

Expert advice problem with noisy low rank loss

Yaxiong Liu

Department of Informatics, Kyushu University/ RIKEN AIP

YAXIONG.LIU@INF.KYUSHU-U.AC.JP

Xuanke Jiang

Department of Informatics, Kyushu University

JIANG.XUANKE.290@S.KYUSHU-U.AC.JP

Kohei Hatano

Faculty of Art and Science, Kyushu University/ RIKEN AIP

HATANO@INF.KYUSHU-U.AC.JP

Eiji Takimoto

Department of Informatics, Kyushu University

EIJI@INF.KYUSHU-U.AC.JP

Editors: Vineeth N Balasubramanian and Ivor Tsang

Abstract

We consider the expert advice problem with a low rank but noisy loss sequence, where a loss vector $\mathbf{l}_t \in [-1, 1]^N$ in each round t is of the form $\mathbf{l}_t = U\mathbf{v}_t + \boldsymbol{\epsilon}_t$ for some fixed but unknown $N \times d$ matrix U called the kernel, some d -dimensional seed vector $\mathbf{v}_t \in \mathbb{R}^d$, and some additional noisy term $\boldsymbol{\epsilon}_t \in \mathbb{R}^N$ whose norm is bounded by ϵ . This is a generalization of the works of Hazan et al. and Barman et al., where the former only treats noiseless loss and the latter assumes that the kernel is known in advance. In this paper, we propose an algorithm, where we re-construct the kernel under the assumptions, that the low rank loss is noised and there is no prior information about kernel. In this algorithm, we approximate the kernel by choosing a set of loss vectors with a high degree of independence from each other, and we give a regret bound of $O(d\sqrt{T} + d^{4/3}(N\epsilon)^{1/3}\sqrt{T})$. Moreover, even if in experiment, the proposed algorithm performs better than Hazan's algorithm and Hedge algorithm.

Keywords: Online expert advice, Low rank loss, Noise, Minimum volume ellipsoid enclosing

1. Introduction

The expert advice problem with low rank loss (Hazan et al., 2016) is an extension of the standard expert problem by considering a latent structure in losses. In this problem, we model the expert advice with d -rank loss as follows: the environment chooses a full rank $N \times d$ matrix U , called kernel. On each round $t \in [T]$, the algorithm picks a prediction \mathbf{w}_t in an N -dimensional simplex over the set of N experts. Then the environment gives a d -rank loss vector $\mathbf{l}_t \in [-1, 1]^N$, where $\mathbf{l}_t = U\mathbf{v}_t$ to the algorithm. At last the algorithm suffers the loss $\mathbf{w}_t \cdot \mathbf{l}_t$. We measure the performance of this algorithm based difference between the cumulative loss and the best expert strategy in hindsight, which is defined as regret in the following equation:

$$\text{Regret}_T = \sum_{t=1}^T \mathbf{w}_t \cdot \mathbf{l}_t - \min_{i \in [N]} \sum_{t=1}^T \mathbf{l}_t(i). \quad (1)$$

This low rank loss setting is popular in recommendation system, especially when experts give losses based on a latent structure (Koren et al., 2009). For instance, these experts share some common information, or their prediction methods are similar and depended on only few factors. As a consequence, the loss vectors given by experts are in fact in a relatively lower dimensional space, which implies that $d \leq N$.

Compared with the original expert advice problem, whose regret bound is $\Theta(\sqrt{T \ln N})$ with Hedge algorithm (Freund and Schapire, 1997; Cesa-Bianchi and Lugosi, 2006), the expert advice problem with d -rank loss incurs an upper bound as $O(d\sqrt{T})$ (Hazan et al., 2016), and $O(\sqrt{dT \ln T})$ (Koren and Livni, 2017), respectively. Moreover, the optimal bound is shown to be $\Theta(\sqrt{dT})$ if the kernel is known to the algorithm (Hazan et al., 2016). The work of Hazan et al. (2016) combines two parts, recovering the kernel U and predicting the \mathbf{w}_t . Hence, this algorithm not only provides a more precise upper bound, if $d \leq o(\sqrt{\ln N})$, but also re-constructs the kernel while running the algorithm. Equivalently, kernel U is explored and exploited in the learning process. Actually, in the algorithm of Hazan et al. (2016), Online mirror descent (OMD) is designed to predict \mathbf{w}_t utilising the recovered U . Without U , OMD can give only a loose bound as $O(\sqrt{NT})$, even if the loss is d -ranked. In contrast, if the kernel is acquired in advance, OMD achieves the optimal bound $O(\sqrt{dT})$.

After this pioneering work (Hazan et al., 2016), Barman et al. (2018) considered a case that, the d -rank loss \mathbf{l}_t is corrupted by a L_2 -norm bounded $\boldsymbol{\epsilon}_t$ noise, i.e., the environment gives the loss vector $\tilde{\mathbf{l}}_t = \mathbf{l}_t + \boldsymbol{\epsilon}_t = U\mathbf{v}_t + \boldsymbol{\epsilon}_t$, which is near to the d -dimensional space. In their work, Barman et al. (2018) gave a regret bound $O(\sqrt{(d + \epsilon)T})$ and a lower bound as $\Omega(\sqrt{T(d + \frac{\epsilon}{2d})})$, where $\|\boldsymbol{\epsilon}_t\|_2^2 \leq \epsilon$. However, there is a strong assumption in their work that the algorithm needs know both U and ϵ .

In this paper, we release the strong assumption in Barman et al. (2018), and assume that the kernel is unknown to the algorithm, while the low rank loss is corrupted by $\boldsymbol{\epsilon}_t$. This problem setting is more realistic in two points: One is the loss vectors are not always strictly well-structured. Although the experts share some common prediction methods, there is still some slight disturbance, which damages the low rank structure loss. The other is the specific structure is unknown to the algorithm before it receives the loss vectors. It is natural that the recommendation system needs to confront new users who are recently registered without any prior information. Thus, the system is required to explore the characteristics of users while recommending their goods.

From the above works, we can see how the low rank structure plays the important role in the algorithms. Either, the kernel is known to the algorithm, then the algorithm can overcome the noised loss vectors easily; or, the loss is precisely low ranked, then the algorithm can extract the kernel from the received losses, even if the kernel is unknown. However, if we apply the algorithms of Hazan et al. (2016), and Koren and Livni (2017) in our problem setting directly, it is equivalent to run OMD on expert advice directly, and the regret bound is $O(\sqrt{NT})$, since the low rank loss no longer exists. Hence, one of the essential problems is how to deal with the ‘‘noisy low rank loss’’ and recover the underlying kernel U . Although recovering a low rank matrix is not a novel concept for machine learning (Goldberg et al., 2010; Foygel and Srebro, 2011; Balcan and Zhang, 2016; Zhang et al., 2018), to the best of the authors knowledge, there is no such research for online expert advice problem with noise attached loss vectors. The most tricky obstacle is

that, it is impossible to re-construct the kernel U exactly with the corrupted loss vectors. Thus, in our proposed algorithm, we attempt to approximate this kernel from the received loss vectors.

In our algorithm, we assume that there is no prior information about the kernel. Thus, we select at most d highly-independent loss vectors (see details in latter section) to construct a pseudo-kernel \tilde{U} in our learning process. Unlike the selection criteria in work (Balcan and Zhang, 2016), which recovers the kernel with a randomised algorithm stochastically, in our paper, each selected vector is supposed to be not only “long” enough, but also far enough to the current sub-space, spanned by all previously selected vectors. Therefore, our algorithm can approximate the kernel of the low rank loss deterministically and obtains a worst case guarantee. Instead of the unknown U in (Hazan et al., 2016), we next utilise the pseudo-kernel, which is spanned by the selected vectors, in OMD to update the predictions. Our algorithm obtains an upper bound with respect to $\tilde{\mathbf{l}}_t$ as $O(\sqrt{T}(d + d^{4/3}(N\epsilon)^{1/3}))$. In practice, we can consider our algorithm and the Hedge algorithm as two meta-experts and run another Hedge algorithm on top of them, which performs nearly as well as our algorithm and the Hedge algorithm with additional negligible ($O(\sqrt{T})$) regret.

