

---

# One-Pass Feature Evolvable Learning with Theoretical Guarantees

---

Cun-Yuan Xing<sup>\*1</sup> Meng-Zhang Qian<sup>\*1</sup> Wu-Yang Chen<sup>1</sup> Wei Gao<sup>1</sup> Zhi-Hua Zhou<sup>1</sup>

## Abstract

Feature evolvable learning studies the scenario where old features will vanish and new features will emerge when learning with data streams, and various methods have been developed by utilizing some useful relationships from old features to new features, rather than re-training from scratch. In this work, we focus on two fundamental problems: How to characterize the relationships between two different feature spaces, and how to exploit those relationships for feature evolvable learning. We introduce the Kernel Ortho-Mapping (KOM) discrepancy to characterize relationships between two different feature spaces via kernel functions, and correlate with the optimal classifiers learned from different feature spaces. Based on this discrepancy, we develop the one-pass algorithm for feature evolvable learning, which requires going through all instances only once without storing the entire or partial training data. Our basic idea is to take online kernel learning with the random Fourier features and incorporate some feature and label relationships via the KOM discrepancy for feature evolvable learning. We finally validate the effectiveness of our proposed method both theoretically and empirically.

## 1. Introduction

Conventional machine learning generally works with the assumption that the data comes from a fixed feature space (Valiant, 1984; Shalev-Shwartz & Ben-David, 2014; Goodfellow et al., 2016; Mohri et al., 2018; Alpaydin, 2021). In some real-world applications, however, we may face more open scenarios; for example, we deploy sensors to collect data in an environmental monitoring task, and each sensor corresponds to a feature. Due to finite lifespan of sensors,

---

<sup>\*</sup>Equal contribution <sup>1</sup>National Key Laboratory for Novel Software Technology, Nanjing University, China; School of Artificial Intelligence, Nanjing University, China. Correspondence to: Wei Gao <gaow@nju.edu.cn>.

we need to deploy new sensors since old sensors will wear out, i.e., features corresponding to old sensors vanish but features corresponding to new sensors emerge.

Feature evolvable learning has been proposed to study the scenarios, where old features will vanish and new features will emerge when learning with streaming data. Recent years have witnessed increasing attention on this direction (Zhang et al., 2016; Hou et al., 2019; Beyazit et al., 2019; Gu et al., 2022; Hou et al., 2023; Schreckenberger et al., 2023; Chen & Liu, 2024). For feature evolvable learning, it is crucial to update model adaptively to accommodate new feature space but retain information of old feature space.

Various methods have been developed for feature evolvable learning by utilizing useful relationships from old features to new features, rather than re-training from scratch (Hou & Zhou, 2018; Zhang et al., 2020; Dong et al., 2022; Hou et al., 2022; Lian et al., 2023; Sajedi & Razzazi, 2024). This is helpful to prevent unnecessary wastes of computational resources and useful information from previous models over old features, and sometimes we may not collect sufficient training data to learn a stable model over new feature space.

Despite successes on the designs of practical algorithms, there are still some fundamental problems unsolved for feature evolvable learning. For example, how to present a formalization on the relationship characterization between different feature spaces, as well as correlations with model performance. Another problem is how to utilize useful relationships and information to improve the performance for feature evolvable learning from a theoretical view.

This work focuses on the one-pass algorithm for feature evolvable learning with theoretical guarantees, and the main contributions can be summarized as follows:

- We introduce the Kernel Ortho-Mapping (KOM) discrepancy to characterize the relationships between two different feature spaces via kernel functions, which essentially reflects kernels' gap under the rotational invariance. We compare our KOM discrepancy with prior characterizations such as kernel alignment and  $\ell_2$  distance (Cristianini et al., 2001; Romero et al., 2015).
- Based on the KOM discrepancy, we develop the one-pass algorithm for feature evolvable data streams, which requires going through all instances only once

without storing the entire or partial training data. Our basic idea is to take online kernel learning with the random Fourier features and incorporate feature and label relationships via the KOM discrepancy<sup>1</sup>.

- Theoretically, we establish the intrinsic relationship between the KOM discrepancy and optimal classifiers learned from different feature spaces, and present the convergence analysis to show better regret bounds via some good model initializations and relationships from old feature space and models.
- We finally conduct extensive experiments to validate the effectiveness of our OPFES method in comparison with the state-of-the-art methods on feature evolvable learning, i.e., our method achieves better performance and the fastest convergence simultaneously.

The rest of this work is constructed as follows: Section 2 presents some preliminaries. Section 3 characterizes the relationship between two different feature spaces. Section 4 develops the OPFES method. Section 5 conducts extensive experiments. Section 6 concludes with future work.

## 2. Preliminaries

**Feature evolvable learning** studies evolvable feature spaces over time, where old features will vanish and new features will emerge. Let  $\mathcal{X}^{[1]} \subseteq \mathbb{R}^{d^{[1]}}$  and  $\mathcal{X}^{[2]} \subseteq \mathbb{R}^{d^{[2]}}$  be the old and new feature spaces, respectively. Feature evolvable learning includes three stages as follows:

- **Previous stage:** receive instances  $\mathbf{x}_t^{[1]}$  from the old space  $\mathcal{X}^{[1]}$  for  $t = 1, \dots, T_1$ ;
- **Evolving stage:** receive instances  $\mathbf{x}_t^{[1]}$  and  $\mathbf{x}_t^{[2]}$  from the old space  $\mathcal{X}^{[1]}$  and new space  $\mathcal{X}^{[2]}$ , respectively, for  $t = T_1 + 1, \dots, T_1 + T_e$  with small positive  $T_e$ ;
- **Current stage:** receive instances from the new space  $\mathcal{X}^{[2]}$  for  $t = T_1 + T_e + 1, \dots, T_1 + T_e + T_2$ .

Figure 1 presents an illustration of single feature evolvable learning (Hou et al., 2017; 2022; Lian et al., 2023), and we could make a similar analysis for multiple cases.

Let  $\mathcal{K}(\cdot, \cdot)$  be a positive-definite and symmetric kernel with a mapping  $\varphi : \mathcal{X} \rightarrow \mathbb{H}$  from a feature space  $\mathcal{X}$  to an RKHS  $\mathbb{H}$ . This work focuses on the shift-invariant kernels  $\mathcal{K}(\mathbf{x}_1, \mathbf{x}_2) = \kappa(\mathbf{x}_1 - \mathbf{x}_2)$ , such as Gaussian kernel and Laplacian kernel (Schölkopf & Smola, 2002).

**Online kernel learning** trains classifiers  $h_1, \dots, h_T \in \mathbb{H}$  from a streaming sample  $S_T = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)\}$

<sup>1</sup>The code is available at [github.com/WeltXing/opfes](https://github.com/WeltXing/opfes)

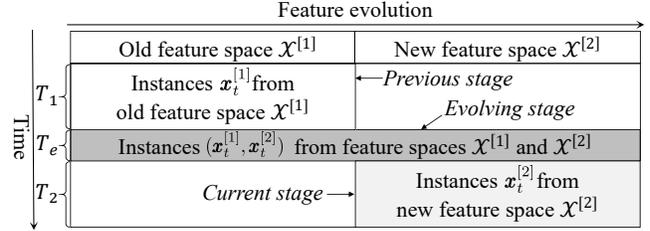


Figure 1. An illustration of feature evolvable stream.

with  $y_i \in \{-1, +1\}$ , by minimizing the following loss

$$h_t \in \arg \min_{h \in \mathbb{H}} \left\{ \frac{1}{t} \sum_{i=1}^t \ell(h, (\mathbf{x}_i, y_i)) + \frac{\lambda}{2} \|h\|_{\mathbb{H}}^2 \right\},$$

where  $\ell(h_t, (\mathbf{x}, y)) = \max\{0, 1 - yh_t(\mathbf{x})\}$ . Based on the representer theorem (Schölkopf & Smola, 2002), we have

$$h_t(\mathbf{x}) = \sum_{i=1}^t \alpha_i \mathcal{K}(\mathbf{x}, \mathbf{x}_i) \text{ and } \|h_t\|_{\mathbb{H}}^2 = \sum_{i,j=1}^t \alpha_i \alpha_j \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j).$$

This shows that online kernel learning requires storing the entire training sample  $S_T$ , which makes it difficult for the large-scale datasets (Shen et al., 2019; Hong & Chae, 2021).

**Online kernel learning with random Fourier features** has been an efficient way for large-scale online kernel learning (Rahimi & Recht, 2007; Lu et al., 2016), which approximates high (or infinite) dimensional mapping  $\varphi(\cdot)$  with the finite-dimensional random Fourier features. Specifically, we approximate the kernel function of  $\mathcal{K}$  as

$$\begin{aligned} \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) &= \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle \\ &\approx \sum_{k=1}^d p(\mathbf{u}_k) \phi(\mathbf{x}_i, \mathbf{u}_k, b_k) \phi(\mathbf{x}_j, \mathbf{u}_k, b_k), \end{aligned} \quad (1)$$

where  $\phi(\mathbf{x}, \mathbf{u}, b) = \sqrt{2} \cos(\langle \mathbf{x}, \mathbf{u} \rangle + b)$ , and  $p(\cdot)$  is the spectral density function of  $\mathcal{K}$ . Here, random vectors  $\mathbf{u}_k$  are sampled i.i.d. from standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and  $b_k$  are randomly selected independently and uniformly over  $[0, 2\pi]$ . For shift-invariant kernel  $\mathcal{K}(\mathbf{x}_1, \mathbf{x}_2)$ , we have

$$p(\mathbf{u}) = \int_{\mathbb{R}^d} \kappa(\mathbf{x}) \exp(-\mathbf{i}\langle \mathbf{x}, \mathbf{u} \rangle) / (2\pi)^d d\mathbf{x},$$

where  $\mathbf{i}$  is the imaginary unit. By random Fourier features, we can approximate a kernel classifier

$$h(\mathbf{x}) = \langle \tilde{\mathbf{w}}, \varphi(\mathbf{x}) \rangle \approx \langle \mathbf{w}, \mathbf{z}(\mathbf{x}) \rangle,$$

where the finite-dimensional approximated vector

$$\mathbf{z}(\mathbf{x}) = (\sqrt{p(\mathbf{u}_1)} \phi(\mathbf{x}, \mathbf{u}_1, b_1), \dots, \sqrt{p(\mathbf{u}_d)} \phi(\mathbf{x}, \mathbf{u}_d, b_d)).$$

We update the classifier according to Fourier online gradient descent (Lu et al., 2016) as follows:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \tau_t \nabla \ell_t(\mathbf{w}_{t-1}), \quad (2)$$

where  $\tau_t$  is the stepsize and loss function

$$\ell_t(\mathbf{w}_{t-1}) = \ell(\mathbf{w}_{t-1}, (z(\mathbf{x}_t), y_t)) + \frac{\lambda}{2} \|\mathbf{w}_{t-1}\|_2^2.$$

This work focuses on random Fourier feature technique over shift-invariant kernels, and it is natural to make similar approximation and online algorithm for other kernels such as polynomial kernel (Pennington et al., 2015) and linear kernel (Wacker et al., 2024).

We introduce some notations throughout this work. Bold uppercase and lowercase letters denote matrices and vectors, respectively. We denote by  $\|\cdot\|_p$  the  $\ell_p$ -norm of a vector, and  $\|\cdot\|_F$  and  $\|\cdot\|_*$  denote the Frobenius norm and nuclear norm of a matrix, respectively. Denote by  $\|\cdot\|_{\text{HS}}$  the Hilbert-Schmidt norm of an operator, which is an extension of the Frobenius norm in Hilbert space. Let  $(\mathbf{v}_i)_{i=1}^n$  be the  $d \times n$  concatenation matrix of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^d$ .

Let  $\mathbf{I}_d$  be an  $d \times d$  identity matrix, and  $\text{diag}(\mathbf{v})$  is a diagonal matrix with diagonal elements  $\mathbf{v}$ . Denote by  $\mathbf{1}_d$  and  $\mathbf{0}_d$   $d$ -dimensional vectors with all-one and all-zero elements, respectively. Let  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  be a Gaussian distribution with parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Let  $\mathcal{U}_d = \{\mathbf{U} \in \mathbb{R}^{d \times d}; \mathbf{U}\mathbf{U}^\top = \mathbf{I}_d\}$  be the set of  $d \times d$  orthogonal matrices, and  $\sqrt{\mathbf{A}}$  is the square root of positive-semidefinite matrix  $\mathbf{A}$ , i.e.,  $\sqrt{\mathbf{A}}\sqrt{\mathbf{A}} = \mathbf{A}$ .

### 3. On the Exploration of Feature Relationship

In this section, we introduce a general framework on the characterization of relationship between different feature spaces based on kernel functions, while previous studies can be viewed as some special selections of different kernels (Hou et al., 2017; 2021; Zhou et al., 2024). We further correlate it with the distance between classifiers trained from different feature spaces and then develop the one-pass learning algorithm for optimization.

#### 3.1. Charactering Relationship between Feature Spaces

Our basic idea is to map the original raw feature spaces into Reproducing Kernel Hilbert Spaces (RKHSs) with kernel functions, which could provide plentiful and implicit non-linear representations for original feature spaces.

We focus on positive-definite kernel functions  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$  over old feature space  $\mathcal{X}^{[1]}$  and new feature space  $\mathcal{X}^{[2]}$ , respectively. For  $S_n = \{(\mathbf{x}_1^{[1]}, \mathbf{x}_1^{[2]}), \dots, (\mathbf{x}_n^{[1]}, \mathbf{x}_n^{[2]})\}$  with  $\mathbf{x}_i^{[1]} \in \mathcal{X}^{[1]}$  and  $\mathbf{x}_i^{[2]} \in \mathcal{X}^{[2]}$ , we define their Gram matrices

$$\mathbf{K}^{[k]} = \left[ \mathcal{K}^{[k]}(\mathbf{x}_i^{[k]}, \mathbf{x}_j^{[k]}) \right]_{n \times n} \quad \text{for } k = 1, 2.$$

For  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$ , we introduce a new distance to measure the difference between two feature spaces as follows.

**Definition 3.1.** We define Kernel Ortho-Mapping (KOM) discrepancy between  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$  over sample  $S_n$  as

$$\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) = \min_{\mathbf{U} \in \mathcal{U}_n} \left\{ \|\mathbf{U}\sqrt{\mathbf{K}^{[1]}} - \sqrt{\mathbf{K}^{[2]}}\|_F / \sqrt{n} \right\}.$$

In this definition, the empirical kernel mapping is introduced to deal with different dimensionalities of kernel mappings as done by Schölkopf & Smola (2002) and Marukatat (2016), and the minimum is taken for the uniqueness of kernel mapping from rotational invariance.

**Lemma 3.2.** We have the closed-form solution for the KOM discrepancy (in Definition 3.1) as

$$\begin{aligned} & \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \\ &= \left( \text{Tr}(\mathbf{K}^{[1]} + \mathbf{K}^{[2]})/n - 2\|\sqrt{\mathbf{K}^{[1]}}\sqrt{\mathbf{K}^{[2]}}\|_*/n \right)^{1/2}. \end{aligned}$$

The detailed proof is given in Appendix A.1, and the basic idea is to take the polar decomposition and upper bound the trace of an orthogonal matrix over  $\mathcal{U}_n$ .

For loss function  $\ell(h, (\mathbf{x}, y)) = \max\{0, 1 - yh_t(\mathbf{x})\}$ , we define the optimal kernel classifiers over sample  $S_n$  in the old and new feature spaces as follows:

$$h_*^{[k]} \in \arg \min_{h^{[k]} \in \mathbb{H}^{[k]}} \sum_{i=1}^n \frac{\ell(h^{[k]}, (\mathbf{x}_i^{[k]}, y_i))}{n} + \frac{\lambda}{2} \|h^{[k]}\|_{\mathbb{H}^{[k]}}^2, \quad (3)$$

where  $k = 1$  and  $k = 2$  correspond to the old and new feature spaces, respectively. We measure the gap between two optimal classifiers  $h_*^{[1]}$  and  $h_*^{[2]}$  over sample  $S_n$  as follows:

$$\hat{\rho}_{S_n}(h_*^{[1]}, h_*^{[2]}) = \frac{1}{n} \sum_{i=1}^n \left| h_*^{[1]}(\mathbf{x}_i^{[1]}) - h_*^{[2]}(\mathbf{x}_i^{[2]}) \right|. \quad (4)$$

We now present the first main result to correlate our KOM discrepancy with two optimal classifiers as follows:

**Theorem 3.3.** Given sample  $S_n$ , we have

$$\hat{\rho}_{S_n}(h_*^{[1]}, h_*^{[2]}) \leq \frac{r}{\lambda} \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) + \frac{r}{\lambda} (2r \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}))^{\frac{1}{2}}$$

for two kernels  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$  bounded by  $r^2$ , where  $\lambda$  is the regularization parameter in Eqn. (3).

In this theorem, we upper bound the distance between two optimal classifiers with our KOM discrepancy, and this may shed some new insights to develop feature evolvable learning algorithms based on our KOM discrepancy, which essentially measures the relationships between two different feature spaces. The detailed proof of Theorem 3.3 is given in Appendix A.2, which linearizes the kernel classifiers via empirical kernel mapping and constructs KOM discrepancy.

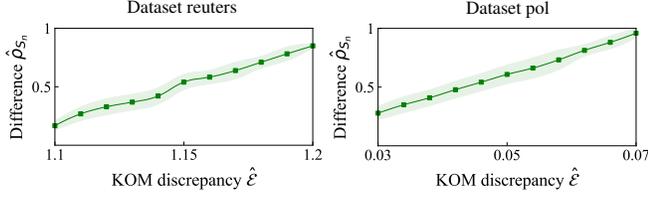


Figure 2. Illustrations of the relationship between  $\hat{\mathcal{E}}$  and  $\hat{\rho}_{S_n}$ .

Theorem 3.3 is limited to binary classification, while it is easy to make a similar analysis for multi-class learning (Crammer & Singer, 2002) and regression (Murphy, 2012).

Figure 2 presents an intuitive illustration on the relationship between KOM discrepancy  $\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})$  and classifiers' difference  $\hat{\rho}_{S_n}(h_*^{[1]}, h_*^{[2]})$  via Gaussian kernels over datasets *pol* and *reuters*. As we can see,  $\hat{\rho}_{S_n}(h_*^{[1]}, h_*^{[2]})$  is positively relevant to KOM discrepancy  $\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})$ , i.e., the bigger the KOM discrepancy  $\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})$ , the larger the distance  $\hat{\rho}_{S_n}(h_*^{[1]}, h_*^{[2]})$ . This is nicely in accordance with our Theorem 3.3 empirically.

Our KOM discrepancy is defined over the sample  $S_n$  in Definition 3.1. We can also define the KOM discrepancy w.r.t. distribution  $\mathcal{D}$  over  $\mathcal{X}^{[1]} \times \mathcal{X}^{[2]}$  as follows

$$\begin{aligned} \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \\ = \min_{\mathbf{U} \in \mathcal{U}} \left\{ \sqrt{\mathbb{E}_{\mathcal{D}} [\|\mathbf{U}\varphi^{[1]}(\mathbf{x}^{[1]}) - \varphi^{[2]}(\mathbf{x}^{[2]})\|_{\text{HS}}^2]} \right\}, \end{aligned}$$

where  $\mathcal{U}$  is a unitary operator set on a real Hilbert space. We present the following convergence analysis.

**Theorem 3.4.** *Let  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$  be two kernels bounded by  $r^2$ . For  $\delta \in (0, 1)$  and for some constant  $c_1 > 0$ , we have, with probability at least  $1 - \delta$  over sample  $S_n$*

$$\left| \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) - \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right| \leq c_1 r \sqrt{\frac{1}{n} \ln \frac{2}{\delta}}.$$

The detailed proof is presented in Appendix A.3, in which the main techniques include McDiarmid's inequality (McDiarmid et al., 1989) and operator Khintchine's inequality in non-commutative probability (Vershynin, 2018).

