \ln m. \quad (29)$$

PROOF. From Lemma A.2, we have

$$\begin{aligned} \Psi(a, \mathbf{z}) - \Psi(b, \mathbf{z}) &\leq \Psi(a, \tilde{h}(\mathbf{z})) - \Psi(b, \tilde{h}(\mathbf{z})) \\ &\leq \Psi(a, \tilde{h}(\tilde{h}(\mathbf{z}))) - \Psi(b, \tilde{h}(\tilde{h}(\mathbf{z}))) \\ &\leq \dots \\ &\leq \Psi(a, \mathbf{z}_{\min}) - \Psi(b, \mathbf{z}_{\min}) \\ &= \left(\frac{1}{b} - \frac{1}{a} \right) \ln m, \end{aligned}$$

where $\mathbf{z}_{\min} = (z_m, z_m, \dots, z_m)$. This completes the proof. \square

We can sort all $\bar{C}_{l+1,t}(\mathbf{i}_{l+1})$'s where $\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))$ in a descending order and construct a non-increasing non-negative vector $\mathbf{z} \in \mathbb{R}^{|\mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))|}$. Then based on Lemma A.2, we have

$$\begin{aligned} \Phi_t(\eta_t) - \Phi_t(\eta_{t+1}) &= \Psi(\eta_t, \mathbf{z}) - \Psi(\eta_{t+1}, \mathbf{z}) \\ &\leq \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln |\mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))| \\ &\leq \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) n \ln 2, \end{aligned} \quad (30)$$

where the last inequality follows from $|\mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))| \leq 2^n$.

This completes the proof for inequality (11).

A.3 Proof of Inequality (13)

First, we have

$$\exp \left(\eta_T \max_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathbf{i}_l)} [-\bar{C}_{l+1,T}(\mathbf{i}_{l+1})] \right) \leq \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathbf{i}_l)} \exp \left(-\eta_T \bar{C}_{l+1,T}(\mathbf{i}_{l+1}) \right). \quad (31)$$

Taking logarithm and dividing by $-\frac{1}{\eta_T} < 0$ in both sides of (31), we can get that

$$\begin{aligned}
\Phi_T(\eta_T) &= -\frac{1}{\eta_T} \ln \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_T \bar{C}_{l+1,T}(\mathbf{i}_{l+1})) \\
&\leq \min_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \bar{C}_{l+1,T}(\mathbf{i}_{l+1}) \\
&\leq \bar{C}_{l+1,T}(\mathbf{i}_{l+1}'') \quad (\forall \mathbf{i}_{l+1}'' \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))) \\
&= \sum_{t=1}^T \bar{c}_{l+1,t}(\mathbf{i}_{l+1}'') \\
&= \sum_{t=1}^T \mathbb{E} \left[\frac{c_{m,t}(\mathbf{I}_{m,t}) - \min_{\mathbf{i}_m \in \mathcal{E}_m(\mathcal{K}_{l+1}(\mathbf{v}_{l-1}(\mathbf{i}_l)))} c_{m,t}(\mathbf{i}_m)}{2^{m-l+1} L_m} | \mathbf{I}_{l+1,t} = \mathbf{i}_{l+1}'' \right], \tag{32}
\end{aligned}$$

where the second last equality follows from the definition of $\bar{C}_{l,t}(\mathbf{i}_l)$ in (8) and the last equality follows from the definition of $\bar{c}_{l,t}(\mathbf{i}_l)$ in (4) and $\mathbf{v}_l(\mathbf{i}_{l+1}'') = \mathbf{i}_l$.

This completes the proof.

A.4 Proof of Lemma 5.2

PROOF. Analogous to Equation (9), we have Equation (33) which characterizes the regret loss at each layer. In Equation (33), equality (E1) is according to the definition for the *normalized expected cost* shown in Equation (17). Equality (E2) is simply by the choosing probability for subsets at each layer. Inequality (E3) has been proved in Section (A.1). Equality (E4) follows from the update law for the *cumulative expected cost*. Inequality (E5) is based on the Equation (11) and (12), as well as the following inequality

$$\begin{aligned}
-\Phi_{l_{l+1}-1}(\eta_{l_{l+1}}) &= \frac{1}{\eta_{l_{l+1}}} \ln \sum_{\mathbf{i}_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))} \exp(-\eta_{l_{l+1}} \bar{C}_{l+1,t_{l+1}-1}(\mathbf{i}_{l+1})) \\
&= \frac{1}{\eta_{l_{l+1}}} \ln |\mathcal{E}_{l+1}(\mathcal{K}_l(\mathbf{i}_l))| \leq \sqrt{nt_{l+1}}.
\end{aligned}$$

Inequality (E6) is proved in Section A.3.

