# A Multilayer Framework for Online Metric Learning

Wenbin Li†, Yanfang Liu†, Jing Huo, Yinghuan Shi, Yang Gao*, *Senior Member, IEEE*, Lei Wang, *Senior Member, IEEE*, and Jiebo Luo, *Fellow, IEEE*

*Abstract*—Online metric learning has been widely applied in classification and retrieval. It can automatically learn a suitable metric from data by restricting similar instances to be separated from dissimilar instances with a given margin. However, the existing online metric learning algorithms have limited performance in real-world classifications, especially when data distributions are complex. To this end, this paper proposes a multilayer framework for online metric learning to capture the nonlinear similarities among instances. Different from the traditional online metric learning, which can only learn one metric space, the proposed *Multi-Layer Online Metric Learning (MLOML)* takes an online metric learning algorithm as a metric layer and learns multiple hierarchical metric spaces, where each metric layer follows a nonlinear layers for the complicated data distribution. Moreover, the forward propagation (FP) strategy and backward propagation (BP) strategy are employed to train the hierarchical metric layers. To build a metric layer of the proposed MLOML, a new *Mahalanobis-based Online Metric Learning (MOML)* algorithm is presented based on the passive-aggressive strategy and one-pass triplet construction strategy. Furthermore, in a progressively and nonlinearly learning way, MLOML has a stronger learning ability than traditional online metric learning in the case of limited available training data. To make the learning process more explainable and theoretically guaranteed, theoretical analysis is provided. The proposed MLOML enjoys several nice properties, indeed learns a metric progressively, and performs better on the benchmark datasets. Extensive experiments with different settings have been conducted to verify these properties of the proposed MLOML.

*Index Terms*—Online Metric Learning, Metric Layer, Passive-Aggressive Strategy, Nonlinearity, Interpretability

## I. INTRODUCTION

Learning a meaningful and quality metric on the original instances is crucial to many classification and retrieval applications. In recent decades, many metric learning methods based on Mahalanobis distance function and bilinear similarity function have been proposed. Mahalanobis distance-based

W. Li, J. Huo and Y. Gao are with the State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China. (e-mail: liwenbin@nju.edu.cn; huojing@nju.edu.cn; gaoy@nju.edu.cn).

Y. Liu is with the State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China, and the College of Mathematics and Information Engineering, Longyan University, Longyan 364012, China. (e-mail: liuyanfang003@163.com).

Y. Shi is with the State Key Laboratory for Novel Software Technology and National Institute of Healthcare Data Science, Nanjing University, Nanjing 210023, China. (e-mail: syh@nju.edu.cn).

L. Wang is with the School of Computing and Information Technology, University of Wollongong, Australia (e-mail: leiw@uow.edu.au).

J. Luo is with the Department of Computer Science, University of Rochester, Rochester, NY 14611, USA (e-mail: jluo@cs.rochester.edu).

† Wenbin Li and Yanfang Liu contributed equally as co-first authors.
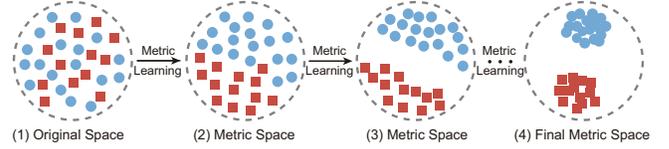
* Corresponding author: Yang Gao.

Fig. 1. An illustration of online hierarchical metric learning by learning new metric spaces progressively.

methods [1]–[8] refer to learning a real-valued distance matrix with a symmetric positive semi-definite (PSD) constraint. Bilinear similarity-based methods [9]–[11] aim to learn a form of bilinear similarity matrix without the PSD constraint. Moreover, there are two kinds of constraints, *i.e.*, pairwise and triplet constraints, that have been widely used in these metric learning methods. A pairwise constraint consists of two similar or dissimilar instances, while a triplet constraint is of the form $\langle x, x^+, x^- \rangle$, where instance $x$ is similar to instance $x^+$, but is dissimilar to instance $x^-$.

In many real-world applications, a lot of data is streaming data which is continuously produced in time, such as wind power data [12], credit data [13], and ADs click data [14]. Online learning algorithm investigates how to learn in a streaming setting [15]–[17]. Therefore, metric learning algorithms should be able to learn metric in an online manner, *i.e., online metric learning (OML)*. In fact, multiple OML algorithms have been proposed [9], [18]–[22]. However, the existing OML algorithms mainly pay attention to rapid constraints construction [9], [21] or low update complexity [19]–[21], while rarely consider the learning ability in the case where all the labeled streaming data cannot be observed. In addition, most of these OML algorithms only learn one linear metric space which cannot learn well refined metrics for a complicated nonlinear data distribution.

To tackle the above limitations of the existing OML algorithms, we propose a *Multi-Layer Online Metric Learning (MLOML)* framework, which is nonlinear and explainable. In this framework, we attempt to design a *metric-algorithm-based layer*, which is stacked by several OML algorithms along with the corresponding nonlinear layers (*e.g.,* ReLU, Sigmoid, tanh). In this way, our proposed MLOML is able to learn a progressively refined metric space by learning another new metric in the former learnt feature space (see Fig. 1). Specifically, in MLOML, one OML algorithm is taken as a metric layer, followed by a nonlinear layer (*i.e.,* ReLU, Sigmoid, or tanh). These two layers are repeatedly stacked multiple times. It is worth noting that each metric layer in MLOML

is a relatively independent OML algorithm, as a result, the parameters of each metric layer can be innovatively updated according to its own local loss during forward propagation (FP). It means that it is possible to train such a metric-algorithm-based layer by only using the FP strategy. The advantages of this FP updating are: (1) the parameter updating is immediate, unlike the delayed updating of the commonly used backward propagation (BP); (2) when additional BP is adopted, FP updating can vastly accelerate the convergence. Note that the second advantage has a similar effect of layer-wise unsupervised pre-training [23]–[25]. However, there are fundamental differences. The existing layer-wise training is unsupervised and only acts as a pre-training operation (or a regularizer [25]), which is not end-to-end. In contrast, the FP updating in the proposed MLOML is supervised and serves the primary training mode rather than a pre-training role (elaborated in Section IV-D), which is end-to-end. In fact, these two updating strategies (*i.e.,* FP and BP) can be combined to train this metric-algorithm-based layer. Ideally, FP updating can explore new feature spaces sequentially, while BP updating can amend the exploration in further.

Furthermore, to achieve a low computational cost when performing MLOML, a new general *Mahalanobis-based Online Metric Learning (MOML)* algorithm is proposed as the metric layer of MLOML. Since all the labeled streaming data cannot be observed in online manner, MOML uses the one-pass triplet construction [21] instead of triple constraints obtained in advance. Simultaneously, MOML has a convex objective function inspired by passive-aggressive learning and enjoys a closed-form solution at each step. We also derive a theoretical regret bound for MOML to prove its convergence. Through stacking MOML hierarchically, the ability of learning feature representation progressively can be explainable and guaranteed.

Our main contributions can be summarized as follows:

- A *Multi-Layer Online Metric Learning (MLOML)* framework is developed for streaming data through forward propagation (FP) strategy or backward propagation (BP) strategy, such that a metric space is learned progressively and deeply, *i.e.,* exploring and learning a new metric in a nonlinear transformation space sequentially.
- Taking *Mahalanobis-based Online Metric Learning (MOML)* as a metric layer, MLOML has theoretical guarantees so that the classification performance will be improved or at least well maintained as the depth of the layers increases.
- MLOML is simple yet effective, as verified by extensive experiments.

## II. RELATED WORK

Online metric learning enjoys several practical and theoretical advantages, making it widely studied and applied in data mining tasks, which can be roughly divided into two categories: Mahalanobis distance-based and bilinear similarity-based methods. In bilinear similarity-based methods, Online Algorithm for Scalable Image Similarity (OASIS) [9] is proposed based on Passive-Aggressive (PA) algorithm, aiming to learn a similarity metric without PSD constraint. Following a similar setting as OASIS, Sparse Online Metric Learning (SOML) [10] learns a diagonal matrix instead of a full matrix to deal with the high-dimensional data. Online Multiple Kernel Similarity (OMKS) [11] has been proposed to handle the multi-modal data. Through adopting an off-diagonal $\ell_1$ norm to the similarity matrix, Sparse Online Relative Similarity (SORS) [26] can obtain a sparse result. Online Similarity Learning via Low Rank and Group Sparsity (OSLLR-GS) [27] is designed to address the over-fitting problem for big data by detecting the feature redundant in the metric matrix and constraining the remaining matrix to a low rank space.

In the second kind of Mahalanobis distance-based methods, Pseudo-Metric Online Learning Algorithm (POLA) [18] introduces the successive projection operation to update a pseudo-metric and map it onto a positive semi-definite cone. As an extended version of Information Theoretic Metric Learning-Online (ITML-Online) [4], LogDet Exact Gradient Online (LEGO) [19] updates a manalanobis metric based on LogDet regularization and gradient descent. Regularized Distance Metric Learning (RDML) [20] with appropriate constraints has a provable regret bound. Mirror Descent for Metric Learning (MDML) [28] is an unified approach which updates a manalanobis metric by composite objective mirror descent. Bellet and Habrard [29] utilize an adaptation of the notion of algorithmic robustness [30] to derive generalization bounds for metric learning. Low-Rank Similarity Metric Learning (LRSM) [31] uses SVD-based projection to solve the challenging high-dimensional learning task, and then employs Alternating Direction Method of Multipliers (ADMM) [32] to optimize the model. Based on the PA algorithm [15], Scalable Large Margin Online Metric Learning (SLMOML) [33] adopts the LogDet divergence to maintain the closeness between two successively learned Mahalanobis matrices, and utilizes the hinge loss to enforce a large margin between relatively dissimilar samples. Fast Low-Rank Metric Learning (FLRML) [34] is an unconstrained optimization on the Stiefel manifold to handle datasets with both high dimensions and large numbers of instances. Large-Margin Distance Metric Learning (LMDML) [35] employs the principle of margin maximization and stochastic gradient descent method to learn the distance metric with PSD constraint. These methods almost assume that the pairwise or triplet constraints can be obtained in advance except RDML, which exactly receives two adjacent samples as a pairwise constraints at each time. In view of adapting the pairwise and triplet constraints to streaming data, Li *et al.* [21] present a one-pass triplet construction strategy and design OPML and COPML algorithms with low time complexity. By incorporating with a smoothed Wasserstein metric distance, Evolving Metric Learning (EML) [36] can handle the instance and feature evolutions simultaneously.