The following table shows the relationship between our work and previous work.

Table 1: Comparison with noisy and non-noisy low rank loss

	known kernel	unknown kernel
without noise	$\Theta(\sqrt{dT})$ (Hazan et al., 2016)	$O(d\sqrt{T})$ (Hazan et al., 2016) $O(\sqrt{dT} \ln T)$ (Koren and Livni, 2017)
ϵ -noise case	$O(\sqrt{(d + \epsilon)T})$ $\Omega\left(\sqrt{T(d + \frac{\epsilon}{2d})}\right)$ (Barman et al., 2018)	$O((1 + \epsilon)\sqrt{T}(d + d^{4/3}(N\epsilon)^{1/3}))$ (Our work)

This paper is composed as follows. In the second section, we define our expert advice with noisy d -rank loss formally and some necessary notations. In the third section, the proposed algorithm is given and we show the upper bound. In section 4, we compare the proposed algorithm and previous work under synthetic environments¹.

2. Preliminaries

We denote the set $\{1, \dots, T\}$ by $[T]$. We denote the N -dimensional vector with all 1 elements by $\mathbb{1}_N$ and $N \times N$ identity matrix by I_N . We define a norm $\|\mathbf{x}\|_H$ with respect to a positive definite matrix H , as $\|\mathbf{x}\|_H = \sqrt{\mathbf{x}^T H \mathbf{x}}$, for a vector \mathbf{x} , and the dual norm of $\|\cdot\|_H$ is defined as $\|\mathbf{x}\|_H^* = \sqrt{\mathbf{x}^T H^{-1} \mathbf{x}}$. We define the angle between two vectors \mathbf{u} and \mathbf{v} as $\theta(\mathbf{u}, \mathbf{v}) = \arccos \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}$. If \mathbf{V} is a subspace then $\theta(\mathbf{u}, \mathbf{V}) = \min_{\mathbf{v} \in \mathbf{V} \setminus \{0\}} \theta(\mathbf{u}, \mathbf{v})$, and $\theta(\mathbf{U}, \mathbf{V}) = \max_{\mathbf{u} \in \mathbf{U}} \theta(\mathbf{u}, \mathbf{V})$. At last, we denote the orthogonal projection from a vector \mathbf{l} to a subspace \mathbf{U} as $\mathcal{P}_{\mathbf{U}}\mathbf{l}$. Moreover, the Euclidean distance from a vector \mathbf{l} to a sub-space \mathbf{U} is defined as $\|\mathbf{l} - \mathcal{P}_{\mathbf{U}}\mathbf{l}\|_2 = \|\mathbf{l}\|_2 \sin \theta(\mathbf{l}, \mathbf{U}) = \|\mathbf{l}\|_2 \sin \theta(\mathbf{l}, \mathcal{P}_{\mathbf{U}}\mathbf{l})$. A linear space spanned by all columns

1. The code is available at <https://github.com/2015211217/LowRankStructureC-.git>

from a matrix $V \in \mathbb{R}^{N \times k}$ can be represented as $\text{span}(V)$. In the following part, we simplify $\text{span}(V)$ as \mathbf{V} , when it leads no ambiguity.

2.1. Problem Setting

Firstly we define the d -rank loss: For a sequence of $\mathbf{l}_t \in [-1, 1]^N$, for $t \in [T]$ we define $L_T \in \mathbb{R}^{N \times T}$ as follows:

$$L_T = [\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_T], \quad (2)$$

where each loss vector \mathbf{l}_t is a column in L_T . We define that the sequence $\{\mathbf{l}_1, \dots, \mathbf{l}_T\}$ is d -rank loss if and only if $\text{rank}(L_T) = d$. Equivalently, for the sequence $\{\mathbf{l}_1, \dots, \mathbf{l}_T\}$, there exists a kernel $U \in \mathbb{R}^{N \times d}$ such that $\mathbf{l}_t = U\mathbf{v}_t$, where $\text{rank}(U) = d$. If $d \ll N$, we can call \mathbf{l}_t as low rank loss vector.

From d -rank loss \mathbf{l}_t , we define the noisy d -rank loss $\tilde{\mathbf{l}}_t : \tilde{\mathbf{l}}_t = \mathbf{l}_t + \boldsymbol{\epsilon}_t$, where $\|\boldsymbol{\epsilon}_t\|_2 \leq \epsilon, \forall t \in [T]$. We call $\boldsymbol{\epsilon}_t$ as the noise vector.

In this paper we consider the following learning problem of online expert advice with noisy d -rank loss. On round $t = 1, \dots, T$, an algorithm gives a prediction $\mathbf{w}_t \in \Delta(N)$; Then an environment gives a loss vector $\tilde{\mathbf{l}}_t \in [-1 - \epsilon, 1 + \epsilon]^N$, note that $\tilde{\mathbf{l}}_t = \mathbf{l}_t + \boldsymbol{\epsilon}_t$. It implies that there exists an underlying $\mathbf{l}_t \in [-1, 1]^N$ as d -rank loss and $\boldsymbol{\epsilon}_t$ as ϵ -noise. However, \mathbf{l}_t, U and $\boldsymbol{\epsilon}_t$ are unknown to the algorithm. At last the algorithm suffers the loss as $\mathbf{w}_t \cdot \tilde{\mathbf{l}}_t$. Next we define the regret as follows:

$$\text{Regret}_T = \sum_{t=1}^T \mathbf{w}_t \cdot \tilde{\mathbf{l}}_t - \min_{i \in [N]} \sum_{t=1}^T \tilde{\mathbf{l}}_t(i). \quad (3)$$

2.2. Online mirror descent

Online mirror descent (OMD) is a basic algorithm utilised in both [Hazan et al. \(2016\)](#) and [Barman et al. \(2018\)](#), as well as in this paper. We give OMD with the time-varying matrix norms as follows:

Algorithm 1 Online Mirror descent

Initialization: $H_0, \dots, H_{T-1} \succ 0, \{\eta_t\}_{t=1}^T$, and $\mathbf{w}_1 \in \Delta(N)$.

for $t = 1, \dots, T$ **do**

 Receive \mathbf{l}_t ,

 Suffer cost $\mathbf{w}_t \cdot \mathbf{l}_t$,

 Update $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \Delta(N)} \mathbf{l}_t \cdot \mathbf{w} + \eta_t^{-1} \|\mathbf{w} - \mathbf{w}_t\|_{H_{t-1}}^2$.

end for

Theorem 1 ([Orabona et al. \(2015\)](#)) *The T -round regret of OMD is bounded as follow:*

$$\sum_{t=1}^T \mathbf{w}_t \cdot \mathbf{l}_t - \sum_{t=1}^T \mathbf{l}_t \cdot \mathbf{w}^* \leq \frac{1}{\eta_T} \|\mathbf{w}_1 - \mathbf{w}^*\|_{H_T}^2 + \frac{1}{2} \sum_{t=1}^T (\eta_t \|\mathbf{l}_t\|_{H_t}^*)^2 + \sum_{t=1}^T \eta_t^{-1} (\|\mathbf{w}_t\|_{H_{t-1}}^2 - \|\mathbf{w}_t\|_{H_t}^2).$$

2.3. Ellipsoid approximation of convex bodies

Given a positive semi-definite matrix $M \succeq 0$, we define the ellipsoid with respect to M as follows:

$$\mathcal{E}(M) = \{\mathbf{x} : \mathbf{x}^T M^\dagger \mathbf{x} \leq 1\}, \quad (4)$$

where M^\dagger is the Moore-Penrose pseudo-inverse matrix of M . In the following theorem we give a result for minimum volume enclosing ellipsoid (MVEE) to a central symmetric convex body.