### Relevant to previous relationship characterizations

Kernel alignment has been used to characterize relationship between two kernels (Cristianini et al., 2001; Cortes et al., 2012; Zhou et al., 2024), which essentially calculates the cosine similarity between two Gram matrices  $\mathbf{K}^{[1]}$  and  $\mathbf{K}^{[2]}$

$$\hat{A}(\mathbf{K}^{[1]}, \mathbf{K}^{[2]}) = \frac{\text{Tr}(\mathbf{K}^{[1]}\mathbf{K}^{[2]})}{\|\mathbf{K}^{[1]}\|_F \|\mathbf{K}^{[2]}\|_F}.$$

We could present the following relationship between our KOM discrepancy and kernel alignment, and the detailed proof is given in Appendix A.4.

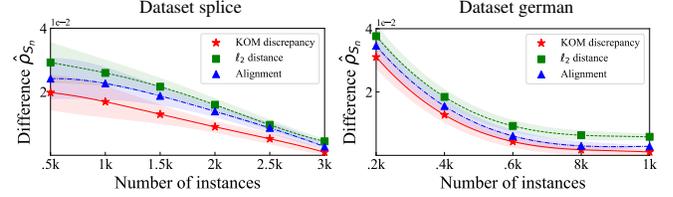


Figure 3. An illustration of the difference between two optimal classifiers on two feature spaces, by optimizing KOM discrepancy and previous kernel alignment and  $\ell_2$  distance, respectively.

**Lemma 3.5.** *For two normalized kernel matrices  $\mathbf{K}^{[1]}$  and  $\mathbf{K}^{[2]}$  with  $\|\mathbf{K}^{[1]}\|_F = \|\mathbf{K}^{[2]}\|_F \leq r^2$ , we have*

$$\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \leq r \sqrt[4]{2(1 - \hat{A}(\mathbf{K}^{[1]}, \mathbf{K}^{[2]}))}.$$

The  $\ell_2$  distance has also been used to align the features of two kernel mappings (Romero et al., 2015; Heo et al., 2019), which solves a finite-dimensional kernel mapping  $\hat{\varphi}^{[2]}$  from the following optimization problem:

$$\hat{\varphi}^{[2]} \in \arg \min_{\varphi^{[2]} \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^n \|\varphi^{[1]}(\mathbf{x}_i^{[1]}) - \varphi^{[2]}(\mathbf{x}_i^{[2]})\|_2^2 \right\}.$$

Here, function space  $\mathcal{F} \subseteq \{\varphi^{[2]} : \mathbb{R}^{d^{[2]}} \rightarrow \mathbb{R}^{\dim(\varphi^{[1]})}\}$ , and  $\varphi^{[1]}$  is the finite-dimensional mapping of kernel  $\mathcal{K}^{[1]}$ .

The  $\ell_2$  distance has been successfully applied for feature evolvable learning. For example, Hou et al. (2017; 2022) selected linear kernels  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$ , while Chen & Liu (2024) considered Gaussian kernel  $\mathcal{K}^{[1]}$  and Mahalanobis kernel  $\mathcal{K}^{[2]}$ . We can also present the following relationship between the KOM discrepancy and  $\ell_2$  distance, and the detailed proof is given in Appendix A.5.

**Lemma 3.6.** *For kernel  $\mathcal{K}^{[1]}$  with mapping  $\varphi^{[1]}$ , we have*

$$\left( \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \hat{\mathcal{K}}^{[2]}) \right)^2 \leq \frac{1}{n} \sum_{i=1}^n \left\| \varphi^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2^2,$$

where kernel  $\hat{\mathcal{K}}^{[2]}(\mathbf{x}^{[2]}, \mathbf{x}^{[2]'}) = \langle \hat{\varphi}^{[2]}(\mathbf{x}^{[2]}), \hat{\varphi}^{[2]}(\mathbf{x}^{[2]'}) \rangle$ .

Figure 3 presents an illustration of the difference between two optimal classifiers over two feature spaces by optimizing our KOM discrepancy, previous  $\ell_2$  distance and kernel alignment, respectively. Here, we focus on simple linear mapping spaces on two datasets *splice* and *german*, and the trends are similar for other datasets.

From Figure 3, it is observable that we could get a smaller difference between two optimal classifiers by optimizing the KOM discrepancy, rather than kernel alignment and  $\ell_2$  distance. Therefore, our KOM discrepancy presents a better characterization of relationships between different feature spaces via kernel functions, and this is partially consistent with Lemma 3.5 and Lemma 3.6.

---

**Algorithm 1** One-pass optimization of Eqn. (6)
 

---

**Input:** Streaming sample  $S_{T_e}^{[e]}$ , number of iterations  $T_M$ , and stepsize  $\eta_M$ 
**Output:** Spectral density  $\mathbf{p}^{(T_M)}$ 

- 1: Initialize  $\mathbf{M}^{(T_1)} = \mathbf{0}$ ,  $\mathbf{v}^{(T_1)} = \mathbf{0}$  and  $\mathbf{p}^{(0)} = \mathbf{1}/d_2$
  - 2: **for**  $t = T_1 + 1, \dots, T_1 + T_e$  **do**
  - 3:   Update  $\mathbf{M}^{(t)}$  and  $\mathbf{v}^{(t)}$  according to Eqns. (7)-(8)
  - 4: **end for**
  - 5: **for**  $i = 1, \dots, T_M$  **do**
  - 6:   Calculate  $\mathbf{q} = \mathbf{p}^{(i-1)} \exp(-\eta_M \nabla f(\mathbf{p}^{(i-1)}))$
  - 7:   Update  $\mathbf{p}^{(i)} = \mathbf{q} / \|\mathbf{q}\|_1$
  - 8: **end for**
  - 9: **return:** Spectral density  $\mathbf{p}^{(T_M)}$
- 

### 3.2. One-Pass Optimization for our KOM discrepancy

During the evolving stage, we reach a streaming sample  $S_{T_e}^{[e]} = \{(\mathbf{x}_{T_1+1}^{[1]}, \mathbf{x}_{T_1+1}^{[2]}, \dots, (\mathbf{x}_{T_1+T_e}^{[1]}, \mathbf{x}_{T_1+T_e}^{[2]})\}$ , and get the kernel  $\mathcal{K}^{[1]}$  learned from the previous stage. Motivated by Theorem 3.3, we learn the kernel  $\mathcal{K}^{[2]}$  by minimizing the KOM discrepancy as follows

$$\min_{\mathcal{K}^{[2]}} \left\{ \hat{\mathcal{E}}(S_{T_e}^{[e]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right\}. \quad (5)$$

For steaming sample  $S_{T_e}^{[e]}$ , it is not allowed to store kernel Gram matrices in memory, and we can not directly optimize the above optimization as in (Cortes et al., 2012; Liu, 2024).

We present one-pass optimization for Eqn. (5) via random Fourier features. The basic idea is to transform the original optimization into a convex problem on a simplex, and then solve it w.r.t. streaming data  $S_{T_e}^{[e]}$ . For the spectral density  $\mathbf{p} = (p^{[2]}(\mathbf{u}_1^{[2]}), \dots, p^{[2]}(\mathbf{u}_{d_2}^{[2]}))$  of  $\mathcal{K}^{[2]}$ , we approximate the KOM discrepancy by random Fourier features as

$$\begin{aligned} & \hat{\mathcal{E}}(S_{T_e}^{[e]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \\ & \approx \sqrt{(\text{Tr}(\mathbf{K}^{[1]}) + \langle \mathbf{p}, \mathbf{v} \rangle - 2\|\mathbf{M}\text{diag}(\sqrt{\mathbf{p}})\|_*) / T_e}, \end{aligned}$$

where  $\mathbf{M} = [m_{kl}]_{d_1 \times d_2}$  and  $\mathbf{v} = (v_1, v_2, \dots, v_{d_2})$  with

$$m_{kl} = \sum_{i=T_1+1}^{T_1+T_e} p^{[1]}(\mathbf{u}_k^{[1]}) \Phi_{ik}^{[1]} \Phi_{il}^{[2]} \quad \text{and} \quad v_k = \sum_{i=T_1+1}^{T_1+T_e} (\Phi_{ik}^{[2]})^2.$$

Here,  $\Phi_{ij}^{[1]} = \phi(\mathbf{x}_i^{[1]}, \mathbf{u}_j^{[1]}, b_j^{[1]})$ ,  $\Phi_{ij}^{[2]} = \phi(\mathbf{x}_i^{[2]}, \mathbf{u}_j^{[2]}, b_j^{[2]})$ , and  $d_1$  and  $d_2$  are the numbers of random Fourier features of kernels  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$ , respectively. From random features approximation, Eqn. (5) is essentially equivalent to

$$\min_{\mathbf{p} \in \Delta} \left\{ f(\mathbf{p}) = \frac{1}{2} \langle \mathbf{p}, \mathbf{v} \rangle - \|\mathbf{M}\text{diag}(\sqrt{\mathbf{p}})\|_* \right\}, \quad (6)$$

where  $\Delta = \{\mathbf{p}: \mathbf{p} \geq \mathbf{0}, \|\mathbf{p}\|_1 = 1\}$ . For Eqn. (6), we initialize  $\mathbf{M}^{(T_1)} = [\mathbf{0}]_{d_1 \times d_2}$  and  $\mathbf{v}^{(T_1)} = \mathbf{0}_{d_2}$ , and in the

$t$ -th round ( $t \geq T_1 + 1$ ), we update  $\mathbf{M}^{(t)}$  and  $\mathbf{v}^{(t)}$  w.r.t. instances  $\mathbf{x}_t^{[1]}$  and  $\mathbf{x}_t^{[2]}$ , respectively, as

$$m_{kl}^{(t)} = m_{kl}^{(t-1)} + p^{[1]}(\mathbf{u}_k^{[1]}) \Phi_{tk}^{[1]} \Phi_{tl}^{[2]} \quad (7)$$

$$v_k^{(t)} = v_k^{(t-1)} + (\Phi_{tk}^{[2]})^2, \quad (8)$$

with  $\Phi_{ij}^{[1]} = \phi(\mathbf{x}_i^{[1]}, \mathbf{u}_j^{[1]}, b_j^{[1]})$  and  $\Phi_{ij}^{[2]} = \phi(\mathbf{x}_i^{[2]}, \mathbf{u}_j^{[2]}, b_j^{[2]})$ .

We finally take the mirror descent method (Bubeck, 2015; Hazan et al., 2016) to solve Eqn. (6).

Algorithm 1 presents the details of our proposed method, and we have the convergence analysis as follows.

**Theorem 3.7.** For Algorithm 1, we have

$$\frac{1}{T_M} \sum_{t=1}^{T_M} f(\mathbf{p}^{(t)}) - f(\mathbf{p}^*) \leq O\left(\frac{1}{\sqrt{T_M}}\right),$$

by setting stepsize  $\eta_M = \Theta(\sqrt{\ln d_2 / T_M})$ , where  $f(\cdot)$  is defined by Eqn. (6) and  $\mathbf{p}^* \in \arg \min_{\mathbf{p} \in \Delta} f(\mathbf{p})$ .

The detailed proof is given in Appendix B, which presents an operation to preserve convexity and then derives the convergence analysis of mirror descent on a simplex.

## 4. The OPFES Approach

This section presents the one-pass optimization for feature evolvable learning in the current stage as shown in Figure 1. Our idea is to incorporate feature and label information from previous relationships and stages, and reuse prior models, rather than re-training a new model from scratch.

### i) Incorporation of feature information via kernel $\mathcal{K}^{[2]}$

Notice that the kernel  $\mathcal{K}^{[2]}$  is learned by minimizing the KOM discrepancy in Section 3. We could train an online kernel model based on  $\mathcal{K}^{[2]}$ , which incorporate implicitly some information from old feature space.

Specifically, we reach an example  $(\mathbf{x}_t^{[2]}, y_t)$  in the  $t$ -th round for  $T_1 + T_e + 1 \leq t \leq T_1 + T_e + T_2$ , and learn an online model  $h_t^{[2]}(\mathbf{x}^{[2]})$  based on kernel  $\mathcal{K}^{[2]}$  as follows

$$h_t^{[2]}(\mathbf{x}^{[2]}) \approx \langle \mathbf{w}_t^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}^{[2]}) \rangle,$$

where  $\mathbf{z}^{[2]}(\mathbf{x}^{[2]}) = ((p^{[2]}(\mathbf{u}_k^{[2]}))^{1/2} \phi(\mathbf{x}^{[2]}, \mathbf{u}_k^{[2]}, b_k^{[2]}))_{k=1}^{d_2}$ . Here, we take the representer theorem (Schölkopf & Smola, 2002) and random features approximation from Eqn. (1). We take online gradient descent with stepsize  $\tau_t^{[2]}$  as

$$\mathbf{w}_t^{[2]} = \mathbf{w}_{t-1}^{[2]} - \tau_t^{[2]} \nabla \ell_t^{[2]}(\mathbf{w}_{t-1}^{[2]}), \quad (9)$$

where the loss function  $\ell_t^{[2]}(\mathbf{w}_{t-1}^{[2]})$  is given by

$$\max \left\{ 0, 1 - y_t \langle \mathbf{w}_{t-1}^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle \right\} + \frac{\lambda}{2} \|\mathbf{w}_{t-1}^{[2]}\|_2^2.$$

## ii) Incorporation of label information via ideal kernel

For label information, we exploit the *ideal kernel* from kernel selections (Cristianini et al., 2001; Kwok et al., 2003). Essentially, the ideal kernel is helpful to learn models with correct predictions for training data (Cristianini et al., 2001). We define the ideal kernel  $\mathcal{K}^*$  over feature space  $\mathcal{X}^{[2]}$  as

$$\mathcal{K}^*(\mathbf{x}_i^{[2]}, \mathbf{x}_j^{[2]}) = y_i y_j \text{ for } i, j \in \{T_1 + 1, \dots, T_1 + T_e\}.$$

We can not obtain the ideal kernel matrix owing to streaming data. Our idea is to learn a new kernel  $\mathcal{K}^l$  aligning with the ideal kernel  $\mathcal{K}^*$  to incorporate label information, i.e.,

$$\mathcal{K}^l = \arg \min_{\mathcal{K}^l} \left\{ \hat{\mathcal{E}} \left( S_{T_e}^{[e]}, \mathcal{K}^*, \mathcal{K}^l \right) \right\}.$$

In the  $t$ -th round, we learn a new online model by representer theorem and random features approximation again, i.e.,

$$h_t^l(\mathbf{x}^{[2]}) \approx \langle \mathbf{w}_t^l, \mathbf{z}^l(\mathbf{x}^{[2]}) \rangle,$$

where  $\mathbf{z}^l(\mathbf{x}^{[2]}) = ((p^l(\mathbf{u}_k^l))^{1/2} \phi(\mathbf{x}^{[2]}, \mathbf{u}_k^l, b_k^l))_{k=1}^{d_2}$  with the density  $p^l$  of kernel  $\mathcal{K}^l$ . Given stepsize  $\tau_t^l$ , we update

$$\mathbf{w}_t^l = \mathbf{w}_{t-1}^l - \tau_t^l \nabla \ell_t^l(\mathbf{w}_{t-1}^l), \quad (10)$$

where the loss function  $\ell_t^l(\mathbf{w}_{t-1}^l)$  is given by

$$\max \{0, 1 - y_t \langle \mathbf{w}_{t-1}^l, \mathbf{z}^l(\mathbf{x}_t^{[2]}) \rangle\} + \frac{\lambda}{2} \|\mathbf{w}_{t-1}^l\|_2^2.$$

## iii) Previous model reuse

We exploit some good model initializations in the current stage, rather than re-training from scratch. Our basic idea is to obtain a new model on the space spanned by  $\mathbf{z}^{[2]}(\mathbf{x}^{[2]})$ , which takes similar predictions with the previous model.

Specifically, we consider the model for the new feature space via an orthogonal transformation

$$\mathbf{w}_{T_1+T_e}^{[2]} = \mathbf{U}_*^\top \mathbf{w}_{T_1}^{[1]}, \quad (11)$$

where  $\mathbf{U}_*$  is in the set of

$$\arg \min_{\mathbf{U} \in \mathcal{U}_{d_1}} \left\{ \left\| \mathbf{U} \left( \mathbf{z}^{[1]}(\mathbf{x}_{T_1+i}^{[1]}) \right)_{i=1}^{T_e} - \left( \mathbf{z}^{[2]}(\mathbf{x}_{T_1+i}^{[2]}) \right)_{i=1}^{T_e} \right\|_F \right\}.$$

In the following, we present an effective solution for  $\mathbf{U}_*$ , and the detailed proof is given in Appendix C.1.

**Proposition 4.1.** *We have  $\mathbf{U}_* = \mathbf{V}\mathbf{W}^\top$  for the optimal solution in Eqn. (11), where  $\mathbf{V}$  and  $\mathbf{W}$  are left and right singular vectors of  $\mathbf{M}^{(T_1+T_e)} \text{diag}(\sqrt{\mathbf{p}^{(T_M)}})$  in Algorithm 1.*

Based on previous analysis, we should learn and update two online kernel models  $\mathbf{w}_t^{[2]}$  and  $\mathbf{w}_t^l$  from Eqns. (9) and (10)

## Algorithm 2 The OPFES method

**Input:** Feature evolvable stream sample  $S_{T_1+T_e+T_2}$ , kernel  $\mathcal{K}^{[1]}$ , stepsize  $\tau_t^{[1]}$ ,  $\tau_t^{[2]}$  and  $\tau_t^l$ , sensitivity parameter  $\gamma$

**Initialize:**  $\mathbf{w}_0^{[1]} = \mathbf{0}$

**Output:** classifier  $h_{T_1+T_e+T_2}$

- 1: Obtain random Fourier features  $(\mathbf{u}_k^{[1]}, b_k^{[1]})_{k=1}^{d_1}$  and  $(\mathbf{u}_k^{[2]}, b_k^{[2]}, \mathbf{u}_k^l, b_k^l)_{k=1}^{d_2}$  via Eqn. (1)
- 2: **for**  $t = 1, \dots, T_1$  **do**
- 3:   Update  $\mathbf{w}_t^{[1]}$  by online gradient descent in Eqn. (2)
- 4: **end for**
- 5: Obtain  $\mathbf{p}^{[2]}$  and  $\mathbf{p}^l$  from Algorithm 1
- 6: Compute  $\mathbf{w}_{T_1+T_e}^{[2]}$  by Eqn. (11)
- 7: **for**  $t = T_1 + T_e + 1, \dots, T_1 + T_e + T_2$  **do**
- 8:   Update  $\mathbf{w}_t^{[2]}$  and  $\mathbf{w}_t^l$  by Eqns. (9)-(10), respectively
- 9:   Update the combined classifier  $h_t$  by Eqn. (12)
- 10: **end for**
- 11: **return:** classifier  $h_{T_1+T_e+T_2}$

in the current stage, respectively. From (Hou et al., 2021), we combine two models, for  $t \geq T_1 + T_e + 1$ ,

$$h_t(\mathbf{x}_t^{[2]}) = \omega_t \langle \mathbf{w}_t^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle + (1 - \omega_t) \langle \mathbf{w}_t^l, \mathbf{z}^l(\mathbf{x}_t^{[2]}) \rangle, \quad (12)$$

where  $\omega_t$  is relevant to a sensitivity parameter  $\gamma > 0$ , i.e.,

$$\omega_t = \frac{\omega_{t-1} e^{-\gamma \ell_t^{[2]}(\mathbf{w}_{t-1}^{[2]})}}{\omega_{t-1} e^{-\gamma \ell_t^{[2]}(\mathbf{w}_{t-1}^{[2]})} + (1 - \omega_{t-1}) e^{-\gamma \ell_t^l(\mathbf{w}_{t-1}^l)}}.$$

Algorithm 2 presents the detailed description of our OPFES approach, which goes through all instances only once without storing the entire or partial training data, while previous methods require storing the entire dataset or partial dataset (Orabona et al., 2008; Jin et al., 2010; Hou et al., 2021; 2022; He et al., 2021a; Wu et al., 2023).

