

Evolving Metric Learning for Incremental and Decremental Features

Jiahua Dong, Yang Cong, *Senior Member, IEEE*, Gan Sun, Tao Zhang, Xu Tang and Xiaowei Xu

Abstract—Online metric learning has been widely exploited for large-scale data classification due to the low computational cost. However, amongst online practical scenarios where the features are evolving (*e.g.*, some features are vanished and some new features are augmented), most metric learning models cannot be successfully applied to these scenarios, although they can tackle the evolving instances efficiently. To address the challenge, we develop a new online **Evolving Metric Learning (EML)** model for incremental and decremental features, which can handle the instance and feature evolutions simultaneously by incorporating with a smoothed Wasserstein metric distance. Specifically, our model contains two essential stages: a Transforming stage (T-stage) and a Inheriting stage (I-stage). For the T-stage, we propose to extract important information from vanished features while neglecting non-informative knowledge, and forward it into survived features by transforming them into a low-rank discriminative metric space. It further explores the intrinsic low-rank structure of heterogeneous samples to reduce the computation and memory burden especially for highly-dimensional large-scale data. For the I-stage, we inherit the metric performance of survived features from the T-stage and then expand to include the new augmented features. Moreover, a smoothed Wasserstein distance is utilized to characterize the similarity relationships among the heterogeneous and complex samples, since the evolving features are not strictly aligned in the different stages. In addition to tackling the challenges in one-shot case, we also extend our model into multi-shot scenario. After deriving an efficient optimization strategy for both T-stage and I-stage, extensive experiments on several datasets verify the superior performance of our EML model.

Index Terms—Online metric learning, instance and feature evolutions, smoothed Wasserstein distance, low-rank constraint.

I. INTRODUCTION

Metric learning has been successfully extended into many fields, *e.g.*, face identification [1], object recognition [2] and medical diagnosis [3]. To efficiently solve the large-scale streaming data problem, learning an online discriminative metric (*i.e.*, online metric learning [4], [5]) attracts

This work is supported by the National Key Research and Development Program of China (2019YFB1310300) and National Nature Science Foundation of China under Grant (61722311, 61821005, 62003336). (*Corresponding author: Yang Cong.*)

Jiahua Dong and Tao Zhang are with the State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences, Shenyang 110016, China, and also with the Institutes for Robotics and Intelligent Manufacturing, Chinese Academy of Sciences, Shenyang 110016, China, and also with the University of Chinese Academy of Sciences, Beijing 100049, China (e-mail: dongjiahua1995@gmail.com, zhangtao2@sia.cn).

Yang Cong, Gan Sun and Xu Tang are with the State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences, Shenyang 110016, China, and also with the Institutes for Robotics and Intelligent Manufacturing, Chinese Academy of Sciences, Shenyang 110016, China (e-mail: congyang81@gmail.com, sungan1412@gmail.com, tangxu@sia.cn).

Xiaowei Xu is with the Department of Information Science, University of Arkansas at Little Rock, Arkansas 72204, USA (e-mail: xwxu@ualr.edu).

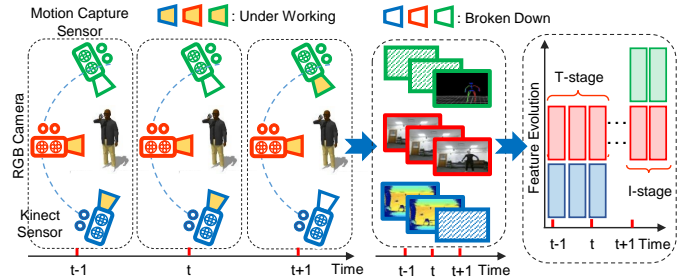


Fig. 1. Illustration example of feature evolution on human motion recognition task, where the blue, red and green colors respectively indicate the vanished features collected from Kinect sensor, the survived features collected from RGB camera and the augmented features collected from motion capture sensor with different lifespans. The vanished features collected from Kinect sensors are decremental in the T-stage, and the augmented features collected from motion capture sensor are incremental in the I-stage. The survived features collected from RGB camera exist in both T-stage and I-stage.

lots of appealing attentions. Generally, most online metric learning models pay attention to the fast metric updating mechanisms [6]–[9] or fast similarity searching strategies [5], [8], [10] for large-scale streaming data, where the streaming data indicate the continuous data flow that the data samples arrive consecutively in a real-time manner.

However, these existing online metric learning methods [5], [6], [10]–[12] only focus on instance evolution, and ignore the feature evolution in many real-world applications, where some features are vanished and some new features are augmented. Take the human motion recognition [13] as an example, as depicted in Fig. 1, the sudden damage of Kinect sensor results in the absence of depth information of human motion, while the emerging of new motion capture sensor could obtain the auxiliary human skeleton knowledge for motion recognition. It leads to a corresponding decrease and increase in the feature dimensionality of the input data, which are considered as the vanished features and augmented features, respectively. The features collected from RGB camera that has been working are regarded as survived features. Such feature evolution setting heavily cripples the human motion recognition performance of the pre-trained model [13]. Another interesting example is that different sensors (*e.g.*, radioisotope, trace metal and biological sensors [14]) are deployed to monitor the dynamic environment change in full aspects. Some sensors expire (vanished features) whereas some new sensors are deployed (augmented features) when different electrochemical conditions and lifespans occur. A fixed or static online metric learning model will fail to take advantage of sensors evolved in this way. Therefore, how to establish a novel metric learning model

to simultaneously handle both instance and feature evolutions amongst these online practical systems is our main focus in this paper.

To address the challenges above, as illustrated in Fig. 1, we develop a new online Evolving Metric Learning (EML) model for incremental and decremental features, which can exploit streaming data with both instance and feature evolutions in an online manner. To be specific, the proposed EML model consists of two significant stages, *i.e.*, a Transforming stage (T-stage) and a Inheriting stage (I-stage). 1) In the T-stage where features are decremental, we propose to explore the important information and data structure from vanished features, and transform them into a low-rank discriminative metric space of survived features, which could be utilized to promote the learning process of the I-stage. Moreover, it explores the intrinsic low-rank structure of the streaming data, which efficiently reduces both memory and computation costs especially for large-scale samples with high dimensional feature. 2) For the I-stage where features are incremental, based on the learned discriminative metric space in the T-stage, we inherit the metric performance of survived features from T-stage, and then expand to consider new augmented features. Furthermore, to better explore the similarity relations amongst the heterogeneous data, a smoothed Wasserstein distance is applied to both T-stage and I-stage where the evolving features are strictly unaligned and heterogeneous in different stages. For the model optimization, we derive an efficient optimization strategy to solve the formulations of T-stage and I-stage. Besides, our EML model could be successfully extended from one-shot scenario into multi-shot scenario, where one-shot scenario indicates that the features of streaming data would only be incremental and decremental by one time (as shown in Fig. 2), while multi-shot scenario denotes that the representations of streaming data would be incremental and decremental multiple times (as shown in Fig. 3). Comprehensive experimental results on several datasets strongly support the effectiveness of our proposed EML model.

The main contributions of this paper are summarized as follows:

- We propose an online Evolving Metric Learning (EML) model for incremental and decremental features to tackle both instance and feature evolutions simultaneously. To our best knowledge, this is the first exploration to tackle this crucial, but rarely-researched challenge in the metric learning field.
- We present two stages for both feature and instance evolutions, *i.e.*, a Transforming stage (T-stage) and a Inheriting stage (I-stage), which can not only make full use of the vanished features in the T-stage, but also take advantage of streaming data with new augmented features in the I-stage.
- A smoothed Wasserstein distance is incorporated into metric learning to characterize the similarity relations of heterogeneous evolving features among different stages. After deriving an alternating direction optimization algorithm to optimize our EML model, extensive experiments on representative datasets validate the superior perfor-

mance of our proposed EML model.

II. RELATED WORK

This section provides a brief overview about metric learning, followed by some representative methods about feature evolution.

A. Metric Learning

Online metric learning has been widely explored for instance evolution to learn large-scale streaming data, which is mainly composed of Mahalanobis distance-based and bilinear similarity-based methods. For the Mahalanobis distance-based methods, POLA [15] is the first attempt to learn the optimal metric in an online manner. Then several variants [5], [10], [16] extend this idea by the fast similarity searching strategies, *e.g.*, [8] proposes a regularized online metric learning model with the provable regret bound. Besides, pairwise constraint [8] and triplet constraint [9] are adopted to learn a discriminative metric function. Generally, triplet constraints perform better than pairwise constraints to learn a discriminative metric function [9], [17]. For the bilinear similarity-based models, OASIS [4] is developed to explore a similarity metric for recognition tasks, and SOML [18] aims to learn a diagonal matrix for high dimensional cases with the similar setting as OASIS [4]. [19] presents an online multiple kernel similarity to tackle multi-modal tasks.