However, the above mentioned metric algorithms only learn one linear metric space which cannot learn well refined metrics for a complicated nonlinear data distribution. In order to solve this issue, we propose MLOML which is developed based on a newly designed Mahalanobis-based online metric learning (MOML). Compared with the above OML algorithms, MLOML has the following advantages: (1) MLOML is hi-
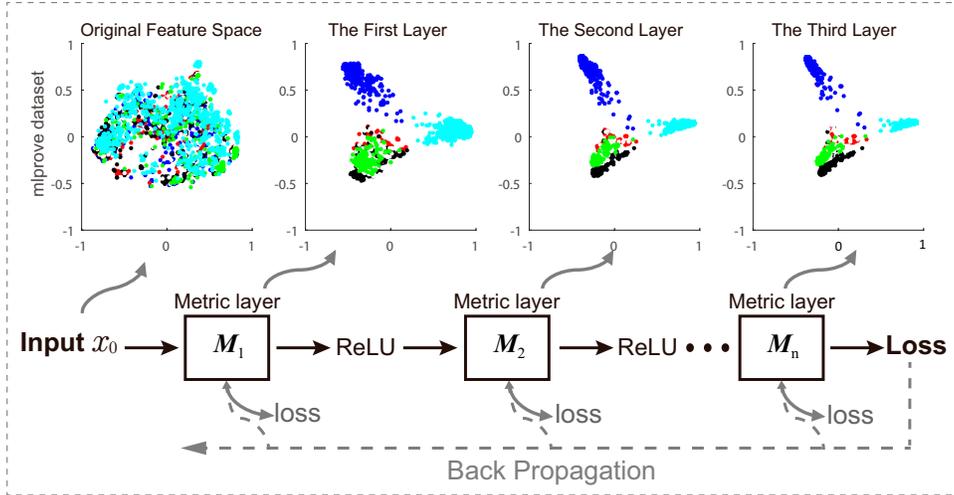
Fig. 2. Framework of the proposed multi-layer online metric learning (MLOML), where each metric layer $M_n$ is an online metric learning algorithm. Here we take three metric layers and two ReLU layers as an example, *i.e.*, $n = 3$.

erarchical and can learn feature representation progressively (*i.e.,* better and better) through FP and BP strategies; (2) MLOML not only has theoretical guarantees by stacking MOML algorithm as its metric layer, but also is nonlinear by employing nonlinear functions; (3) MLOML enjoys a stronger learning ability than traditional OML algorithms with the same amount of data.

## III. OUR FRAMEWORK

Our goal is to design a novel multi-layer online metric learning framework (MLOML) for streaming data, which is stacked by metric-algorithm-based layers along with the corresponding nonlinear layers (*e.g.,* ReLU, Sigmoid, tanh). The framework is illustrated in Fig. 2.

### A. Multi-Layer Online Metric Learning

In this section, we propose and explain our MLOML in detail. MLOML is made up of multiple metric layers and nonlinear layers, in which one metric layer is an OML algorithm and one nonlinear layer is ReLU, Sigmoid or tanh. To ensure the progressively learning ability of MLOML, we should guarantee the convexity of each metric layer, which can easily guarantee the convergence of each layer. Therefore, a new *Mahalanobis-based OML algorithm (MOML)* is designed specifically. MOML has a convex objective function and enjoys a closed-form solution. Moreover, a tight regret bound of MOML is also proved (see Theorem 2).

Specifically, MOML is built on triplet-based constraints $\langle x, x^+, x^- \rangle$, where instance $x$ is similar to instance $x^+$, but is dissimilar to instance $x^-$, and these triplets can encode the proximity comparison information. Therefore, MLOML is also learnt from triplet constraints. For computational efficiency, a one-pass triplet construction strategy presented by OPML [21] is also employed to construct triplets rapidly which can solve the inability to observe all the streaming data and its labels. In brief, for each new sample, two latest samples from both the same and different classes in the past samples are selected.

By using this strategy, triplets can be constructed in an online manner. There are two types of layers in MLOML, that are OML layer and non-linear layer, where MLOML-r, MLOML-s, MLOML-t correspond to MLOML with the ReLU, Sigmoid, tanh layers respectively. If we design a three-layer MLOML-r model, there should be three OML layers in this model. Moreover, each OML layer is followed by a ReLU layer except the last OML layer (*i.e.*, the third OML layer). This principle is also satisfied for the MLOML-s and MLOML-t models.

A loss layer can also be added, which can give a global adjustment of the entire metric-algorithm-based model via backward propagation. To adequately use the effect of each local metric layer, the local loss is also utilized to update all the former layers (*i.e.,* the loss of the $i$-th metric layer can be used to update the 1-st to the $(i-1)$-th layers). In this way, vanishing gradient problem can also be alleviated. The novel loss function can be formulated as follows:

$$\Gamma = \frac{1}{2}\Gamma_{triplet} + \sum_{i=1}^{n} w_i \Gamma_{local}^i + \frac{\lambda}{2} \sum_{i=1}^{n} \|\boldsymbol{L}^i\|_F^2, \qquad (1)$$

where $\Gamma_{triplet} = [\|\boldsymbol{x}_t^{(n)} - \boldsymbol{x}_p^{(n)}\|_2^2 + 1 - \|\boldsymbol{x}_t^{(n)} - \boldsymbol{x}_q^{(n)}\|_2^2]_+$ indicates the triplet loss of the final output of the model (where $[z]_+ = \max(0, z)$), $\Gamma_{local}^i$ denotes the local loss of the $i$-th OML layer (*i.e.,* Eq.(4)), and $\|\boldsymbol{L}^i\|_F^2$ represents the Frobenius norm of parameter matrix $\boldsymbol{L}^i$, *i.e.,* the transformation matrix learnt in the $i$-th OML layer. Moreover, $\lambda$ is the predefined hyper-parameter. While $w_i$, the weight of the $i$-th metric layer can be learnt by SGD during training phase, indicating the importance of each metric layer.

**A New Mahalanobis-based OML (MOML):** To build a metric layer of the proposed MLOML, a new OML algorithm named MOML is presented, which can act as a representative of Mahalanobis-based algorithms. Note that, in essence, MLOML can be constructed by other Mahalanobis-based algorithms. However, with MOML as a building component, MLOML enjoys better theoretical properties. The goal of MOML, learnt from triplet constraints, is to learn a Maha-

lanobis distance function $D$ that satisfies the following large margin constraint:

$$D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_q) > D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_p) + r, \forall \boldsymbol{x}, \boldsymbol{x}_p, \boldsymbol{x}_q \in \mathbb{R}^d, \quad (2)$$

where $\boldsymbol{x}$ and $\boldsymbol{x}_p$ belong to the same class, while $\boldsymbol{x}$ and $\boldsymbol{x}_q$ come from different classes. $D_{\boldsymbol{M}}(\boldsymbol{x}_1, \boldsymbol{x}_2) = (\boldsymbol{x}_1 - \boldsymbol{x}_2)^\top \boldsymbol{M}(\boldsymbol{x}_1 - \boldsymbol{x}_2)$, where $\boldsymbol{M} \in \mathbb{R}^{d \times d}$ is a positive semi-definite parameter matrix. Also, $r$ is the margin. Naturally, the hinge loss (*i.e.*, $r = 1$) can be employed as below,

$$\ell(\boldsymbol{M}, \langle \boldsymbol{x}, \boldsymbol{x}_p, \boldsymbol{x}_q \rangle) = \max(0, 1 + D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_p) - D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_q)). \quad (3)$$

In a sequential manner, given a triplet $\langle \boldsymbol{x}_t, \boldsymbol{x}_p, \boldsymbol{x}_q \rangle$ at the $t$-th time step. Inspired by the Passive-Aggressive (PA) algorithms (*i.e.*, a family of margin based online learning algorithms) [15], we design a convex objective function at each time step as follows,

$$\Gamma = \arg\min_{\boldsymbol{M} \succeq 0} \frac{1}{2} \|\boldsymbol{M} - \boldsymbol{M}_{t-1}\|_F^2 + \gamma \Big[1 + D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_p) - D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_q)\Big]_+$$
$$(D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_p) = (\boldsymbol{x} - \boldsymbol{x}_p)^\top \boldsymbol{M}(\boldsymbol{x} - \boldsymbol{x}_p), \; D_{\boldsymbol{M}}(\boldsymbol{x}, \boldsymbol{x}_q) = (\boldsymbol{x} - \boldsymbol{x}_q)^\top \boldsymbol{M}(\boldsymbol{x} - \boldsymbol{x}_q))$$
$$= \arg\min_{\boldsymbol{M} \succeq 0} \frac{1}{2} \|\boldsymbol{M} - \boldsymbol{M}_{t-1}\|_F^2 + \gamma \Big[1 + \mathrm{Tr}(\boldsymbol{M} \boldsymbol{A}_t)\Big]_+, \quad (4)$$

where $\| \cdot \|_F$ is Frobenius norm, $[z]_+ = \max(0, z)$ is the hinge loss, $\mathrm{Tr}(\cdot)$ denotes the trace operation, $\gamma$ is the regularization parameter and $\boldsymbol{A}_t = (\boldsymbol{x}_t - \boldsymbol{x}_p)(\boldsymbol{x}_t - \boldsymbol{x}_p)^\top - (\boldsymbol{x}_t - \boldsymbol{x}_q)(\boldsymbol{x}_t - \boldsymbol{x}_q)^\top$. We can easily get that $\Gamma$ is a convex function for $\boldsymbol{M}$, because $\mathrm{Tr}(\boldsymbol{M} \boldsymbol{A}_t)$ is a linear function of $\boldsymbol{M}$ which is convex, the hinge loss function $[1 + z]_+$ is convex (not continuous at $z = -1$), and $\| \cdot \|_F$ and the domain $\boldsymbol{M} \succeq 0$ are convex too. It can be shown that an optimal solution can be found within the domain $\boldsymbol{M} \succeq 0$ by properly setting the value of $\gamma$. Thus, we can get the optimal solution of Eq. (4) by calculating the gradient $\frac{\partial \Gamma(\boldsymbol{M})}{\partial \boldsymbol{M}} = 0$:

$$\frac{\partial \Gamma(\boldsymbol{M})}{\partial \boldsymbol{M}} = \begin{cases} \boldsymbol{M} - \boldsymbol{M}_{t-1} + \gamma \boldsymbol{A}_t = 0 & [z]_+ > 0 \\ \boldsymbol{M} - \boldsymbol{M}_{t-1} = 0 & [z]_+ = 0. \end{cases} \quad (5)$$
$$s.t. \quad \boldsymbol{M} \succeq 0$$