## Theoretical guarantee

We begin with the upper bounds for prediction difference between  $\mathbf{w}_{T_1}^{[1]}$  and  $\mathbf{w}_{T_1+T_e}^{[2]}$  via our KOM discrepancy, and the detailed proof is presented in Appendix C.2.

**Lemma 4.2.** *For bounded kernels, we have, for previous model  $\mathbf{w}_{T_1}^{[1]}$  and reused model  $\mathbf{w}_{T_1+T_e}^{[2]}$  from Eqn. (11),*

$$\begin{aligned} & \frac{1}{T_e} \sum_{t=T_1+1}^{T_1+T_e} \left| \langle \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) \rangle - \langle \mathbf{w}_{T_1+T_e}^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle \right| \\ & \leq \sqrt{2} \hat{\mathcal{E}}(S_{T_e}^{[e]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) / \lambda, \end{aligned}$$

where  $\lambda$  is the regularization parameter in Eqn. (3).

Table 1. Details of datasets

Dataset	# Inst.	# Feat.	Dataset	# Inst.	# Feat.	Dataset	# Inst.	# Feat.	Dataset	# Inst.	# Feat.	Dataset	# Inst.	# Feat.
jungle	2351	87	svmguidel	7089	4	elevators	16599	18	nomao	34465	118	higgs	98049	28
splice	3175	60	usps	9298	25	magic	19020	10	adult	48842	108	miniboone	130064	50
bioresponse	3751	1776	aileron	13750	40	letter	20000	16	acoustic	78823	50	ijcnn1	141691	22
christine	5418	1636	pol	15000	44	house	22784	16	runwalk	88588	6	covtype	581012	54

Denote by the optimal model in the current stage

$$\mathbf{w}_*^{[2]} \in \arg \min_{\mathbf{w}} \left\{ \mathcal{L}_{T_2}^{[2]}(\mathbf{w}) = \frac{1}{T_2} \sum_{t=T_1+T_e+1}^{T_1+T_e+T_2} \ell_t^{[2]}(\mathbf{w}) \right\},$$

and the cumulative loss

$$\hat{\mathcal{L}}_{T_2}^{[2]} = \frac{1}{T_2} \sum_{t=T_1+T_e+1}^{T_1+T_e+T_2} \ell_t^{[2]}(\mathbf{w}_t^{[2]}),$$

with  $\ell_t^{[2]}(\mathbf{w}) = \max\{0, 1 - y_t \langle \mathbf{w}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle\} + \lambda \|\mathbf{w}\|_2^2 / 2$ .

Let  $S_{T_2} = \{(\mathbf{x}_i^{[2]}, y_i)\}_{i=T_1+T_e+1}^{T_1+T_e+T_2}$  be a streaming sample drawn i.i.d. from a distribution. For Algorithm 2, we have

**Theorem 4.3.** *For kernels  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$  with bound  $r^2$ , and for  $\delta \in (0, 1)$ , the following holds with probability at least  $1 - \delta$  over sample  $S_{T_2}$*

$$\begin{aligned} \hat{\mathcal{L}}_{T_2}^{[2]} - \mathcal{L}_{T_2}^{[2]}(\mathbf{w}_*^{[2]}) &\leq \frac{4r^2}{\lambda\sqrt{T_2}} \left( \frac{\mathcal{E}}{r} + \sqrt{\frac{\mathcal{E}}{r}} \right)^{1/2} \\ &+ \frac{c_2 r^2}{\lambda\sqrt{T_2}} \left[ \left( \frac{1}{\sqrt{T_1}} + \frac{1}{\sqrt{T_e}} + \frac{1}{\sqrt{T_2}} \right) \sqrt{\ln \frac{6}{\delta}} \right]^{1/2}, \end{aligned}$$

where  $\mathcal{E} = \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})$  and  $c_2 > 0$  is some constant.

Theorem 4.3 gives the convergence analysis for Algorithm 2. We obtain tighter bound as for smaller KOM discrepancy, i.e., a closer relationship between old and new feature spaces. It is also useful to exploit information and model from old feature space theoretically, because of tighter bounds as for larger  $T_1$  and  $T_e$ . The detailed proof is given in Appendix C.3, which is motivated from regret analysis (Hazan et al., 2016), generalization bounds via Rademacher complexity (Bartlett & Mendelson, 2002), some online-to-batch conversion techniques (Cesa-Bianchi et al., 2004).

Denote by the cumulative loss for  $\mathbf{w}^l$

$$\hat{\mathcal{L}}_{T_2}^l = \frac{1}{T_2} \sum_{t=T_1+T_e+1}^{T_1+T_e+T_2} \ell_t^l(\mathbf{w}_t^l),$$

with  $\ell_t^l(\mathbf{w}) = \max\{0, 1 - y_t \langle \mathbf{w}, \mathbf{z}^l(\mathbf{x}_t^l) \rangle\} + \lambda \|\mathbf{w}\|_2^2 / 2$ . We also analyze the cumulative error rate for Algorithm 2.

**Theorem 4.4.** *For kernels  $\mathcal{K}^{[2]}$  and  $\mathcal{K}^l$  with bound  $r^2$  and for parameter  $\gamma = \sqrt{\ln 2 / ((1 + 3r^2 / 2\lambda) T_2)}$ , we have*

$$\sum_{t=T_1+T_e+1}^{T_1+T_e+T_2} \frac{\mathbb{I}[y_t h_t(\mathbf{x}_t^{[2]}) \leq 0]}{T_2} \leq \min \{ \hat{\mathcal{L}}_{T_2}^{[2]}, \hat{\mathcal{L}}_{T_2}^l \} + \sqrt{\frac{2 \ln 2}{T_2}}$$

where  $\mathbb{I}[\cdot]$  is the indicator function, which returns 1 if the argument is true, and 0 otherwise.

Theorem 4.4 shows that the cumulative error rate of our OPFES method converges to the minimum of the cumulative loss of two classifiers. The detailed proof is presented in Appendix C.4, and the basic idea is to construct a potential function and apply Hoeffding’s lemma (Hoeffding, 1963).

## 5. Empirical Study

We conduct experiments on 20 datasets<sup>2</sup>, and the details are summarized in Table 1. Most of the datasets have been well-studied for previous feature evolvable learning, and all features have been scaled to  $[0, 1]$ . We compare our OPFES with state-of-the-art methods on feature evolvable learning.

- align-FESL: Feature evolvable method of random features and kernel alignment for feature and label relationships (Sinha & Duchi, 2016).
- rff-ROGD: Feature evolvable method of random features and  $\ell_2$  distance for feature relationships (Lu et al., 2016);
- rff-FESL: Online ensemble of rff-ROGD and random feature models learned from scratch (Hou et al., 2021);
- ker-ROGD: Feature evolvable method with kernel model and  $\ell_2$  distance for feature relationships (Hou et al., 2021);
- ker-FESL: Online ensemble of ker-ROGD and kernel models learned from scratch (Hou et al., 2021);
- lin-ROGD: Feature evolvable method with linear models and  $\ell_2$  distance for feature relationships (Hou et al., 2017);
- lin-FESL: Online ensemble of lin-ROGD and a linear model learned from scratch (Hou et al., 2017);
- OCDS: Capricious streaming method with a linear model via generative graphical model (He et al., 2021b).

For each dataset, we randomly split the feature space into old feature space  $\mathcal{X}^{[1]}$  and new feature space  $\mathcal{X}^{[2]}$  with almost equal number of features, following (Gu et al., 2022; Ni et al., 2024). We set  $T_e = 1000$  for datasets with a size larger than 10000; otherwise, set  $T_e$  as 10% of the dataset’s size. We also set  $T_1$  and  $T_2$  as half of the amount of dataset

<sup>2</sup>Downloaded from [OpenML](#) and [UCI datasets repository](#)

Table 2. Cumulative error rate (CER) evaluation of our OPFES and compared methods (mean $\pm$ std). ●/○ indicates that our OPFES is significantly better/worse than the corresponding algorithms (pairwise t-tests at 95% significance level).

Dataset	Our OPFES	align-FESL	rff-FESL	rff-ROGD	ker-FESL	ker-ROGD	lin-FESL	lin-ROGD	OCDS
jungle	.0097 $\pm$ .0047	.0099 $\pm$ .0028	.0161 $\pm$ .0035●	.0246 $\pm$ .0069●	.0276 $\pm$ .0055●	.0329 $\pm$ .0061●	.1084 $\pm$ .0152●	.1471 $\pm$ .0144●	.1106 $\pm$ .0138●
splice	.3070 $\pm$ .0079	.3126 $\pm$ .0136	.3234 $\pm$ .0087●	.3662 $\pm$ .0215●	.4192 $\pm$ .0160●	.4240 $\pm$ .0188●	.3447 $\pm$ .0097●	.4307 $\pm$ .0213●	.3547 $\pm$ .0156●
bioresponse	.2763 $\pm$ .0117	.2921 $\pm$ .0106●	.3051 $\pm$ .0137●	.4285 $\pm$ .0192●	.3690 $\pm$ .0095●	.4454 $\pm$ .0116●	.2938 $\pm$ .0093●	.3684 $\pm$ .0102●	.2951 $\pm$ .0112●
christine	.3192 $\pm$ .0095	.3205 $\pm$ .0090	.3316 $\pm$ .0108●	.3503 $\pm$ .0092●	.3858 $\pm$ .0096●	.4506 $\pm$ .0117●	.3443 $\pm$ .0098●	.3439 $\pm$ .0098●	.3663 $\pm$ .0116●
svmguide1	.1614 $\pm$ .0052	.1617 $\pm$ .0056	.1632 $\pm$ .0054●	.2295 $\pm$ .0107●	.1900 $\pm$ .0061●	.2316 $\pm$ .0050●	.2399 $\pm$ .0062●	.2451 $\pm$ .0102●	.2442 $\pm$ .0070●
usps	.1684 $\pm$ .0044	.1875 $\pm$ .0073●	.1658 $\pm$ .0051	.2184 $\pm$ .0061●	.2267 $\pm$ .0084●	.2857 $\pm$ .0063●	.2654 $\pm$ .0081●	.2839 $\pm$ .0073●	.2746 $\pm$ .0085●
aileron	.1963 $\pm$ .0034	.2144 $\pm$ .0081●	.2139 $\pm$ .0059●	.2344 $\pm$ .0098●	.2531 $\pm$ .0076●	.2523 $\pm$ .0076●	.2466 $\pm$ .0066●	.2465 $\pm$ .0066●	.3026 $\pm$ .0047●
pol	.0654 $\pm$ .0036	.0692 $\pm$ .0044●	.0686 $\pm$ .0023●	.0807 $\pm$ .0028●	.0865 $\pm$ .0023●	.0956 $\pm$ .0034●	.1484 $\pm$ .0035●	.1655 $\pm$ .0041●	.1515 $\pm$ .0038●
elevators	.2422 $\pm$ .0039	.2419 $\pm$ .0038	.2467 $\pm$ .0037●	.2619 $\pm$ .0045●	.2963 $\pm$ .0042●	.3003 $\pm$ .0051●	.3073 $\pm$ .0043●	.3045 $\pm$ .0040●	.3073 $\pm$ .0042●
magic	.2154 $\pm$ .0039	.2206 $\pm$ .0039●	.2121 $\pm$ .0046	.2434 $\pm$ .0057●	.2656 $\pm$ .0033●	.3119 $\pm$ .0074●	.2535 $\pm$ .0040●	.2988 $\pm$ .0073●	.2554 $\pm$ .0045●
letter	.1354 $\pm$ .0043	.1568 $\pm$ .0086●	.1557 $\pm$ .0037●	.2311 $\pm$ .0067●	.3139 $\pm$ .0060●	.3373 $\pm$ .0076●	.3380 $\pm$ .0038●	.3565 $\pm$ .0071●	.3390 $\pm$ .0034●
house	.1849 $\pm$ .0040	.1927 $\pm$ .0030●	.1894 $\pm$ .0043●	.2001 $\pm$ .0093●	.2598 $\pm$ .0112●	.2597 $\pm$ .0113●	.2623 $\pm$ .0084●	.2658 $\pm$ .0120●	.2853 $\pm$ .0037●
nomao	.0646 $\pm$ .0026	.0680 $\pm$ .0024●	.0778 $\pm$ .0017●	.0845 $\pm$ .0041●	.1302 $\pm$ .0034●	.1355 $\pm$ .0032●	.0860 $\pm$ .0019●	.1107 $\pm$ .0039●	.0882 $\pm$ .0023●
adult	.1875 $\pm$ .0023	.1942 $\pm$ .0027●	.1906 $\pm$ .0021●	.1932 $\pm$ .0029●	.2277 $\pm$ .0019●	.2218 $\pm$ .0031●	.2050 $\pm$ .0033●	.2036 $\pm$ .0042●	.2303 $\pm$ .0026●
acoustic	.3074 $\pm$ .0024	.3227 $\pm$ .0036●	.2967 $\pm$ .0045○	.2977 $\pm$ .0043○	.4168 $\pm$ .0073●	.4107 $\pm$ .0072●	.4321 $\pm$ .0079●	.4317 $\pm$ .0075●	.4668 $\pm$ .0022●
runwalk	.2602 $\pm$ .0033	.2890 $\pm$ .0055●	.2578 $\pm$ .0016○	.3496 $\pm$ .0130●	.3558 $\pm$ .0021●	.4355 $\pm$ .0061●	.4945 $\pm$ .0021●	.4963 $\pm$ .0033●	.4972 $\pm$ .0027●
higgs	.3946 $\pm$ .0045	.4135 $\pm$ .0061●	.3803 $\pm$ .0074○	.3807 $\pm$ .0080○	.4577 $\pm$ .0054●	.4570 $\pm$ .0055●	.4309 $\pm$ .0028●	.4481 $\pm$ .0139●	.4366 $\pm$ .0021●
miniboone	.1602 $\pm$ .0036	.2804 $\pm$ .0011●	.1729 $\pm$ .0029●	.1603 $\pm$ .0029	.2488 $\pm$ .0047●	.2484 $\pm$ .0047●	.2384 $\pm$ .0039●	.2384 $\pm$ .0039●	.2803 $\pm$ .0011●
ijcnn1	.0616 $\pm$ .0115	.0746 $\pm$ .0038●	.0673 $\pm$ .0028●	.0747 $\pm$ .0083●	.0957 $\pm$ .0009●	.0957 $\pm$ .0009●	.0951 $\pm$ .0007●	.0957 $\pm$ .0009●	.0957 $\pm$ .0009●
covtype	.3782 $\pm$ .0008	.3813 $\pm$ .0025●	.3783 $\pm$ .0009	.3795 $\pm$ .0012●	.4095 $\pm$ .0025●	.4093 $\pm$ .0025●	.3790 $\pm$ .0007●	.3920 $\pm$ .0034●	.3792 $\pm$ .0007●
Win/Tie/Loss		15/5/0	14/3/3	17/1/2	20/0/0	20/0/0	20/0/0	20/0/0	20/0/0

size subtracting  $T_e$ . For ker-ROGD and ker-FESL, we set the buffer size to 10% of dataset size and consider reservoir sampling as done by Hou et al. (2021).

For rff-ROGD, rff-FESL, align-FESL and our OPFES, we fix the dimensionality of random Fourier feature as 1000. In the previous stage, we employ Gaussian kernels with widths in  $2^{[-6:6]}$  for all methods. For OPFES, we set  $T_M = 1000$  of the optimal stepsize from Theorem 3.7. The stepsize  $\tau_t$  is constrained within  $10^{[-4:2]}/\sqrt{t}$ , and the regularization parameter  $\lambda$  is selected from  $10^{[-10:1]}$ . For OCDS,  $\alpha$  and  $\beta$  are chosen from  $10^{[-5:0]}$  by cross validations.

The performance of the compared methods is evaluated by 50 times on each dataset with random partitions and random ordering in the previous, evolving and current stages, where the cumulative error rate (CER) is obtained by averaging over these 50 runs, as summarized in Table 2.

It is observable that, from Table 2, our OPFES method takes significantly better performance than three linear methods lin-ROGD, lin-FESL, and OCDS, since these methods rely on simple linear classifiers and  $\ell_2$  distance to characterize feature relationships. Our OPFES outperforms ker-ROGD, ker-FESL, rff-ROGD most times because of exploration on feature relationships via our KOM discrepancy, while other methods fix Gaussian or Mahalanobis kernels.

Our OPFES is also better than align-FESL, since the KOM discrepancy is more effective in capturing some feature relationships than kernel alignment, as shown in Lemma 3.5.

Table 3. Ablation studies for our OPFES method (mean $\pm$ std): (i) OPFES without KOM discrepancy; (ii) OPFES without ideal kernel; (iii) OPFES without initialization from previous model.

Datasets	OPFES	(i)	(ii)	(iii)
Pol	.0654 $\pm$ .0036	.1290 $\pm$ .0044●	.0674 $\pm$ .0059	.1418 $\pm$ .0072●
House	.1849 $\pm$ .0033	.1976 $\pm$ .0032●	.1957 $\pm$ .0034●	.2350 $\pm$ .0033●
Nomao	.0646 $\pm$ .0026	.0933 $\pm$ .0028●	.0750 $\pm$ .0047●	.1875 $\pm$ .0056●
Adult	.1875 $\pm$ .0023	.2013 $\pm$ .0035●	.1988 $\pm$ .0043●	.2133 $\pm$ .0064●

Our OPFES achieves better and comparable performance in contrast to rff-FESL, except for datasets acoustic, runwalk, and higgs. This is partially because of the class-imbalance problem on the three datasets, which results in the hardness of learning ideal kernel  $\mathcal{K}^l$  for label information and the degrade of learning performance of base learners  $w^l$ .

Table 3 presents the ablation study of OPFES to verify the effectiveness of KOM discrepancy, ideal kernel and initializations from previous models. Due to pages limit, we consider four datasets Pol, House, Nomao and Adult, while the trends are similar for other datasets. It is clear that the performance of our OPFES will decrease drastically without the consideration of KOM discrepancy and good initializations. Ideal kernel takes limited improvement from label correlation, in particularly for dataset Pol, where  $\mathcal{K}^{[2]}$  possibly takes comparable performance to  $\mathcal{K}^l$ .

Figure 4 presents the convergence analysis of cumulative error rate for our OPFES. It is evident that our OPFES takes faster convergence than other methods in the current stage,

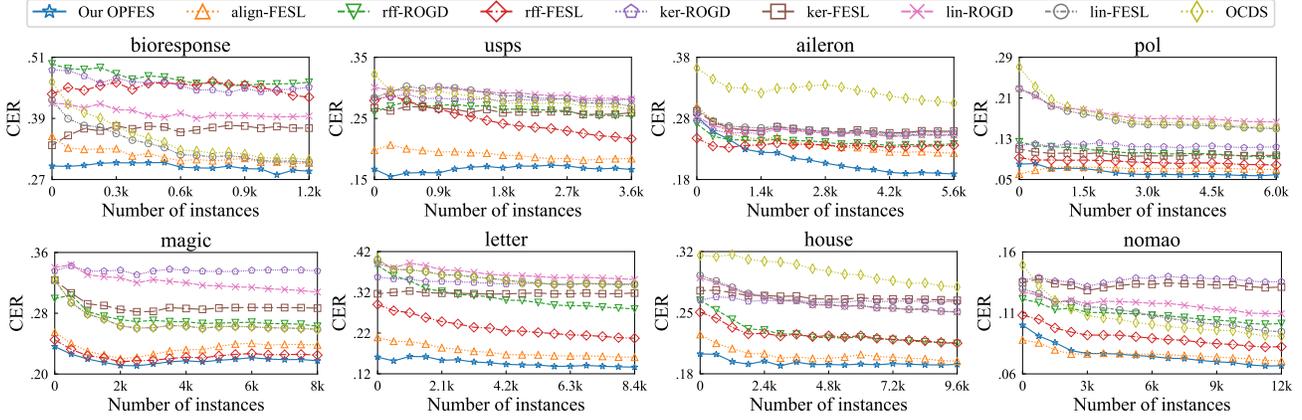


Figure 4. Cumulative error rate (CER) versus the number of instances in the current stage for our OPFES and compared methods. The lower the curve, the faster the convergence.