Unfortunately, these recently-proposed online metric learning methods cannot exploit the discriminative similarity relations for the strictly unaligned heterogeneous data in different evolution stages. To explore heterogeneous relationships among different data samples, [11] focuses on learning a nonlinear metric to distinguish the foreground boundary and background for robust visual tracking. Duan *et al.* [12] design fine-grained localized distance metrics to learn hierarchical nonlinear transformations between heterogeneous samples. Ding *et al.* [20] introduce the fast low-rank learning mechanism and representation denoising strategy to explore a more robust metric learning framework. Furthermore, [21] proposes a multi-modal distance metric method for image ranking by incorporating both click and visual representations in distance metric learning. [22] presents a multi-view stochastic learning model with high-order distance metric to explore modality-specific statistical information. However, above-mentioned metric methods cannot be successfully applied to the challenging online scenarios, where the features are evolving due to the different sensor lifespans (*e.g.*, some features are vanished and some new features are augmented).

B. Feature Evolution

For the feature evolution, with the assumption that there exists samples from both vanished feature space and augmented feature space in an overlapping period, [23] develops an evolvable feature learning model by reconstructing the vanished features and exploiting it along with new emerging features for large-scale streaming data. [24] proposes an one-pass incremental and decremental learning model for

streaming data, which consists of a compressing stage and an expanding stage. Different from [23], [24] assumes that there are overlapping features instead of overlapping period. Similar to [24], [25] focuses on learning the mapping function from two different feature spaces by using optimal transport technique. Furthermore, [26], [27] intend to classify trapezoidal data stream with feature and instance increasing doubly. However, the new emerging samples often have overlapping features with the previously existing samples. [28] develops an incremental feature learning model to tackle the emergence of new activity recognition sensors, which encourages the proposed model to well generalize the sudden emergence of incremental features.

Amongst the discussion above, there are no any feature evolution models highly related to our work except for OPID (OPIDe) [24]. However, there are several key differences between [24] and our EML model: 1) Our work is the first attempt to explore both instance and feature evolutions simultaneously via T-stage and I-stage in the metric learning field, when compared with [24]. 2) Due to the strictly unaligned evolving features in the different stages, we utilize the smoothed Wasserstein distance to explore the distance relationships among the heterogeneous and complex data, rather than the Euclidean distance in [24]. 3) Compared with [24], the low-rank regularizer for distance matrix could effectively learn a discriminative low-rank metric space, while neglecting non-informative knowledge for heterogeneous data in different feature evolution stages.

III. EVOLVING METRIC LEARNING (EML)

This section first reviews online metric learning, and then detailedly introduces how to tackle both instance and feature evolutions via our proposed EML model.

A. Revisit Online Metric Learning

Metric learning focuses on exploring an optimal distance metric matrix, in the light of different measure functions, e.g., Mahalanobis distance function: $d_M(x_p, x_q) = \sqrt{(x_p - x_q)^\top M (x_p - x_q)}$, where $x_p \in \mathbb{R}^d$ and $x_q \in \mathbb{R}^d$ are the p -th and q -th samples, respectively. $M \in \mathbb{R}^{d \times d}$ is the symmetric positive semi-definite matrix, which can be formulated as $L^\top L$ [5], where $L \in \mathbb{R}^{r \times d}$ (r denotes the rank of M) is the transformation matrix. The Mahalanobis distance function between x_p and x_q can be rewritten as $d_L(x_p, x_q) = \|L(x_p - x_q)\|_2$. Given an online constructed triplet (x_p, x_q, x_k) , L could be updated in an online manner via the Passive-Aggressive algorithm [29], i.e.,

$$L_t = \arg \min_L \frac{1}{2} \|L - L_{t-1}\|_F^2 + \frac{\gamma}{2} \ell_L(x_p, x_q, x_k), \quad (1)$$

where $\ell_L(x_p, x_q, x_k) = [1 + d_L(x_p, x_q) - d_L(x_p, x_k)]_+$ is a hinge loss function. $[z]_+ = \max(0, z)$. x_p and x_q belong to the same class, and x_p and x_k belong to different classes. $\gamma \geq 0$ is the regularization parameter.

However, most existing online metric learning models only focus on instance evolution with a fixed feature dimensionality, which cannot be utilized in the feature evolution scenario,

i.e., streaming data with incremental and decremental features. Furthermore, they mainly aim to promote the discrimination of the learned distance matrix L by minimizing the squared Mahalanobis distance from similar sample pairs. Especially, they assume that the feature descriptors of the sample pairs they focus on addressing are often aligned well in advance. Unfortunately, due to some unavoidable factors like non-linear lighting changes, heavily intensity noise and geometrical deformation, such assumption is heavily violated in the real-world tasks, especially for the feature evolution tasks. Therefore, the learned distance matrix L in Eq. (1) is not applicable and discriminative to explore similarity relationships between the heterogeneous and complex samples, whose evolving feature descriptors are not strictly aligned in different evolution stages [30].

B. The Proposed EML Model

This subsection first introduces how to integrate a smoothed Wasserstein distance into online metric formulation (i.e., Eq. (1)) to characterize the similarity relations of heterogeneous data with feature evolution in the different stages. Then the details about how to tackle feature evolution via Transforming stage (T-stage) and Inheriting stage (I-stage) in one-shot scenario are elaborated, followed by the extension of multi-shot scenario.

1) *Online Wasserstein Metric Learning*: Wasserstein distance [31] is an optimal transportation to transport all the earth from the source to target destination, while requiring the minimum amount of efforts. Formally, given two signatures $P = \{(x_{pi}, \mu_{pi})\}_{i=1}^m$ and $Q = \{(x_{qj}, \mu_{qj})\}_{j=1}^n$, the smoothed Wasserstein distance [32] between P and Q is:

$$W_\sigma(P, Q) = \min_{F \in \mathbb{F}(P, Q)} \langle D(P, Q), F \rangle - \sigma h(F), \quad (2)$$

$$s.t. \mathbb{F}(P, Q) = \{F | F \mathbf{1}_n = \mu_p, F^\top \mathbf{1}_m = \mu_q, F \geq 0\},$$

where $D(P, Q) = \{d_L(i, j)\}_{i,j=1}^{m,n} \in \mathbb{R}^{m \times n}$, and $d_L(i, j)$ denotes the cost of transporting one unit of earth from the source sample x_{pi} to the target sample x_{qj} . $F = \{f(i, j)\}_{i,j=1}^{m,n}$ indicates the flow network matrix, and $f(i, j)$ represents the amount of earth that is transported from x_{pi} to x_{qj} . $\mu_p = [\mu_{p1}, \dots, \mu_{pm}] \in \mathbb{R}^m$ and $\mu_q = [\mu_{q1}, \dots, \mu_{qn}] \in \mathbb{R}^n$ are normalized marginal probability mass vectors, and they satisfy $\sum_i \mu_{pi} = 1$ and $\sum_j \mu_{qj} = 1$. $\sigma \geq 0$ is a balance parameter, and $h(F) = -\langle F, \log(F) \rangle$ is the strictly concave entropic function.

In Eq. (2), the Mahalanobis distance is employed as ground distance to construct smoothed Wasserstein distance. Thus, each element $d_L(i, j)$ of $D(P, Q)$ in Eq. (2) represents the squared Mahalanobis distance between the source sample x_{pi} of P and the target sample x_{qj} of Q , i.e., $d_L(i, j) = \|L(x_{pi} - x_{qj})\|_2^2$. Given the online constructed triplet (P, Q, K) via [33], where the samples of P and Q belong to the same class, and the samples of P and K belong to different classes. After substituting Mahalanobis distance in Eq. (1) with the smooth Wasserstein distance

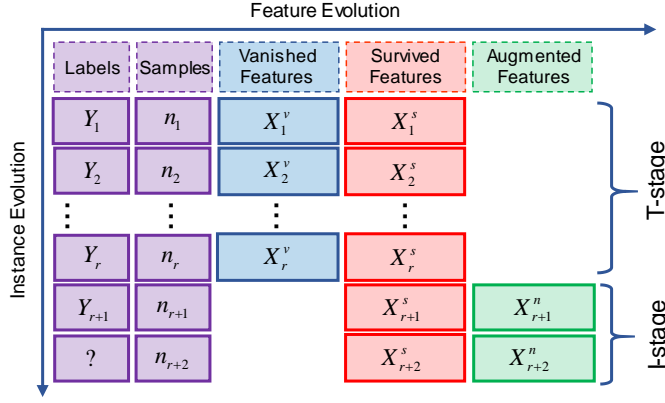


Fig. 2. The illustration of our EML model in one-shot scenario, which evolves instances and features simultaneously via T-stage and I-stage. Different colors denote different kinds of features, e.g., blue, red and green colors denote the vanished, survived and augmented features, respectively. The purple color indicates labels and the number of corresponding samples.

defined in Eq. (2), online Wasserstein metric learning could be formulated as follows:

$$\min_{L,F} \mathcal{L}_L(P, Q, K) = \frac{1}{2} \|L - L_{t-1}\|_F^2 + \frac{\gamma}{2} \ell_L(P, Q, K), \quad (3)$$

where $\ell_L(P, Q, K) = [1 + W_\sigma(P, Q) - W_\sigma(P, K)]_+$. When compared with the triplet (x_p, x_q, x_k) , each signature in (P, Q, K) consists of several samples belonging to same class rather than only one sample.