According to Theorem 1 (presented below), with a proper $\gamma$, the semi-positive definitiveness of $\boldsymbol{M}$ can be guaranteed. Thus, at the $t$-th time step, the parameter matrix $\boldsymbol{M}_t$ can be updated as below,

$$\boldsymbol{M}_t = \begin{cases} \boldsymbol{M}_{t-1} - \gamma \boldsymbol{A}_t & [z]_+ > 0 \\ \boldsymbol{M}_{t-1} & [z]_+ = 0. \end{cases} \quad (6)$$

From Eq. (6), we can see that the time complexity of MOML is $O(d^2)$ at each time step. Using MOML as the base metric layer of MLOML has the following advantages: (1) the objective function of MOML is convex and enjoys a closed-form solution, which is beneficial to theoretical analysis; (2) without loss of generality, MOML can act as a representative of Mahalanobis-based OML algorithms.

**Theoretical Guarantee:** Several theoretical guarantees are given for the proposed algorithms. Theorem 1 is a positive-definite guarantee of the parameter matrix $\boldsymbol{M}$ in MOML. Moreover, Theorem 2 presents a regret bound of MOML. Proposition 1 tries to analyze and explain the effectiveness

of the proposed framework *i.e.*, MLOML. All the details of the proofs can be found in the appendix.

**Theorem 1.** *Suppose $\boldsymbol{M}_t$ is positive-definite, then $\boldsymbol{M}_{t+1}$ given by the MOML update, i.e., $\boldsymbol{M}_{t+1} = \boldsymbol{M}_t - \gamma \boldsymbol{A}_{t+1}$ is positive definite by properly setting $\gamma$.*

**Theorem 2.** *Let $\langle \boldsymbol{x}_1, \boldsymbol{x}_p, \boldsymbol{x}_q \rangle, \cdots, \langle \boldsymbol{x}_T, \boldsymbol{x}_p, \boldsymbol{x}_q \rangle$ be a sequence of triplet constraints where each sample $\boldsymbol{x}_t|_{t=1}^T \in \mathbb{R}^d$ has $\|\boldsymbol{x}_t\|_2 = 1$ for all $t$. Let $\boldsymbol{M}_t \in \mathbb{R}^{d \times d}$ be the solution of MOML at the $t$-th time step, and $\boldsymbol{U} \in \mathbb{R}^{d \times d}$ denotes an arbitrary parameter matrix. By setting $\gamma = \frac{1}{\Phi \sqrt{T}}$ (where $\Phi \in \mathbb{R}^+$), the regret bound is*

$$R(\boldsymbol{U}, T) = \sum_{t=1}^T \ell(\boldsymbol{M}_t) - \sum_{t=1}^T \ell(\boldsymbol{U}) \leq \frac{1}{2} \|\boldsymbol{I} - \boldsymbol{U}\|_F^2 + \frac{32}{\Phi^2}. \quad (7)$$

**Proposition 1.** *Let $\boldsymbol{M}_1, \cdots, \boldsymbol{M}_n$ be the parameter matrixes learnt by each metric layer of MLOML. The subsequent metric layer can learn a feature space that is at least as good as the one learnt by the former metric layer. That is, the composite feature space learnt by both $\boldsymbol{M}_1$ and $\boldsymbol{M}_2$ is better than the feature space learnt only by $\boldsymbol{M}_1$ in most cases (i.e., the feature space is more discriminating for classification).*

**Other OML Algorithms** In addition to MOML, other OML algorithms such as LEGO [19], RDML [20] and OPML [21] *etc.*, can also be adapted into the proposed multi-layer framework (namely LEGO-Multi, RDML-Multi and OPML-Multi). It is worth mentioning that both LEGO and RDML learn a Mahalanobis parameter matrix $\boldsymbol{M}$, while OPML just learns a transformation matrix $\boldsymbol{L}$. Hence, OPML doesn't need an additional matrix decomposition operation (*i.e.*, $\boldsymbol{M} = \boldsymbol{L}^\top \boldsymbol{L}$). The experimental results of LEGO-Multi, RDML-Multi and OPML-Multi will be discussed in Section IV-G.

### B. Forward and Backward Propagation

The proposed MLOML (*i.e.*, MLOML-r, MLOML-s, MLOML-t) is made up of a series of OML algorithms (*i.e.*, MOML metric layer) and nonlinear functions (*i.e.*, ReLU, Sigmoid, tanh). Then, MLOML attempts to explore a new way to train the metric layer by introducing forward propagation (FP) updating. In fact, MLOML can not only be learnt by forward propagation, but also be learnt by backward propagation. Moreover, these two strategies can be adopted simultaneously too. During forward propagation, each metric layer can be learnt immediately, through this way, new feature space can be explored sequentially. When backward propagation, the return gradients can be used to fine-tune all the metric layers, amending the feature spaces learnt by the forward propagation.

Therefore, MLOML can be trained with three different propagation strategies as follows: (1) **MLOML-FP**, which is only trained by employing forward propagation strategy. (2) **MLOML-FBP**, which utilizes forward and backward propagation strategies simultaneously. Specifically, a loss layer is added as the last layer to calculate the final loss, where the loss function (*i.e.*, Eq. (1)) is adopted. (3) **MLOML-BP** is similar to MLOML-FBP, while MLOML-BP only utilizes the
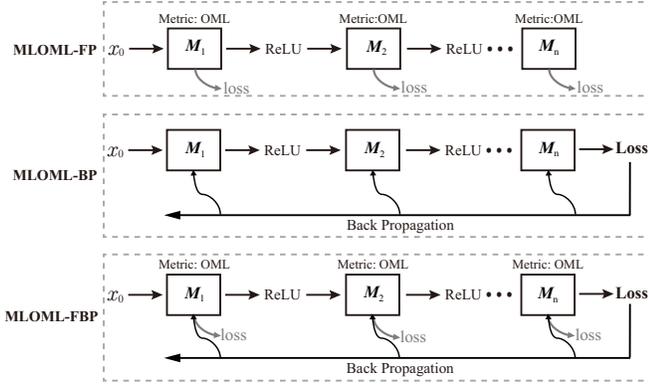
Fig. 3. Flowcharts of MLOML-FP, MLOML-BP and MLOML-FBP, respectively, where MLOML is the MLOML-r model.

TABLE I
TWELVE UCI DATASETS WITH DIFFERENT SCALES (*i.e.,* #INST) AND
FEATURE DIMENSIONS (*i.e.,* #FEAT).

| Datasets | #inst | #feat | #class | Datasets | #inst | #feat | #class |
|---|---|---|---|---|---|---|---|
| lsvt | 126 | 310 | 2 | balance | 625 | 4 | 3 |
| iris | 150 | 4 | 3 | breast | 683 | 9 | 2 |
| wine | 178 | 13 | 3 | pima | 768 | 8 | 2 |
| spect | 267 | 22 | 2 | diabetic | 1151 | 19 | 2 |
| ionophere | 351 | 34 | 2 | waveform | 5000 | 21 | 3 |
| pems | 440 | 137710 | 7 | mlprove | 6118 | 57 | 6 |

final loss to train the entire model without the local losses and without forward updating. The flowcharts of these three variations can be seen in Fig. 3. The comparison between these variations will be shown in Section IV-D.

## IV. EXPERIMENTS

To verify the effectiveness and applicability of the proposed MLOML, we conduct various experiments on the UCI datasets, which include multiple real-world machine learning tasks for which only vectorized features can be accessed, to analyze and interpret the properties of MLOML. First, we introduce the training process of MLOML.

### A. Training

We will describe how to train MLOML in detail in this section. Note that, MLOML is trained from scratch in an end-to-end manner, which is totally different from the traditional layer-by-layer training.

**Initialization:** Parameter matrix $M_i$ $(i = 1, 2, \ldots, n)$ is initialized as an identity matrix. The hyper-parameter $\gamma$ in MOML and the $\lambda$ in loss layer need to be chosen by cross-validation according to the specific task. All $w_i$ $(i = 1, 2, \ldots, n)$ is initialized as 1. The number of layers in MLOML is also a hyper-parameter, which can be chosen according to a specific task (3 or 5 layers are usually enough).

**Forward Propagation:** At the $t$-th time step, one triplet $\langle x_t^{(0)}, x_p^{(0)}, x_q^{(0)} \rangle$ is constructed. Then the triplet is fed into the first OML layer, and the current local triplet loss (*i.e.,* Eq. (4)) is calculated by using the current metric matrix $M_1$.

According to the updating strategy of MOML (*i.e.,* Eq. (6)), the metric matrix $M_1$ is updated for the first time. Then, $M_1$ is mathematically decomposed as $L_1^\top L_1$. After transformation by using $L_1$, the new triplet $\langle x_t^{(1)} = L_1 x_t^{(0)}, x_p^{(1)} = L_1 x_p^{(0)}, x_q^{(1)} = L_1 x_q^{(0)} \rangle$ is fed into the next ReLU, Sigmoid or tanh layer. In a serial manner, the final output of the last layer is $\langle x_t^{(n)}, x_p^{(n)}, x_q^{(n)} \rangle$. Through the linear (*i.e.,* OML layer) and nonlinear transformation (*i.e.,* ReLU, Sigmoid, or tanh layer), new feature spaces are sequentially explored. At the same time, the metric matrix of each OML layer is also learnt.