Theorem 2 (John's Theorem (Ball et al., 1997)) *Let K be a convex body in \mathbb{R}^d that is symmetric around zero. Let \mathcal{E} be an ellipsoid with minimum volume enclosing K . Then:*

$$\frac{1}{\sqrt{d}}\mathcal{E} \subseteq K \subseteq \mathcal{E}. \quad (5)$$

Moreover for a polytope defined as $P_A = \{\mathbf{x} : \|A\mathbf{x}\|_\infty \leq 1, \mathbf{x} \in \mathbb{R}^d\}$, corresponding to a given matrix $A \in \mathbb{R}^{N \times d}$, we have the following theorem.

Theorem 3 (Grötschel et al. (2012)) *There exists a poly-time procedure $\text{MVEE}(A)$ that receives as input a matrix $A \in \mathbb{R}^{N \times d}$ and returns a matrix M such that*

$$\frac{1}{\sqrt{2d}}\mathcal{E}(M) \subseteq P_A \subseteq \mathcal{E}(M).$$

3. Algorithm for no prior information about kernel

In this section, we consider the case that there is no prior information about the kernel in the online expert advice with noisy d -rank loss problem. Under this assumption, we construct a pseudo kernel, \tilde{U} , in our learning process, when the d -rank structure is corrupted by ϵ_t . Note that in this algorithm, we construct this pseudo kernel cumulatively. We denote it as $\tilde{U}^k \in \mathbb{R}^{N \times k}$ for convenience.

Concisely, our algorithm will select some loss vectors to construct the pseudo kernel and run OMD with respect to from \tilde{U}^k processed matrix H^k . The criteria to select $\tilde{\mathbf{l}}_t$ is twofold: firstly, the L_2 -norm of $\tilde{\mathbf{l}}_t$ is supposed to be large enough; then, the Euclidean distance from $\tilde{\mathbf{l}}_t$ to the current pseudo kernel spanned space $\text{span}(\tilde{U}^k)$ should be far enough. Thus, for some $k \in [d]$ and corresponding parameters s_k, γ_k , the basic idea of our algorithm is as follows:

1. $\|\tilde{\mathbf{l}}_t\|_2 \leq 2s_k \rightarrow$ Do OMD with respect to matrix H^k .
2. $\|\tilde{\mathbf{l}}_t\|_2 \geq 2s_k$ and $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\text{span}(\tilde{U}^k)} \tilde{\mathbf{l}}_t\|_2 \leq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon) \rightarrow$ Do OMD with respect to H^k .
3. $\|\tilde{\mathbf{l}}_t\|_2 \geq 2s_k$ and $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\text{span}(\tilde{U}^k)} \tilde{\mathbf{l}}_t\|_2 \geq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon) \rightarrow$ Adding $\tilde{\mathbf{l}}_t$ as a column to \tilde{U}^k , and reset k as $k + 1$, then do OMD with respect to H^k .

Remark 4 *If we have that $\|\tilde{\mathbf{l}}_t\|_2 \geq 2s_k > 2\epsilon$, then we obtain the following result: Since $\|\tilde{\mathbf{l}}_t - \mathbf{l}_t\|_2 \leq \epsilon$, $\theta(\tilde{\mathbf{l}}_t, \mathbf{l}_t) \leq 2 \sin \theta(\tilde{\mathbf{l}}_t, \mathbf{l}_t) \leq \frac{\epsilon}{s_k}$*

According to this criteria, we can observe that we need to guarantee that $K = \max k \leq d$, to ensure that our algorithm effective, since the algorithm from (Hazan et al., 2016) updates N -times for low rank noisy loss vector. It implies that the previous algorithm in (Hazan et al., 2016) is OMD and obtains $O(\sqrt{(N + \epsilon)T})$ actually.

We confirm s_k, γ_k according to the following proposition.

Proposition 5 Define a sequence s_i and non-decreasing sequence γ_i such that $\gamma_i \in (0, \frac{\pi}{2})$ and

$$\frac{\gamma_i}{\gamma_{i-1} - (i-1)\frac{\beta\epsilon}{s_{i-1}\gamma_{i-1}}} \cdot \frac{\pi}{2\beta} + \frac{s_i\gamma_i}{s_i\gamma_{i-1}}(i-1) \leq i \quad \forall i \in \{2, \dots, k\}, \quad (6)$$

and

$$\gamma_i - \frac{\beta\epsilon}{s_i\gamma_i}i > 0 \quad \forall i \in \{1, \dots, k\}. \quad (7)$$

Given any space $\mathbf{W} \subseteq \mathbb{R}^N$: Let $\mathbf{U}^k = \text{span}\{\mathbf{W}, \mathbf{l}_1, \dots, \mathbf{l}_k\}$ and $\tilde{\mathbf{U}}^k = \text{span}\{\mathbf{W}, \tilde{\mathbf{l}}_1, \dots, \tilde{\mathbf{l}}_k\}$ be two subspace such that $\theta(\mathbf{l}_i, \tilde{\mathbf{l}}_i) \leq \epsilon/s_i$ for all $i \in \{2, \dots, k\}$.

If $\theta(\tilde{\mathbf{l}}_i, \mathbf{U}^{i-1}) > \gamma_i$, and $\beta \geq \frac{\pi}{2}$ then we have that

$$\theta(\mathbf{U}^k, \tilde{\mathbf{U}}^k) < \beta k \frac{\epsilon}{s_k \gamma_k}. \quad (8)$$

Proof For $i = 1$, due to Lemma 10, by setting that $\theta(l_1, \emptyset) = \frac{\pi}{2}$, it is trivial that

$$\theta(\mathbf{U}^1, \tilde{\mathbf{U}}^1) \leq \frac{\pi\epsilon}{2s_1\gamma_1} \leq \beta \frac{\epsilon}{s_1\gamma_1}. \quad (9)$$

Here $\theta(\tilde{\mathbf{l}}_1, \mathbf{l}_1) \leq \frac{\epsilon}{s_1}$, and $\theta(\tilde{\mathbf{l}}_1, \tilde{\mathbf{U}}^0) = \theta(\tilde{\mathbf{l}}_1, \mathbf{W}) \geq \gamma_1$. If $\mathbf{W} = \emptyset$ then we have that

$$\theta(\mathbf{U}^1, \tilde{\mathbf{U}}^1) = \theta(\mathbf{l}_1, \tilde{\mathbf{l}}_1) \leq \frac{\epsilon}{s_1} \leq \frac{\pi}{2} \times \frac{\epsilon}{\gamma_1 s_1},$$

if $\gamma_1 \leq \frac{\pi}{2}$.

Then if it holds for $i - 1$: Let us involve $\mathbf{U}_0^i = \text{span}\{\mathbf{U}^{i-1}, \tilde{\mathbf{l}}_i\}$, note that

$$\theta(\mathbf{U}_0^i, \tilde{\mathbf{U}}^i) = \theta(\text{span}\{\mathbf{U}^{i-1}, \tilde{\mathbf{l}}_i\}, \text{span}\{\tilde{\mathbf{U}}^{i-1}, \tilde{\mathbf{l}}_i\}) \quad (10)$$

Since the induction hypothesis so we have that

$$\theta(\text{span}\{\mathbf{W}, \mathbf{l}_1, \dots, \mathbf{l}_{i-1}\}, \text{span}\{\mathbf{W}, \tilde{\mathbf{l}}_1, \dots, \tilde{\mathbf{l}}_{i-1}\}) \leq (i-1) \frac{\beta\epsilon}{s_{i-1}\gamma_{i-1}}. \quad (11)$$

then we have that

$$\begin{aligned} \theta(\mathbf{U}^k, \tilde{\mathbf{U}}^k) &\leq \theta(\mathbf{U}^k, \mathbf{U}_0^k) + \theta(\mathbf{U}_0^k, \tilde{\mathbf{U}}^k) \\ &\leq \frac{\pi}{2} \frac{\theta(\tilde{\mathbf{l}}_k, \mathbf{l}_k)}{\theta(\tilde{\mathbf{l}}_k, \mathbf{U}^{k-1})} + \theta(\mathbf{U}_0^i, \tilde{\mathbf{U}}^i) \end{aligned}$$

The first inequality is from Remark 13. The second inequality is due to Lemma 10.