2) *Transforming Stage (T-stage) & Inheriting Stage (I-stage)*: In one-shot scenario where the features of streaming data would only be incremental and decremental by one time, two essential stages (*i.e.*, T-stage and I-stage) of our proposed EML model for steaming data with feature evolution are elaborated below.

I. Transforming Stage (T-stage): As shown in Fig. 2, suppose that $\{X_i, Y_i\}_{i=1}^r$ denotes the streaming data in the T-stage, where $X_i = [X_i^v, X_i^s] \in \mathbb{R}^{n_i \times (d_v + d_s)}$ and $Y_i \in \mathbb{R}^{n_i}$ denote the samples and labels in the i -th batch, respectively. r is the total batches in T-stage and n_i indicates the sample number in the i -th batch. Obviously, each instance of X_i consists of vanished and survived features, and d_v and d_s indicate the corresponding dimensions of vanished features $X_i^v \in \mathbb{R}^{n_i \times d_v}$ and survived features $X_i^s \in \mathbb{R}^{n_i \times d_s}$.

If we directly combine both vanished and survived features to learn a unified metric function, it fails to be utilized in I-stage where some features are vanished and some other new features are augmented. We thus propose to extract important information from vanished features and forward it into survived features by exploring a common discriminative metric space. In other words, we aim to train a model using only survived features to characterize the effective information extracted from both vanished and survived features.

In the i -th batch of T-stage, inspired by [33], the triplet (P_i^s, Q_i^s, K_i^s) for survived features is constructed in an online manner, where the samples of $P_i^s \in \mathbb{R}^{n_p \times d_s}$ and $Q_i^s \in \mathbb{R}^{n_q \times d_s}$ belong to same class while the samples of P_i^s and $K_i^s \in \mathbb{R}^{n_k \times d_s}$ belong to different classes. n_p, n_q and n_k are the numbers of samples in each signature. Likewise, we can construct the triplet (P_i^a, Q_i^a, K_i^a) for all features (containing

both vanished and survived features) in T-stage, where the samples of $P_i^a \in \mathbb{R}^{n_p \times (d_v + d_s)}$ and $Q_i^a \in \mathbb{R}^{n_q \times (d_v + d_s)}$ belong to same class while the samples of P_i^a and $K_i^a \in \mathbb{R}^{n_k \times (d_v + d_s)}$ belong to different classes.

Let $L^s \in \mathbb{R}^{k \times d_s}$ and $L^a \in \mathbb{R}^{k \times (d_v + d_s)}$ denote the distance matrices trained on survived features and all features (containing both vanished and survived features) in T-stage. Since the dimensions of L^s and L^a are different, it is reasonable to add some essential consistency constraints on the optimal distance matrices L^s and L^a to extract important information from vanished features, and forward it into survived features. Generally, based on the smoothed Wasserstein metric learning in Eq. (3), the formulation of the i -th batch in the T-stage could be expressed as follows:

$$\min_{L^s, L^a, F} \mathcal{L}_{L^s}(P_i^s, Q_i^s, K_i^s) + \mathcal{L}_{L^a}(P_i^a, Q_i^a, K_i^a) + \rho \mathcal{C}_{L^s, L^a}(P_i^s, Q_i^s, K_i^s; P_i^a, Q_i^a, K_i^a) + \lambda \text{rank}(L^s, L^a), \quad (4)$$

where $\mathcal{L}_{L^s}(P_i^s, Q_i^s, K_i^s)$ and $\mathcal{L}_{L^a}(P_i^a, Q_i^a, K_i^a)$ denote the triplet losses of smoothed Wasserstein metric learning on survived features and all features (containing both vanished and survived features), respectively. $\text{rank}(\cdot) = \text{rank}(L^s) + \text{rank}(L^a)$ denotes the regularization term, which learns the underlying low-rank property of heterogeneous samples. $\rho \geq 0$ and $\lambda \geq 0$ are the balance parameters. $\mathcal{C}_{L^s, L^a}(\cdot; \cdot)$ in Eq. (4) is designed to enable the consistence constraint for L^s and L^a , which aims to use only survived features to characterize the efficient information extracted from both vanished and survived features.

Specifically, $\mathcal{C}_{L^s, L^a}(\cdot; \cdot)$ constructs an essential triplet loss by incorporating smoothed Wasserstein metric learning on different feature spaces, *i.e.*, survived features and all features (containing both vanished and survived features). We attempt to compute the smoothed Wasserstein distance between different heterogeneous distributions based on vanished features and all features. For example, $W_\sigma(P_i^a, Q_i^s) = \{d_L(u, v)\}_{u,v=1}^{n_p, n_q} \in \mathbb{R}^{n_p \times n_q}$ denotes the smoothed Wasserstein distance between P_i^a from all features and Q_i^s from survived features, where $d_L(u, v) = \|L^a x_{pu}^a - L^s x_{qv}^s\|_2^2$ indicates the Mahalanobis distance between the u -th source sample x_{pu}^a of P_i^a and the v -th target sample x_{qv}^s of Q_i^s . Likewise, $W_\sigma(P_i^a, K_i^s)$, $W_\sigma(P_i^s, Q_i^a)$ and $W_\sigma(P_i^s, K_i^a)$ have similar definitions with $W_\sigma(P_i^a, Q_i^s)$. Formally, the consistence constraint $\mathcal{C}_{L^s, L^a}(\cdot; \cdot)$ is concretely expressed as follows:

$$\mathcal{C}_{L^s, L^a}(\cdot; \cdot) = [W_\sigma(P_i^a, Q_i^s) - W_\sigma(P_i^a, K_i^s) + 1]_+ + [W_\sigma(P_i^s, Q_i^a) - W_\sigma(P_i^s, K_i^a) + 1]_+. \quad (5)$$

II. Inheriting Stage (I-stage): Suppose that $\{X_{r+1}, Y_{r+1}\}$ denotes the data samples in the $r+1$ -th batch of I-stage, where $X_{r+1} = [X_{r+1}^s, X_{r+1}^n] \in \mathbb{R}^{n_{r+1} \times (d_s + d_n)}$ indicates the samples and $Y_{r+1} \in \mathbb{R}^{n_{r+1}}$ is the corresponding labels, as shown in Fig. 2. X_{r+1}^s and X_{r+1}^n represent the survived features and new augmented features in the $r+1$ -th batch. d_n and n_{r+1} are the dimension of the new augmented features and the number of samples. Thus, the goal of I-stage is to use $\{X_{r+1}, Y_{r+1}\}$ for training and make the prediction for the $r+2$ -th batch data $X_{r+2} = [X_{r+2}^s, X_{r+2}^n] \in \mathbb{R}^{n_{r+2} \times (d_s + d_n)}$ whose number of samples is same as that of $X_{r+1} \in \mathbb{R}^{n_{r+1} \times (d_s + d_n)}$.

To classify the $r+2$ -th batch data, we propose to inherit the metric performance of optimal distance matrix L^s learned on survived features in T-stage, since a set of common survived features exist in both T-stage and I-stage. Although we could construct the triplets directly from the $r+1$ -th batch for training, this trivial strategy has two significant shortcomings: 1) the trained metric model is difficult to be extended into multi-shot scenario; 2) the metric model learned only with the $r+1$ -th batch data would have worse prediction performance due to the lack of full usage of data in T-stage.

To this end, we utilize a similar stacking strategy with [34], [35], where [34], [35] focus on forming linear combinations of different predictors to train a unified classifier and achieve improved prediction accuracy. However, we propose to concatenate all feature descriptors as in stacking and train a unified predictor on the stacked features. It could better inherit the metric performance learned in T-stage. Concretely, let $Z_{r+1}^s = X_{r+1}^s (L^s)^\top \in \mathbb{R}^{n_{r+1} \times k}$ as the transformed discriminative metric space, which can be regarded as the new representation of X_{r+1}^s for stacking. X_{r+1} could then be represented as $Z_{r+1} = [Z_{r+1}^s, X_{r+1}^n] \in \mathbb{R}^{n_{r+1} \times (k+d_n)}$. Likewise, X_{r+2} is characterized as Z_{r+2} . Furthermore, we learn an optimal distance matrix $L^z \in \mathbb{R}^{k \times (k+d_n)}$ on Z_{r+1} with online constructed triplet $(P_{r+1}^z, Q_{r+1}^z, K_{r+1}^z)$, and evaluate the performance on Z_{r+2} , where the samples of P_{r+1}^z and Q_{r+1}^z belong to same class while the samples of P_{r+1}^z and K_{r+1}^z belong to different classes. Formally, at the t -th iterative step, the objective function of learning $L^z \in \mathbb{R}^{k \times (k+d_n)}$ in I-stage can be formulated as:

$$\begin{aligned} \min_{L^z, F} \frac{1}{2} \|L^z - L_{t-1}^z\|_F^2 + \lambda \text{rank}(L^z) \\ + \frac{\gamma}{2} [W_\sigma(P_{r+1}^z, Q_{r+1}^z) - W_\sigma(P_{r+1}^z, K_{r+1}^z) + 1], \end{aligned} \quad (6)$$

where $\gamma \geq 0$ and $\lambda \geq 0$ are the balance parameters. In our experiments, λ and γ in both Eq. (4) and Eq. (6) are set as the same value for simplification. $\text{rank}(L^z)$ denotes the regularization term, which aims to explore the intrinsic low-rank structure of heterogeneous samples in I-stage.