**Backward Propagation (optional):** The final loss is calculated according to Eq. (1) by using the output of the last OML layer. By using chain rule, SGD is adopted to update all the decomposed transformation matrix $L_i$ $(i = 1, 2, \ldots, n)$. Then each metric matrix $M_i$ $(i = 1, 2, \ldots, n)$ can be obtained naturally by $M_i = L_i^\top L_i$. Note that all these three samples in a triplet are used to calculate the gradients. Ideally, forward updating can explore new feature spaces, while backward updating can amend the exploration. In this way, that is, exploration with amendment, a much better feature space can be found. In practice, the backward propagation indeed can further slightly improve the feature space learnt by the forward propagation in some cases, but this could also bring additional computation load. As a trade-off between time and performance, if not specified, we will train the proposed MLOML only by forward propagation, similar to the Deep forest [37]. More details can be seen in Section IV-D.

### B. Datasets

We pick twelve commonly used datasets from UCI Machine Learning Repository [38], which vary in the dimensionality and size. The details of these twelve datasets can be seen in Table I. The reason of choosing these datasets is that they are all vectorized data and can be representative data for the real-world applications. For example, *lsvt* is a real voice rehabilitation treatment dataset. *pems* contains 15 months worth of daily data that describes the occupancy rate of different car lanes of the San Francisco bay area freeways. Also, *ionophere* is real radar data, which is collected by a system in Goose Bay, Labrador.

Classification task will be conducted on these datasets. For each dataset, 50% samples are randomly sampled as training set, and the rest is taken as test set. Each dataset will be resampled 30 times, and each algorithm will be tested on all these sampled datasets. When the feature dimensionality $d \geq 200$, the $d$-dimensional feature will be reduced to a 100-dimensional feature by principal component analysis (PCA) for easier handling. All datasets are normalized by $\ell_2$ normalization. Error rate is adopted as the evaluation criterion.

### C. Comparison with the State of the Art

To evaluate the effectiveness of the family of MLOML (MLOML-r, MLOML-s and MLOML-t), six state-of-the-art online metric learning (OML) algorithms, *i.e.,* RDML [20], LEGO [19], OASIS [9], OPML [21], SLMOML [33] and the new designed MOML are employed as comparisons. Note that all of these compared OML algorithms are single layer

TABLE II
Error rates (mean ± std. deviation) on the UCI datasets. ●/○ indicates that MLOML-r, MLOML-s, MLOML-t are significantly better/worse than the respective algorithm according to the $t$-tests at 95% significance level. The statistics of win/tie/loss between MLOML-r and other algorithms is also counted.

| Datasets | Euclidean | Batch | | | | Online | |
|---|---|---|---|---|---|---|---|
| | | LMNN | KISSME | LMDML | ML-CC | RDML | LEGO |
| Isvt | .369±.051● | .387±.057● | .403±.098● | .374±.064● | .384±.047● | .400±.055● | .369±.051● |
| iris | .038±.016● | .040±.018● | .039±.020● | .028±.015 | .058±.034● | .028±.019 | .037±.016● |
| wine | .218±.039 | .170±.044○ | **.069±.021**○ | .229±.041 | .112±.031○ | .350±.028● | .231±.041● |
| spect | .354±.031● | .357±.032● | .365±.031● | .342±.036● | .409±.044● | .347±.035● | .326±.035 |
| ionophere | .180±.017● | .157±.016● | .156±.023● | .115±.012● | .104±.022● | .096±.015● | .129±.019● |
| pems | .498±.033● | .402±.038 | **.188±.028**○ | .352±.028○ | .227±.025○ | .421±.030● | .461±.033● |
| balance | .108±.013● | .088±.013● | .101±.011● | .075±.010● | .066±.011 | .070±.011 | .091±.011● |
| breast | .106±.012 | .107±.012 | .106±.014 | .107±.013 | .113±.012● | .115±.015● | **.104±.016** |
| pima | .324±.018 | .326±.020 | .333±.021● | .330±.019● | .322±.022 | .357±.022● | .322±.020 |
| diabetic | .343±.018 | .335±.017○ | **.288±.018**○ | .316±.014○ | .338±.011 | .353±.015● | .322±.006○ |
| waveform | .195±.006● | .190±.005● | **.158±.006**○ | .176±.007○ | .208±.006● | .166±.006○ | .198±.006● |
| mlprove | .084±.005● | .037±.004● | .234±.273● | .007±.002● | .032±.025● | .027±.011● | .024±.003● |
| **win/tie/loss** | **8/4/0** | **7/3/2** | **7/1/4** | **6/3/3** | **7/3/2** | **9/2/1** | **8/3/1** |

| Datasets | Online | | | | | | |
|---|---|---|---|---|---|---|---|
| | OASIS | OPML | SLMOML | MOML | MLOML-t | MLOML-s | MLOML-r |
| Isvt | .333±.000 | .370±.051● | .369±.051● | .369±.051● | .369±.054● | .369±.052● | **.326±.053** |
| iris | .333±.000● | .035±.016● | .038±.016● | .028±.018 | **.027±.015** | **.025±.017** | **.026±.017** |
| wine | .586±.061● | .220±.040 | .225±.040 | .226±.041 | .214±.038○ | .216±.039 | .219±.039 |
| spect | .385±.033● | .321±.029 | .355±.032● | .331±.032● | .323±.028 | **.317±.024** | **.320±.025** |
| ionophere | .183±.017● | .107±.020● | .107±.020● | .108±.032● | .102±.021● | **.086±.013●** | **.081±.016** |
| pems | .651±.036● | .331±.029○ | .495±.033● | .416±.057● | .319±.032○ | .407±.030 | .397±.036 |
| balance | .125±.010● | .073±.012● | .077±.012● | .070±.010● | **.064±.013** | .066±.011 | .066±.012 |
| breast | .175±.050● | .109±.014● | .106±.012 | .112±.014● | .106±.013 | **.104±.013** | .105±.012 |
| pima | .349±.003● | .324±.022 | .334±.022● | .323±.020 | .322±.019 | **.321±.017** | .323±.018 |
| diabetic | .451±.021● | .322±.019○ | .348±.016● | .342±.016 | .341±.015 | .341±.014 | .342±.017 |
| waveform | .298±.049● | .175±.006○ | .161±.005○ | .173±.006○ | .168±.006○ | .162±.005○ | .187±.006 |
| mlprove | .002±.001○ | .006±.002● | .011±.004● | **.001±.001**○ | .002±.001○ | .002±.001○ | .004±.001 |
| **win/tie/loss** | **10/1/1** | **6/3/3** | **9/2/1** | **6/4/2** | **2/6/4** | **2/8/2** | |

algorithms, while the proposed famliy of MLOML (MLOML-r, MLOML-s and MLOML-t) are built on MOML is multi-layer algorithm. Euclidean distance is adopted as the baseline algorithm. Besides, four batch metric learning algorithms *i.e.,* LMNN [2], KISSME [39], LMDML [35] and ML-CC [40] are also employed for reference. Note that these three algorithms (*i.e.,* Euclidean, LMNN and KISSME) are offline.

Cross-validation is used for hyper-parameter selection for all algorithms. Specifically, the regularization parameter $\gamma$ for the family of MLOML (*i.e.,* the $\gamma$ in MOML metric layer, $\gamma \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$), the learning rate $\lambda$ for RDML ($\lambda \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$), the regularization parameter $\eta$ for LEGO ($\eta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$), the regularization parameter $\gamma$ for OPML ($\gamma \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$), the weighting parameter $\mu$ for LMNN ($\mu \in \{0.125, 0.25, 0.5\}$), the parameters $K$ and $\mu$ for ML-CC ($K \in \{2, 4, 8\}$ and $\mu \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$, and the aggressiveness parameter $C$ for SLMOML ($C \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$) are all set up in this

way.

For fair comparison, all OML algorithms adopt the same triplet construction strategy introduced by OPML to construct the pairwise or triplet constraints. The difference is that, in OPML the triplet construction strategy is one-pass, while here multiple-scan strategy is employed to construct more constraints for adequately training (the scanning number is set to 20). Note that, all OML algorithms are still trained in an online manner. Moreover, three metric layers MLOML is adopted in this experiment. A $k$-NN classifier (*i.e.,* $k = 5$) is used to get the final classification results. The results are summarized in Table II. For each dataset, the mean and standard deviation of error rate are calculated, and pairwise $t$-tests between MLOML and other algorithms at 95% significance level are also performed. Then the win/tie/loss is counted according to the $t$-test. From this table, we can see that the family of MLOML can not only achieve superior performance compared with other state-of-the-art OML algorithms, but also better than batch metric learning algorithms except KISSME

TABLE III
ERROR RATES ON TWELVE UCI DATASETS BY EMPLOYING DIFFERENT PROPAGATION STRATEGIES FOR MLOML.