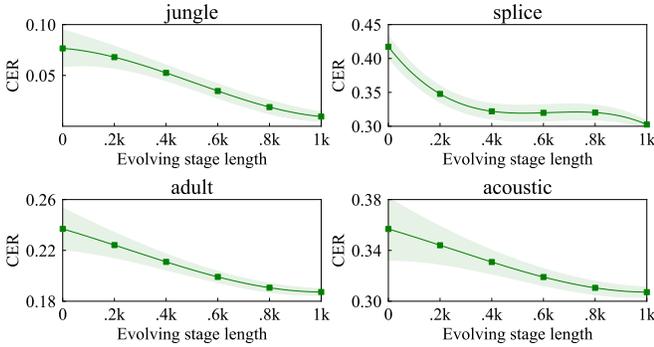


Figure 5. Cumulative error rate (CER) versus the length of evolving stage for our OPFES.

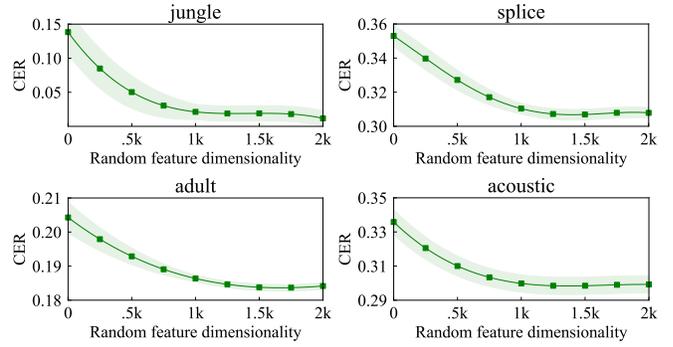


Figure 6. Cumulative error rate (CER) versus the dimensionality of random Fourier features for our OPFES.

partially because our OPFES can exploit feature relationships between two feature spaces and label information. It is observable that our method obtains lower cumulative error rates in the beginning of the current stage, indicating the importance of good model initializations from the evolving stage, which is consistent with Theorem 4.3.

Figure 5 empirically illustrates the influences of the evolving stage length  $T_e$  in Algorithm 2. It is evident that larger  $T_e$  could yield better performance, because of more available instances in training to capture feature relationships and label information, which is consistent with our theoretical analysis in Theorem 4.3. We finally present the influence of dimensionality of random Fourier features for our OPFES on 4 datasets in Figure 6, and trends are similar on other datasets. It is clear that our OPFES obtains stable performance if we set the dimensionality  $d$  larger than 1000, while smaller dimensionality could yield heavy information loss.

## 6. Conclusion

This work focuses on two fundamental problems on feature evolvable learning. We propose the Kernel Ortho-Mapping (KOM) discrepancy to characterize intrinsic relationships

between two feature spaces via kernel functions, and then theoretically correlate it with optimal classifiers learned from different feature spaces. Based on this discrepancy, we develop one-pass algorithm for feature evolvable learning without storing the entire or partial training data. We verify the effectiveness of our proposed OPFES both theoretically and empirically. An interesting future work is to apply our KOM discrepancy to deep learning via neural tangent kernel, and exploit more effective tools to characterize the feature relationships for feature evolvable learning.

## Acknowledgements

The authors want to thank the reviewers for their helpful comments and suggestions. This research was supported by National Key R&D Program of China (2021ZD0112802) and NSFC (62376119).

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning especially Feature Evolvable Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

- Alpaydin, E. *Machine Learning*. MIT Press, 2021.
- Banerjee, A., Guo, X., and Wang, H. On the optimality of conditional expectation as a bregman predictor. *IEEE Transactions on Information Theory*, 51(7):2664–2669, 2005.
- Bartlett, P. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- Beyazit, E., Alagurajah, J., and Wu, X. Online learning from data streams with varying feature spaces. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, pp. 3232–3239, Honolulu, HI, 2019.
- Bhatia, R. *Matrix Analysis*. Springer Science & Business Media, 2013.
- Boyd, S. and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, 2004.
- Bubeck, S. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8: 231–357, 2015.
- Cesa-Bianchi, N., Conconi, A., and Gentile, C. On the generalization ability of on-line learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050–2057, 2004.
- Chen, Y. and Liu, S. A novel learning method for feature evolvable streams. *Evolving Systems*, 15:1–19, 2024.
- Cortes, C., Mohri, M., and Rostamizadeh, A. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13(1):795–828, 2012.
- Cramer, K. and Singer, Y. On the learnability and design of output codes for multiclass problems. *Machine learning*, 47:201–233, 2002.
- Cristianini, N., Kandola, J., Elisseeff, A., and Shawe-Taylor, J. On kernel-target alignment. In *Advances in Neural Information Processing Systems 14*, pp. 367–373, Vancouver, Canada, 2001.
- Dong, J., Cong, Y., Sun, G., Zhang, T., Tang, X., and Xu, X. Evolving metric learning for incremental and decremental features. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(4):2290–2302, 2022.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016.
- Gu, S., Qian, Y., and Hou, C. Incremental feature spaces learning with label scarcity. *ACM Transactions on Knowledge Discovery from Data*, 16(6):1–26, 2022.
- Hazan, E. et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- He, Y., Dong, J., Hou, B.-J., Wang, Y., and Wang, F. Online learning in variable feature spaces with mixed data. In *Proceedings of IEEE International Conference on Data Mining*, pp. 181–190, Auckland, New Zealand, 2021a.
- He, Y., Wu, B., Wu, D., Beyazit, E., Chen, S., and Wu, X. Toward mining capricious data streams: A generative approach. *IEEE Transactions on Neural Networks and Learning Systems*, 32(3):1228–1240, 2021b.
- Heo, B., Kim, J., Yun, S., Park, H., Kwak, N., and Choi, J. Y. A comprehensive overhaul of feature distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1921–1930, Seoul, Korea, 2019.
- Hoeffding, W. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- Hong, S. and Chae, J. Distributed online learning with multiple kernels. *IEEE Transactions on Neural Networks and Learning Systems*, 34(3):1263–1277, 2021.
- Hou, B.-J., Zhang, L., and Zhou, Z.-H. Learning with feature evolvable streams. In *Advances in Neural Information Processing Systems 30*, pp. 1417–1427, Long Beach, CA, 2017.
- Hou, B.-J., Yan, Y.-H., Zhao, P., and Zhou, Z.-H. Storage fit learning with feature evolvable streams. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, pp. 7729–7736, Virtual Event, 2021.
- Hou, B.-J., Zhang, L., and Zhou, Z.-H. Prediction with unpredictable feature evolution. *IEEE Transactions on Neural Networks and Learning Systems*, 33(10):5706–5715, 2022.
- Hou, C. and Zhou, Z.-H. One-pass learning with incremental and decremental features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(11): 2776–2792, 2018.
- Hou, C., Zeng, L.-L., and Hu, D. Safe classification with augmented features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(9):2176–2192, 2019.
- Hou, C., Fan, R., Zeng, L.-L., and Hu, D. Adaptive feature selection with augmented attributes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):9306–9324, 2023.
- Jin, R., Hoi, S. C., and Yang, T. Online multiple kernel learning: Algorithms and mistake bounds. In *Proceedings of*

- the 21st International Conference on Algorithmic Learning Theory, pp. 390–404, Canberra, Australia, 2010.
- Kwok, J. T. et al. Learning with idealized kernels. In *Proceedings of the 20th International Conference on Machine Learning*, pp. 400–407, Washington, DC, 2003.
- Lian, H., Wu, D., Hou, B.-J., Wu, J., and He, Y. Online learning from evolving feature spaces with deep variational models. *IEEE Transactions on Knowledge and Data Engineering*, 36(8):4144–4162, 2023.
- Liu, X. Incomplete multiple kernel alignment maximization for clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(3):1412–1424, 2024.
- Lu, J., Hoi, S. C., Wang, J., Zhao, P., and Liu, Z.-Y. Large scale online kernel learning. *Journal of Machine Learning Research*, 17(47):1–43, 2016.
- Marukatat, S. Kernel matrix decomposition via empirical kernel map. *Pattern Recognition Letters*, 77:50–57, 2016.
- McDiarmid, C. et al. On the method of bounded differences. *Surveys in Combinatorics*, 141(1):148–188, 1989.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Foundations of Machine Learning*. MIT Press, 2018.
- Murphy, K. P. *Machine Learning: A Probabilistic Perspective*. MIT Press, 2012.
- Ni, H., Gu, S., Fan, R., and Hou, C. Feature incremental learning with causality. *Pattern Recognition*, 146:110033, 2024.
- Orabona, F., Keshet, J., and Caputo, B. The projectron: A bounded kernel-based perceptron. In *Proceedings of the 25th International Conference on Machine Learning*, pp. 720–727, Helsinki, Finland, 2008.
- Pennington, J., Yu, F. X. X., and Kumar, S. Spherical random features for polynomial kernels. In *Advances in Neural Information Processing Systems 28*, pp. 1846–1854, Montreal, Canada, 2015.
- Rahimi, A. and Recht, B. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems 20*, pp. 1177–1184, Vancouver, Canada, 2007.
- Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *Proceedings of 3rd International Conference on Learning Representations*, pp. 1–14, San Diego, CA, 2015.
- Sajedi, R. and Razzazi, M. Data stream classification in dynamic feature space using feature mapping. *The Journal of Supercomputing*, 80(9):12043–12061, 2024.
- Schölkopf, B. and Smola, A. J. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- Schreckenberger, C., He, Y., Lütke, S., Bartelt, C., and Stuckenschmidt, H. Online random feature forests for learning in varying feature spaces. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence*, pp. 4587–4595, Washington, DC, 2023.
- Shalev-Shwartz, S. and Ben-David, S. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- Shen, Y., Chen, T., and Giannakis, G. B. Random feature-based online multi-kernel learning in environments with unknown dynamics. *Journal of Machine Learning Research*, 20(22):1–36, 2019.
- Sinha, A. and Duchi, J. C. Learning kernels with random features. In *Advances in Neural Information Processing Systems 29*, pp. 1298–1306, Barcelona, Spain, 2016.
- Valiant, L. G. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- Vershynin, R. *High-dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- Wacker, J., Kanagawa, M., and Filippone, M. Improved random features for dot product kernels. *Journal of Machine Learning Research*, 25(235):1–75, 2024.
- Watson, G. A. Characterization of subdifferential of some matrix norms. *Linear Algebra Appl*, 170(1):33–45, 1992.
- Wu, D., Zhuo, S., Wang, Y., Chen, Z., and He, Y. Online semi-supervised learning with mix-typed streaming features. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence*, pp. 4720–4728, Washington, DC, 2023.
- Zhang, Q., Zhang, P., Long, G., Ding, W., Zhang, C., and Wu, X. Online learning from trapezoidal data streams. *IEEE Transactions on Knowledge and Data Engineering*, 28(10):2709–2723, 2016.
- Zhang, Z.-Y., Zhao, P., Jiang, Y., and Zhou, Z.-H. Learning with feature and distribution evolvable streams. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 11317–11327, Virtual Event, 2020.
- Zhou, Z., Shen, Y., Shao, S., Gong, L., and Lin, S. Rethinking centered kernel alignment in knowledge distillation. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pp. 5680–5688, Jeju, South Korea, 2024.

## A. Proofs for Kernel Ortho-Mapping Discrepancy

### A.1. Proof of Lemma 3.2

For Frobenius norm and orthogonal matrix  $\mathbf{U} \in \mathcal{U}_n$ , we have

$$\begin{aligned} \left\| \sqrt{\mathbf{K}^{[1]}}\mathbf{U} - \sqrt{\mathbf{K}^{[2]}} \right\|_F^2 &= \text{Tr} \left( \sqrt{\mathbf{K}^{[1]}}\mathbf{U} - \sqrt{\mathbf{K}^{[2]}} \right) \left( \sqrt{\mathbf{K}^{[1]}}\mathbf{U} - \sqrt{\mathbf{K}^{[2]}} \right)^\top \\ &= \text{Tr}(\mathbf{K}^{[1]}) + \text{Tr}(\mathbf{K}^{[2]}) - 2 \text{Tr} \left( \mathbf{U}^\top \sqrt{\mathbf{K}^{[2]}} \sqrt{\mathbf{K}^{[1]}} \right). \end{aligned} \quad (13)$$

The minimization of Kernel Ortho-Mapping (KOM) (in Definition 3.1) is equivalent to the following optimization:

$$\max_{\mathbf{U} \in \mathcal{U}_n} \left\{ \text{Tr} \left( \mathbf{U}^\top \sqrt{\mathbf{K}^{[2]}} \sqrt{\mathbf{K}^{[1]}} \right) \right\}. \quad (14)$$

Denote by  $\mathbf{X} = \sqrt{\mathbf{K}^{[2]}}\sqrt{\mathbf{K}^{[1]}}$  and  $|\mathbf{X}| = \sqrt{\mathbf{X}^\top \mathbf{X}} = \sqrt{\sqrt{\mathbf{K}^{[1]}}\mathbf{K}^{[2]}\sqrt{\mathbf{K}^{[1]}}}$ . There exists  $\mathbf{V} \in \mathcal{U}_n$  such that  $\mathbf{X} = \mathbf{V}|\mathbf{X}|$ , which is a polar decomposition of  $\mathbf{X}$ . We then have

$$\text{Tr} \left( \mathbf{U}^\top \sqrt{\mathbf{K}^{[2]}} \sqrt{\mathbf{K}^{[1]}} \right) = \text{Tr} \left( \mathbf{U}^\top \mathbf{X} \right) = \text{Tr} \left( \mathbf{U}^\top \mathbf{V} |\mathbf{X}| \right).$$

Denote by  $\mathbf{W} = \mathbf{U}^\top \mathbf{V}$ . We take the eigen-decomposition  $|\mathbf{X}| = \mathbf{P}\mathbf{D}\mathbf{P}^\top$  for orthogonal  $\mathbf{P}$  and diagonal matrix  $\mathbf{D}$  with non-negative elements from the semi-positive definiteness of  $|\mathbf{X}|$ . From the cyclic invariance of matrix trace, we have

$$\text{Tr} \left( \mathbf{U}^\top \sqrt{\mathbf{K}^{[2]}} \sqrt{\mathbf{K}^{[1]}} \right) = \text{Tr} \left( \mathbf{W}\mathbf{P}\mathbf{D}\mathbf{P}^\top \right) = \text{Tr} \left( \mathbf{D}\mathbf{P}^\top \mathbf{W}\mathbf{P} \right),$$

Denote by  $\hat{\mathbf{W}} = \mathbf{P}^\top \mathbf{W}\mathbf{P}$  another unitary matrix in  $\mathcal{U}_n$ . The optimization problem in Eqn. (14) can be solved by

$$\max_{\hat{\mathbf{W}} \in \mathcal{U}_n} \left\{ \text{Tr} \left( \mathbf{D}\hat{\mathbf{W}} \right) \right\} = \max_{\hat{\mathbf{W}} \in \mathcal{U}_n} \left\{ \sum_{i=1}^n \mathbf{D}_{ii} \hat{\mathbf{W}}_{ii} \right\} = \text{Tr}(\mathbf{D}),$$

where the last equality holds from  $\hat{\mathbf{W}} = \mathbf{I}_n$ . From nuclear norm and matrix trace of positive semi-definite matrix, we have

$$\text{Tr}(\mathbf{D}) = \text{Tr} \left( \sqrt{\sqrt{\mathbf{K}^{[1]}}\mathbf{K}^{[2]}\sqrt{\mathbf{K}^{[1]}}} \right) = \sum_{i=1}^n \sqrt{\sigma_i \left( \sqrt{\mathbf{K}^{[1]}}\mathbf{K}^{[2]}\sqrt{\mathbf{K}^{[1]}} \right)} = \sum_{i=1}^n \sigma_i \left( \sqrt{\mathbf{K}^{[1]}}\sqrt{\mathbf{K}^{[2]}} \right) = \left\| \sqrt{\mathbf{K}^{[1]}}\sqrt{\mathbf{K}^{[2]}} \right\|_*,$$

which completes the proof by combining with Eqn. (13).  $\square$

### A.2. Proof for Theorem 3.3.

We begin with the empirical feature mapping as follows:

**Lemma A.1.** *Let  $\mathbf{K}$  be the Gram matrix w.r.t. kernel  $\mathcal{K}$  and sample  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , and consider eigen decomposition  $\mathbf{K} = \mathbf{S}\mathbf{D}\mathbf{S}^\top$ . For kernel classifier  $h(\mathbf{x}) = \sum_{j=1}^n \alpha_j \mathcal{K}(\mathbf{x}_j, \mathbf{x})$ , we have*

$$h(\mathbf{x}_i) = \left\langle \mathbf{w}, \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \mathbf{e}_i \right\rangle \quad \text{with} \quad \mathbf{w} = \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \sum_{j=1}^n \alpha_j \mathbf{e}_j,$$

where  $\mathbf{e}_i$  denotes a unit vector of the  $i$ -th element being 1.