3) *Multi-shot Scenario*: Different from one-shot scenario, the features of streaming data in multi-shot scenario would be incremental and decremental M times. This subsection extends our model from one-shot case into multi-shot scenario, and the illustration example of multi-shot scenario when $M=2$ is depicted in Fig. 3. Specifically, $\{X_i, Y_i\}_{i=1}^{R/2}$ denotes the streaming data in Stage 1, where $X_i = [X_i^v, X_i^s] \in \mathbb{R}^{n_i \times (d_v+d_s)}$ and $Y_i \in \mathbb{R}^{n_i}$ respectively represent the samples and labels in the i -th batch. n_i indicates the sample number in the i -th batch, and $R/2$ denotes the total batches in Stage 1. When the streaming data $\{X_i, Y_i\}_{i=R/2+1}^R$ in Stage 2 arriving, it performs features evolution for the first time (*i.e.*, some features are vanished and some new features are augmented), where $X_i = [X_i^s, X_i^n] \in \mathbb{R}^{n_i \times (d_s+d_n)}$. Moreover, in Stage 3, the streaming data $\{X_{R+1}, Y_{R+1}\}$ performs features evolution for the second time, and we predict the results of our proposed EML model on the $R+2$ -th batch data X_{R+2} , where $X_{R+1} = [X_{R+1}^s, X_{R+1}^n]$ and $X_{R+2} = [X_{R+2}^s, X_{R+2}^n]$. Note that there are overlapped feature representations between

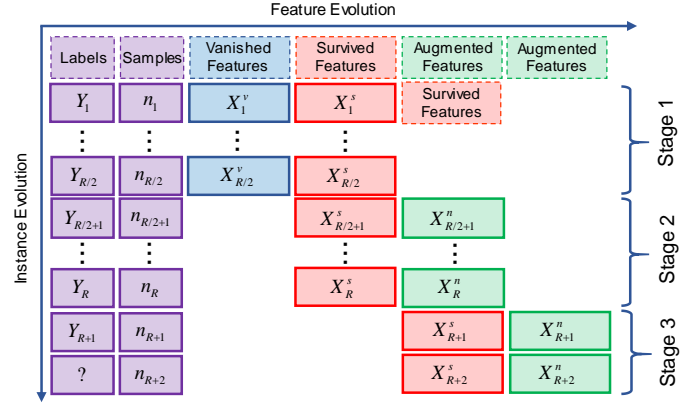


Fig. 3. The illustration of our EML model in multi-shot scenario when $M=2$, where Stage 1 and Stage 2 share the survived features, and Stage 2 and Stage 3 share the new augmented features. Specifically, our proposed model respectively regards Stage 1 and Stage 2 as T-stage and I-stage for the first feature evolution, and considers Stage 2 and Stage 3 as T-stage and I-stage for the second feature evolution.

any two adjacent stages. For example, as presented in Fig. 3, the survived features in Stage 1 are regarded as the vanished features in Stage 2, and the augmented feature in Stage 2 are considered as the survived features in Stage 3. Therefore, there are multiple Transforming stages (T-stage) and Inheriting stages (I-stage) in multi-shot scenario. To be specific, our proposed model first regards Stage 1 and Stage 2 as T-stage and I-stage for the first feature evolution. Then, it considers Stage 2 and Stage 3 as T-stage and I-stage for the second feature evolution. Generally, in multi-shot scenario, we have two essential learning tasks:

- Task I: Similar to the prediction task in one-shot case, we aim to classify testing data X_{R+2} in Stage 3 by training our proposed model on previous $R+1$ batch streaming data $\{X_i, Y_i\}_{i=1}^{R+1}$.
- Task II: Different from the prediction task in one-shot scenario, we attempt to make predictions for all stages (*i.e.*, Stage 1, Stage 2 and Stage 3 when $M=2$) by training our proposed model on the streaming data $\{X_i, Y_i\}_{i=1}^{R+1}$ in all stages.

IV. MODEL OPTIMIZATION

This section presents an alternating optimization strategy to update our proposed EML model amongst two stages, *i.e.*, T-stage and I-stage, followed by the computational complexity analysis of our model. The whole optimization strategy of our proposed EML model is introduced in **Algorithm 1**.

Note that the low-rank minimization in Eq. (4) and Eq. (6) is a well-known NP hard problem. Take L^z as an example, $\text{rank}(L^z)$ in Eq. (6) can be effectively surrogated by trace norm $\|L^z\|_*$. Different from traditional Singular Value Thresholding (SVT) [36], we employ a regularization term to guarantee the low-rank property, *i.e.*, $\|L^z\|_* = \text{tr}((L^z{}^\top L^z)^{1/2}) = \text{tr}(L^z{}^\top (L^z L^z{}^\top)^{-1/2} L^z)$. As a result, $\text{rank}(L^z)$ in Eq. (6) could be formulated as $\text{tr}(L^z{}^\top H^z L^z)$, where $H^z = (L^z L^z{}^\top)^{-1/2}$. Likewise, the low rank optimization of L^a and L^s shares the same strategy with L^z . $\text{rank}(L^a)$ and $\text{rank}(L^s)$ are respectively surrogated by $\text{tr}(L^a H^a L^a{}^\top)$

Algorithm 1 The Optimization of Our Proposed EML Model**Input:** The data $\{X_i, Y_i\}_{i=1}^{r+1}$, the parameters γ, λ, ρ ;**Output:** L^s and L^z ;

- 1: **Initialize:** L^s, L^a, L^z, F ;
- 2: **Transforming stage (T-stage):**
- 3: **for** $i = 1, \dots, r$ **do**
- 4: Calculate the smoothed Wasserstein distance for data X_i , and construct the triplets for training;
- 5: **repeat**
- 6: Solve F when fixing L^a and L^s ;
- 7: Update L^a via Eq. (8);
- 8: Update L^s via Eq. (10);
- 9: Update H^a and H^s via $H^a = (L^a L^{a\top})^{-1/2}$ and $H^s = (L^s L^{s\top})^{-1/2}$;
- 10: **until** Converge
- 11: **end for**
- 12: **Inheriting stage (I-stage):**
- 13: Transform X_{r+1} as Z_{r+1} to calculate smoothed Wasserstein distance, and construct the training triplets;
- 14: **repeat**
- 15: Solve the distance flow-network F when fixing L^z ;
- 16: Update L^z via Eq. (12);
- 17: Update H^z via $H^z = (L^z L^{z\top})^{-1/2}$;
- 18: **until** Converge

and $\text{tr}(L^s H^s L^{s\top})$, where $H^a = (L^a L^{a\top})^{-1/2}$ and $H^s = (L^s L^{s\top})^{-1/2}$.

A. Optimizing T-stage via an Alternating Strategy

1) *Updating L^a by fixing $\{L^s, H^a, F\}$:* When fixing the variables L^s, H^a and F , the optimization problem in Eq. (4) for solving variable L^a can be concretely expressed as:

$$\begin{aligned}
L_t^a = \arg \min_{L^a} & \frac{1}{2} \|L^a - L_{t-1}^a\|_F^2 + \lambda \text{tr}(L^a H^a L^{a\top}) \\
& + \frac{\gamma}{2} [\text{tr}(D(P_i^a, Q_i^a)F) - \text{tr}(D(P_i^a, K_i^a)F) + 1]_+ \\
& + \frac{\rho}{2} [\text{tr}(D(P_i^a, Q_i^s)F) - \text{tr}(D(P_i^a, K_i^s)F) + 1]_+ \\
& + \frac{\rho}{2} [\text{tr}(D(P_i^s, Q_i^a)F) - \text{tr}(D(P_i^s, K_i^a)F) + 1]_+.
\end{aligned} \quad (7)$$

The optimal solution of L_t^a could be relaxedly achieved by nulling the gradient of Eq. (7):

$$L_t^a = (L_{t-1}^a - \rho L_t^s (G_3 + G_4)) (I + \lambda H^a + \gamma G_1 + \rho G_2)^{-1}, \quad (8)$$

where $G_1 = Q_i^{a\top} \text{diag}(\mathbf{1}^\top F) Q_i^a - K_i^{a\top} \text{diag}(\mathbf{1}^\top F) K_i^a - P_i^{a\top} F Q_i^a - Q_i^{a\top} F^\top P_i^a + P_i^{a\top} F K_i^a + K_i^{a\top} F^\top P_i^a$, $G_2 = P_i^{a\top} \text{diag}(F \mathbf{1}) P_i^a - K_i^{a\top} \text{diag}(\mathbf{1}^\top F) K_i^a$, $G_3 = K_i^{s\top} F P_i^a - Q_i^{s\top} F^\top P_i^a$, $G_4 = P_i^{s\top} F K_i^a - K_i^{s\top} F^\top P_i^a$.