Now due to the Lemma 12 we have that

$$\theta(\mathbf{U}^k, \tilde{\mathbf{U}}^k) \leq \frac{\pi}{2} \frac{\theta(\mathbf{l}_i, \tilde{\mathbf{l}}_i)}{\theta(\tilde{\mathbf{l}}_i, \tilde{\mathbf{U}}^{i-1}) - \theta(\mathbf{U}^{i-1}, \tilde{\mathbf{U}}^{i-1})} + \theta(\mathbf{U}_0^i, \tilde{\mathbf{U}}^i) \quad (12)$$

Since the fact that $\theta(\tilde{\mathbf{l}}_k, \tilde{\mathbf{U}}^{k-1}) > \gamma_k$, and the induction hypothesis we obtain that

$$\begin{aligned} \theta(\mathbf{U}^i, \tilde{\mathbf{U}}^i) &< \frac{\pi}{2} \frac{\frac{\epsilon}{s_i}}{\gamma_{i-1} - \beta(i-1) \frac{\epsilon}{s_{i-1}\gamma_{i-1}}} + (i-1) \frac{\beta\epsilon}{s_{i-1}\gamma_{i-1}} \\ &= \frac{\epsilon\beta}{s_i\gamma_i} \left(\frac{\gamma_i}{\gamma_{i-1} - (i-1) \frac{\beta\epsilon}{s_{i-1}\gamma_{i-1}}} \times \frac{\pi}{2\beta} + (i-1) \frac{s_i\gamma_i}{s_{i-1}\gamma_{i-1}} \right) \end{aligned} \quad (13)$$

Since Equation 6 we have that

$$\frac{\epsilon\beta}{s_i\gamma_i} \left(\frac{\gamma_i}{\gamma_{i-1} - (i-1) \frac{\beta\epsilon}{s_{i-1}\gamma_{i-1}}} \times \frac{\pi}{2\beta} + (i-1) \frac{s_i\gamma_i}{s_{i-1}\gamma_{i-1}} \right) \leq \frac{\beta\epsilon}{s_i\gamma_i} i \quad (14)$$

Thus we have our conclusion. ■

The description of the algorithm is shown in Algorithm 2.

Algorithm 2 Mirror descent with l_2 noise

Initialization: $\mathbf{w}_1 = \frac{1}{N}\mathbb{I}_N$, $\tau = 0$, $k = 0$, $\tilde{\mathbf{U}} = \{\}$, a sequence of $\{s_k\}$ and a non decreasing sequence $\{\gamma_k\}$ for all $\gamma_k \in (0, \frac{\pi}{2})$ such that γ_k and s_k satisfies the conditions in Equation (6) and (7). Setting $m_k = \max\{2s_k, 6\epsilon + \gamma_k\sqrt{N}\}$.

for $t = 1, \dots, T$ **do**

Observe $\tilde{\mathbf{l}}_t$, suffer loss $\mathbf{w}_t \cdot \tilde{\mathbf{l}}_t$.

if $\|\tilde{\mathbf{l}}_t\|_2 \geq 2s_k$ **then**

if $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{\mathbf{U}}^k}\tilde{\mathbf{l}}_t\|_2 \geq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)$ **then**

Add $\tilde{\mathbf{l}}_t$ as a new column of $\tilde{\mathbf{U}}^k$, reset $\tau = 0$ and set $k \leftarrow k + 1$.

Compute $M = \text{MVEE}(\tilde{\mathbf{U}}^k)$ and $H^k = I_N + \tilde{\mathbf{U}}^k M (\tilde{\mathbf{U}}^k)^T$.

end if

end if

let $\tau \leftarrow \tau + 1$ and $\eta_t = \sqrt{8k/(1 + m_k)^2\tau}$ and set:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \Delta(N)} \tilde{\mathbf{l}}_t \cdot \mathbf{w} + \eta_t^{-1} \|\mathbf{w} - \mathbf{w}_t\|_{H^k}^2.$$

end for

According to our proposed algorithm, we see that our updating rule divides the T rounds into K epoches, where the final size of pseudo kernel, $K \leq d$ (will be show in the following Theorem). In each epoch, OMD procees with respect to H^k , a fixed matrix. Before we give the regret bound, we need some useful lemmata.

Lemma 6 Let $\|\tilde{\mathbf{l}}_t - \mathbf{l}_t\|_2 \leq \epsilon$, if $\theta(\mathbf{l}_t, \text{span}(\tilde{U}^k)) \leq \gamma_k$ then we have

$$\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2 \leq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon). \quad (15)$$

Proof Firstly we obtain that

$$\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2 \leq \|\tilde{\mathbf{l}}_t - \mathbf{l}_t\|_2 + \|\mathbf{l}_t - \mathcal{P}_{\tilde{U}^k} \mathbf{l}_t\|_2 + \|\mathcal{P}_{\tilde{U}^k} \mathbf{l}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2.$$

Since the fact that $\|\tilde{\mathbf{l}}_t - \mathbf{l}_t\|_2 \leq \epsilon$, and $\|\mathbf{l}_t - \mathcal{P}_{\tilde{U}^k} \mathbf{l}_t\|_2 \leq \gamma_k \|\mathbf{l}_t\|_2 \leq \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)$, we have our conclusion. \blacksquare

Lemma 6 states the criteria of distance between $\tilde{\mathbf{l}}_t$ and current space \tilde{U}^k .

Lemma 7 Given $\tilde{U}^k \in \mathbb{R}^{N \times k}$, and $\text{span}(\tilde{U}^k) \subset \mathbb{R}^N$, for any $\tilde{\mathbf{l}}_t$ such that $\|\tilde{\mathbf{l}}_t\|_2 \leq 2s_k$ or $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2 \leq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)$, there exists a unique $\mathbf{v}_t \in \mathbb{R}^k$, and \mathbf{e}_t such that $\tilde{\mathbf{l}}_t = \tilde{U}^k \mathbf{v}_t + \mathbf{e}_t$, where $\theta(\mathbf{e}_t, \tilde{U}^k) = \frac{\pi}{2}$. Therefore we have that

$$\|\mathbf{e}_t\|_2 \leq \max\{2s_k, 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)\}. \quad (16)$$

In particular, we set that $\forall t \in [T]$

$$m_k = \max\{2s_k, 6\epsilon + \gamma_k \sqrt{N}\} \geq \max\{2s_k, 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)\}. \quad (17)$$

Proof Since the fact that $\|\mathbf{e}_t\|_2 \leq \|\tilde{\mathbf{l}}_t\|_2$ and the updating rule, we have that

$$\|\mathbf{e}_t\|_2 \leq \max\{2s_k, 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)\}. \quad (18)$$

Due to our setting that $\gamma_k \leq \frac{\pi}{2} \leq 2$, and $\|\tilde{\mathbf{l}}_t\|_2 \leq \sqrt{N} + \epsilon$, thus we obtain Equation (17). \blacksquare

Theorem 8 Running Algorithm 2 on T -rounds, denoting that K is the final size of pseudo kernel, and $K \leq d$, we have

$$\text{Regret}_T \leq O \left((1 + \epsilon) \sqrt{T} \left(K + \sqrt{\sum_{k=1}^K k m_k^2} \right) \right). \quad (19)$$

Proof Firstly let us prove that $K \leq d$. Compared with \tilde{U}^k with k loss vectors of $\tilde{\mathbf{l}}_t$, we define a matrix U^k with respect to \mathbf{l}_t corresponding to $\tilde{\mathbf{l}}_t$.