*Proof.* From eigen decomposition  $\mathbf{K} = \mathbf{S}\mathbf{D}\mathbf{S}^\top$ , we have

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{K}_{ij} = (\mathbf{S}\mathbf{D}\mathbf{S}^\top)_{ij} = \left\langle \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \mathbf{e}_i, \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \mathbf{e}_j \right\rangle,$$

and this gives one data-dependent feature mapping of  $\mathcal{K}$  as  $\hat{\varphi} : \mathbf{x}_i \mapsto \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \mathbf{e}_i$  (Schölkopf & Smola, 2002). For kernel classifier  $h(\mathbf{x}) = \sum_{j=1}^n \alpha_{1,j} \mathcal{K}(\mathbf{x}_j, \mathbf{x})$ , we have

$$h(\mathbf{x}_i) = \sum_{j=1}^n \alpha_{1,j} \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{j=1}^n \alpha_{1,j} \langle \hat{\varphi}(\mathbf{x}_j), \hat{\varphi}(\mathbf{x}_i) \rangle = \left\langle \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \sum_{j=1}^n \alpha_{1,j} \mathbf{e}_j, \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \mathbf{e}_i \right\rangle = \left\langle \mathbf{w}, \mathbf{S}\sqrt{\mathbf{D}}\mathbf{S}^\top \mathbf{e}_i \right\rangle,$$

which completes the proof.  $\square$

**Proof of Theorem 3.3.** From Lemma A.1, the kernel learning problem in Eqn. (3) can be equivalently linearized as

$$\mathbf{w}_*^{[k]} \in \arg \min_{\mathbf{w}^{[k]} \in \mathbb{R}^n} \left\{ \hat{R}^{[k]}(\mathbf{w}^{[k]}) + \frac{\lambda}{2} \|\mathbf{w}^{[k]}\|_2^2 \right\}, \quad \text{for } k = 1, 2, \quad (15)$$

where

$$\hat{R}^{[k]}(\mathbf{w}^{[k]}) = \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - y_i \langle \mathbf{w}^{[k]}, \hat{\varphi}^{[k]}(\mathbf{x}_i^{[k]}) \rangle \right\},$$

and  $\hat{\varphi}^{[k]}$  are the empirical kernel mappings of  $\mathcal{K}^{[k]}$ . For two classifiers and from strong convexity of Eqn. (15), we have

$$\begin{aligned} \frac{\lambda}{2} \|\mathbf{w}_*^{[1]}\|_2^2 + \hat{R}^{[1]}(\mathbf{w}_*^{[1]}) + \frac{\lambda}{2} \|\mathbf{w}_*^{[1]} - \mathbf{w}_*^{[2]}\|_2^2 &\leq \hat{R}^{[1]}(\mathbf{w}_*^{[2]}) + \frac{\lambda}{2} \|\mathbf{w}_*^{[2]}\|_2^2, \\ \frac{\lambda}{2} \|\mathbf{w}_*^{[2]}\|_2^2 + \hat{R}^{[2]}(\mathbf{w}_*^{[2]}) + \frac{\lambda}{2} \|\mathbf{w}_*^{[2]} - \mathbf{w}_*^{[1]}\|_2^2 &\leq \hat{R}^{[2]}(\mathbf{w}_*^{[1]}) + \frac{\lambda}{2} \|\mathbf{w}_*^{[1]}\|_2^2. \end{aligned}$$

This holds that, from 1-Lipschitz continuous hinge loss and Cauchy-Schwarz inequality,

$$\begin{aligned} \|\mathbf{w}_*^{[1]} - \mathbf{w}_*^{[2]}\|_2^2 &\leq \frac{1}{\lambda} \left( (\hat{R}^{[1]}(\mathbf{w}_*^{[2]}) - \hat{R}^{[2]}(\mathbf{w}_*^{[2]})) + (\hat{R}^{[2]}(\mathbf{w}_*^{[1]}) - \hat{R}^{[1]}(\mathbf{w}_*^{[1]})) \right) \\ &= \frac{1}{n\lambda} \sum_{i=1}^n \left( \ell \left( \langle \mathbf{w}_*^{[2]}, \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \rangle, y_i \right) - \ell \left( \langle \mathbf{w}_*^{[2]}, \hat{\varphi}_2(\mathbf{x}_i^{[2]}) \rangle, y_i \right) \right) \\ &\quad + \frac{1}{n\lambda} \sum_{i=1}^n \left( \ell \left( \langle \mathbf{w}_*^{[1]}, \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \rangle, y_i \right) - \ell \left( \langle \mathbf{w}_*^{[1]}, \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \rangle, y_i \right) \right) \\ &\leq \frac{1}{n\lambda} \left( \|\mathbf{w}_*^{[1]}\|_2 + \|\mathbf{w}_*^{[2]}\|_2 \right) \sum_{i=1}^n \left\| \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2. \end{aligned} \quad (16)$$

For  $\|\mathbf{w}_*^{[1]}\|_2 + \|\mathbf{w}_*^{[2]}\|_2$ , we have the following constraints on  $\mathbf{w}_*^{[k]}$  from the KKT condition of Eqn. (15)

$$\mathbf{w}_*^{[k]} \in \left\{ -\frac{1}{n\lambda} \sum_{i=1}^n \mathbf{g}_i : \mathbf{g}_i \in \partial \max \left\{ 0, 1 - y_i \langle \mathbf{w}_*^{[k]}, \hat{\varphi}^{[k]}(\mathbf{x}_i^{[k]}) \rangle \right\}, i \in [n] \right\} \quad \text{for } k = 1, 2,$$

where  $\partial(\cdot)$  is the sub-gradient operator. We can upper bound  $\|\mathbf{w}^{[1]*}\|_2 + \|\mathbf{w}^{[2]*}\|_2 \leq 2r/\lambda$  from

$$\|\mathbf{w}_*^{[k]}\|_2 \leq \max \left\{ \left\| -\frac{\sum_{i=1}^n \mathbf{g}_i}{n\lambda} \right\|_2 : \mathbf{g}_i \in \partial \max \left\{ 0, 1 - y_i \langle \mathbf{w}_*^{[k]}, \hat{\varphi}_k(\mathbf{x}_i^{[k]}) \rangle \right\} \right\} \leq \frac{1}{n\lambda} \sum_{i=1}^n \left\| y_i \hat{\varphi}^{[k]}(\mathbf{x}_i^{[k]}) \right\|_2 \leq \frac{r}{\lambda}.$$

It remains to bound  $\sum_{i=1}^n \|\hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]})\|_2$ . We observe that the empirical kernel mapping is unitarily invariant in Lemma A.1, i.e., we have  $\hat{\varphi}(\mathbf{x}_i)^\top \hat{\varphi}(\mathbf{x}_j) = (\mathbf{U}\hat{\varphi}(\mathbf{x}_i))^\top (\mathbf{U}\hat{\varphi}(\mathbf{x}_j))$  for  $\mathbf{U} \in \mathcal{U}_n$  and  $i, j \in [n]$ , and the  $i$ -th column of  $\mathbf{U}\sqrt{\mathbf{K}}$  is also a legal empirical kernel mapping of each  $\mathbf{x}_i$ . This follows that, for KOM discrepancy,

$$\sum_{i=1}^n \left\| \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2 \leq \sqrt{n} \min_{\mathbf{U}_1, \mathbf{U}_2 \in \mathcal{U}_n} \left\{ \left\| \sqrt{\mathbf{K}^{[1]}} \mathbf{U}_1 - \sqrt{\mathbf{K}^{[2]}} \mathbf{U}_2 \right\|_F \right\} = \sqrt{n} \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}),$$

from the unitary invariance of Frobenius norm. From Eqn. (16), we have

$$\left\| \mathbf{w}_*^{[1]} - \mathbf{w}_*^{[2]} \right\|_2^2 \leq \frac{2r \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})}{\lambda^2},$$

and we bound the difference between two optimal classifiers

$$\begin{aligned}
 & \left| h_*^{[1]}(\mathbf{x}_i^{[1]}) - h_*^{[2]}(\mathbf{x}_i^{[2]}) \right| \\
 &= \left| \langle \mathbf{w}_*^{[1]}, \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \rangle - \langle \mathbf{w}_*^{[2]}, \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \rangle \right| \\
 &\leq \left| \langle \mathbf{w}_*^{[1]}, \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \rangle - \langle \mathbf{w}_*^{[2]}, \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \rangle \right| + \left| \langle \mathbf{w}_*^{[2]}, \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \rangle - \langle \mathbf{w}_*^{[2]}, \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \rangle \right| \\
 &\leq \left\| \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) \right\|_2 \cdot \left\| \mathbf{w}_*^{[1]} - \mathbf{w}_*^{[2]} \right\|_2 + \left\| \mathbf{w}_*^{[2]} \right\|_2 \cdot \left\| \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2 \\
 &\leq \frac{r}{\lambda} \sqrt{2r\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})} + \frac{r}{\lambda} \left\| \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2,
 \end{aligned}$$

from equivalent linearization of kernel classifier and bounded norm for optimal classifier. We finally have

$$\begin{aligned}
 \hat{\rho}_{S_n}(h_*^{[1]}, h_*^{[2]}) &\leq \frac{r}{\lambda} \sqrt{2r\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})} + \frac{r}{n\lambda} \sum_{i=1}^n \left\| \hat{\varphi}^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2 \\
 &\leq \frac{r}{\lambda} \left( \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) + \sqrt{2r\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})} \right),
 \end{aligned}$$

which completes the proof.  $\square$

### A.3. Proof of Theorem 3.4

We begin with some useful lemmas as follows:

**Lemma A.2.** For a distribution  $\mathcal{D} \sim \mathcal{X}^{[1]} \times \mathcal{X}^{[2]}$ , we have

$$\begin{aligned}
 \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) &= \min_{\mathbf{U} \in \mathcal{U}} \left\{ \sqrt{\mathbb{E}_{\mathcal{D}} \left[ \left\| \mathbf{U} \varphi^{[1]}(\mathbf{x}^{[1]}) - \varphi^{[2]}(\mathbf{x}^{[2]}) \right\|_{HS}^2 \right]} \right\} \\
 &= \left( \mathbb{E}_{\mathcal{D}} \left[ \mathcal{K}^{[1]}(\mathbf{x}^{[1]}, \mathbf{x}^{[1]}) + \mathcal{K}^{[2]}(\mathbf{x}^{[2]}, \mathbf{x}^{[2]}) \right] - 2 \left\| \mathbb{E}_{\mathcal{D}} \left[ \varphi^{[1]}(\mathbf{x}^{[1]}) \varphi^{[2]}(\mathbf{x}^{[2]})^\top \right] \right\|_* \right)^{1/2}.
 \end{aligned}$$

*Proof.* For Hilbert-Schmidt norm, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathcal{D}} \left[ \left\| \mathbf{U} \varphi^{[1]}(\mathbf{x}^{[1]}) - \varphi^{[2]}(\mathbf{x}^{[2]}) \right\|_{HS}^2 \right] \\
 &= \mathbb{E}_{\mathcal{D}} \left[ \mathcal{K}^{[1]}(\mathbf{x}^{[1]}, \mathbf{x}^{[1]}) + \mathcal{K}^{[2]}(\mathbf{x}^{[2]}, \mathbf{x}^{[2]}) \right] - 2 \mathbb{E}_{\mathcal{D}} \left[ \varphi^{[2]}(\mathbf{x}^{[2]})^\top \mathbf{U} \varphi^{[1]} \right] \\
 &= \mathbb{E}_{\mathcal{D}} \left[ \mathcal{K}^{[1]}(\mathbf{x}^{[1]}, \mathbf{x}^{[1]}) + \mathcal{K}^{[2]}(\mathbf{x}^{[2]}, \mathbf{x}^{[2]}) \right] - 2 \mathbb{E}_{\mathcal{D}} \left[ \text{Tr} \left( \varphi^{[2]}(\mathbf{x}^{[2]})^\top \mathbf{U} \varphi^{[1]} \right) \right] \\
 &= \mathbb{E}_{\mathcal{D}} \left[ \mathcal{K}^{[1]}(\mathbf{x}^{[1]}, \mathbf{x}^{[1]}) + \mathcal{K}^{[2]}(\mathbf{x}^{[2]}, \mathbf{x}^{[2]}) \right] - 2 \text{Tr} \left( \mathbf{U} \mathbb{E}_{\mathcal{D}} \left[ \varphi^{[1]}(\mathbf{x}^{[1]}) \varphi^{[2]}(\mathbf{x}^{[2]})^\top \right] \right),
 \end{aligned}$$

where the last inequality holds from the linearity of operator trace w.r.t expectation, and independence between  $\mathbf{U}$  and  $\mathcal{D}$ . We complete the proof from the similar derivations as in the proof of Lemma 3.2.  $\square$

**Lemma A.3** (Perturbation bound for singular values (Bhatia, 2013)). For  $n \times n$  real matrices  $\mathbf{A}$  and  $\mathbf{B}$ , we have

$$\sum_{i=1}^n |\sigma_i(\mathbf{A}) - \sigma_i(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_*,$$

where  $\sigma_i(\mathbf{A})$  and  $\sigma_i(\mathbf{B})$  are their respective  $i$ -th singular values, i.e.,  $\sigma_1(\mathbf{A}) \geq \dots \geq \sigma_n(\mathbf{A})$  and  $\sigma_1(\mathbf{B}) \geq \dots \geq \sigma_n(\mathbf{B})$ .

**Lemma A.4** (McDiarmid's inequality (McDiarmid et al., 1989)). Let  $X_1, X_2, \dots, X_n$  be independent random variables taking values in a set  $A$ , and  $f : A^n \rightarrow \mathbb{R}$  satisfies

$$\sup_{X_1, X_2, \dots, X_n, X'_i \in A} |f(X_1, X_2, \dots, X_n) - f(X_1, X_2, \dots, X_{i-1}, X'_i, X_{i+1}, X_n)| \leq c_i$$

for every  $i \in [n]$ . Then, for  $t > 0$ , we have

$$\Pr [f(X_1, \dots, X_i, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X'_i, \dots, X_n)] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

**Lemma A.5.** Let  $\mathcal{K}^{[1]}$  and  $\mathcal{K}^{[2]}$  be two kernels bounded by  $r^2$ . For samples  $S_n = \{(\mathbf{x}_1^{[1]}, \mathbf{x}_1^{[2]}), \dots, (\mathbf{x}_n^{[1]}, \mathbf{x}_n^{[2]})\}$  and  $S'_n = S_n \setminus \{(\mathbf{x}_k^{[1]}, \mathbf{x}_k^{[2]})\} \cup \{(\mathbf{x}_k^{[1]'}, \mathbf{x}_k^{[2]'})\}$  ( $k \in [n]$ ), we have

$$\left| \left\| \sqrt{\mathbf{K}^{[1]}} \sqrt{\mathbf{K}^{[2]}} \right\|_* - \left\| \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]'}} \right\|_* \right| \leq 2r^2,$$

where  $\mathbf{K}^{[1]}$  and  $\mathbf{K}^{[1]'}$  are Gram matrices w.r.t kernel  $\mathcal{K}^{[1]}$  over  $S_n$  and  $S'_n$ , respectively, and define  $\mathbf{K}^{[2]}$  and  $\mathbf{K}^{[2]'}$  similarly.

*Proof.* From Lemma A.3, we have

$$\begin{aligned} & \left| \left\| \sqrt{\mathbf{K}^{[1]}} \sqrt{\mathbf{K}^{[2]}} \right\|_* - \left\| \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]'}} \right\|_* \right| \\ & \leq \left| \left\| \sqrt{\mathbf{K}^{[1]}} \sqrt{\mathbf{K}^{[2]}} \right\|_* - \left\| \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]}} \right\|_* \right| + \left| \left\| \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]}} \right\|_* - \left\| \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]'}} \right\|_* \right| \\ & = \left| \sum_{i=1}^n \left| \sigma_i(\sqrt{\mathbf{K}^{[1]}} \sqrt{\mathbf{K}^{[2]}}) \right| - \left| \sigma_i(\sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]}}) \right| \right| + \left| \sum_{i=1}^n \left| \sigma_i(\sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]}}) \right| - \left| \sigma_i(\sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]'}}) \right| \right| \\ & \leq \min_{\mathbf{U}, \mathbf{V} \in \mathcal{U}_n} \left\{ \sum_{i=1}^n \left| \sigma_i(\sqrt{\mathbf{K}^{[1]}} \sqrt{\mathbf{K}^{[2]}}) - \sigma_i(\mathbf{U} \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]}}) \right| + \left| \sigma_i(\sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]}}) - \sigma_i(\mathbf{V} \sqrt{\mathbf{K}^{[1]'}} \sqrt{\mathbf{K}^{[2]'}}) \right| \right\} \\ & \leq \min_{\mathbf{U} \in \mathcal{U}_n} \left\{ \left\| (\sqrt{\mathbf{K}^{[1]}} - \mathbf{U} \sqrt{\mathbf{K}^{[1]'}}) \sqrt{\mathbf{K}^{[2]}} \right\|_* \right\} + \min_{\mathbf{V} \in \mathcal{U}_n} \left\{ \left\| (\sqrt{\mathbf{K}^{[2]}} - \mathbf{V} \sqrt{\mathbf{K}^{[2]'}}) \sqrt{\mathbf{K}^{[1]'}} \right\|_* \right\}. \end{aligned} \quad (17)$$

We now prove that there is an  $\hat{\mathbf{U}} \in \mathcal{U}_n$  such that the difference between  $\sqrt{\mathbf{K}^{[1]}}$  and  $\hat{\mathbf{U}} \sqrt{\mathbf{K}^{[1]'}}$  lies only in the  $k$ -th column. Denote by  $\mathbf{v} = [\hat{\mathbf{U}} \sqrt{\mathbf{K}^{[1]'}}]_k$ , i.e., the  $k$ -th column of  $\hat{\mathbf{U}} \sqrt{\mathbf{K}^{[1]'}}$ . The existence of  $\hat{\mathbf{U}} \in \mathcal{U}_n$  is equivalent to solving the following underdetermined system of  $n - 1$  equations with  $n$  variables

$$\langle \mathbf{v}, (\sqrt{\mathbf{K}^{[1]}})_i \rangle = \frac{\mathcal{K}^{[1]}(\mathbf{x}_i^{[1]}, \mathbf{x}_k^{[1]'})}{(\mathcal{K}^{[1]}(\mathbf{x}_k^{[1]'}, \mathbf{x}_k^{[1]'})^{1/2}} \quad \text{for } i \in [n] \setminus \{k\}. \quad (18)$$

From the full-rank Gram matrix  $\mathbf{K}^{[1]}$ , there exists a solution  $\mathbf{v}_0$  for the system in Eqn. (18) because of  $n - 1$  equations with  $n$  variables. By setting  $\mathbf{v} = \mathbf{v}_0$ , we have an  $\hat{\mathbf{U}} \in \mathcal{U}_n$  such that the difference between  $\sqrt{\mathbf{K}^{[1]}}$  and  $\hat{\mathbf{U}} \sqrt{\mathbf{K}^{[1]'}}$  lies only in the  $k$ -th column. This follows that, from the nuclear norm of the rank-1 matrix,

$$\min_{\mathbf{U} \in \mathcal{U}_n} \left\{ \left\| (\sqrt{\mathbf{K}^{[1]}} - \mathbf{U} \sqrt{\mathbf{K}^{[1]'}}) \sqrt{\mathbf{K}^{[2]}} \right\|_* \right\} \leq \left\| ([\sqrt{\mathbf{K}^{[1]}}]_k - [\hat{\mathbf{U}} \sqrt{\mathbf{K}^{[1]'}}]_k) [\sqrt{\mathbf{K}^{[2]}}]_k^\top \right\|_* \leq 2r^2.$$

This completes the proof from Eqn. (17) and similar analysis for  $\min_{\mathbf{V} \in \mathcal{U}_n} \left\{ \left\| (\sqrt{\mathbf{K}^{[2]}} - \mathbf{V} \sqrt{\mathbf{K}^{[2]'}}) \sqrt{\mathbf{K}^{[1]'}} \right\|_* \right\}$ .  $\square$

**Lemma A.6** (Non-commutative Khintchine inequality, (Vershynin, 2018)). For independent Rademacher random variables  $\epsilon_1, \dots, \epsilon_n$  and for real matrices  $\mathbf{X}_1, \dots, \mathbf{X}_n$  of the same size, we have

$$\mathbb{E} \left[ \left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_* \right] \leq C \max \left\{ \left\| \sqrt{\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top} \right\|_*, \left\| \sqrt{\sum_{i=1}^n \mathbf{X}_i^\top \mathbf{X}_i} \right\|_* \right\} \text{ for some positive constant } C.$$

**Lemma A.7.** For independent Rademacher random variables  $\epsilon_1, \dots, \epsilon_n$  and for real matrices  $\mathbf{X}_1, \dots, \mathbf{X}_n$  of the same size with  $\mathbb{E}[\|\mathbf{X}_i\|_*] < \infty$ , we have

$$\mathbb{E} \left[ \left\| \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i]) \right\|_* \right] \leq 2\mathbb{E} \left[ \left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_* \right].$$

*Proof.* Let  $\{\mathbf{X}'_1, \dots, \mathbf{X}'_n\}$  denote an independent copy of the sequence  $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ . From i.i.d. assumption and Jensen's inequality, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} \left[ \left\| \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} [\mathbf{X}_i]) \right\|_* \right] &= \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} \left[ \left\| \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} [\mathbf{X}_i]) - \mathbb{E}' \left[ \mathbf{X}'_i - \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} [\mathbf{X}_i] \right] \right\|_* \right] \\ &\leq \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} \left[ \mathbb{E}_{\mathbf{X}'_1, \dots, \mathbf{X}'_n} \left[ \left\| \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i]) - (\mathbf{X}'_i - \mathbb{E}[\mathbf{X}'_i]) \right\|_* \right] \right] = \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{X}'_1, \dots, \mathbf{X}'_n} \left[ \left\| \sum_{i=1}^n (\mathbf{X}_i - \mathbf{X}'_i) \right\|_* \right]. \end{aligned}$$