2) *Updating L^s by fixing $\{L^a, H^s, F\}$:* With the obtained distance matrix L^a and flow matrix F , the optimization

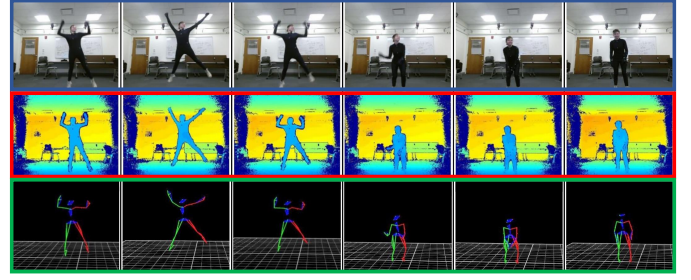


Fig. 4. The examples of human motions in EV-Action dataset, where the first, second and third rows denote the samples collected from RGB camera, Kinect sensor and motion capture sensor, respectively.

problem for variable L^s in Eq. (4) could be formulated as:

$$\begin{aligned}
L_t^s = \arg \min_{L^s} & \frac{1}{2} \|L^s - L_{t-1}^s\|_F^2 + \lambda \text{tr}(L^s H^s L^{s\top}) \\
& + \frac{\gamma}{2} [\text{tr}(D(P_i^s, Q_i^s)F) - \text{tr}(D(P_i^s, K_i^s)F) + 1]_+ \\
& + \frac{\rho}{2} [\text{tr}(D(P_i^a, Q_i^s)F) - \text{tr}(D(P_i^a, K_i^s)F) + 1]_+ \\
& + \frac{\rho}{2} [\text{tr}(D(P_i^s, Q_i^a)F) - \text{tr}(D(P_i^s, K_i^a)F) + 1]_+.
\end{aligned} \quad (9)$$

Concretely, the updating operator for L_t^s could be given as:

$$L_t^s = (L_{t-1}^s - \rho L_t^a (G_6 + G_8)) (I + \lambda H^s + \gamma G_5 + \rho G_7)^{-1}, \quad (10)$$

where $G_5 = Q_i^{s\top} \text{diag}(\mathbf{1}^\top F) Q_i^s - K_i^{s\top} \text{diag}(\mathbf{1}^\top F) K_i^s + P_i^{s\top} F K_i^s + K_i^{s\top} F^\top P_i^s - P_i^{s\top} F Q_i^s - Q_i^{s\top} F^\top P_i^s$, $G_6 = P_i^{a\top} F K_i^s - P_i^{a\top} F Q_i^s$, $G_7 = Q_i^{s\top} \text{diag}(F \mathbf{1}) Q_i^s - K_i^{s\top} \text{diag}(\mathbf{1}^\top F) K_i^s$, $G_8 = K_i^{a\top} F^\top P_i^s - Q_i^{a\top} F^\top P_i^s$.

3) *Updating F by fixing $\{L^a, L^s\}$:* When the distance matrices L^a and L^s are fixed, we split Eq. (4) into some independent traditional smoothed Wasserstein distance subproblems, which could be solved by [33]. We omit the detailed process of solving smoothed Wasserstein distance subproblems for simplification.

B. Optimizing I-stage via an Alternating Strategy

1) *Updating L^z by fixing $\{H^z, F\}$:* Given the fixed H^z and F , the formulation for L^z in Eq. (6) is rewritten as:

$$\begin{aligned}
L_t^z = \arg \min_{L^z} & \frac{1}{2} \|L^z - L_{t-1}^z\|_F^2 + \lambda \text{tr}(L^z H^z L^{z\top}) + \\
& \frac{\gamma}{2} [\text{tr}(D(P_{r+1}^z, Q_{r+1}^z)F) - \text{tr}(D(P_{r+1}^z, K_{r+1}^z)F) + 1]_+.
\end{aligned} \quad (11)$$

By nulling the gradient of Eq. (11), the optimization solution of Eq. (11) for L^z could be given as:

$$L_t^z = L_{t-1}^z (I + \lambda H^z + \gamma G_9)^{-1}, \quad (12)$$

where $G_9 = Q_{r+1}^{z\top} \text{diag}(\mathbf{1}^\top F) Q_{r+1}^z + P_{r+1}^{z\top} F K_{r+1}^z - K_{r+1}^{z\top} \text{diag}(\mathbf{1}^\top F) K_{r+1}^z + K_{r+1}^{z\top} F^\top P_{r+1}^z - P_{r+1}^{z\top} F Q_{r+1}^z - Q_{r+1}^{z\top} F^\top P_{r+1}^z$.

2) *Updating F by fixing L^z :* The optimization procedure of variable F in I-stage is same as that in T-stage: with the fixed L^z , the formulation Eq. (6) is split into some independent traditional smoothed Wasserstein distance subproblems, and we solve the variable F via [33].

TABLE I
THE EXPERIMENTAL SETTINGS IN ONE-SHOT SCENARIO.

Datasets	c	$\sum_{i=1}^r n_i$	n_i	d_v	d_s	d_n
EV-Action	20	4200	500, 600, 700	1024	1024	75
Mnist0vs5	2	3200	80, 160, 320	114	228	113
Mnist0vs3vs5	3	4800	120, 240, 480	123	245	121
Splice	2	2240	80, 160, 320	10	40	10
Gisette	2	6000	100, 200, 300	1239	2478	1238
USPS0vs5	2	960	120, 160, 240	64	128	64
USPS0vs3vs5	3	1440	180, 240, 300	64	128	64
Satimage	3	1080	60, 90, 120	10	18	8
ImageNet	1000	1200000	10000, 12000, 14000	512	1024	512
PAMAP2	18	7200	600, 700, 800	81	162	81

C. Computational Complexity Analysis

The main computational cost in our EML model involves the updating operations in both T-stage and I-stage. Specifically, in the T-stage, the computational costs of updating L^s and L^a are $O(kd_s + k(d_v + d_s)^2d_s + d_s^3)$ and $O(k(d_s + d_v) + kd_s^2(d_v + d_s) + (d_v + d_s)^3)$, respectively. For the I-stage, solving the variable L_z in Eq. (6) takes $O((k + d_a)^3)$. Besides, the computational cost of solving F in both T-stage and I-stage is $O(n_p^2n_q^2 + n_p^2n_k^2 + n_q^2n_k^2)$, where $n_p, n_q, n_k \ll n_i$. When compared with the feature dimension and sample number, the value of k is often small, and thus our proposed model is efficient to optimize in an online manner.

V. EXPERIMENTS

This section first presents detailed experimental configurations and competing methods. Then the experimental performance along with some analyses about our EML model in both one-shot and multi-shot cases are provided.

A. Configurations and Competing Methods

The experimental configurations of our EML model in one-shot scenario and some competing methods are detailedly introduced in this subsection.

1) *Experimental Configurations*: As shown in Table I, we conduct extensive comparisons on two real-world human motion recognition datasets (*i.e.*, EV-Action [13] and PAMAP2 [41]), a large-scale visual recognition dataset (*i.e.*, ImageNet [42]) and five synthetic benchmark datasets¹ containing three digit datasets (*i.e.*, Mnist, Gisette and USPS), one DNA dataset (*i.e.*, Splice) and one image dataset (*i.e.*, Satimage). Specifically, EV-Action dataset [13] is a human action dataset with 5300 samples, which consists of 20 common action categories, where 10 actions are finished by single subject and the others are accomplished by the same subjects interacting with other objects. It is a typical application for feature evolution in the real-world, where the features from depth information, RGB image, and human skeleton are respectively regarded as vanished, survived and augmented features. Some example samples about human actions are visualized as Fig. 4. PAMAP2 [41] is composed of 18 activities performed by 9 different subjects wearing three inertial measurement units (IMU) and a heart rate monitor. We only utilize the data information from IMU in our experiments, due to the large missing values

collected from the heart rate monitor. Each IMU contains one gyroscope, two accelerometers, one magnetometer, where the features from them are regarded as vanished, survived and augmented features, respectively. Moreover, ImageNet [42] including 1000 different categories is a large-scale challenging visual recognition dataset, where each of 1000 classes has roughly 1000 samples. We utilize ResNet [43] as feature extractor to obtain 2048-dimension feature representations for ImageNet [42].

For a fair comparison, as presented in Table I, we adopt the same experimental settings with [24] in one-shot and multi-shot cases, which are elaborated as follows:

- The number of streaming data in each batch is same, *i.e.*, $n_i = n_{r+1} = n_{r+2}$ ($i \in \{1, 2, \dots, r\}$), and the sample number in each class is equal for all training and testing batches.
- In T-stage, the total number of training data is fixed and the sample number in each batch is varied. In the light of this, the number of training and evaluation samples also varies in the last evaluation phase.
- We allocate the first d_v features, the next d_s features and the rest of features as vanished features, survived features and new augmented features, respectively. The first and last quarters are corresponding vanished and augmented features in our experiments.
- The experimental performance in each run may have slightly difference due to the influence of computer system and simulation environment, even though we run each experiment under the same experimental settings. To circumvent the randomness effect of experimental performance, all experimental results are the averaged results over fifty random runs, which is more convincing to illustrate the superiority of our EML model.