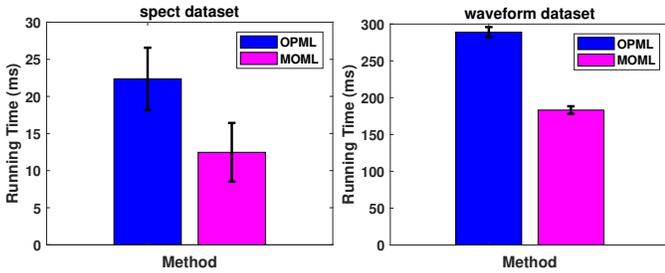| Datasets | MLOML-r | | | MLOML-s | | | MLOML-t | | |
|---|---|---|---|---|---|---|---|---|---|
| | BP | FBP | FP | BP | FBP | FP | BP | FBP | FP |
| lsvt | .354±.053● | .325±.055 | .326±.053 | .369±.051 | .369±.052 | .369±.052 | .369±.051 | .369±.051 | .369±.054 |
| iris | .032±.015● | .026±.017 | .026±.017 | .040±.017● | .025±.017 | .025±.017 | .039±.019● | .027±.015 | .027±.015 |
| wine | .220±.040 | .220±.040 | .219±.039 | .216±.039 | .216±.039 | .216±.039 | .216±.041 | .214±.041 | .214±.038 |
| spect | .358±.029● | .319±.026 | .320±.025 | .352±.030● | .314±.025● | .317±.024 | .346±.030● | .321±.025 | .323±.028 |
| ionophere | .128±017● | .081±.017 | .081±.016 | .175±.016● | .088±.013 | .086±.013 | .150±.016● | .102±.021 | .102±.021 |
| pems | .466±.036● | .396±.033 | .397±.036 | .500±.032● | .408±.030 | .407±.030 | .496±.032● | .319±.031 | .319±.032 |
| balance | .070±.013● | .066±.011 | .066±.012 | .109±.012● | .066±.012 | .066±.011 | .066±.011● | .064±.013 | .064±.013 |
| breast | .109±.015● | .107±.014 | .105±.012 | .106±.012 | .104±.013 | .104±.013 | .106±.016 | .104±.014 | .106±.013 |
| pima | .323±.017 | .324±.017 | .323±.018 | .322±.017 | .321±.017 | .321±.017 | .323±.018 | .323±.020 | .322±.019 |
| diabetic | .340±.017 | .340±.014 | .342±.017 | .342±.016 | .340±.014 | .341±.014 | .341±.015 | .342±.015 | .341±.015 |
| waveform | .180±.006○ | .175±.006○ | .187±.006 | .194±.006● | .163±.004 | .162±.005 | .169±.005 | .166±.005● | .168±.006 |
| mlprove | .006±.002● | .003±.001○ | .004±.001 | .084±.005● | .002±.001 | .002±.001 | .007±.002● | .001±.001● | .002±.001 |
| **win/tie/loss** | **8/3/1** | **0/10/2** | | **7/5/0** | **0/11/1** | | **6/6/0** | **0/10/2** | |



Fig. 4. Running time of OPML and MOML on the spect and waveform datasets.

algorithm. A possible reason is that KISSME learns a distance metric from equivalence constraints which is easier to specify labels. We can also see that MLOML is robust on small datasets, *e.g.,* lsvt, iris, spect and ionophere, which means that MLOML can handle small-scale data very well.

The traditional online metric learning algorithms, RDML, LEGO, OASIS, OPML and SLMOML algorithms, like our proposed MOML algorithm, can be used as metric layers for multi-layer online metric learning algorithms. However, in the literature of the RDML, LEGO, OASIS and SLMOML algorithms, they need to first sample triplets, and then apply these triplets to the training process. Because of this, their running time are relatively long. In order to guarantee the fairness of the comparison, Fig. 4 shows the running time of OPML and MOML on the spect and waveform datasets, where OPML and MOML use the one-pass triplet construction strategy. We can see that the proposed MOML has a shorter running time than OPML.

### D. Forward and Backward Propagation

In this section, we analyze the learning ability of MLOML by adopting different propagation strategies, *i.e.,* MLOML-FP, MLOML-BP and MLOML-FBP. Specifically, we conduct classification task on the twelve UCI datasets to compare these three variations of MLOML, each of which contains three metric layers. The results are exhibited in Table III. From the results, we can see that MLOML-FP performs better than

MLOML-BP. The reason is not difficult to perceive, because BP may suffer from the vanishing gradient problem. Taking advantage of the fact that each metric layer of MLOML is a MOML algorithm, it can learn a good metric in each layer during FP. We can also observe that MLOML-FP performs similarly to MLOML-FBP. The reason may be that MLOML-FP has achieved quite good classification performance on some datasets, so additional BP updating cannot further improve the performance. However, on other datasets, MLOML-FBP indeed achieves the best classification performance as expected, such as iris, spect and mlprove *etc.* It is worth mentioning that MLOML-FP is the fastest one among these three variations with a time complexity of $O(nd^2)$, where $n$ is the number of metric layers. Overall, for the proposed MLOML, the FP training strategy is the best one when considering both training performance and training efficiency. It should be noted that in the following chapters, MLOML refers specifically to the MLOML-r model, and the three-layer and five-layer MLOML networks are denoted by MLOML-3L and MLOML-5L.

### E. Progressive Feature Representation

In this section, we will analyze the progressive feature representation ability of each metric layer in MLOML and verify that the metric space can become better and better by adding metric layer gradually. Particularly, an MLOML-5L model is employed. To test the feature representation ability of each metric layer, we perform classification task on the output features of each metric layer respectively. We choose nine UCI datasets and take Euclidean distance, MOML and LMNN as the baseline algorithms. Note that only the test sets of these datasets are used to perform this experiment. From Fig. 5, we can see that the classification performance of MLOML-5L becomes better with the increase of the number of metric layers. Besides, in some datasets, the curve of error rate can converge smoothly. Moreover, we apply PCA to four UCI datasets to obtain the new space with 2-dimensional features, and then visualize the feature space learnt by each metric layer for more intuition (shown in Fig. 6). The four UCI datasets are picked and entered into one learnt MLOML-3L model. Next, all output samples of each metric layer are $\ell_2$
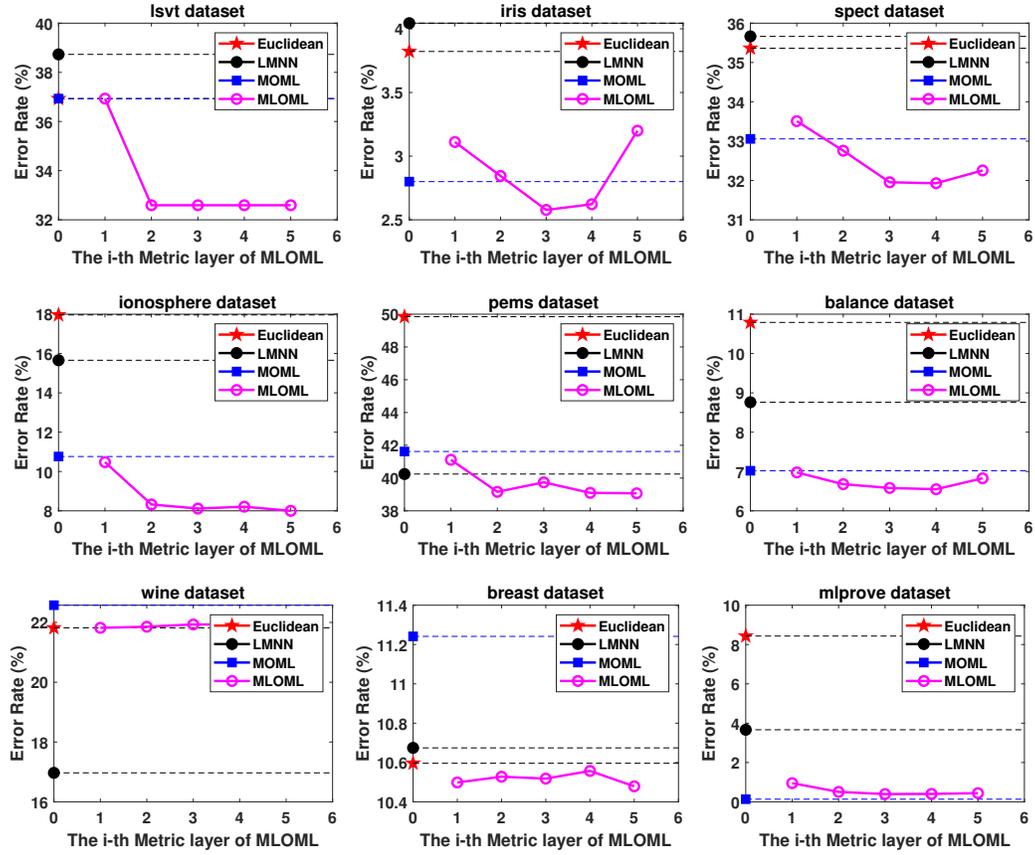
Fig. 5. Results of different metric layers of MLOML. Moreover, Euclidean, MOML and LMNN are taken as the baseline algorithms.
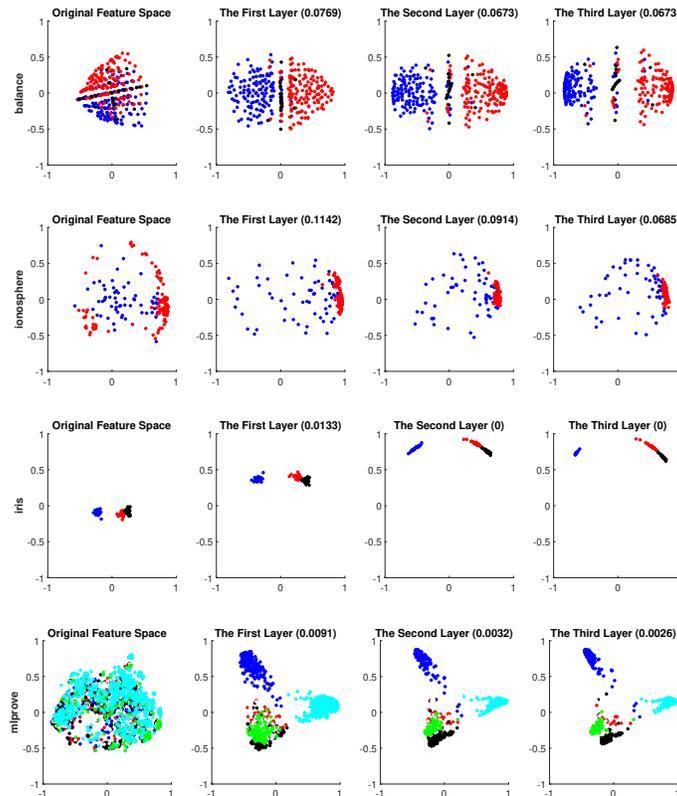