Due to our updating rule in algorithm, we have that

$$\theta(\tilde{\mathbf{l}}_t, \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t) \geq \gamma_k,$$

since if $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2 \geq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)$, then we have that

$$\theta(\tilde{\mathbf{l}}_t, \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t) \geq \sin \theta(\tilde{\mathbf{l}}_t, \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t) \geq \frac{2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_t\|_2 + \epsilon)}{\|\tilde{\mathbf{l}}_t\|_2} \geq \gamma_k.$$

Therefore by Proposition 5 we have that $\theta(\text{span}(U^k), \text{span}(\tilde{U}^k)) \leq \frac{\beta k \epsilon}{\gamma_k s_k}$. So we have that

$$\theta(\mathbf{l}_t, \text{span}(U^k)) \geq \theta(\mathbf{l}_t, \text{span}(\tilde{U}^k)) - \theta(\text{span}(U^k), \text{span}(\tilde{U}^k)) \geq \gamma_k - \frac{\beta \epsilon}{s_k \gamma_k} k > 0.$$

The first inequality is from Lemma 12. Since Lemma 6 we have that $\theta(\mathbf{l}_t, \text{span}(\tilde{U}^k)) \geq \gamma_k$. It implies that if $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2 \geq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_2\|_2 + \epsilon)$, then we have $\mathbf{l}_t \notin \text{span}(U^k)$. So we have that $\text{rank}(\tilde{U}^K) \leq \text{rank}(L_T) = d$.

Now let us prove the second part. Due to above result we can divide the total learning rounds T into K episodes. $T = \sum_{k=1}^K |T_k|$ and $T_k \cup T_{k'} = \emptyset$, if $k \neq k'$. For a given k , if $t \in T_k$, we obtain that $\|\tilde{\mathbf{l}}_t - \mathcal{P}_{\tilde{U}^k} \tilde{\mathbf{l}}_t\|_2 \leq 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_2\|_2 + \epsilon)$ or $\|\tilde{\mathbf{l}}_t\|_2 \leq 2s_k$. Thus we need give upper bound of $(\|\tilde{\mathbf{l}}_t\|_{H^k}^*)^2$, and $\|\mathbf{w}_1 - \mathbf{w}^*\|_{H^k}$.

Given a k -dimensional polytope as

$$P = \{\mathbf{v} \in \mathbb{R}^k : \|\tilde{U}^k \mathbf{v}\|_\infty \leq 1 + \epsilon\},$$

we are able to have

$$\mathcal{E}\left(\frac{1}{2k}M\right) \subseteq P \subseteq \mathcal{E}(M), \quad (20)$$

due to the MVEE procedure, and Theorem 3 in previous section. For each $\tilde{\mathbf{l}}_t$ we give a unique orthogonal decomposition of $\tilde{\mathbf{l}}_t$ upon $\text{span}(\tilde{U}^k)$. Thus we obtain $\tilde{\mathbf{l}}_t = \tilde{U}^k \mathbf{v}_t + \mathbf{e}_t$. Since $\tilde{\mathbf{l}}_t \in [0, 1]^N$, and we have that $\mathbf{v}_t \in P$, and $\theta(\mathbf{e}_t, \tilde{U}^k) = \frac{\pi}{2}$. Therefore we have following inequality:

$$(\|\tilde{\mathbf{l}}_t\|_{H^k}^*)^2 = (\|\tilde{U}^k \mathbf{v}_t + \mathbf{e}_t\|_{H^k}^*)^2 \leq (\|\tilde{U}^k \mathbf{v}_t\|_{H^k}^* + \|\mathbf{e}_t\|_{H^k}^*)^2 \leq 2(\|\tilde{U}^k \mathbf{v}_t\|_{H^k}^*)^2 + 2(\|\mathbf{e}_t\|_{H^k}^*)^2.$$

Firstly we give upper bound of $(\|\tilde{U}^k \mathbf{v}_t\|_{H^k}^*)^2$.

$$\begin{aligned} (\|\tilde{U}^k \mathbf{v}_t\|_{H^k}^*)^2 &= (\tilde{U}^k \mathbf{v}_t)^T (I_N + \tilde{U}^k M (\tilde{U}^k)^T)^{-1} (\tilde{U}^k \mathbf{v}_t) \\ &\leq (\tilde{U}^k \mathbf{v}_t)^T (\tilde{U}^k M (\tilde{U}^k)^T)^+ (\tilde{U}^k \mathbf{v}_t) = \mathbf{v}_t^T M^{-1} \mathbf{v}_t \leq 1. \end{aligned}$$

The last equality is due to Lemma 11 in Appendix.

Next we try to bound $(\|\mathbf{e}_t\|_{H^k}^*)^2$. We obtain that

$$(\|\mathbf{e}_t\|_{H^k}^*)^2 = \mathbf{e}_t^T H^{-1} \mathbf{e}_t = \mathbf{e}_t^T (I_N + \tilde{U}^k M (\tilde{U}^k)^T)^{-1} \mathbf{e}_t \leq \mathbf{e}_t^T I_n \mathbf{e}_t \leq \|\mathbf{e}_t\|_2^2. \quad (21)$$

Based on previous discussion we have that $\|\mathbf{e}_t\|_2 \leq \max\{2s_k, 2\epsilon + \gamma_k(\|\tilde{\mathbf{l}}_2\|_2 + \epsilon)\}$. Therefore we have that

$$\|\mathbf{e}_t\|_2 \leq \max\{2s_k, 2\epsilon + \gamma_k(\sqrt{N} + \epsilon)\}. \quad (22)$$

Due to above Lemma 7, we have that $\|\mathbf{e}\|_2 \leq m_k$. In conclusion we have:

$$(\|\tilde{\mathbf{l}}_t\|_{H^k}^*)^2 \leq 2(1 + m_k)^2. \quad (23)$$

Now we show the bound of $\|\mathbf{w}_1 - \mathbf{w}^*\|_{H^k}$. Since $\|\mathbf{w}_1 - \mathbf{w}^*\|_{H^k} \leq 2 \max_{\mathbf{w} \in \Delta(N)} \|\mathbf{w}\|_{H^k}$, it suffices to bound that $\max_{\mathbf{w} \in \Delta(N)} \|\mathbf{w}\|_{H^k}$. Since $\|\mathbf{w}\|_{H^k}^2 \leq 1 + 2k \|\mathbf{w}\|_{H(k)'}^2$ with $H(k)' = \frac{1}{2k} \tilde{U}^k M (\tilde{U}^k)^T$.