2) *Competing Methods*: We validate the superior performance of our EML model by comparing it with the following competing methods: One-pass **Pegasos** [37] assumes that the vanished and augmented features are available in different feature evolution stages; **OPMV** [38] regards the features in T-stage and I-stage as the first and second views; **TCA** [39] assumes that the streaming samples in T-stage and I-stage are drawn from the source and target distributions; **BDML** [2], **OPML** [7] and **CDML** [40] are the representative metric learning methods, which only utilize the samples with the augmented features, and ignore the previous vanished features; As for the feature evolution approaches, **OPID** and **OPIDe** [24] propose an one-pass incremental and decremental model for feature evolution. **FIRF** [28] designs a feature incremental random forest framework to tackle the emergence of new sensors (*i.e.*, new augmented features) in a dynamic environment.

B. Experiments in One-shot Scenario

In this subsection, we introduce the comprehensive experimental analysis, ablation studies, effects of hyper-parameters and convergence investigation of our proposed EML model in one-shot scenario, followed by computational costs of optimization complexity.

¹<http://archive.ics.uci.edu/ml/>

TABLE II

COMPARISONS BETWEEN OUR MODEL AND STATE-OF-THE-ART METHODS IN TERMS OF ACCURACY (%) ON TEN DATASETS: MEAN AND STANDARD ERRORS AVERAGED OVER FIFTY RANDOM RUNS IN ONE-SHOT SCENARIO. MODELS WITH THE BEST PERFORMANCE ARE BOLDDED.

Dataset	n_i	Pegasos [37]	OPMV [38]	TCA [39]	BDML [2]	OPML [7]	CDML [40]	OPIDe [24]	OPID [24]	FIRF [28]	Ours
EV-Action	500	57.38±1.51	56.37±1.91	53.88±2.04	56.42±0.71	54.10±1.71	55.08±0.83	57.84±1.06	57.57±1.08	57.13±0.84	58.87±0.68
	600	57.46±1.60	56.94±1.82	54.61±1.73	56.81±0.65	55.37±1.64	55.92±1.03	57.22±0.95	56.71±1.40	56.92±1.25	58.65±0.84
	700	57.22±1.34	56.68±1.87	54.37±1.69	56.63±0.77	55.82±1.62	56.22±0.71	57.09±1.13	56.85±1.27	57.23±1.16	58.32±0.82
Mnist Ovs5	80	97.74±0.73	97.39±0.92	96.53±1.75	97.00±1.66	96.45±1.72	96.75±1.32	98.68±0.88	98.88±0.99	98.14±0.87	99.85±0.91
	160	98.11±1.03	95.82±1.84	93.08±2.94	98.25±0.80	96.83±1.38	97.04±0.58	97.94±0.97	98.75±0.90	96.79±1.52	99.78±0.57
	320	97.68±0.79	96.47±1.79	92.43±3.82	98.24±0.75	96.98±1.03	97.16±0.85	97.38±0.58	97.21±0.66	96.83±1.37	99.27±0.37
Mnist Ovs3vs5	120	91.47±3.92	95.87±1.82	91.26±3.87	92.23±2.86	92.42±2.22	92.66±1.49	94.58±1.78	94.97±1.30	95.03±0.83	96.91±1.38
	240	89.95±3.08	93.96±1.18	90.85±1.74	92.87±1.40	91.99±1.64	92.47±1.31	93.45±1.41	93.48±1.35	94.24±1.13	95.37±0.92
	480	90.12±1.93	93.28±1.69	91.14±3.95	93.21±1.06	92.74±1.17	93.04±0.96	93.30±0.86	93.37±0.79	93.85±0.95	95.54±0.87
Splice	80	79.65±4.13	80.13±3.86	76.93±4.52	65.65±5.53	69.60±4.38	68.85±2.27	81.22±3.73	80.50±3.53	79.83±2.55	82.65±3.32
	160	82.25±3.26	81.95±2.84	80.93±3.47	71.55±4.07	78.21±2.53	75.85±2.65	84.00±2.03	83.91±2.05	82.06±1.91	85.25±2.06
	320	82.32±3.18	78.72±4.37	81.53±3.38	72.16±3.40	80.86±2.01	78.93±1.17	85.55±1.32	85.94±1.38	83.69±1.73	87.03±1.52
Gisette	100	97.53±1.33	95.27±2.85	94.11±3.35	90.25±3.13	94.17±3.02	93.71±2.39	97.14±1.28	97.56±1.26	94.21±0.96	97.29±1.25
	200	95.14±2.97	94.05±3.36	93.03±3.16	91.50±1.25	93.61±3.19	92.68±1.72	95.59±0.95	95.39±1.06	93.76±0.79	96.82±0.91
	300	96.84±1.35	93.71±3.11	94.37±3.72	93.83±2.12	93.77±2.96	93.24±1.56	96.36±0.69	95.33±0.93	94.18±1.06	97.89±0.43
USPS Ovs5	120	98.52±1.67	95.27±2.67	96.42±1.81	95.90±1.65	93.72±2.32	94.74±2.46	96.17±1.44	96.51±1.25	95.85±1.33	97.23±1.64
	160	97.84±0.82	95.65±1.72	95.46±2.13	96.38±1.23	93.04±4.05	95.21±1.57	96.78±1.31	96.93±1.00	95.75±1.12	98.91±0.67
	240	97.93±0.72	96.17±1.28	95.85±2.07	96.78±1.18	93.62±3.11	95.62±1.83	94.93±1.28	95.06±1.10	93.72±0.93	98.94±0.70
USPS Ovs3vs5	180	94.68±1.20	92.46±1.07	93.88±1.37	90.62±2.48	92.06±1.64	91.85±1.62	94.47±1.77	94.13±1.92	94.63±1.45	95.73±0.88
	240	94.39±1.09	91.69±2.31	92.94±1.58	91.48±1.68	91.23±1.73	91.73±1.24	92.08±1.93	92.50±1.66	93.36±2.07	95.52±1.26
	300	95.47±0.94	92.25±1.60	93.26±1.44	92.13±1.09	91.60±1.71	92.07±1.36	92.95±1.12	92.67±1.46	93.18±1.54	94.05±1.46
Satimage	60	94.25±2.56	96.48±1.47	97.25±1.08	97.14±1.59	97.47±1.59	97.39±1.46	98.17±2.19	97.60±2.31	97.92±2.05	99.20±0.91
	90	96.49±1.49	96.83±1.18	96.52±1.32	97.62±1.52	97.69±1.16	97.84±1.31	98.58±1.12	97.29±2.08	98.16±1.85	99.71±1.06
	120	98.03±1.13	97.38±1.94	97.12±1.87	97.12±1.48	97.15±1.49	97.22±1.63	98.45±1.14	96.85±1.94	97.24±1.36	99.52±1.07
ImageNet	10000	55.28±1.83	51.03±2.58	50.44±3.15	52.49±3.14	52.74±2.54	52.15±2.71	55.63±1.22	55.70±2.03	53.94±2.05	56.47±1.57
	12000	56.37±1.75	50.24±2.39	50.83±2.96	52.68±2.33	52.94±1.87	52.06±2.64	55.94±1.83	56.31±2.33	54.82±1.77	57.83±1.93
	14000	58.04±2.38	51.61±3.52	50.62±2.74	53.02±3.14	53.73±2.19	52.64±2.37	56.85±1.52	57.06±1.84	54.79±2.29	59.17±1.84
PAMAP2	600	91.64±1.08	89.85±1.33	85.73±1.84	87.67±1.74	86.23±1.81	87.92±2.16	92.17±0.93	92.64±1.05	93.56±0.84	95.27±0.71
	700	91.85±1.15	90.14±1.29	86.04±2.03	88.06±2.20	87.94±1.57	88.73±1.91	91.85±1.28	92.39±1.04	93.28±1.13	95.46±0.85
	800	91.57±0.89	90.25±1.56	85.49±2.75	88.75±2.06	88.13±1.90	89.37±1.68	93.05±0.88	93.62±0.93	93.84±0.79	95.66±0.94

TABLE III

ABLATION STUDY OF OUR EML MODEL IN ONE-SHOT SCENARIO.