We also have, from the triangle inequality of nuclear norm and symmetry of Rademacher random variables,

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{X}'_1, \dots, \mathbf{X}'_n} \left[ \left\| \sum_{i=1}^n (\mathbf{X}_i - \mathbf{X}'_i) \right\|_* \right] &= \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{X}'_1, \dots, \mathbf{X}'_n} \left[ \left\| \sum_{i=1}^n \epsilon_i (\mathbf{X}_i - \mathbf{X}'_i) \right\|_* \right] \\ &\leq \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} \left[ \left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_* \right] + \mathbb{E}_{\mathbf{X}'_1, \dots, \mathbf{X}'_n} \left[ \left\| \sum_{i=1}^n -\epsilon_i \mathbf{X}'_i \right\|_* \right] = 2 \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_n} \left[ \left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_* \right], \end{aligned}$$

which completes the proof.  $\square$

**Proof of Theorem 3.4.** By triangle inequality, we have

$$\begin{aligned} &\left| \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) - \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right| \\ &\leq \underbrace{\left| \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) - \mathbb{E}_{S_n} \left[ \hat{\mathcal{E}}(S_n, \mathcal{K}_1, \mathcal{K}_2) \right] \right|}_{\text{Concentration analysis}} + \underbrace{\left| \mathbb{E}_{S_n} \left[ \hat{\mathcal{E}}(S_n, \mathcal{K}_1, \mathcal{K}_2) \right] - \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right|}_{\text{Random matrix analysis}}. \quad (19) \end{aligned}$$

For sample  $S'_n$  with the  $k$ -th instance replaced by  $(\mathbf{x}_k^{[1]'}, \mathbf{x}_k^{[2]'})$  from  $S_n$ , we have

$$\begin{aligned} &n \left| \left( \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right)^2 - \left( \hat{\mathcal{E}}(S'_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right)^2 \right| \\ &\leq \left| \Delta \left( \left\| \sqrt{\mathbf{K}^{[1]}} \sqrt{\mathbf{K}^{[2]}} \right\|_* \right) \right| + \left| \mathcal{K}^{[1]}(\mathbf{x}_k^{[1]}, \mathbf{x}_k^{[1]}) - \mathcal{K}^{[1]}(\mathbf{x}_k^{[1]'}, \mathbf{x}_k^{[1]'}) \right| + \left| \mathcal{K}^{[2]}(\mathbf{x}_k^{[2]}, \mathbf{x}_k^{[2]}) - \mathcal{K}^{[2]}(\mathbf{x}_k^{[2]'}, \mathbf{x}_k^{[2]'}) \right| \leq 6r^2, \end{aligned}$$

from bounded kernel and Lemma A.5. Based on the McDiarmid's inequality (Lemma A.4), the following holds with probability at least  $1 - \delta$ ,

$$\left| \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) - \mathbb{E}_{S_n} \left[ \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right] \right| \leq 3r \sqrt{\frac{2}{n} \ln \frac{2}{\delta}}. \quad (20)$$

For the second term in Eqn. (19), we consider the random matrix analysis with nuclear norm (Vershynin, 2018). Denote by  $\Phi = \mathbb{E}_{\mathcal{D}}[\varphi^{[1]}(\mathbf{x}^{[1]})\varphi^{[2]}(\mathbf{x}^{[2]})^\top]$  and operator  $\mathbf{X}_i = (\varphi^{[1]}(\mathbf{x}_i^{[1]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top - \Phi)/n$ , and we have

$$\mathbb{E}_{S_n} \left[ \frac{1}{n} \text{Tr}(\mathbf{K}^{[k]}) \right] = \mathbb{E}_{\mathcal{D}}[\mathcal{K}^{[k]}(\mathbf{x}^{[k]}, \mathbf{x}^{[k]})] \quad \text{for } k \in [2].$$

This follows that, from the linearity of expectation and Lemma A.3,

$$\begin{aligned} &\left| \left( \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right)^2 - \mathbb{E}_{S_n} \left[ \left( \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right)^2 \right] \right| = 2 \left| \mathbb{E}_{S_n} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \varphi^{[1]}(\mathbf{x}_i^{[1]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top - \Phi \right\|_* \right] \right| \\ &\leq 2 \mathbb{E}_{S_n} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \varphi^{[1]}(\mathbf{x}_i^{[1]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top - \Phi \right\|_* \right] \leq 4 \mathbb{E} \left[ \left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_* \right], \end{aligned}$$

where  $\epsilon_1, \dots, \epsilon_n$  are independent Rademacher random variables. From  $\mathbf{X}_i = (\varphi^{[1]}(\mathbf{x}_i^{[1]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top - \Phi)/n$ , we have

$$\mathbf{X}_i \mathbf{X}_i^\top = \frac{\mathcal{K}^{[1]}(\mathbf{x}_i^{[1]}, \mathbf{x}_i^{[1]}) \cdot \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top - \varphi^{[1]}(\mathbf{x}_i^{[1]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top \Phi^\top - \Phi \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[1]}(\mathbf{x}_i^{[1]})^\top + \Phi \Phi^\top}{n^2}.$$

This follows that, by the sub-additivity of square root w.r.t nuclear norm and some algebraic calculations,

$$\begin{aligned} \left\| \sqrt{\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top} \right\|_* &\leq \frac{1}{n} \left\| \sqrt{\sum_{i=1}^n \mathcal{K}^{[1]}(\mathbf{x}_i^{[1]}, \mathbf{x}_i^{[1]}) \cdot \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top} \right\|_* \\ &\quad + \frac{\sqrt{2}}{n} \left\| \sqrt{\sum_{i=1}^n \Phi \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[1]}(\mathbf{x}_i^{[1]})^\top} \right\|_* + \frac{1}{n} \left\| \sqrt{\Phi \Phi^\top} \right\|_* \leq \left( \frac{1 + \sqrt{2}}{\sqrt{n}} + \frac{1}{n} \right) r^2, \end{aligned}$$

where we use

$$\left\| \sqrt{\Phi \Phi^\top} \right\|_* = \|\Phi\|_* = \left\| \mathbb{E}_{\mathcal{D}} \left[ \varphi^{[1]}(\mathbf{x}^{[1]})\varphi^{[2]}(\mathbf{x}^{[2]})^\top \right] \right\|_* \leq \mathbb{E}_{\mathcal{D}} \left[ \left\| \varphi^{[1]}(\mathbf{x}^{[1]})\varphi^{[2]}(\mathbf{x}^{[2]})^\top \right\|_* \right] \leq r^2;$$

and

$$\left\| \sqrt{\sum_{i=1}^n \Phi \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[1]}(\mathbf{x}_i^{[1]})^\top} \right\|_* \leq r \left\| \sqrt{\sum_{i=1}^n \varphi^{[1]}(\mathbf{x}_i^{[1]})\mathbb{E}_{\mathcal{D}} \left[ \varphi^{[1]}(\mathbf{x}^{[1]})^\top \right]} \right\|_* \leq \sqrt{nr}^{\frac{3}{2}} \left\| \sqrt{\mathbb{E}_{\mathcal{D}} \left[ \varphi^{[1]}(\mathbf{x}^{[1]}) \right]} \right\|_* \leq \sqrt{nr}^2;$$

and, from linearity of expectation and relationship between nuclear norm and operator norm for rank-1 operator,

$$\left\| \sqrt{\sum_{i=1}^n \mathcal{K}^{[1]}(\mathbf{x}_i^{[1]}, \mathbf{x}_i^{[1]}) \cdot \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top} \right\|_* \leq r \left\| \sqrt{\sum_{i=1}^n \varphi^{[2]}(\mathbf{x}_i^{[2]})\varphi^{[2]}(\mathbf{x}_i^{[2]})^\top} \right\|_* \leq \sqrt{nr}^2.$$

Similarly, we can present the same upper bounds for  $\sum_{i=1}^n \mathbf{X}_i^\top \mathbf{X}_i$ . This follows that, from Lemmas A.6 and A.7,

$$\left| \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) - \mathbb{E}_{S_n} \left[ \hat{\mathcal{E}}(S_n, \mathcal{K}_1, \mathcal{K}_2) \right] \right| \leq C \left( \frac{1 + \sqrt{2}}{\sqrt{n}} + \frac{1}{n} \right) r. \quad (21)$$

We complete the proof from Eqns. (20) and (21), and triangle inequality.  $\square$

#### A.4. Proof of Lemma 3.5

From (Bhatia, 2013), we have  $\|\sqrt{\mathbf{A}} - \sqrt{\mathbf{B}}\|_F \leq \sqrt{\|\mathbf{A} - \mathbf{B}\|_*}$ , and this follows that, for KOM discrepancy,

$$\hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \leq \frac{1}{\sqrt{n}} \left\| \sqrt{\mathbf{K}^{[1]}} - \sqrt{\mathbf{K}^{[2]}} \right\|_F \leq \sqrt{\frac{1}{n} \|\mathbf{K}^{[1]} - \mathbf{K}^{[2]}\|_*}.$$

We further have, for normalized kernel matrices,

$$\begin{aligned} \sqrt{\frac{1}{n} \|\mathbf{K}^{[1]} - \mathbf{K}^{[2]}\|_*} &\leq \sqrt{\frac{1}{\sqrt{n}} \|\mathbf{K}^{[1]} - \mathbf{K}^{[2]}\|_F} \\ &= \sqrt[4]{\frac{\|\mathbf{K}^{[1]}\|_F \|\mathbf{K}^{[2]}\|_F}{n} \left( \frac{\|\mathbf{K}^{[1]}\|_F}{\|\mathbf{K}^{[2]}\|_F} + \frac{\|\mathbf{K}^{[2]}\|_F}{\|\mathbf{K}^{[1]}\|_F} - \frac{2 \text{Tr}(\mathbf{K}^{[1]}\mathbf{K}^{[2]})}{\|\mathbf{K}^{[1]}\|_F \|\mathbf{K}^{[2]}\|_F} \right)} \leq r \sqrt[4]{2(1 - \hat{A}(\mathbf{K}^{[1]}, \mathbf{K}^{[2]}))}, \end{aligned}$$

which completes the proof.  $\square$

### A.5. Proof of Lemma 3.6

Denote by  $d = \dim(\varphi^{[1]})$  for simplicity. For sample  $S_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , we write

$$\mathbf{P}^{[1]} = \left( \varphi^{[1]}(\mathbf{x}_i^{[1]}) \right)_{i=1}^n \quad \text{and} \quad \hat{\mathbf{P}}^{[2]} = \left( \hat{\varphi}^{[2]}(\mathbf{x}_i^{[1]}) \right)_{i=1}^n,$$

and we also have their respective Gram matrices  $\mathbf{K}^{[1]} = \mathbf{P}^{[1]\top} \mathbf{P}^{[1]}$  and  $\hat{\mathbf{K}}^{[2]} = \hat{\mathbf{P}}^{[2]\top} \hat{\mathbf{P}}^{[2]}$ . It remains to prove

$$\min_{\mathbf{U} \in \mathcal{U}_n} \left\{ \left\| \sqrt{\mathbf{K}^{[1]}} \mathbf{U} - \sqrt{\hat{\mathbf{K}}^{[2]}} \right\|_F \right\} \leq \left\| \mathbf{P}^{[1]} - \hat{\mathbf{P}}^{[2]} \right\|_F.$$

If  $n \leq d$ , then there exists a matrix  $\mathbf{V} \in \mathcal{U}_n$ , from Lemma A.1 and the unitary invariance of Frobenius norm, such that

$$\left\| \mathbf{P}^{[1]} - \hat{\mathbf{P}}^{[2]} \right\|_F = \left\| \mathbf{P}^{[1]} \mathbf{V} - \hat{\mathbf{P}}^{[2]} \mathbf{V} \right\|_F = \left\| \mathbf{P}^{[1]\top} \mathbf{V} - \sqrt{\hat{\mathbf{K}}^{[2]}} \right\|_F,$$

We also have

$$\left\| \mathbf{P}^{[1]} - \hat{\mathbf{P}}^{[2]} \right\|_F \geq \min_{\mathbf{W} \in \mathcal{U}_n} \left\{ \left\| \sqrt{\mathbf{K}^{[1]}} \mathbf{W} - \sqrt{\hat{\mathbf{K}}^{[2]}} \right\|_F \right\},$$

for some  $\mathbf{W} \in \mathcal{U}_n$  with  $\sqrt{\mathbf{K}^{[1]}} \mathbf{W} = \mathbf{P}^{[1]\top} \mathbf{V}$ .

If  $n > d$ , then we have, similarly to the proof of Theorem A.2,

$$\left\| \mathbf{P}^{[1]} - \hat{\mathbf{P}}^{[2]} \right\|_F \geq \min_{\mathbf{W} \in \mathcal{U}_d} \left\{ \left\| \mathbf{P}^{[1]\top} \mathbf{W} - \hat{\mathbf{P}}^{[2]\top} \right\|_F \right\} = \sqrt{\text{Tr}(\mathbf{K}^{[1]}) + \text{Tr}(\hat{\mathbf{K}}^{[2]}) - 2 \left\| \mathbf{P}^{[1]} \hat{\mathbf{P}}^{[2]\top} \right\|_*},$$

and this follows that, from the unitary invariance of the nuclear norm,

$$\left\| \mathbf{P}^{[1]} \hat{\mathbf{P}}^{[2]\top} \right\|_* = \left\| \begin{bmatrix} \mathbf{P}^{[1]} \\ \mathbf{0}_{(n-d) \times n} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{P}}^{[2]\top} & \mathbf{0}_{n \times (n-d)} \end{bmatrix} \right\|_* = \left\| \sqrt{\mathbf{K}^{[1]}} \sqrt{\hat{\mathbf{K}}^{[2]}} \right\|_*.$$

We finally have, from Lemma 3.2,

$$\min_{\mathbf{W} \in \mathcal{U}_d} \left\{ \left\| \mathbf{P}^{[1]\top} \mathbf{W} - \hat{\mathbf{P}}^{[2]\top} \right\|_F^2 \right\} = \left( \hat{\mathcal{E}}(S_n, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) \right)^2 \leq \frac{1}{n} \sum_{i=1}^n \left\| \varphi^{[1]}(\mathbf{x}_i^{[1]}) - \hat{\varphi}^{[2]}(\mathbf{x}_i^{[1]}) \right\|_2^2,$$

which completes the proof.  $\square$

### B. Proof of Theorem 3.7

We begin with two useful lemmas as follows:

**Lemma B.1.** *The sub-gradient of  $f(\mathbf{p})$  in Eqn. (6) is given by*

$$\mathbf{v} - \text{diag}(\mathbf{M}^\top \mathbf{U} \mathbf{V}^\top) / \sqrt{\mathbf{p}} \in \partial f(\mathbf{p}),$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are left and right singular vectors matrices of  $\mathbf{M} \text{diag}(\sqrt{\mathbf{p}})$ .

*Proof.* We set  $\mathbf{X} = \mathbf{M} \sqrt{\mathbf{p}}$  with the singular value decomposition  $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^\top$ . Following (Watson, 1992), the sub-gradient of  $\|\mathbf{X}\|_*$  is given by

$$\partial \|\mathbf{X}\|_* = \left\{ \mathbf{U} \mathbf{V}^\top + \mathbf{W} \mid \mathbf{W} \in \mathbb{R}^{d_1 \times d_2}, \mathbf{U}^\top \mathbf{W} = \mathbf{0}, \mathbf{W} \mathbf{V} = \mathbf{0}, \|\mathbf{W}\|_2 \leq 1 \right\}.$$

Obviously,  $\mathbf{U} \mathbf{V}^\top$  is a sub-gradient by taking  $\mathbf{W} = \mathbf{0}$ , and we have

$$\frac{\partial \|\mathbf{X}\|_*}{\partial \mathbf{p}} = \frac{\partial \|\mathbf{X}\|_*}{\partial \mathbf{X}} \frac{\partial \mathbf{X}}{\partial \text{diag}(\sqrt{\mathbf{p}})} \frac{\partial \text{diag}(\sqrt{\mathbf{p}})}{\partial \mathbf{p}} = \text{diag}(\mathbf{M}^\top \mathbf{U} \mathbf{V}^\top) / \sqrt{\mathbf{p}},$$

which completes the proof.  $\square$

This proposition shows that a singular value decomposition is required to compute the sub-gradient of  $f(\mathbf{p})$  in each iteration, which is computationally expensive. However,  $\mathbf{M}\text{diag}(\sqrt{\mathbf{p}})$  has a low rank structure with the rank no more than  $\min\{T_e, d_1, d_2\}$  because evolving stage  $T_e$  is usually small compared to  $d_1$  and  $d_2$  in our setting. This could drastically reduce the computational cost on singular value decomposition.

**Proof of Theorem 3.7.** From Eqn. (6), we have

$$f(\mathbf{p}) = \min_{\mathbf{W} \in \mathcal{U}_{T_e}} \{g(\mathbf{p}, \mathbf{W})\} \quad \text{with} \quad g(\mathbf{p}, \mathbf{W}) = \frac{1}{2} \left\| \sqrt{\tilde{\mathbf{K}}^{[1]}} - \sqrt{\Phi^{[2]} \text{diag}(\mathbf{p}) \Phi^{[2]\top}} \right\|_F^2 - \frac{1}{2} \text{Tr} \left( \tilde{\mathbf{K}}^{[1]} \right),$$

where  $\tilde{\mathbf{K}}^{[1]} = [\sum_{k=1}^{d_1} \Phi_{i,k}^{[1]} \Phi_{j,k}^{[1]} / d_1]_{T_e \times T_e}$  and  $\Phi^{[2]} = [\Phi_{i,j}^{[2]}]_{T_e \times d_2}$ . Hence,  $f(\mathbf{p})$  is a convex function w.r.t.  $\mathbf{p}$  from the convexity of  $\min_{\mathbf{W} \in \mathcal{U}_{T_e}} g(\mathbf{p}, \mathbf{W})$  in (Boyd & Vandenberghe, 2004).