Given a convex set P in \mathbb{R}^k , so the dual set P^* is defined as

$$P^* = \{\mathbf{x} : \sup_{\mathbf{p} \in P} |\mathbf{x} \cdot \mathbf{p}| \leq 1\}.$$

The dual of an ellipsoid $\mathcal{E}(M)$ is given by $(\mathcal{E}(M))^* = \mathcal{E}(M^{-1})$ and since Equation (20) it is standard to show that

$$(\mathcal{E}(M))^* \subseteq P^* \subseteq (\mathcal{E}(\frac{1}{2k}M))^*. \quad (24)$$

It implies that $P^* \subseteq \mathcal{E}(2kM^{-1})$. Note that due to the definition of P^* we have that for all $i \in [N]$, $(1 + \epsilon)^{-1} \tilde{\mathbf{u}}_i \cdot \mathbf{v} \leq 1$, where $\tilde{\mathbf{u}}_i$ is the i -th row of \tilde{U}^k . So each row of \tilde{U}^k are in P^* , thus we have for each $\tilde{\mathbf{u}}_i$:

$$\begin{aligned} (1 + \epsilon)^{-1} \tilde{\mathbf{u}}_i^T (2kM^{-1})^{-1} (1 + \epsilon) \tilde{\mathbf{u}}_i \leq 1 &\Leftrightarrow (1 + \epsilon)^{-1} \tilde{\mathbf{u}}_i^T M (1 + \epsilon)^{-1} \tilde{\mathbf{u}}_i \leq 2k \\ &\Leftrightarrow \|(1 + \epsilon)^{-1} \tilde{\mathbf{u}}_i\|_M^2 \leq 2k \end{aligned} \quad (25)$$

Since $\mathbf{w} \in \Delta(N)$, we have that

$$\|\mathbf{w}\|_{H(k)'}^2 = \frac{1}{2k} \|\tilde{U}\mathbf{w}\|_M^2 \leq \frac{1}{2k} \max_i \|\mathbf{u}_i\|_M^2 \leq (1 + \epsilon)^2.$$

Therefore we have that $\|\mathbf{w}\|_{H^k}^2 \leq 1 + (1 + \epsilon)^2 2k$. Moreover we obtain that $\|\mathbf{w}\|_{H^k} \leq 2(1 + \epsilon)\sqrt{k}$. Given $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{t=1}^T \mathbf{w} \cdot \tilde{\mathbf{l}}_t$, then we have that

$$\begin{aligned} \text{Regret}_{T_k} &= \sum_{t \in T_k} \mathbf{w}_t \cdot \tilde{\mathbf{l}}_t - \sum_{t \in T_k} \tilde{\mathbf{l}}_t \cdot \mathbf{w}^* \\ &\leq \frac{1}{\eta_{T_k}} \|\mathbf{w}_{T_{k-1}} - \mathbf{w}^*\|_{H^k}^2 + \frac{1}{2} \sum_{t \in T_k} \eta_t (\|\tilde{\mathbf{l}}_t\|_{H^k}^*)^2 \\ &\leq \frac{4}{\eta_{T_k}} \max_{\mathbf{w} \in \Delta(N)} \|\mathbf{w}\|_{H^k}^2 + \frac{1}{2} \sum_{t \in T_k} \eta_t (\|\tilde{\mathbf{l}}_t\|_{H^k}^*)^2 \\ &\leq \frac{8k(1 + \epsilon)^2}{\eta_{T_k}} + \sum_{t \in T_k} \eta_t (1 + m_k)^2, \end{aligned}$$

where we set that $\mathbf{w}_{T_{k-1}}$ as last term in episode T_{k-1} , the first inequality is due to Theorem 1.

Particularly, we can choose $\eta_t = \sqrt{8k(1 + \epsilon)^2 / (1 + m_k^2)t}$, then we get the regret bound as

$$\text{Regret}_{T_k} \leq O(\sqrt{(8k(1 + \epsilon)^2)(1 + m_k)^2 T_k}).$$

Thus, we have the regret as

$$\text{Regret}_T \leq O\left((1 + \epsilon)\sqrt{T} \left(K + \sqrt{\sum_{k=1}^K km_k^2}\right)\right) \quad (26)$$

■

3.1. Parameter optimization

In this sub-section we are going to give an optimal setting of parameter γ_i , m_i and s_i . Since Equation (17), we can simplify

$$m_i = \max\{2s_i, 6\epsilon + \gamma_i\sqrt{N}\} \leq 2s_i + 2\gamma_i\sqrt{N}, \quad (27)$$

if we assume that $6\epsilon \leq \gamma_i\sqrt{N}, \forall i \in [K]$. Then we need to solve the following optimal problem:

$$\begin{aligned} & \min s_i + \gamma_i\sqrt{N} \\ & \text{s.t. } \gamma_i - \frac{\beta\epsilon}{s_i\gamma_i}i > 0 \quad \wedge \quad \beta \geq \frac{\pi}{2} \geq \gamma_i > 0 \quad \forall i \in \{1, \dots, K\}. \end{aligned} \quad (28)$$

where the latter inequality is due to Proposition 5.

It is equivalent to solve:

$$\begin{aligned} & \min s_i + \gamma_i\sqrt{N} \\ & \text{s.t. } s_i \geq (1+x)\frac{\beta\epsilon i}{\gamma_i} \quad \wedge \quad \beta \geq \frac{\pi}{2} \geq \gamma_i > 0 \quad \forall i \in \{1, \dots, K\}, \forall x > 0. \end{aligned} \quad (29)$$

With simple calculation we have that

$$\begin{cases} s_i = (N(1+x)\beta i\epsilon)^{1/3}, \\ \gamma_i = ((1+x)\beta i\epsilon)^{1/3} N^{-1/6}, \end{cases} \quad (30)$$

for any $x > 0$.

Let us re-consider the constraint in Equation (6)

$$\frac{\gamma_i}{\gamma_{i-1} - (i-1)\frac{\beta\epsilon}{s_{i-1}\gamma_{i-1}}} \cdot \frac{\pi}{2\beta} + \frac{s_i\gamma_i}{s_i\gamma_{i-1}}(i-1) \leq i \quad \forall i \in \{2, \dots, K\},$$

with above s_i and γ_i .

We can simplify this equation as

$$\frac{\pi}{2\beta} \leq i^{1/3}(i-1)^{1/3}(i^{1/3} - (i-1)^{1/3}) \cdot x, \forall i \in \{2, \dots, K\}. \quad (31)$$

Without loss the generality, we can set that $x = 1$ and $\beta = 10$, according to Lemma 14.

Therefore if we set that $s_k = (20kN\epsilon)^{1/3}$ and $\gamma_k = \sqrt{20k\epsilon/s_k}$, in Algorithm 2, we have the following bound:

$$\text{Regret}_T \leq O\left((1+\epsilon)\sqrt{T}(K + K^{4/3}(N\epsilon)^{1/3})\right), \quad (32)$$

with the constraint that $6\epsilon \leq \gamma_k\sqrt{N} \Leftrightarrow \epsilon^2 \leq \frac{20kN}{216}$, $\gamma_k \leq \frac{\pi}{2} \Leftrightarrow \epsilon^2 \leq \frac{64\pi^6 N}{400k^2}$ and $s_k \geq \epsilon \Leftrightarrow \epsilon \leq 20kN$ for all $k \in [K]$, where the last constraint is due to Remark 4.

Corollary 9 *If $\epsilon^2 \leq \min\{\frac{64\pi^6 N}{400k^2}, 20kN, \frac{20kN}{216}\}, \forall k \in [d]$, setting $s_k = (20kN\epsilon)^{1/3}$ and $\gamma_k = \sqrt{20k\epsilon/s_k}$, running our algorithm for T times we have the regret bound as*

$$\text{Regret}_T \leq O\left((1+\epsilon)\sqrt{T}(K + K^{4/3}(N\epsilon)^{1/3})\right) \leq O\left((1+\epsilon)\sqrt{T}(d + d^{4/3}(N\epsilon)^{1/3})\right). \quad (33)$$

4. Experiments

We perform preliminary experiments using a synthetic environment. In our experiments, we construct the environment as follows: Firstly, we produce a $N \times d$ matrix U , if the rank of U is not d , then produce another one, until we obtain a d -rank matrix U as kernel.