Dataset	n_i	Ours-woT	Ours-woI	Ours-woW	Ours
EV-Action	500	56.68±1.74	54.36±1.61	57.93±0.85	58.33±0.76
	600	56.23±1.81	55.70±1.49	57.70±1.04	57.94±0.88
	700	57.02±1.56	55.93±1.76	57.83±0.92	58.12±0.86
Mnist Ovs5	80	97.85±1.24	96.70±1.71	98.90±0.97	99.07±0.94
	160	97.54±1.46	96.84±1.85	98.87±1.06	99.22±0.61
	320	97.23±3.34	96.88±0.96	98.95±0.83	99.27±0.37
Mnist Ovs3vs5	120	94.55±1.48	92.78±2.11	96.02±1.85	96.53±1.49
	240	93.49±1.07	92.88±1.31	94.88±1.37	95.37±0.92
	480	94.32±0.81	93.37±1.13	95.13±1.22	95.54±0.87
Splice	80	81.58±3.10	70.83±4.47	82.45±3.38	82.65±3.32
	160	84.07±2.51	78.87±3.01	84.87±2.19	85.25±2.06
	320	84.85±2.38	81.56±1.99	85.94±1.61	86.40±1.59
Gisette	100	95.22±1.30	92.47±1.68	96.84±1.40	97.29±1.25
	200	94.38±1.52	92.96±1.75	96.27±1.53	96.82±0.91
	300	96.11±0.95	95.08±1.19	97.14±0.87	97.79±0.46
USPS Ovs5	120	95.42±1.82	94.82±2.02	96.26±1.33	97.23±1.64
	160	96.04±1.33	94.95±1.70	97.03±1.47	98.31±0.82
	240	96.35±1.06	95.17±1.16	97.24±0.96	98.87±0.74
USPS Ovs3vs5	180	93.36±1.77	91.97±2.00	94.86±1.17	95.28±0.96
	240	93.13±1.38	92.01±1.45	94.33±1.54	94.96±1.37
	300	92.99±1.35	91.81±1.67	93.47±1.83	94.05±1.46
Satimage	60	96.50±1.59	97.43±1.36	98.31±1.10	98.97±0.95
	90	96.78±2.72	97.31±1.10	98.19±1.16	98.71±1.13
	120	96.22±1.91	97.23±1.22	98.02±1.22	98.53±1.20

1) *Experimental Analysis*: The experimental results for one-shot scenario are presented in Table II. From the presented performance, we have the following observations: 1) Although our proposed model has no access to the vanished features in T-stage, both transforming and inheriting strategies could efficiently exploit useful information of vanished feature and expand it into new augmented features in I-stage. 2) Our proposed EML model could be successfully applied to both high-dimensional (e.g., EV-Action, Gisette and ImageNet) and low-dimensional (e.g., Satimage and Splice) feature evolution, which are the challenging tasks to explore the intrinsic data

structure and informative knowledge using the existing features; 3) When we utilize the learned distance matrix in T-stage to assist the training procedure in I-stage, the evaluation performance of our proposed model increases significantly, even though the training samples in I-stage are relatively rare, i.e., n_i contains a small number of samples in I-stage. 4) Our EML model performs better than OPID and OPIDe [24], since T-stage could explore important information from vanished features, and I-stage efficiently inherits the metric performance from T-stage to take advantage of new augmented features.

2) *Ablation Studies*: To verify the effectiveness of our EML model, we intend to research the effects of different components of our model, i.e., training without T-stage (denoted as Ours-woT), training without I-stage (denoted as Ours-woI) and training without the Wasserstein distance metric (denoted as Ours-woW). The performance of Ours-woW is evaluated under the metric of Mahalanobis distance. From the presented results in Table III, our proposed EML model has the best performance when both transforming and inheriting strategies work together to tackle incremental and decremental features via the Wasserstein distance metric, which validates the reasonable design of our proposed model. Compared with other metric distances (e.g., Mahalanobis distance), the smoothed Wasserstein distance could better mine the similarity relationships between the heterogeneous and complex streaming samples, since the evolving features are not strictly aligned in different stages. Both T-stage and I-stage play an essential role in tackling instance and feature evolutions simultaneously.

3) *Effects of Hyper-Parameters*: In this subsection, as shown in Fig. 5, we introduce extensive parameter experiments on several representative datasets (MnistOvs3vs5, USPSOvs5,

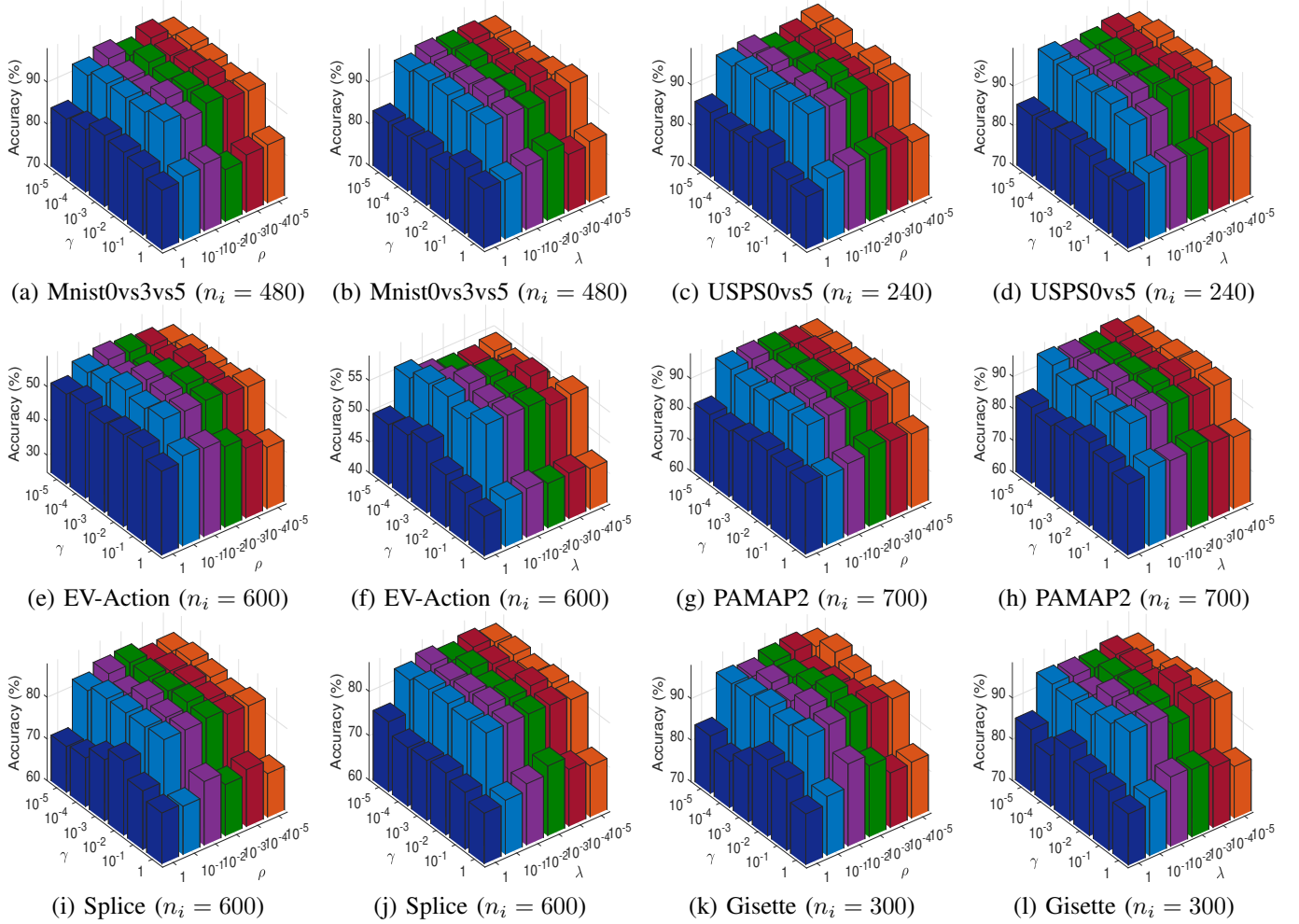


Fig. 5. Effect Investigations of hyper-parameters $\{\gamma, \rho\}$ when $\lambda = 10^{-4}$ and $\{\gamma, \lambda\}$ when $\rho = 10^{-3}$ on Mnist0vs3vs5 (a)(b), USPS0vs5 (c)(d), EV-Action (e)(f), PAMAP2 (g)(h), Splice (i)(j) and Gisette (k)(l) datasets in one-shot scenario.

TABLE IV

COMPUTATIONAL TIME IN TERMS OF MINUTES: MEAN AND STANDARD ERRORS AVERAGED OVER FIFTY RANDOM RUNS IN ONE-SHOT SCENARIO.