In Algorithm 1, we select

$$h(\mathbf{p}) = \sum_{i=1}^{d_2} p_i (\ln p_i - 1),$$

and we have the corresponding Bregman divergence

$$D_h(\mathbf{p} \parallel \mathbf{q}) = h(\mathbf{p}) - h(\mathbf{q}) - \langle \nabla h(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle \quad \text{for} \quad \mathbf{p}, \mathbf{q} \in \Delta.$$

From the Fenchel conjugate, we also have

$$h^*(\boldsymbol{\theta}) = \sup_{\mathbf{p} \in \mathbb{R}^{d_2}} \{\boldsymbol{\theta}^\top \mathbf{p} - h(\mathbf{p})\} \quad \text{and} \quad \nabla h^*(\boldsymbol{\theta}) = \arg \max_{\mathbf{p} \in \mathbb{R}^{d_2}} \{\boldsymbol{\theta}^\top \mathbf{p} - h(\mathbf{p})\},$$

and define  $D_{h^*}(\mathbf{p} \parallel \mathbf{q}) = h^*(\mathbf{p}) - h^*(\mathbf{q}) - \langle \nabla h^*(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$  similarly. In Algorithm 1, we rewrite mirror descent iteration as

$$\mathbf{p}^{(t)} = \nabla h^* \left( \nabla h(\mathbf{p}^{(t-1)}) - \tau_{k-1} \mathbf{g}^{(t-1)} \right) \quad \text{with} \quad \mathbf{g}^{(t-1)} \in \partial f(\mathbf{p}^{(t-1)}).$$

Let  $\mathbf{p}^* \in \arg \min_{\mathbf{p} \in \Delta} f(\mathbf{p})$  and  $\boldsymbol{\theta}^* = \nabla h(\mathbf{p}^*)$ . For Bregman divergence (Banerjee et al., 2005), we have

$$D_{h^*}(\boldsymbol{\theta}^{(t)} \parallel \boldsymbol{\theta}^*) = D_{h^*}(\boldsymbol{\theta}^{(t-1)} \parallel \boldsymbol{\theta}^*) + \left( \boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(t-1)} \right)^\top \left( \nabla h^*(\boldsymbol{\theta}^{(t-1)}) - \nabla h^*(\boldsymbol{\theta}^*) \right) + D_{h^*}(\boldsymbol{\theta}^{(t)} \parallel \boldsymbol{\theta}^{(t-1)}),$$

and

$$\left( \boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(t-1)} \right)^\top \left( \nabla h^*(\boldsymbol{\theta}^{(t-1)}) - \nabla h^*(\boldsymbol{\theta}^*) \right) = -\tau_{t-1} \mathbf{g}^{(t-1)\top} \left( \mathbf{p}^{(t-1)} - \mathbf{p}^* \right).$$

For convex  $f$  and  $\mathbf{g}^{(t-1)} \in \partial f(\mathbf{p}^{(t-1)})$ , we have  $f(\mathbf{p}^{(t)}) - f(\mathbf{p}^*) \leq \mathbf{g}^{(t)\top} (\mathbf{p}^{(t)} - \mathbf{p}^*)$ , and

$$\tau_{t-1} \left[ f(\mathbf{p}^{(t-1)}) - f(\mathbf{p}^*) \right] \leq D_{h^*}(\boldsymbol{\theta}^{(t-1)} \parallel \boldsymbol{\theta}^*) - D_{h^*}(\boldsymbol{\theta}^{(t)} \parallel \boldsymbol{\theta}^*) + D_{h^*}(\boldsymbol{\theta}^{(t-1)} \parallel \boldsymbol{\theta}^{(t)}).$$

Summing from  $t = 0$  to  $T_m$ , we have

$$\sum_{t=1}^{T_m} \tau_t \left[ f(\mathbf{p}^{(t)}) - f(\mathbf{p}^*) \right] \leq D_{h^*}(\boldsymbol{\theta}^{(1)} \parallel \boldsymbol{\theta}^*) + \frac{1}{2} \sum_{t=1}^{T_m} \tau_t^2 \left\| \mathbf{g}^{(t)} \right\|_*^2 \leq \ln d_2 + \frac{1}{2} \sum_{t=1}^{T_m} \tau_t^2 \left\| \mathbf{g}^{(t)} \right\|_*^2, \quad (22)$$

where the norm  $\|\cdot\|_*$  is defined on  $h^*$  and  $D_{h^*}(\boldsymbol{\theta}^{(1)} \parallel \boldsymbol{\theta}^*) = D_h(\mathbf{p}^{(1)} \parallel \mathbf{p}^*)$ .

For each  $t \in [T_m]$ , we project  $\mathbf{p}_t$  onto  $\mathcal{P} = \{\mathbf{p} \in \Delta : \|\mathbf{p}\|_\infty \geq \epsilon\}$  for a small  $\epsilon > 0$ . We consider the singular value decomposition  $\mathbf{M}\text{diag}(\sqrt{\mathbf{p}^{(t)}}) = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^\top$ , and denote by  $\mathbf{u}_t = \text{diag}(\mathbf{M}^\top \mathbf{U}_t \mathbf{V}_t^\top)$ . This follows that

$$\|\mathbf{u}_t\|_\infty = \max_{i \in [d_2]} \{(\mathbf{M} \mathbf{U}_t \mathbf{V}_t^\top)_{i,i}\} \leq \max_{i \in [d_2]} \{\|\mathbf{M}_i\|_2\} \leq T_e \sqrt{d_2},$$

from Lemma B.1 and orthogonality of  $\mathbf{U}_t \mathbf{V}_t^\top$ . Hence, we have, from the conjugation between  $\ell_1$  and  $\ell_\infty$  norm,

$$\left\| \mathbf{g}^{(t)} \right\|_* \leq \|\mathbf{v}\|_\infty + \frac{\|\mathbf{u}_t\|_\infty}{\epsilon} \leq T_e \left( 1 + \frac{\sqrt{d_2}}{\epsilon} \right), \quad (23)$$

and this follows that, from Eqns. (22)-(23) and by selecting stepsize  $\tau_t = \tau$ ,

$$\frac{1}{T_m} \sum_{t=1}^{T_m} f(\mathbf{p}^{(t)}) - f(\mathbf{p}^*) \leq \frac{\ln d_2}{\tau T_m} + \frac{\tau T_e}{2} \left(1 + \frac{\sqrt{d_2}}{\epsilon}\right).$$

This completes the proof by setting  $\tau = \sqrt{2 \ln d_2 / (T_m T_e (1 + \sqrt{d_2}/\epsilon))}$ .  $\square$

## C. Analysis for OPFES

### C.1. Proof of Proposition 4.1

Let  $d_1$  and  $d_2$  be the dimensionalities of random features of  $\mathcal{K}^1$  and  $\mathcal{K}^2$ , respectively. Denote by

$$\mathbf{Z}^{[1]} = (\mathbf{z}_{T_1+t}^{[k]})_{t=1}^{T_e} \in \mathbb{R}^{d_k \times T_e} \quad \text{and} \quad \mathbf{Z}^{[2]} = (\mathbf{z}_{T_1+t}^{[k]})_{t=1}^{T_e} \in \mathbb{R}^{d_k \times T_e}.$$

For  $d_1 \geq d_2$ , we can rewrite the optimization for Eqn. (11) as

$$\mathbf{U}^* \in \arg \min_{\mathbf{U} \in \mathcal{U}_{d_2 \times d_1}} \left\{ \left\| \mathbf{U} \mathbf{Z}^{[1]} - \mathbf{Z}^{[2]} \right\|_F \right\},$$

where  $\mathcal{U}_{d_2 \times d_1} = \{\mathbf{U} \in \mathbb{R}^{d_2 \times d_1} : \mathbf{U} \mathbf{U}^\top = \mathbf{I}_{d_2}\}$  is the set of semi-orthogonal matrices. From the proof of Lemma 3.2, we have an equivalent optimization as

$$\max_{\mathbf{U} \in \mathcal{U}_{d_2 \times d_1}} \left\{ \text{Tr} \left( \mathbf{U} \mathbf{M}^{(T_1+T_e)} \text{diag}(\sqrt{\mathbf{p}^{(T_M)}}) \right) \right\}.$$

Let  $\mathbf{M}^{(T_1+T_e)} \text{diag}(\sqrt{\mathbf{p}^{(T_M)}}) = \mathbf{V} \mathbf{\Sigma} \mathbf{W}^\top$  be the singular value decomposition with left and right singular vector matrices  $\mathbf{V}^\top \in \mathcal{U}_{d_2 \times d_1}$  and  $\mathbf{W} \in \mathcal{U}_{d_2}$ , respectively. Denote by  $\mathbf{S} = \mathbf{W}^\top \mathbf{U} \mathbf{V}$ , and we have, from cyclic invariance of matrix trace,

$$\max_{\mathbf{U} \in \mathcal{U}_{d_2 \times d_1}} \left\{ \text{Tr} \left( \mathbf{U} \mathbf{V} \mathbf{\Sigma} \mathbf{W}^\top \right) \right\} = \max_{\mathbf{S} \in \mathcal{U}_{d_2}} \left\{ \text{Tr} \left( \mathbf{S} \mathbf{\Sigma} \right) \right\}.$$

We get the maximum when  $\mathbf{S} = \mathbf{I}_{d_2}$  and the optimal solution set for  $\mathbf{U}$  is given by

$$\mathbf{U}_* \in \left\{ \mathbf{W} \mathbf{V}^\top + \mathbf{Q} (\mathbf{I} - \mathbf{V} \mathbf{V}^\top) : \mathbf{Q} \in \mathbb{R}^{d_1 \times d_2} \right\}. \quad (24)$$

For  $d_1 < d_2$ , we have the equivalent optimization of Eqn. (11) as

$$\mathbf{U}^* \in \arg \min_{\mathbf{U}^\top \in \mathcal{U}_{d_1 \times d_2}} \left\{ \left\| \mathbf{U} \mathbf{Z}^{[1]} - \mathbf{Z}^{[2]} \right\|_F \right\}.$$

We also take the singular value decomposition of  $\mathbf{M}^{(T_1+T_e)} \text{diag}(\sqrt{\mathbf{p}^{(T_M)}})$  with  $\mathbf{V} \in \mathcal{U}_{d_1}$  and  $\mathbf{W}^\top \in \mathcal{U}_{d_1 \times d_2}$ , respectively. For  $\mathbf{S} = \mathbf{W}^\top \mathbf{U} \mathbf{V} \in \mathbb{R}^{d_1 \times d_1}$ , we similarly have the optimal solution set for  $\mathbf{U}$  as

$$\mathbf{U}_* \in \left\{ \mathbf{W} \mathbf{V}^\top + (\mathbf{I} - \mathbf{W} \mathbf{W}^\top) \mathbf{Q} : \mathbf{Q} \in \mathbb{R}^{d_2 \times d_1} \right\}. \quad (25)$$

From Eqns. (24) and (25),  $\mathbf{U} = \mathbf{W} \mathbf{V}^\top$  is always an optimal solution for Eqn. (11) by setting  $\mathbf{Q} = \mathbf{0}$ .  $\square$

### C.2. Proof of Lemma 4.2

For  $t \in [T_1 + 1, T_1 + T_e]$ , we have, by simple calculations and Cauchy-Schwarz inequality,

$$\begin{aligned} & \left| \langle \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) \rangle - \langle \mathbf{w}_{T_1+T_e}^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle \right| = \left| \langle \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) \rangle - \langle \mathbf{U}_*^\top \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle \right| \\ & \leq \left| \langle \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) - \mathbf{U}_* \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle \right| \leq \left\| \mathbf{w}_{T_1}^{[1]} \right\|_2 \left\| \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) - \mathbf{U}_* \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \right\|_2. \end{aligned}$$

This follows that, by summing  $t$  from  $T_1 + 1$  to  $T_1 + T_e$ ,

$$\begin{aligned} \frac{1}{T_e} \left| \langle \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[1]}(\mathbf{x}_i^{[1]}) \rangle - \langle \mathbf{w}_{T_1+T_e}^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}_i^{[2]}) \rangle \right| &\leq \frac{\|\mathbf{w}_{T_1}^{[1]}\|_2}{T_e} \sum_{i=T_1+1}^{T_1+T_e} \left\| \mathbf{U}_* \mathbf{z}^{[1]}(\mathbf{x}_i^{[1]}) - \mathbf{z}^{[2]}(\mathbf{x}_i^{[2]}) \right\|_2 \\ &\leq \frac{\sqrt{2}}{\lambda \sqrt{T_e}} \min_{\mathbf{U} \in \mathcal{U}_{d_2 \times d_1}} \left\{ \left\| \mathbf{U}_* \left( \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) \right)_{t=T_1+1}^{T_1+T_e} - \left( \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \right)_{t=T_1+1}^{T_1+T_e} \right\|_F \right\} = \sqrt{2} \hat{\mathcal{E}}(S_{T_e}^{[e]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) / \lambda, \end{aligned}$$

where the second inequality holds from Lemma C.5 and Cauchy-Schwarz inequality, and the equality holds from the unitary invariance of Frobenius norm.  $\square$

### C.3. Proof of Theorem 4.3

We first introduce some useful lemmas.

**Lemma C.1** (Hoeffding's bounds (Hoeffding, 1963)). *Let  $X_1, X_2, \dots, X_n$  be independent random variables in  $[a, b]$ , and  $\bar{X} = \sum_{i=1}^n X_i/n$ . For  $t > 0$ , we have*

$$\Pr [\bar{X} - \mathbb{E}[\bar{X}] \geq t] \leq \exp\left(-\frac{2nt^2}{(b-a)^2}\right).$$

**Lemma C.2** (Generalization bound of kernel methods (Mohri et al., 2018)). *Given  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  drawn i.i.d from the distribution  $\mathcal{D}$ . Let  $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a kernel bounded by  $r^2$ , and  $\varphi(\cdot)$  is the feature mapping of  $\mathcal{K}$  and let  $\mathcal{H} = \{\mathbf{x} \rightarrow \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle : \|\mathbf{w}\|_{\mathbb{H}_{\mathcal{K}}} \leq \Lambda\}$  for some  $\Lambda \geq 0$ . For loss function  $|\ell(\mathbf{w}, (\mathbf{x}, y))| \leq M$  and  $\delta \in (0, 1)$ , the following holds with probability at least  $1 - \delta$  for  $h \in \mathcal{H}$ ,*

$$R(h) \leq \hat{R}_S(h) + 2r\Lambda \sqrt{\frac{1}{n}} + M \sqrt{\frac{1}{2n} \ln \frac{1}{\delta}},$$

where  $R(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(\mathbf{w}, (\mathbf{x}, y))]$  and  $\hat{R}_S(h) = \sum_{i=1}^n \ell(\mathbf{w}, (\mathbf{x}_i, y_i))/n$ .

**Lemma C.3** (Online to batch conversion (Cesa-Bianchi et al., 2004)). *Let  $S_T = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)\}$  be a sample drawn i.i.d. from  $\mathcal{D}$ ,  $\ell$  a loss bounded by  $M$  and  $h_1, \dots, h_T$  the sequence of hypotheses generated by an online algorithm. For  $\delta \in (0, 1)$ , the following holds with a probability at least  $1 - \delta$ ,*

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ \ell \left( \frac{1}{T} \sum_{t=1}^T h_t(x_t), y_t \right) \right] \leq \frac{1}{T} \sum_{i=1}^T \ell(h_i(x_i), y_i) + M \sqrt{\frac{2 \ln(1/\delta)}{T}}.$$

**Lemma C.4.** *Let  $S_T = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)\}$  be a sample and  $B$  be a closed convex set with projection  $\Pi_B(\mathbf{w}) = \arg \min_{\mathbf{w}' \in B} \|\mathbf{w} - \mathbf{w}'\|$ . For a strongly convex loss function  $\ell$  with bounded gradient w.r.t.  $\mathbf{w}$ , i.e.,  $\|\nabla_{\mathbf{w}} \ell(\mathbf{w}, (x_t, y_t))\| \leq G$  for  $\mathbf{w} \in B$  and  $t \in [T]$ . For the update rule  $\mathbf{w}_t = \Pi_B(\mathbf{w}_{t-1} - \nabla \ell(\mathbf{w}_t, (x_t, y_t))/\lambda t)$  with  $\mathbf{w}_0, \mathbf{w} \in B$ , we have*

$$\frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}_t, (x_t, y_t)) - \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}, (x_t, y_t)) \leq \frac{G^2(1 + \ln T)}{2\lambda T}.$$

*Proof.* Denote by  $\nabla_t = \nabla_{\mathbf{w}} \ell(\mathbf{w}_t, (x_t, y_t))$  for simplicity. For  $\lambda$ -strongly convex functions, we have

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 &\leq \|\mathbf{w}_t - \mathbf{w}\|^2 - 2\tau_t \langle \nabla_t, \mathbf{w}_t - \mathbf{w} \rangle + \tau_t^2 \|\nabla_t\|^2 \\ &\leq \|\mathbf{w}_t - \mathbf{w}\|^2 - 2\tau_t \left( \ell(\mathbf{w}_t, (x_t, y_t)) - \ell(\mathbf{w}, (x_t, y_t)) + \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{w}\|^2 \right) + \tau_t^2 \|\nabla_t\|^2 \\ &\leq (1 - \lambda\tau_t) \|\mathbf{w}_t - \mathbf{w}\|^2 - 2\tau_t (\ell(\mathbf{w}_t, (x_t, y_t)) - \ell(\mathbf{w}, (x_t, y_t))) + \tau_t^2 \|\nabla_t\|^2. \end{aligned}$$

This follows that

$$\ell(\mathbf{w}_t, (x_t, y_t)) - \ell(\mathbf{w}, (x_t, y_t)) \leq \frac{\tau_t^{-1} - \lambda}{2} \|\mathbf{w}_t - \mathbf{w}\|^2 - \frac{1}{2\tau_t} \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 + \frac{\tau_t G^2}{2}.$$

We have, by summing  $t = 1$  to  $T$ , and setting  $\tau_t = 1/(\lambda t)$  with  $1/\tau_0 = 0$ ,

$$2 \sum_{t=1}^T (\ell(\mathbf{w}_t, (x_t, y_t)) - \ell(\mathbf{w}, (x_t, y_t))) \leq \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}\|^2 \left( \frac{1}{\tau_t} - \frac{1}{\tau_{t-1}} - \lambda \right) + G^2 \sum_{t=1}^T \tau_t = G^2 \sum_{t=1}^T \frac{1}{\lambda t} \leq \frac{G^2}{\lambda} (1 + \ln T),$$

which completes the proof.  $\square$

**Lemma C.5.** *Let  $S_T = \{(x_1, y_1), \dots, (x_T, y_T)\}$  be a sample. For kernel function  $\mathcal{K}(\mathbf{x}, \mathbf{x}) \leq r^2$ . we have  $T$  classifiers  $h_1, \dots, h_T$  generated by online kernel learning with  $\ell_t(h) = \max\{1 - y_t h(\mathbf{x}_t), 0\} + \lambda \|h\|_{\mathbb{H}}^2/2$ . We have*

$$\frac{1}{T} \sum_{t=1}^T \ell_t(h_t) - \frac{1}{T} \sum_{t=1}^T \ell_t(h_*) \leq \frac{4r^2(1 + \ln T)}{\lambda T},$$

where  $h_* = \arg \min_{h \in \mathbb{H}} \{\sum_{t=1}^T \ell_t(h)/T\}$ .

*Proof.* By setting  $\tau_t = 1/(\lambda t)$  and  $h_0 = 0$ , we rewrite the update rule as

$$h_t = \left(1 - \frac{1}{t}\right) h_{t-1} - \frac{1}{\lambda t} g_t \quad \text{with} \quad g_t = \mathbb{I}_{[y_t h(\mathbf{x}_t) < 1]} y_t \varphi(\mathbf{x}_t), \quad (26)$$

where  $\varphi(\mathbf{x}_t)$  is the feature mapping of  $\mathcal{K}$ . Hence, we have

$$h_t = \frac{1}{\lambda t} \sum_{i=1}^t g_i \quad \text{and} \quad \|h_t\|_{\mathbb{H}} \leq \frac{r}{\lambda},$$

from  $\prod_{j=i+1}^t (1 - 1/j) = i/t$  for  $i \leq t - 1$ . We complete the proof from Eqn. (26) and Lemma C.4.  $\square$

From Lemma C.5, we have the following corollary.