To approach the maximal regret in our experiment, before the algorithms start, we randomly produce $T \times M$ seed vectors \mathbf{v}_t^j , where $\mathbf{v}_t^j(i) \in [-1, 1], \forall i \in [d], t \in [T]$ and $j \in [M]$. Next we denote $\mathbf{l}_t^j = (U\mathbf{v}_t^j)/\|U\mathbf{v}_t^j\|_\infty$, and the noise vector $\epsilon_t = s\tilde{\epsilon}_t/(\|\tilde{\epsilon}_t\|_2)$, where $\tilde{\epsilon}_t(i)$ is randomly produced between $[-1, +1]$. We define $\tilde{\mathbf{l}}_t^j = \mathbf{l}_t^j + \epsilon_t$. Then, we process Hedge, Hazan’s and our algorithm(Algorithm 2) with the same input sequences accordingly, and by going this procedure for M sequence respectively, we can obtain M regret results for each algorithm, then choosing the maximum value among those M as follows:

On round t , for any algorithm \mathcal{A} , we record maximum of the regret with respect to M loss sequence as $\text{Regret}_{\mathcal{A}}(t)$

$$\text{Regret}_{\mathcal{A}}(t) = \max_{j \in [M]} \left\{ \sum_{s=1}^t \mathbf{w}_s^j \cdot \tilde{\mathbf{l}}_s^j - \min_{i \in [N]} \sum_{s=1}^t \tilde{\mathbf{l}}_s^i \right\}, \tag{34}$$

where we denote \mathbf{w}_s^j as output of \mathcal{A} with j -th loss sequence on round s .

Firstly, we set $T = 1000, M = 5, N = 300, d = 2$. and plot $\text{Regret}_{\mathcal{A}}(T)$ of each algorithm as a function of ϵ for $T = 1000$ in Figure 1.

As can be seen in Figure 1, our algorithm performs more robustly than others for different choices of noise ϵ . In particular, we set ϵ (i.e., setting s) as $\{0, 0.2, 0.4, 0.6, 0.8, 1\}$, respectively, and other parameters remain, and detailed graphs for each ϵ is in the supplementary material.

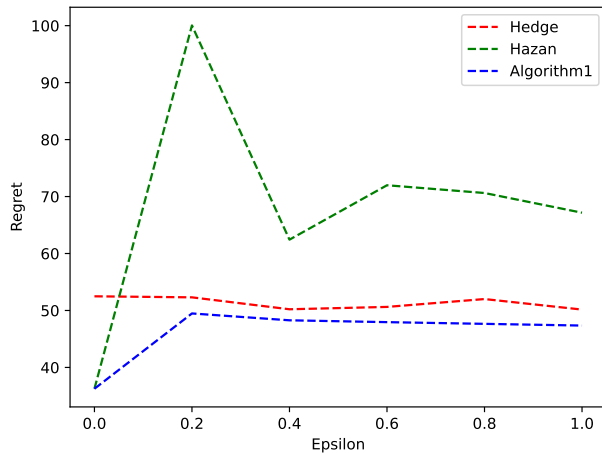


Figure 1: $\text{Regret}_{\mathcal{A}}(T)$ as a function of ϵ for $T = 1000$

Secondly, in order to see the results based on the different dimensions, we construct the experiments and set $N = 50, 100, 150, \dots, 500$, and $\epsilon = 0.2, 0.4, 0.6, 0.8$, since the unsatisfying

performance of Hazan’s algorithm shown in Figure 1, and rapidly increased running time with respect to dimension, we discard it to make sure the experiments can be done in an acceptable amount of time.

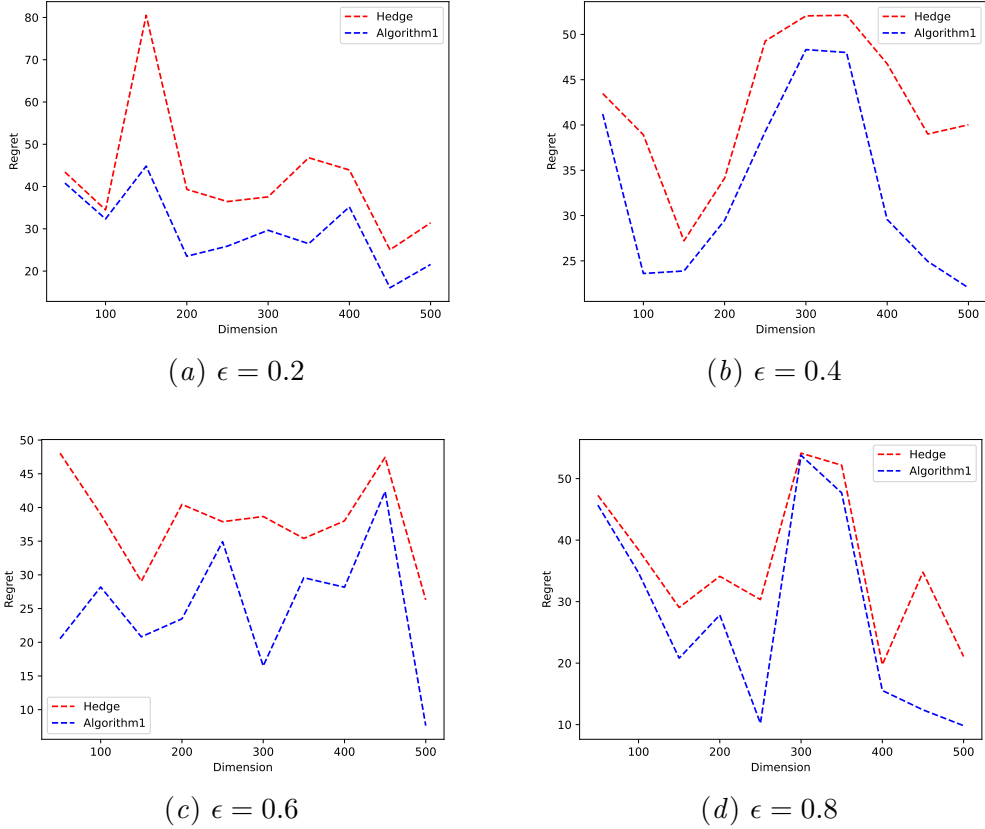


Figure 2: Results for different choices of N and ϵ

5. Concluding remarks

In this paper, we construct an algorithm for expert advice with noisy low rank loss. This algorithm is designed for the problem where the algorithm obtains no prior information about the low rank structure but only the noise bound ϵ . Theoretically, we achieve a regret bound as $O(\sqrt{T}(d + d^{4/3}(N\epsilon)^{1/3}))$, however, in the experiments, our algorithm performs better even if $\epsilon \geq \Omega(\frac{1}{N})$, which indicates that there might be a gap between our bound and the optimal one.

6. Acknowledgement

This work was supported by JSPS KAKENHI Grant Numbers JP19H04174 and JP19H04067, respectively. We thank all reviewers for the helpful comments and suggestions. Liu thank Guangsheng Ma for the beneficial discussion about some mathematical details.

References

- Maria-Florina F Balcan and Hongyang Zhang. Noise-tolerant life-long matrix completion via adaptive sampling. In *Advances in Neural Information Processing Systems*, pages 2955–2963, 2016.
- Keith Ball et al. An elementary introduction to modern convex geometry. *Flavors of geometry*, 31:1–58, 1997.
- Siddharth Barman, Aditya Gopalan, and Aadirupa Saha. Online learning for structured loss spaces. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Rina Foygel and Nathan Srebro. Concentration-based guarantees for low-rank matrix reconstruction. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 315–340, 2011.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- Andrew Goldberg, Ben Recht, Junming Xu, Robert Nowak, and Jerry Zhu. Transduction with matrix completion: Three birds with one stone. In *Advances in neural information processing systems*, pages 757–765, 2010.
- Martin Grötschel, László Lovász, and Alexander Schrijver. *Geometric algorithms and combinatorial optimization*, volume 2. Springer Science & Business Media, 2012.
- Elad Hazan, Tomer Koren, Roi Livni, and Yishay Mansour. Online learning with low rank experts. In *Conference on Learning Theory*, pages 1096–1114, 2016.
- Tomer Koren and Roi Livni. Affine-invariant online optimization and the low-rank experts problem. *Advances in Neural Information Processing Systems*, 30:4747–4755, 2017.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- Francesco Orabona, Koby Crammer, and Nicolo Cesa-Bianchi. A generalized online mirror descent with applications to classification and regression. *Machine Learning*, 99(3):411–435, 2015.
- Xiao Zhang, Lingxiao Wang Wang, and Quanquan Gu. A unified framework for nonconvex low-rank plus sparse matrix recovery. In *International Conference on Artificial Intelligence and Statistics*, 2018.