Fig. 6. Feature visualization on four UCI datasets, demonstrating the feature representation learnt by each metric layer in MLOML-3L.
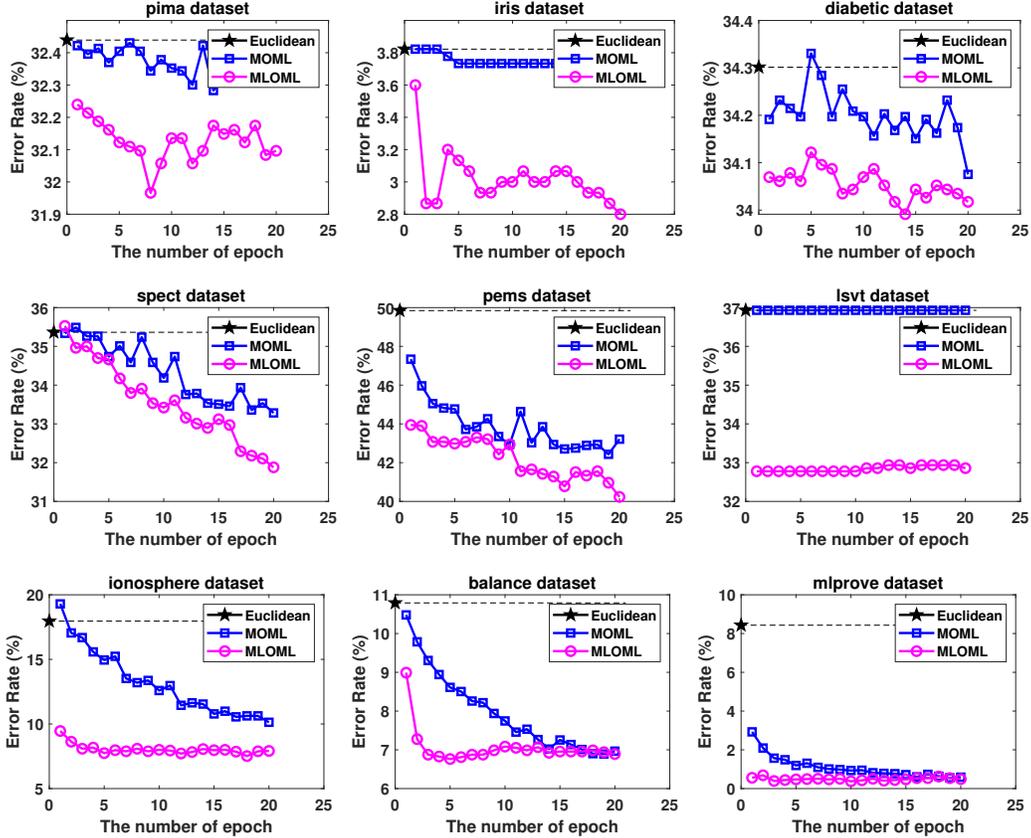
Fig. 7. Error rates on nine UCI datasets by changing the number of scans for MOML and MLOML.

normalized and reduced to a two-dimensional space by PCA. As seen, in original feature space, the distribution of samples is disordered. As the number of metric layers increases, the intra-class distance becomes smaller, the inter-class distance becomes larger, and the distribution of samples becomes more separable.

### F. Learning Ability of MLOML

Since the multiple-scan strategy is performed in the training phase, it is necessary to test the learning ability of MLOML by setting different numbers of scans. Note that $m$ times scanning will scan the training data $m$ times. Therefore, we set the number of scans from 1 to 20, and compare the classification performance between MLOML and MOML under different scans. Specifically, nine datasets are picked, and Euclidean distance is taken as the baseline. The results are presented in Fig. 7. From the figure, we can see that as the number of scans increases, the classification performance of MLOML is significantly improved and then converge, which can reflect the ability of MLOML for reusing data. Compared with MOML, with the same amount of data (*i.e.,* the same scan), MLOML can learn better feature representation (*i.e.,* lower error rate). In other words, the learning ability of MLOML is stronger than MOML, which means that MLOML can gain more learning ability from the multi-layer architecture.

### G. Extendability of MLOML

In order to verify the extendability of the proposed framework, we take the other three OML algorithms (*e.g.,* LEGO, RDML and OPML) as the base OML layer followed by the ReLU layer and construct their corresponding multi-layer versions, respectively (*i.e.,* LEGO-multi, RDML-multi and OPML-multi). Note that these three algorithms are all Mahalanobis-based OML algorithms. For simplicity, FP strategy is employed for these three algorithms. Other settings are similar to the ones in Section IV-E. From Fig. 8, we can see that LEGO-multi, RDML-multi and OPML-multi have similar characteristic to MLOML. In most cases, multi-layer versions of these algorithms perform better than their corresponding shallow versions. Moreover, the progressive learning ability of feature representation is demonstrated. Therefore, the effectiveness and extendability of the proposed framework are verified.

### V. DISCUSSIONS AND CONCLUSIONS

In this study, we propose a multi-layer framework for online metric learning. Specifically, we implement *multi-layer online metric learning (MLOML)* by stacking a set of OML algorithms. Extensive experiments have been conducted to analyze and verify the properties of MLOML. For future work, we will analyze and discuss this framework from three aspects as follows.

- **Extendability:** Although only OML-based algorithms are implemented (*e.g.,* MLOML), the proposed framework is
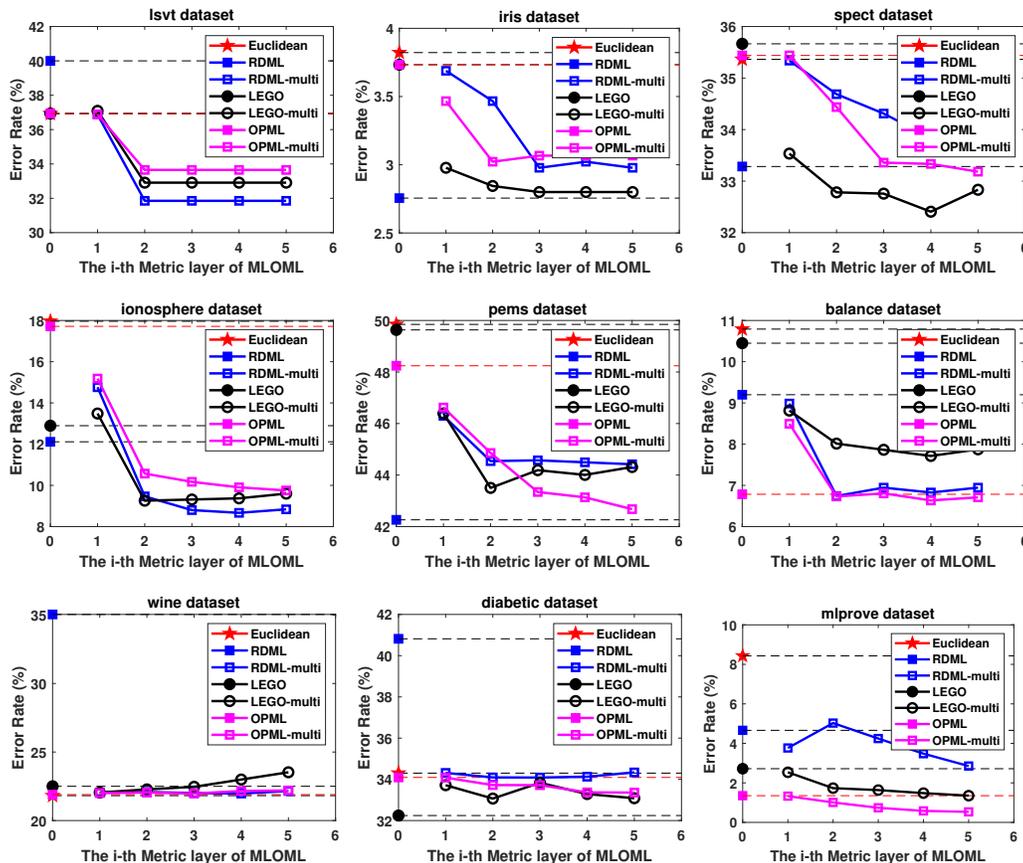
Fig. 8. Results of different metric layers of RDML-multi, LEGO-multi and OPML-multi (marked by hollow shapes), which are stacked by online metric algorithms RDML, LEGO and OPML (marked by corresponding solid shapes) along with the ReLU layers, respectively.

extensible, such as: a) mini-batch or batch metric learning based metric layer can be constructed; b) different metric learning algorithms can be combined as different metric layers.

- **Advantages:** The proposed MLOML has many nice properties: a) it is online; b) it can be trained by either forward or backward propagation; c) it is quite fast and effective, which can be trained by CPU; d) it can progressively learn feature representation.

- **Drawbacks:** Because MLOML is based on MOML, the performance of MLOML depends on the performance of MOML. Currently, MLOML cannot efficiently handle high dimensional data well due to a full matrix $M$ learned in MOML. This problem can be tackled by learning a diagonal matrix or employing dimensionality reduction through online feature selection, which will be investigated in the next work. Meanwhile, the number of metric layers in MLOML is uniformly specified according to the experimental results. However, the optimal number of layers is often different for different tasks. The another question is how many metric layers is sufficient for a task will be studied in the future.