With the results in Equation (33), we have Equation (34).

In Equation (34), the equalities follow from the fact that

$$\begin{aligned}
&\sum_{t \in \mathcal{T}_l} \mathbb{E}[c'_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l,t} = \mathbf{i}_l^*] - \sum_{t \in \mathcal{T}_l} c'_{m,t}(\mathbf{i}_{m_t}^*) \\
&= \sum_{t \in \mathcal{T}_{l+1}} \mathbb{E}[c'_{m,t}(\mathbf{I}_{m,t}) | \mathbf{I}_{l,t} = \mathbf{i}_l^*] - \sum_{t \in \mathcal{T}_{l+1}} c'_{m,t}(\mathbf{i}_{m_t}^*).
\end{aligned}$$

□

Received February 2018; revised May 2018; accepted June 2018

$$\begin{aligned}
& \sum_{t \in \mathcal{T}_{l+1}} \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_{l, t} = i_l] \\
&= \sum_{t \in \mathcal{T}_{l+1}} \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \Pr[I_{l+1, t} = i_{l+1} | I_{l, t} = i_l] \cdot \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_{l+1, t} = i_{l+1}] \\
&\stackrel{(E1)}{=} \sum_{t \in \mathcal{T}_{l+1}} \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \Pr[I_{l+1, t} = i_{l+1} | I_{l, t} = i_l] \cdot \left[\bar{c}'_{l+1, t}(i_{l+1}) \cdot 2^{m_t-l} L_{m_t} + \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \right] \\
&= \sum_{t \in \mathcal{T}_{l+1}} \left\{ \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \Pr[I_{l+1, t} = i_{l+1} | I_{l, t} = i_l] \cdot \bar{c}'_{l+1, t}(i_{l+1}) \cdot 2^{m_t-l} L_{m_t} + \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \right\} \\
&\stackrel{(E2)}{=} \sum_{t \in \mathcal{T}_{l+1}} \left\{ \left[\sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \frac{\exp(-\eta_t \bar{C}'_{l+1, t-1}(i_{l+1}))}{\sum_{i'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\eta_t \bar{C}'_{l+1, t-1}(i'_{l+1}))} \cdot \bar{c}'_{l+1, t}(i_{l+1}) \right] \cdot 2^{m_t-l} L_{m_t} + \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \right\} \\
&\stackrel{(E3)}{\leq} \sum_{t \in \mathcal{T}_{l+1}} \left\{ \left[-\frac{1}{\eta_t} \ln \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \frac{\exp(-\eta_t \bar{C}'_{l+1, t-1}(i_{l+1})) \cdot \exp(-\eta_t \bar{c}'_{l+1, t}(i_{l+1}))}{\sum_{i'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\eta_t \bar{C}'_{l+1, t-1}(i'_{l+1}))} + \frac{\eta_t}{2} \right] \cdot 2^{m_t-l} L_{m_t} + \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \right\} \\
&\stackrel{(E4)}{=} \sum_{t \in \mathcal{T}_{l+1}} \left[-\frac{1}{\eta_t} \ln \sum_{i_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \frac{\exp(-\eta_t \bar{C}'_{l+1, t}(i_{l+1}))}{\sum_{i'_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))} \exp(-\eta_t \bar{C}'_{l+1, t-1}(i'_{l+1}))} + \frac{\eta_t}{2} \right] \cdot 2^{m_t-l} L_{m_t} + \sum_{t \in \mathcal{T}_{l+1}} \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \\
&= \sum_{t \in \mathcal{T}_{l+1}} \left[\Phi_t(\eta_t) - \Phi_{t-1}(\eta_t) + \frac{\eta_t}{2} \right] \cdot 2^{m_t-l} L_{m_t} + \sum_{t \in \mathcal{T}_{l+1}} \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \\
&= \left\{ \Phi_T(\eta_T) + \sum_{t=l+1}^{T-1} (\Phi_t(\eta_t) - \Phi_t(\eta_{t+1})) - \Phi_{t_{l+1}-1}(\eta_{t_{l+1}}) + \sum_{t=l+1}^T \frac{\eta_t}{2} \right\} \cdot 2^{m_t-l} L_{m_t} + \sum_{t \in \mathcal{T}_{l+1}} \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \\
&\stackrel{(E5)}{\leq} \left\{ \Phi_T(\eta_T) + \ln 2 \cdot n \left(\frac{1}{\eta_T} - \frac{1}{\eta_{t_{l+1}}} \right) - \Phi_{t_{l+1}-1}(\eta_{t_{l+1}}) + \sum_{t=l+1}^T \frac{\eta_t}{2} \right\} \cdot 2^{m_t-l} L_{m_t} + \sum_{t \in \mathcal{T}_{l+1}} \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \\
&\leq \left\{ \Phi_T(\eta_T) + \sqrt{nT} + \sqrt{n t_{l+1}} + \sqrt{nT} \right\} \cdot 2^{m_t-l} L_{m_t} + \sum_{t \in \mathcal{T}_{l+1}} \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \\
&\stackrel{(E6)}{\leq} \left\{ \sum_{t \in \mathcal{T}_{l+1}} \mathbb{E} \left[\frac{c'_{m_t, t}(I_{m_t, t}) - \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t})}{2^{m_t-l} L_{m_t}} | I_{l+1, t} = i'_{l+1} \right] + 3\sqrt{nT} \right\} \cdot 2^{m_t-l} L_{m_t} + \sum_{t \in \mathcal{T}_{l+1}} \min_{i_{m_t} \in \mathcal{E}_{m_t}(\mathcal{K}_l(i_l))} c'_{m_t, t}(i_{m_t}) \\
&\leq \sum_{t \in \mathcal{T}_{l+1}} \mathbb{E} [c'_{m_t, t}(I_{m_t, t}) | I_{l+1, t} = i''_{l+1}] + 3\sqrt{nT} \cdot 2^{m_t-l} L_{m_t} \\
&= \sum_{t \in \mathcal{T}_{l+1}} \mathbb{E} [c'_{m_t, t}(I_{m_t, t}) | I_{l+1, t} = i''_{l+1}] + 3DLn\sqrt{T} \cdot \frac{1}{2}, \quad \forall i''_{l+1} \in \mathcal{E}_{l+1}(\mathcal{K}_l(i_l))
\end{aligned}$$