7. Appendix A. Necessary Lemmata

Lemma 10 ((Balcan and Zhang, 2016)) *Let $\mathbf{W} = \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_{k-1}\}$, $U = \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_{k-1}, \mathbf{u}\}$, and $\tilde{U} = \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_{k-1}, \tilde{\mathbf{u}}\}$ be subspaces spanned by vectors in \mathbb{R}^m . Then*

$$\theta(U, \tilde{U}) \leq \frac{\pi}{2} \frac{\theta(\tilde{\mathbf{u}}, \mathbf{u})}{\theta(\tilde{\mathbf{u}}, \mathbf{W})}$$

Lemma 11 ((Hazan et al., 2016) Lemma 11) *Let $M \in \mathbb{R}^{k \times k}$, $U \in \mathbb{R}^{N \times d}$ such that $M \succ 0$ and U . Then*

$$U^T (UMU^T)^+ U = M^{-1}. \quad (35)$$

Lemma 12 *For any $\mathbf{l}_t \in \mathbb{R}^N$, and two sub-spaces $\tilde{U}^k, U^k \subseteq \mathbb{R}^N$. We have that*

$$\theta(\mathbf{l}_t, U^k) \geq \theta(\mathbf{l}_t, \tilde{U}^k) - \theta(U^k, \tilde{U}^k). \quad (36)$$

Proof Naturally we obtain that $\theta(\mathbf{l}_t, \tilde{x}) \leq \theta(x, \tilde{x}) + \theta(\mathbf{l}_t, x)$, $\forall x, \tilde{x}$.

Given $\tilde{x}^* = \arg \min_{\tilde{x} \in \tilde{U}^k} \theta(x, \tilde{x})$, then we have that

$$\min_{\tilde{x} \in \tilde{U}^k} \theta(\mathbf{l}_t, \tilde{x}) \leq \theta(\mathbf{l}_t, \tilde{x}^*) \leq \theta(x, \tilde{x}^*) + \theta(\mathbf{l}_t, x) \quad (37)$$

Rearranging the terms we obtain that

$$\min_{\tilde{x} \in \tilde{U}^k} \theta(\mathbf{l}_t, \tilde{x}) - \min_{\tilde{x} \in \tilde{U}^k} \theta(x, \tilde{x}) \leq \theta(\mathbf{l}_t, x) \quad \forall x. \quad (38)$$

Since the definition we have that $\max_{x \in U^k} \min_{\tilde{x} \in \tilde{U}^k} \theta(x, \tilde{x}) = \theta(U^k, \tilde{U}^k)$, we have

$$\min_{\tilde{x} \in \tilde{U}^k} \theta(\mathbf{l}_t, \tilde{x}) - \theta(U^k, \tilde{U}^k) \leq \min_{\tilde{x} \in \tilde{U}^k} \theta(\mathbf{l}_t, \tilde{x}) - \min_{\tilde{x} \in \tilde{U}^k} \theta(x, \tilde{x}) \leq \theta(\mathbf{l}_t, x) \quad \forall x. \quad (39)$$

At last, setting $x = \arg \min_{x \in U^k} \theta(\mathbf{l}_t, x)$, we have that

$$\theta(\mathbf{l}_t, \tilde{U}^k) - \theta(U^k, \tilde{U}^k) = \min_{\tilde{x} \in \tilde{U}^k} \theta(\mathbf{l}_t, \tilde{x}) - \theta(U^k, \tilde{U}^k) \leq \min_{x \in U^k} \theta(\mathbf{l}_t, x) = \theta(\mathbf{l}_t, U^k). \quad (40)$$

■

Remark 13 *For any $\mathbf{l}_t \in U^k \subseteq \mathbb{R}^N$ and two spaces $U_0^k, \tilde{U}^k \subseteq \mathbb{R}^N$. We have*

$$\theta(U^k, U_0^k) + \theta(U_0^k, \tilde{U}^k) \geq \theta(U^k, \tilde{U}^k). \quad (41)$$

Proof According to Lemma 12 we have that

$$\theta(\mathbf{l}_t, U_0^k) \geq \theta(\mathbf{l}_t, \tilde{U}^k) - \theta(U_0^k, \tilde{U}^k). \quad (42)$$

Rearranging the equation we have that

$$\theta(\mathbf{l}_t, U_0^k) + \theta(U_0^k, \tilde{U}^k) \geq \theta(\mathbf{l}_t, \tilde{U}^k). \quad (43)$$

Letting $l'_t = \arg \max_{l_t \in U^k} \theta(l_t, \tilde{U}^k)$, we obtain that $\theta(l'_t, U_0^k) + \theta(U_0^k, \tilde{U}^k) \geq \theta(U^k, \tilde{U}^k)$. Since $\theta(U^k, U_0^k) = \max_{l_t \in U^k} \theta(l_t, U_0^k) \geq \theta(l'_t, U_0^k)$, we have our conclusion that

$$\theta(U^k, U_0^k) + \theta(U_0^k, \tilde{U}^k) \geq \theta(U^k, \tilde{U}^k). \quad (44)$$

■

Lemma 14 *For any $k \geq 2$, we have that*

$$\left(\frac{k}{k-1}\right)^{2/3} \left(\pi \cdot \left(\frac{k-1}{k}\right)^{1/3} + 10(k-1)\right) \leq 10k. \quad (45)$$

Proof First we have that:

$$\begin{aligned} & \left(\frac{k}{k-1}\right)^{2/3} \left(\pi \cdot \left(\frac{k-1}{k}\right)^{1/3} + 10(k-1)\right) \leq 10k \\ \Leftrightarrow & \left(\frac{k}{k-1}\right)^{1/3} \pi + 10k^{2/3}(k-1)^{1/3} \leq 10k \\ \Leftrightarrow & \pi \leq (10k - 10k^{2/3}(k-1)^{1/3}) \cdot \left(\frac{k-1}{k}\right)^{1/3} \\ \Leftrightarrow & \pi \leq 10k^{2/3}(k^{1/3} - (k-1)^{1/3}) \cdot k^{-1/3} \cdot (k-1)^{1/3} \\ \Leftrightarrow & \pi \leq 10k^{1/3}(k-1)^{1/3}(k^{1/3} - (k-1)^{1/3}). \end{aligned}$$

We set that $g(k) = k^{1/3}(k-1)^{1/3}(k^{1/3} - (k-1)^{1/3})$, then we have that

$$g'(k) = \frac{-3k(k-1)^{1/3} + 3k^{4/3} - 2k^{1/3} + (k-1)^{1/3}}{3k^{2/3}(k-1)^{2/3}}. \quad (46)$$

Next we can show that

$$\begin{aligned} & -3k(k-1)^{1/3} + 3k^{4/3} - 2k^{1/3} + (k-1)^{1/3} > 0 \\ \Leftrightarrow & k^{1/3}(3k-1) > (k-1)^{1/3}(3k-1) \\ \Leftrightarrow & k(3k-2)^3 > (k-1)(3k-1)^3 \\ \Leftrightarrow & 2k-1 > 0. \end{aligned}$$

Above equation implies that $\forall k \geq 2$, then we have $g'(k) \geq 0$, and $g(k)$ is an increasing function. Meanwhile since $\pi \leq 10 \cdot g(2)$, we obtain our conclusion. ■