## APPENDIX

### A. Proof of Theorem 1

*Proof.* As $\boldsymbol{A}_{t+1} = (\boldsymbol{x}_{t+1}-\boldsymbol{x}_p)(\boldsymbol{x}_{t+1}-\boldsymbol{x}_p)^\top - (\boldsymbol{x}_{t+1}-\boldsymbol{x}_q)(\boldsymbol{x}_{t+1}-\boldsymbol{x}_q)^\top$, whose rank is 1 or 2, it has at most 2 non-zero eigenvalues. That is to say, $\mathrm{Tr}(\boldsymbol{A}_{t+1}) = \lambda_1 + \lambda_2$. Specifically, we can also easily get that,

$$-\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_q\|_2^2 \leq \lambda(\boldsymbol{A}_{t+1}) \leq \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_p\|_2^2, \quad (8)$$

where $\lambda(\boldsymbol{A}_{t+1})$ means the eigenvalue of $\boldsymbol{A}_{t+1}$ (*i.e.,* $\lambda_1$ or $\lambda_2$). For each sample $\boldsymbol{x}$ is $\ell_2$ normalized, the ranges of $\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_p\|_2^2$ and $\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_q\|_2^2$ vary from $[0, 4]$. Thus,

$$\lambda_{\min}(\boldsymbol{M}_t) - 4\gamma \leq \lambda(\boldsymbol{M}_t - \gamma\boldsymbol{A}_{t+1}) \leq \lambda_{\max}(\boldsymbol{M}_t) + 4\gamma. \quad (9)$$

When $\gamma \leq \frac{1}{4}\lambda_{\min}(\boldsymbol{M}_t)$, it is guaranteed that the minimum eigenvalue of $\boldsymbol{M}_t - \gamma\boldsymbol{A}_{t+1}$ is greater than zero. As the initial matrix $\boldsymbol{M}_1 = \boldsymbol{I}$ is positive definite (*i.e.,* $\lambda_{\min}(\boldsymbol{M}_1) = 1$). By properly setting a small $\gamma$, the minimum eigenvalue of $\boldsymbol{M}_t - \gamma\boldsymbol{A}_{t+1}$ is generally large than zero. Thus, the positive definiteness of $\boldsymbol{M}_{t+1} = \boldsymbol{M}_t - \gamma\boldsymbol{A}_{t+1}$ can be guaranteed. Same theoretical guarantee (*i.e.,* the small pertubations of positive

definite matrix) can also be found in the chapter 9.6.12 of [41]. □

### B. Proof of Theorem 2

*Proof.* According to the objective function of MOML, *i.e.,*

$$\Gamma = \arg\min_{\boldsymbol{M} \succcurlyeq 0} \frac{1}{2}\|\boldsymbol{M} - \boldsymbol{M}_{t-1}\|_F^2 + \gamma\Big[1 + \mathrm{Tr}(\boldsymbol{M}\boldsymbol{A}_t)\Big]_+ , \quad (10)$$

we denote $\ell_t$ as the instantaneous loss suffered by MOML at each $t$-time step with the learnt $\boldsymbol{M}_t \in \mathbb{R}^{d\times d}$, and denote by $\ell_t^*$ the loss suffered by an arbitrary parameter matrix $\boldsymbol{U} \in \mathbb{R}^{d\times d}$, which can be formalized as below:

$$\begin{aligned}\ell_t &= \ell(\boldsymbol{M}_t; \langle\boldsymbol{x}_t, \boldsymbol{x}_p, \boldsymbol{x}_q\rangle) = [1 + \mathrm{Tr}\,(\boldsymbol{M}_t\boldsymbol{A}_t)]_+ \\ \ell_t^* &= \ell(\boldsymbol{U}; \langle\boldsymbol{x}_t, \boldsymbol{x}_p, \boldsymbol{x}_q\rangle) = [1 + \mathrm{Tr}\,(\boldsymbol{U}\boldsymbol{A}_t)]_+ ,\end{aligned} \quad (11)$$

where $\boldsymbol{A}_t = (\boldsymbol{x}_t - \boldsymbol{x}_p)(\boldsymbol{x}_t - \boldsymbol{x}_p)^\top - (\boldsymbol{x}_t - \boldsymbol{x}_q)(\boldsymbol{x}_t - \boldsymbol{x}_q)^\top$, Tr denotes trace and $[z]_+ = \max(0, z)$. As $\mathrm{Tr}(\boldsymbol{M}_t\boldsymbol{A}_t)$ is a linear function, it is convex *w.r.t* $\boldsymbol{M}_t$ by natural. Besides, the hinge loss function $[z]_+$ is a convex function (but not continuous at $z = 0$) *w.r.t* $z$. Hence, the resulting composite function $\ell_t(\boldsymbol{M}_t)$ is convex *w.r.t* $\boldsymbol{M}_t$. As $\ell$ is a convex function, we can introduce the first-order condition as follow:

$$\ell(\boldsymbol{Y}) \geq \ell(\boldsymbol{X}) + \mathrm{VEC}(\triangledown\ell(\boldsymbol{X}))^\top \mathrm{VEC}(\boldsymbol{Y} - \boldsymbol{X}), \quad (12)$$

where $\boldsymbol{X}, \boldsymbol{Y} \in \mathbb{R}^{d\times d}$, VEC denotes vectorization of a matrix, and $\triangledown\ell(\boldsymbol{X})$ is the gradient of function $\ell$ at $\boldsymbol{X}$.

Inspired by [15], we define $\Delta_t$ to be $\|\boldsymbol{M}_t - \boldsymbol{U}\|_F^2 - \|\boldsymbol{M}_{t+1} - \boldsymbol{U}\|_F^2$. Then calculating the cumulative sum of $\Delta_t$ over all $t \in \{1, 2, \cdots, T\}$, we can easily obtain $\sum_t \Delta_t$,

$$\begin{aligned}\sum_{t=1}^{T} \Delta_t &= \sum_{t=1}^{T}(\|\boldsymbol{M}_t - \boldsymbol{U}\|_F^2 - \|\boldsymbol{M}_{t+1} - \boldsymbol{U}\|_F^2) \\ &= \|\boldsymbol{M}_1 - \boldsymbol{U}\|_F^2 - \|\boldsymbol{M}_{T+1} - \boldsymbol{U}\|_F^2 \\ &\leq \|\boldsymbol{M}_1 - \boldsymbol{U}\|_F^2 .\end{aligned} \quad (13)$$

For simplicity, we employ stochastic gradient descent (SGD) to update the parameter matrix $\boldsymbol{M}_t$. Hence, according to the definition of SGD, $\boldsymbol{M}_{t+1} = \boldsymbol{M}_t - \eta \triangledown \ell(\boldsymbol{M}_t)$, where $\eta$ is the learning rate, and $\triangledown\ell(\boldsymbol{M}_t) = \gamma\boldsymbol{A}_{t+1}$. Then, we can rewrite the $\Delta_t$ as,

$$\begin{aligned}\Delta_t &= \|\boldsymbol{M}_t - \boldsymbol{U}\|_F^2 - \|\boldsymbol{M}_{t+1} - \boldsymbol{U}\|_F^2 \\ &= \|\boldsymbol{M}_t - \boldsymbol{U}\|_F^2 - \|\boldsymbol{M}_t - \eta \triangledown \ell(\boldsymbol{M}_t) - \boldsymbol{U}\|_F^2 \\ &= \|\boldsymbol{M}_t\|_F^2 - 2\langle\boldsymbol{M}_t, \boldsymbol{U}\rangle_F + \|\boldsymbol{U}\|_F^2 - \|\boldsymbol{M}_t - \boldsymbol{U}\|_F^2 \\ &\quad + 2\langle\boldsymbol{M}_t - \boldsymbol{U}, \eta \triangledown \ell(\boldsymbol{M}_t)\rangle_F - \eta^2\|\triangledown\ell(\boldsymbol{M}_t)\|_F^2 \\ &= 2\eta\,\mathrm{VEC}(\boldsymbol{M}_t - \boldsymbol{U})^\top \mathrm{VEC}(\triangledown\ell(\boldsymbol{M}_t)) - \eta^2\|\triangledown\ell(\boldsymbol{M}_t)\|_F^2 \\ &\scriptstyle(employ\ the\ Eq.\ (12)\ i.e.,\ \ell(\boldsymbol{U}) \geq \ell(\boldsymbol{M}_t) + \mathrm{VEC}(\triangledown\ell(\boldsymbol{M}_t))^\top \mathrm{VEC}(\boldsymbol{U} - \boldsymbol{M}_t)) \\ &\geq 2\eta(\ell_t - \ell_t^*) - \eta^2\|\triangledown\ell(\boldsymbol{M}_t)\|_F^2 .\end{aligned} \quad (14)$$

We can easily get that,

$$\sum_{t=1}^{T}\Big[2\eta(\ell_t - \ell_t^*) - \eta^2\|\triangledown\ell(\boldsymbol{M}_t)\|_F^2\Big] \leq \|\boldsymbol{M}_1 - \boldsymbol{U}\|_F^2 . \quad (15)$$

As all samples are $\ell_2$ normalized, the 2-norm of each sample is 1, namely $\|\boldsymbol{x}_t\|_2 \equiv 1, t \in \{1, 2, \cdots, T\}$. We can easily calculate the Frobenius norm of $\boldsymbol{A}_{t+1}$.

$$\begin{aligned}\|\boldsymbol{A}_{t+1}\|_F &\leq \|(\boldsymbol{x}_{t+1} - \boldsymbol{x}_p)(\boldsymbol{x}_{t+1} - \boldsymbol{x}_p)^\top\|_F + \|(\boldsymbol{x}_{t+1} - \boldsymbol{x}_q)(\boldsymbol{x}_{t+1} - \boldsymbol{x}_q)^\top\|_F \\ &\scriptstyle(employ\ \|\boldsymbol{a}\boldsymbol{b}^\top\|_F^2 = (\sum_{i=1}^{d}|\boldsymbol{a}_i|^2)(\sum_{j=1}^{d}|\boldsymbol{b}_j|^2),\ where\ \boldsymbol{a}, \boldsymbol{b} \in \mathbb{R}^d) \\ &= \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_p\|_2 \cdot \|\boldsymbol{x}_{t+1}^\top - \boldsymbol{x}_p^\top\|_2 + \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_q\|_2 \cdot \|\boldsymbol{x}_{t+1}^\top - \boldsymbol{x}_q^\top\|_2 \\ &= \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_p\|_2^2 + \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_q\|_2^2 \\ &\scriptstyle(for\ \|\boldsymbol{a} - \boldsymbol{b}\|_2^2 \leq (\|\boldsymbol{a}\|_2 + \|\boldsymbol{b}\|_2)^2) \\ &\leq 8 .\end{aligned} \quad (16)$$