(33)

$$\begin{aligned}
& \sum_{t \in \mathcal{T}} \mathbb{E}[c'_{m_t, t}(I_{m_t, t})] - \sum_{t \in \mathcal{T}} c'_{m_t, t}(i_{m_t}^*) = \sum_{t \in \mathcal{T}} \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_0, t = i_0^*] - \sum_{t \in \mathcal{T}} c'_{m_t, t}(i_{m_t}^*) \\
&\stackrel{(E1)}{=} \sum_{t \in \mathcal{T}_1} \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_0, t = i_0^*] - \sum_{t \in \mathcal{T}_1} c'_{m_t, t}(i_{m_t}^*) \leq \sum_{t \in \mathcal{T}_1} \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_1, t = i_1^*] - \sum_{t \in \mathcal{T}_1} c'_{m_t, t}(i_{m_t}^*) + 3DLn\sqrt{T} \cdot \frac{1}{20} \\
&\stackrel{(E2)}{=} \sum_{t \in \mathcal{T}_2} \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_1, t = i_1^*] - \sum_{t \in \mathcal{T}_2} c'_{m_t, t}(i_{m_t}^*) + 3DLn\sqrt{T} \cdot \frac{1}{20} \\
&\leq \sum_{t \in \mathcal{T}_2} \mathbb{E}[c'_{m_t, t}(I_{m_t, t}) | I_1, t = i_2^*] - \sum_{t \in \mathcal{T}_2} c'_{m_t, t}(i_{m_t}^*) + 3DLn\sqrt{T} \cdot \left(\frac{1}{20} + \frac{1}{2^1} \right) \\
&\leq \dots \leq 3DLn\sqrt{T} \cdot \left(\frac{1}{20} + \frac{1}{2^1} + \dots + \frac{1}{2^m} \right) \leq 6DLn\sqrt{T}.
\end{aligned}$$

(34)