**Corollary C.6.** *For online kernel learning in the previous stage (Figure 1), we have*

$$\|\mathbf{w}_{T_1}^{[1]}\|_2 \leq \frac{r}{\lambda} \quad \text{and} \quad \left\| \frac{1}{T_1} \sum_{t=1}^{T_1} \mathbf{w}_t^{[1]} \right\|_2 \leq \frac{r}{\lambda}.$$

**Lemma C.7.** *Given  $S_n = \{(x_1, y_1), \dots, (x_T, y_T)\}$ , and for a kernel  $\mathcal{K}$  bounded by  $r^2$  and a classifier  $h_0 \in \mathbb{H}$ , let  $h_1, \dots, h_T$  be classifiers generated by online kernel learning with  $\ell(h, (\mathbf{x}, y)) = \max\{1 - y h(\mathbf{x}), 0\} + \lambda \|h\|_{\mathbb{H}}^2/2$  and  $\lambda > 0$ . For  $h_* = \arg \min_{h \in \mathbb{H}} \sum_{t=1}^T \ell_t(h)$ , we have*

$$\frac{1}{T} \sum_{t=1}^T \ell_t(h_t) - \frac{1}{T} \sum_{t=1}^T \ell_t(h_*) \leq \frac{2r \|h_0 - h_*\|_{\mathbb{H}}}{\sqrt{T}}.$$

*Proof.* For the norm in RKHS, we have

$$\|h_{t+1} - h_*\|_{\mathbb{H}}^2 = \|h_t - \eta \nabla \ell_t(h_t) - h_*\|_{\mathbb{H}}^2 = \|h_t - h_*\|_{\mathbb{H}}^2 + \eta^2 \|\nabla \ell_t(h_t)\|_{\mathbb{H}}^2 - 2\eta \nabla \ell_t(h_t)^\top (h_t - h_*),$$

and this follows that, from convex loss function  $\ell_t(h_t) - \ell_t(h_*) \leq \nabla \ell_t(h_t)^\top (h_t - h_*)$ ,

$$\ell_t(h_t) - \ell_t(h_*) \leq \frac{\|h_t - h_*\|_{\mathbb{H}}^2 - \|h_{t+1} - h_*\|_{\mathbb{H}}^2}{2\eta} + \frac{\eta}{2} \|\nabla \ell_t(h_t)\|_{\mathbb{H}}^2.$$

We have, by summing from  $t = 0$  to  $T - 1$ ,

$$\sum_{t=1}^T (\ell_t(h_t) - \ell_t(h_*)) \leq \frac{\|h_0 - h_*\|_{\mathbb{H}}^2}{2\eta} + 2\eta r^2 T,$$

which completes the proof by setting  $\eta_t = \|h_0 - h_*\|_{\mathbb{H}} / (r\sqrt{T})$ .  $\square$

**Lemma C.8.** Given samples  $S_1 = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_1}$  and  $S_2 = \{(\mathbf{x}_i, y_i)\}_{i=n_1+1}^{n_1+n_2}$  with  $\|\mathbf{x}\|_2 \leq r$ , denote by

$$\begin{aligned} \mathbf{w}_1^* &\in \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \hat{R}_1(\mathbf{w}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \ell(\mathbf{w}, (\mathbf{x}_i, y_i)) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right\}, \\ \mathbf{w}_2^* &\in \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \hat{R}_2(\mathbf{w}) = \frac{1}{n_2} \sum_{i=n_1+1}^{n_1+n_2} \ell(\mathbf{w}, (\mathbf{x}_i, y_i)) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right\}. \end{aligned}$$

For  $\delta \in (0, 1)$ , the following holds with probability at least  $1 - \delta$ ,

$$\hat{R}_1(\mathbf{w}_1^*) - \hat{R}_2(\mathbf{w}_2^*) \leq \frac{r^2}{\lambda} \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right) \ln\left(\frac{1}{\delta}\right)}.$$

*Proof.* We introduce a new function

$$f(S_1, S_2) = \hat{R}_1(\mathbf{w}_1^*) - \hat{R}_2(\mathbf{w}_2^*),$$

and consider the sample  $S'_1 = S_1 \setminus \{(\mathbf{x}_k, y_k)\} \cup \{(\mathbf{x}'_k, y'_k)\}$  for  $k \in [n_1]$ . From Cauchy-Schwarz inequality and 1-Lipschitz hinge loss, we have

$$|f(S'_1, S_2) - f(S_1, S_2)| = \left| \hat{R}_1(\mathbf{w}_1^*) - \hat{R}'_1(\mathbf{w}'_1^*) \right| \leq r \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2 + \frac{\lambda}{2} \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2 \cdot \|\mathbf{w}_1^* + \mathbf{w}'_1^*\|_2 \leq 2r \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2,$$

where the norm of optimal classifiers satisfy

$$\|\mathbf{w}_1^*\|_2 \leq \frac{r}{\lambda}, \quad \|\mathbf{w}'_1^*\|_2 \leq \frac{r}{\lambda} \quad \text{and} \quad \|\mathbf{w}_1^* + \mathbf{w}'_1^*\|_2 \leq \frac{2r}{\lambda},$$

from the KKT condition as in the proof of Theorem 3.3. From strong convexity, we have

$$\hat{R}_1(\mathbf{w}'_1^*) \geq \hat{R}_1(\mathbf{w}_1^*) + \frac{\lambda}{2} \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2^2 \quad \text{and} \quad \hat{R}'_1(\mathbf{w}_1^*) \geq \hat{R}'_1(\mathbf{w}'_1^*) + \frac{\lambda}{2} \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2^2, \quad (27)$$

and this follows that,

$$\begin{aligned} \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2^2 &\leq \frac{1}{\lambda} \left( \hat{R}_1(\mathbf{w}'_1^*) - \hat{R}_1(\mathbf{w}_1^*) - \hat{R}'_1(\mathbf{w}_1^*) + \hat{R}'_1(\mathbf{w}'_1^*) \right) \\ &= \frac{1}{\lambda n_1} \left( \ell(\mathbf{w}_1^*, (\mathbf{x}_k, y_k)) - \ell(\mathbf{w}'_1^*, (\mathbf{x}_k, y_k)) + \ell(\mathbf{w}_1^*, (\mathbf{x}'_k, y'_k)) - \ell(\mathbf{w}'_1^*, (\mathbf{x}'_k, y'_k)) \right) \leq \frac{r}{\lambda n_1} \|\mathbf{w}_1^* - \mathbf{w}'_1^*\|_2. \end{aligned}$$

Hence, we have, from Eqn. (27),

$$|f(S_1, S_2) - f(S'_1, S_2)| \leq \frac{2r^2}{\lambda n_1}.$$

We could make a similar analysis for  $S_2$ . This completes the proof from Lemma A.4.  $\square$

**Proof of Theorem 4.3.** For  $k = 1$  and  $k = 2$ , we introduce some notations as follows:

$$\begin{aligned} \hat{R}_{T_1}(\mathbf{w}^{[k]}) &= \frac{1}{T_1} \sum_{t=1}^{T_1} \ell(\mathbf{w}^{[k]}, (\mathbf{x}_t^{[k]}, y_t)) + \frac{\lambda}{2} \|\mathbf{w}^{[k]}\|_2^2, \\ \hat{R}_{T_e}(\mathbf{w}^{[k]}) &= \frac{1}{T_e} \sum_{t=T_1+1}^{T_1+T_e} \ell(\mathbf{w}^{[k]}, (\mathbf{x}_t^{[k]}, y_t)) + \frac{\lambda}{2} \|\mathbf{w}^{[k]}\|_2^2, \\ \hat{R}_{T_2}(\mathbf{w}^{[k]}) &= \frac{1}{T_2} \sum_{t=T_1+T_e+1}^{T_1+T_e+T_2} \ell(\mathbf{w}^{[k]}, (\mathbf{x}_t^{[k]}, y_t)) + \frac{\lambda}{2} \|\mathbf{w}^{[k]}\|_2^2. \end{aligned}$$

Denote by  $\ell_t(\mathbf{w}^{[k]}) = \ell(\mathbf{w}^{[k]}, (\mathbf{x}_t^{[k]}, y_t)) + \lambda \|\mathbf{w}^{[k]}\|_2^2$ . From the i.i.d assumption, it is natural to consider samples on new feature space  $\mathcal{X}^{[2]}$  in the previous stage and on old feature space  $\mathcal{X}^{[1]}$  in the current stages respectively.

We have, from Lemma C.5 and strong convexity,

$$\hat{\mathcal{L}}_{T_2}^{[2]} - \mathcal{L}_{T_2}^{[2]}(\mathbf{w}_*^{[2]}) \leq 2r \sqrt{\frac{2}{\lambda T_2} \left( \hat{R}_{T_2}^{[2]}(\mathbf{w}_{T_1+T_e}^{[2]}) - \hat{R}_{T_2}^{[2]}(\mathbf{w}_*^{[2]}) \right)}, \quad (28)$$

and we also have  $\hat{R}_{T_2}(\mathbf{w}_{T_1+T_e}^{[2]}) - \hat{R}_{T_2}(\mathbf{w}_{T_2*}^{[2]}) = Z_1 + Z_2 + Z_3 + Z_4 + Z_5 + Z_6 + Z_7$ , where

$$\begin{aligned} Z_1 &= \hat{R}_{T_2}(\mathbf{w}_{T_1+T_e}^{[2]}) - \hat{R}_{T_e}(\mathbf{w}_{T_1+T_e}^{[2]}), \quad Z_2 = \hat{R}_{T_e}(\mathbf{w}_{T_1+T_e}^{[2]}) - \hat{R}_{T_e}(\mathbf{w}_{T_1}^{[1]}), \quad Z_3 = \hat{R}_{T_e}(\mathbf{w}_{T_1}^{[1]}) - R_{T_1}(\mathbf{w}_{T_1}^{[1]}), \\ Z_4 &= R_{T_1}(\mathbf{w}_{T_1}^{[1]}) - \frac{1}{T_1} \sum_{t=1}^{T_1} \ell_t(\mathbf{w}_t^{[1]}), \quad Z_5 = \frac{1}{T_1} \sum_{t=1}^{T_1} \ell_t(\mathbf{w}_t^{[1]}) - \hat{R}_{T_1}(\mathbf{w}_{T_1*}^{[1]}), \quad Z_6 = \hat{R}_{T_1}(\mathbf{w}_{T_1*}^{[1]}) - \hat{R}_{T_2}(\mathbf{w}_{T_2*}^{[1]}), \\ Z_7 &= \hat{R}_{T_2}(\mathbf{w}_{T_2*}^{[1]}) - \hat{R}_{T_2}(\mathbf{w}_{T_2*}^{[2]}). \end{aligned}$$

From the i.i.d assumption for evolving and current stage, the following holds with probability at least  $1 - \delta/6$ ,

$$\begin{aligned} Z_1 &\leq \mathbb{E}[Z_1] + \max_{(\mathbf{x}^{[2]}, y) \in \mathcal{X}^{[2]} \times \mathcal{Y}} \left\{ \ell(\mathbf{w}_{T_1+T_e}^{[2]}, (\mathbf{x}^{[2]}, y)) \right\} \sqrt{\left( \frac{1}{2T_e} + \frac{1}{2T_2} \right) \ln \frac{6}{\delta}} \\ &= \left( \frac{r^2}{\lambda} + 1 \right) \sqrt{\left( \frac{1}{2T_e} + \frac{1}{2T_2} \right) \ln \frac{6}{\delta}}, \end{aligned} \quad (29)$$

from Lemma C.1 and Corollary C.6. Similarly, the following holds with the probability as least  $1 - \delta/6$ ,

$$Z_3 \leq \mathbb{E}[Z_3] + \max_{(\mathbf{x}^{[2]}, y) \in \mathcal{X}^{[2]} \times \mathcal{Y}} \left\{ \ell(\mathbf{w}_{T_1+T_e}^{[2]}, (\mathbf{x}^{[2]}, y)) \right\} \sqrt{\frac{\ln(1/\delta_2)}{2T_e}} = \left( \frac{r^2}{\lambda} + 1 \right) \sqrt{\frac{\ln(1/\delta_2)}{2T_e}}. \quad (30)$$

From Lemma 4.2 and Theorem 3.4, the following holds with the probability at least  $1 - \delta/6$ ,

$$\begin{aligned} Z_2 &= \frac{1}{T_e} \sum_{t=T_1+1}^{T_1+T_e} \left( \ell(\mathbf{w}_{T_1}^{[1]}, (\mathbf{x}_t^{[1]}, y_t)) - \ell(\mathbf{w}_{T_1+T_e}^{[2]}, (\mathbf{x}_t^{[2]}, y_t)) \right) \\ &\leq \frac{1}{T_e} \sum_{t=T_1+1}^{T_1+T_e} \left| \langle \mathbf{w}_{T_1}^{[1]}, \mathbf{z}^{[1]}(\mathbf{x}_t^{[1]}) \rangle - \langle \mathbf{w}_{T_1+T_e}^{[2]}, \mathbf{z}^{[2]}(\mathbf{x}_t^{[2]}) \rangle \right| \\ &\leq \frac{r \hat{\mathcal{E}}(S_{T_e}^{[e]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})}{\lambda} \leq \frac{r}{\lambda} \left( \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) + c_1 r \sqrt{\frac{1}{T_e} \ln \frac{6}{\delta}} \right). \end{aligned} \quad (31)$$

From Lemma C.3, the following holds with the probability at least  $1 - \delta/6$ ,

$$Z_4 \leq \left( \frac{r^2}{\lambda} + 1 \right) \sqrt{\frac{\ln(6/\delta)}{2T_1}}, \quad (32)$$

and we have, from Lemma C.5,

$$Z_5 \leq \frac{4r^2 \ln(1 + T_1)}{\lambda T_1} \leq \frac{4r^2}{\lambda \sqrt{T_1}}. \quad (33)$$

From Lemma C.8, the following holds with a probability at least  $1 - \delta/6$ ,

$$Z_6 \leq \frac{r^2}{\lambda} \sqrt{\left( \frac{1}{T_1} + \frac{1}{T_2} \right) \ln \frac{6}{\delta}}. \quad (34)$$

From Theorem 3.3 and Theorem 3.4, the following holds with a probability at least  $1 - \delta/6$ ,

$$\begin{aligned}
 Z_7 &\leq \frac{1}{T_2} \sum_{t=T_1+T_e+1}^{T_1+T_e+T_2} \left| \langle \mathbf{w}_{T_2^*}^{[2]}, \mathbf{x}_t^{[2]} \rangle - \langle \mathbf{w}_{T_2^*}^{[1]}, \mathbf{x}_t^{[1]} \rangle \right| + \frac{\lambda}{2} \left\| \mathbf{w}_{T_2^*}^{[1]} - \mathbf{w}_{T_2^*}^{[2]} \right\|_2 \left\| \mathbf{w}_{T_2^*}^{[1]} + \mathbf{w}_{T_2^*}^{[2]} \right\|_2 \\
 &\leq \frac{r}{\lambda} \hat{\mathcal{E}}(S_{T_2}^{[2]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) + \frac{r}{\lambda} \sqrt{2r \hat{\mathcal{E}}(S_{T_2}^{[2]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})} + r \left\| \mathbf{w}_{T_2^*}^{[1]} - \mathbf{w}_{T_2^*}^{[2]} \right\|_2 \\
 &\leq \frac{r}{\lambda} \hat{\mathcal{E}}(S_{T_2}^{[2]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) + \frac{2r}{\lambda} \sqrt{2r \hat{\mathcal{E}}(S_{T_2}^{[2]}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})} \\
 &\leq \frac{r}{\lambda} \left( \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]}) + 2\sqrt{2r \mathcal{E}(\mathcal{D}, \mathcal{K}^{[1]}, \mathcal{K}^{[2]})} \right) + O\left( \sqrt[4]{\frac{\ln(6/\delta)}{T_2}} \right). \tag{35}
 \end{aligned}$$

From Eqns. (29)-(35) and union bounds, the following holds with probability at least  $1 - \delta$ ,

$$\hat{R}_{T_2}(\mathbf{w}_{T_1+T_e}^{[2]}) - \hat{R}_{T_2}(\mathbf{w}_{T_2^*}^{[2]}) \leq \frac{r^2}{\lambda} \left[ c \left( \frac{1}{\sqrt{T_1}} + \frac{1}{\sqrt{T_e}} + \frac{1}{\sqrt[4]{T_2}} \right) \ln\left(\frac{6}{\delta}\right) + \frac{\mathcal{E}}{r} + \sqrt{\frac{2\mathcal{E}}{r}} \right], \tag{36}$$

for some constant  $c > 0$  with  $(1 - \delta/6)^6 \geq 1 - \delta$ . We complete the proof by combining Eqn. (28) and Eqn. (36).  $\square$

#### C.4. Proof of Theorem 4.4

For simplicity, we reindex the time-step of samples in the current stage as  $t = 1, \dots, T_2$ . Denote by  $\ell_{t,i}$  the loss for the  $i$ -th base learner at the  $t$ -th iteration with  $i \in [2]$ , and  $\mathcal{L}_{t,i}$  is the corresponding cumulative loss. We define the potential function

$$\Phi_t = \frac{1}{\gamma} \ln \left( \sum_{i=1}^2 \exp(-\gamma \mathcal{L}_{t,i}) \right),$$

and this follows that, from  $e^{-x} \leq 1 - x + x^2$  and  $\ln(1+x) \leq x$ ,

$$\begin{aligned}
 \Phi_t - \Phi_{t-1} &= \frac{1}{\gamma} \left( \frac{\exp(-\gamma \mathcal{L}_{t,i})}{\sum_{i=1}^2 \exp(-\gamma \mathcal{L}_{t-1,i})} \right) = \frac{1}{\gamma} \left( \sum_{i=1}^2 \omega_{t,i} \exp(-\gamma \ell_{t,i}) \right) \\
 &\leq \frac{1}{\gamma} \left( \sum_{i=1}^2 \omega_{t,i} (1 - \gamma \ell_{t,i} + \gamma^2 \ell_{t,i}^2) \right) = \frac{1}{\gamma} \ln \left( 1 - \gamma \langle \boldsymbol{\omega}_t, \boldsymbol{\ell}_t \rangle + \gamma^2 \sum_{i=1}^2 \omega_{t,i} \ell_{t,i}^2 \right) \leq \langle \boldsymbol{\omega}_t, \boldsymbol{\ell}_t \rangle + \gamma^2 \sum_{i=1}^2 \omega_{t,i} \ell_{t,i}^2,
 \end{aligned}$$

where  $\boldsymbol{\omega}_t = (\omega_{t,1}, \omega_{t,2})$ ,  $\boldsymbol{\ell}_t = (\ell_{t,1}, \ell_{t,2})$ , and the last equality holds from Eqn. (12). We have, by summing over  $t \in [T_2]$ ,

$$\begin{aligned}
 \sum_{t=1}^{T_2} \langle \boldsymbol{\omega}_t, \boldsymbol{\ell}_t \rangle &\leq \Phi_0 - \Phi_{T_2} + \gamma \sum_{t=1}^{T_2} \sum_{i=1}^2 \omega_{t,i} \ell_{t,i}^2 \\
 &\leq \frac{\ln 2}{\gamma} - \frac{1}{\gamma} \ln(\exp(-\gamma \mathcal{L}_{T_2, i^*})) + \gamma \sum_{t=1}^{T_2} \sum_{i=1}^2 \omega_{t,i} \ell_{t,i}^2 \leq \frac{\ln 2}{\gamma} + \mathcal{L}_{T_2, i^*} + \gamma \sum_{t=1}^{T_2} \sum_{i=1}^2 \omega_{t,i} \ell_{t,i}^2,
 \end{aligned}$$

where  $\mathcal{L}_{T_2, i^*} = \min_{i \in \{1,2\}} \mathcal{L}_{T_2, i}$ . We have, by rearranging and from Theorem 4.3,

$$\sum_{t=1}^{T_2} \langle \boldsymbol{\omega}_t, \boldsymbol{\ell}_t \rangle - \min_{i=1,2} \mathcal{L}_{T_2, i} \leq \frac{\ln 2}{\gamma} + \gamma \sum_{t=1}^{T_2} \sum_{i=1}^2 \omega_{t,i} \ell_{t,i}^2 \leq \frac{\ln 2}{\gamma} + \gamma T_2 \left( 1 + \frac{3r^2}{2\lambda} \right),$$

which completes the proof by setting  $\gamma = \sqrt{\ln 2 / ((1 + 3r^2/2\lambda)T_2)}$ .  $\square$