Thus,

$$\begin{aligned}\sum_{t=1}^{T}(\ell_t - \ell_t^*) &\leq \frac{1}{2\eta}\|\boldsymbol{M}_1 - \boldsymbol{U}\|_F^2 + \frac{\eta}{2}\sum_{t=1}^{T}\|\triangledown\ell(\boldsymbol{M}_t)\|_F^2 \\ &= \frac{1}{2\eta}\|\boldsymbol{M}_1 - \boldsymbol{U}\|_F^2 + \frac{\eta}{2}\sum_{t=1}^{T}\|\gamma\boldsymbol{A}_{t+1}\|_F^2 \\ &\leq \frac{1}{2\eta}\|\boldsymbol{M}_1 - \boldsymbol{U}\|_F^2 + 32T\eta\gamma^2 \\ &\scriptstyle(\boldsymbol{M}_1\ is\ initialized\ to\ an\ identity\ matrix\ \boldsymbol{I}) \\ &= \frac{1}{2\eta}\|\boldsymbol{I} - \boldsymbol{U}\|_F^2 + 32T\eta\gamma^2 .\end{aligned} \quad (17)$$

In particular, setting $\eta = \frac{1}{\Phi\sqrt{T}}$ (where $\Phi > 0$ is a constant) yields the regret bound $R(\boldsymbol{U}, T) \leq \big(\frac{\Phi}{2}\|\boldsymbol{I} - \boldsymbol{U}\|_F^2 + \frac{32\gamma^2}{\Phi}\big)\sqrt{T}$. In fact, in this study, as a closed-form solution is employed (*i.e.,* $\eta = 1$), the regret bound is $R(\boldsymbol{U}, T) \leq \frac{1}{2}\|\boldsymbol{I} - \boldsymbol{U}\|_F^2 + 32T\gamma^2$. By setting $\gamma$ in a decreasing way with the iteration number $T$, for example, $\gamma = \frac{1}{\Phi\sqrt{T}}$, we can obtain a regret bound $R(\boldsymbol{U}, T) \leq \frac{1}{2}\|\boldsymbol{I} - \boldsymbol{U}\|_F^2 + \frac{32}{\Phi^2}$. Hence proved. □

### C. Theoretical analysis of Proposition 1

*Proof.* For simplicity, we just consider to analyze and prove this proposition of MLOML-FP that only uses forward propagation strategy. In fact, as MLOML-FP only has forward propagation, each metric layer is a relatively independent MOML algorithm. Thus, Theorem 2 is applicable to each metric layer. In other words, each metric layer (*i.e.,* a MOML algorithm) has its own tight regret bound. As the subsequent metric layer is learnt based on the output of the former metric layer, the metric space should not be worse according to the theoretical guarantee of regret bound. Moreover, ReLU, Sigmoid, tanh activation functions can introduce nonlinear and sparsity into the feature mapping, which is also beneficial to the exploration of feature space. In some cases, if the latter metric layer is in the wrong direction, backward propagation can be chosen to correct and adjust the direction to some extent. □

### REFERENCES

[1] E. Xing, M. Jordan, S. J. Russell, and A. Ng, "Distance metric learning with application to clustering with side-information," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 15, pp. 521–528, 2002.

[2] K. Q. Weinberger, J. Blitzer, and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," in *Conference on Neural Information Processing Systems (NeurIPS)*, 2005, pp. 1473–1480.

[3] B. Nguyen and B. De Baets, "Kernel-based distance metric learning for supervised $k$-means clustering," *IEEE Transactions on Neural Networks and Learning systems*, vol. 30, no. 10, pp. 3084–3095, 2019.

[4] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *International Conference on Machine Learning (ICML)*, 2007, pp. 209–216.

[5] S. Xiang, F. Nie, and C. Zhang, "Learning a mahalanobis distance metric for data clustering and classification," *Pattern Recognition*, vol. 41, no. 12, pp. 3600–3612, 2008.

[6] S. Ying, Z. Wen, J. Shi, Y. Peng, J. Peng, and H. Qiao, "Manifold preserving: An intrinsic approach for semisupervised distance metric learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 7, pp. 2731–2742, 2018.

[7] H.-J. Ye, D.-C. Zhan, and Y. Jiang, "Fast generalization rates for distance metric learning," *Machine Learning*, vol. 108, no. 2, pp. 267–295, 2019.

[8] Y. Liu, B. Du, W. Tu, M. Gong, Y. Guo, and D. Tao, "Logdet metric-based domain adaptation," *IEEE Transactions on Neural Networks and Learning systems*, vol. 31, no. 11, pp. 4673–4687, 2020.

[9] G. Chechik, V. Sharma, U. Shalit, and S. Bengio, "Large scale online learning of image similarity through ranking," *Journal of Machine Learning Research*, vol. 11, pp. 1109–1135, 2010.

[10] X. Gao, S. C. H. Hoi, Y. Zhang, J. Wan, and J. Li, "SOML: sparse online metric learning with application to image retrieval," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2014, pp. 1206–1212.

[11] H. Xia, S. C. H. Hoi, R. Jin, and P. Zhao, "Online multiple kernel similarity learning for visual search," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 3, pp. 536–549, 2014.

[12] A. Kavousi-Fard, A. Khosravi, and S. Nahavandi, "A new fuzzy-based combined prediction interval for wind power forecasting," *IEEE Transactions on Power Systems*, vol. 31, no. 1, pp. 18–26, 2016.

[13] E. M. Duhon, "On-line consumer credit data reporting system," 2001, uS Patent 6,311,169.

[14] X. He, J. Pan, O. Jin, T. Xu, B. Liu, T. Xu, Y. Shi, A. Atallah, R. Herbrich, S. Bowers *et al.*, "Practical lessons from predicting clicks on ads at facebook," in *International Workshop on Data Mining for Online Advertising (ADKDD)*. ACM, 2014, pp. 1–9.

[15] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online passive-aggressive algorithms," *Journal of Machine Learning Research*, vol. 7, pp. 551–585, 2006.

[16] H. Yu, J. Lu, and G. Zhang, "Online topology learning by a gaussian membership-based self-organizing incremental neural network," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 10, pp. 3947–3961, 2020.

[17] Yu, Hang and Lu, Jie and Zhang, Guangquan, "An online robust support vector regression for data streams," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 1, pp. 150–163, 2022.

[18] S. Shalev-Shwartz, Y. Singer, and A. Y. Ng, "Online and batch learning of pseudo-metrics," in *International Conference on Machine Learning (ICML)*, 2004.

[19] P. Jain, B. Kulis, I. S. Dhillon, and K. Grauman, "Online metric learning and fast similarity search," in *Conference on Neural Information Processing Systems (NeurIPS)*, 2008, pp. 761–768.

[20] R. Jin, S. Wang, and Y. Zhou, "Regularized distance metric learning: Theory and algorithm," in *Conference on Neural Information Processing Systems (NeurIPS)*, 2009, pp. 862–870.

[21] W. Li, Y. Gao, L. Wang, L. Zhou, J. Huo, and Y. Shi, "OPML: A one-pass closed-form solution for online metric learning," *Pattern Recognition*, vol. 75, pp. 302–314, 2018.

[22] Y. Gao, Y.-F. Li, S. Chandra, L. Khan, and B. Thuraisingham, "Towards self-adaptive metric learning on the fly," in *The World Wide Web Conference*. ACM, 2019, pp. 503–513.

[23] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

[24] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *Conference on Neural Information Processing Systems (NeurIPS)*, 2006, pp. 153–160.

[25] D. Erhan, Y. Bengio, A. Courville, P.-A. Manzagol, P. Vincent, and S. Bengio, "Why does unsupervised pre-training help deep learning?" *Journal of Machine Learning Research*, vol. 11, no. Feb, pp. 625–660, 2010.

[26] D. Yao, P. Zhao, C. Yu, H. Jin, and B. Li, "Sparse online relative similarity learning," in *IEEE International Conference on Data Mining (ICDM)*, 2015, pp. 529–538.

[27] Y. Cong, J. Liu, B. Fan, P. Zeng, H. Yu, and J. Luo, "Online similarity learning for big data with overfitting," *IEEE Transactions on Big Data*, vol. 4, no. 1, pp. 78–89, 2018.

[28] G. Kunapuli and J. W. Shavlik, "Mirror descent for metric learning: A unified approach," in *ECML-PKDD*, 2012, pp. 859–874.

[29] A. Bellet and A. Habrard, "Robustness and generalization for metric learning," *Neurocomputing*, vol. 151, pp. 259–267, 2015.

[30] H. Xu and S. Mannor, "Robustness and generalization," *Machine Learning*, vol. 86, no. 3, pp. 391–423, 2012.

[31] W. Liu, C. Mu, R. Ji, S. Ma, J. Smith, and S.-F. Chang, "Low-rank similarity metric learning in high dimensions," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2015, pp. 2792–2799.

[32] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.

[33] G. Zhong, Y. Zheng, S. Li, and Y. Fu, "Slmoml: online metric learning with global convergence," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 10, pp. 2460–2472, 2018.

[34] H. Liu, Z. Han, Y.-S. Liu, and M. Gu, "Fast low-rank metric learning for large-scale and high-dimensional data," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 817–827.

[35] B. Nguyen, C. Morell, and B. De Baets, "Scalable large-margin distance metric learning using stochastic gradient descent," *IEEE Transactions on Cybernetics*, vol. 50, no. 3, pp. 1072–1083, 2020.

[36] J. Dong, Y. Cong, G. Sun, T. Zhang, X. Tang, and X. Xu, "Evolving metric learning for incremental and decremental features," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 14, no. 8, pp. 1–13, 2021.

[37] Z. Zhou and J. Feng, "Deep forest: Towards an alternative to deep neural networks," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2017, pp. 3553–3559.

[38] M. Lichman, "UCI machine learning repository," 2013. [Online]. Available: http://archive.ics.uci.edu/ml

[39] M. Köstinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof, "Large scale metric learning from equivalence constraints," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 2288–2295.

[40] X. Guo, C. Dang, J. Liang, W. Wei, and J. Liang, "Metric learning with clustering-based constraints," *International Journal of Machine Learning and Cybernetics*, vol. 12, no. 12, pp. 3597–3605, 2021.

[41] K. B. Petersen, M. S. Pedersen *et al.*, "The matrix cookbook," *Technical University of Denmark*, vol. 7, p. 15, 2008.