Dataset	Pegasus [37]	OPMV [38]	TCA [39]	BDML [2]	OPML [7]	CDML [40]	OPIDc [24]	OPID [24]	FIRF [28]	Ours
EV-Action ($n_i = 500$)	27.11±0.04	28.58±0.04	37.72±0.09	36.18±0.18	22.95±0.07	37.42±0.33	26.47±0.09	26.33±0.14	23.04±0.11	25.48±0.06
Mnist0vs5 ($n_i = 80$)	6.18±0.07	7.45±0.06	14.96±0.12	16.27±0.10	3.84±0.04	16.58±0.19	5.13±0.11	4.95±0.07	3.89±0.07	4.68±0.05
USPS0vs5 ($n_i = 120$)	3.16±0.02	4.75±0.10	11.93±0.04	13.06±0.12	1.24±0.03	13.42±0.21	1.93±0.05	1.87±0.06	1.32±0.05	1.53±0.08
Gisette ($n_i = 100$)	40.52±0.03	41.06±0.16	51.28±0.04	49.73±0.14	35.26±0.03	49.80±0.17	38.24±0.08	39.15±0.05	35.54±0.15	37.57±0.11
Satimage ($n_i = 120$)	2.64±0.05	3.08±0.10	10.23±0.07	11.47±0.09	0.52±0.02	11.62±0.14	0.66±0.04	0.71±0.08	0.55±0.04	0.68±0.03
PAMAP2 ($n_i = 360$)	9.84±0.04	9.27±0.06	16.74±0.15	18.05±0.18	4.69±0.13	18.32±0.09	7.84±0.09	7.63±0.07	4.81±0.07	6.45±0.06

EV-Action, PAMAP2, Splice and Gisette) as the examples to investigate the effects of hyper-parameters $\{\gamma, \lambda, \rho\}$ in one-shot scenario. Specifically, the experimental performances of our proposed model are averaged over fifty random repetitions by empirically tuning $\{\gamma, \lambda, \rho\}$ in a wide selection range of $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ to choose the optimal values of hyper-parameters. When fixing λ as 10^{-4} , we investigate the effects of $\{\gamma, \rho\}$, and introduce the hyper-parameter influence of $\{\gamma, \lambda\}$ when $\rho = 10^{-3}$. From the performance depicted in Fig. 5, we could observe our EML model has stable prediction performance over the wide selection range of different hyper-parameters. Moreover, when $\gamma = 10^{-2}$, $\rho = 10^{-3}$ and $\lambda = 10^{-4}$, our EML model performs the best prediction performance on most benchmark dataset,

except for Mnist0vs3vs5 dataset performing the best when $\gamma = 10^{-2}$, $\rho = 10^{-4}$ and $\lambda = 10^{-4}$.

4) *Convergence Investigations*: The convergence condition of **Algorithm 1** is depending on the little change (we set it as 2.5×10^{-5}) in the consecutive objective function values, and Fig. 6 depicts the convergence curves of our EML model on Mnist and USPS datasets. From the presented results in Fig. 6, we notice that our EML model could converge asymptotically to a stable value with respect to the objective function value after a few iterations. Furthermore, it validates that our proposed optimization algorithm could efficiently achieve stable performance with appropriate convergence condition.

5) *Computational Costs*: Table IV presents the computational costs (*i.e.*, optimization time) of our proposed EML

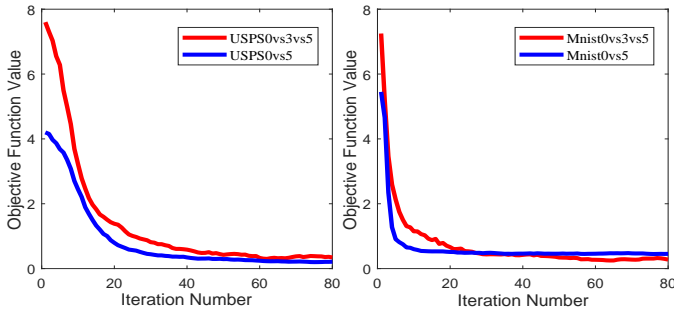


Fig. 6. The convergence analysis of our proposed EML model on the USPS (left) and Mnist (right) dataset in one-shot scenario.

TABLE V
EFFECT INVESTIGATION OF LOW RANK CONSTRAINT OF OUR EML MODEL IN ONE-SHOT SCENARIO.

Dataset	n_i	Ours-woLR	Ours
EV-Action	500	57.61±0.56	58.33±0.76
	600	57.46±0.89	57.94±0.88
	700	57.35±1.34	58.12±0.86
Mnist Ovs3vs5	120	96.18±1.36	96.53±1.49
	240	94.62±1.04	95.37±0.92
	480	95.20±1.01	95.54±0.87
Splice	80	82.16±0.94	82.65±3.32
	160	84.69±1.93	85.25±2.06
	320	85.76±1.45	86.40±1.59
Gisette	100	96.65±1.23	97.29±1.25
	200	96.40±1.16	96.82±0.91
	300	97.08±0.75	97.79±0.46
USPS Ovs3vs5	180	94.59±1.05	95.28±0.96
	240	94.51±1.27	94.96±1.37
	300	93.29±1.18	94.05±1.46
Satimage	60	98.42±1.22	98.97±0.95
	90	98.32±1.30	98.71±1.13
	120	98.11±1.05	98.53±1.20

model and other competing methods. From the reported results, we have the following conclusions: 1) Our model is computationally efficient in an online manner for real-world applications since $n_p, n_q, n_k \ll n_i$ and k is often a small value when compared with the feature dimension and the sample number. 2) The computational time costs (by the minute) of our model are less than other competing methods about 0.67~13.71 minutes on most experimental datasets except for OPML [7], since OPML only takes advantage of training samples in I-stage for optimization procedure.

6) *Effect Investigation of Low Rank Constraint*: This subsection investigates the effectiveness of low rank regularizer in our proposed EML model, as introduced in Table V. We substitute the low rank constraint with Frobenius norm, and denote its classification performance as Ours-woLR. The presented results in Table V clearly demonstrates that the performance of our proposed EML model degrades about 0.42% ~ 0.72% in terms of accuracy, when the low rank constraint is abandoned. It illustrates that our EML model could effectively explore the intrinsic low rank structure of heterogeneous samples for different evolving features by incorporating with the low rank regularizer.

C. Experiments in Multi-shot Scenario

This subsection introduces the experimental configurations and comparison performance of our proposed EML model in

multi-shot scenario.

1) *Experimental Configurations*: In multi-shot scenario, we set $M = 2$, *i.e.*, two-shot scenario with three stages for illustration, as depicted in Fig. 3. The streaming samples used in one-shot scenario are split into three stages. Except for the configurations introduced in one-shot scenario, the additional experimental configurations for multi-shot scenario are summarized as follows:

- All batches of T-stage in one-shot scenario are split into Stage 1 and Stage 2 with equal number of samples, as shown in Fig. 3. Under this setting, the survived features in Stage 2 would be the vanished features in Stage 3, and the new augmented features in Stage 2 would be the survived features in Stage 3. In other words, Stage 1 and Stage 2 are respectively considered as T-stage and I-stage for the first feature evolution. Moreover, Stage 2 and Stage 3 are regarded as T-stage and I-stage for the second feature evolution.
- The features are divided into four equal parts with the same partition order as one-shot scenario. Concretely, the second quarter is the shared part of Stage 1 and Stage 2. The third quarter is the shared part of Stage 2 and Stage 3. The first quarter in Stage 1 and the last quarter in Stage 3 denote the vanished and new augmented features.

2) *Experiments for Task I and II*: To address the Task I in multi-shot case, we directly use the last two adjacent evolution stages and regard it as the one-shot scenario for predictions, since the streaming data in any two adjacent stages share the common features. To be specific, we first utilize the transforming strategy in Eq. (4) on the streaming data in Stage 2 to learn the discriminative distance matrix, and the inheriting strategy in Eq. (6) is then applied to classify samples in Stage 3. To tackle the Task II in multi-shot scenario, we regard two adjacent stages as one-shot scenario (*i.e.*, T-stage and I-stage) and repeat this procedure until the last stage in multi-shot scenario. Specifically, the transforming and inheriting strategies are first integrated into Stage 1 and Stage 2, and then we make predictions on the second batch streaming data in Stage 2. After inheriting the metric performance of Stage 1, we extract the useful information from the survived features in Stage 2 and forward it into the new augmented features via common discriminative space, when new labeled streaming data in Stage 2 arriving. Furthermore, we perform the same inheriting strategy on survived features in Stage 2 to promote the performance predictions in Stage 3.

The experimental results of our proposed EML model averaged over fifty random repetitions for Task I and II are presented in Table VI and Fig. 7. Notice that: 1) Our model significantly outperforms other competing methods (*e.g.*, OPIDe and OPID [24]) especially in Task I, since it could inherit the metric performance of survived features in any two adjacent stages. 2) Compared with Task I, our proposed model performs better for Task II in most cases, since the survived features existing in Stage 1 could effectively promote the predictions for following streaming batches. 3) Our model could be successfully extended from one-shot case into multi-shot scenario to address both Task I and Task II, which further

TABLE VI

COMPARISONS BETWEEN OUR MODEL AND STATE-OF-THE-ART METHODS IN TERMS OF ACCURACY (%) ON SEVEN DATASETS: MEAN AND STANDARD ERRORS AVERAGED OVER FIFTY RANDOM RUNS IN MULTI-SHOT SCENARIO FOR TASK I. MODELS WITH THE BEST PERFORMANCE ARE BOLDED.

Dataset	n_i	Pegasos [37]	OPMV [38]	TCA [39]	BDML [2]	OPML [7]	CDML [40]	OPIDe [24]	OPID [24]	FIRF [28]	Ours
EV-Action	500	54.25±1.42	53.60±1.53	50.63±1.89	53.26±1.18	51.36±1.48	52.77±0.69	54.61±1.53	54.27±1.24	54.16±0.57	55.93±1.04
	600	55.59±1.27	53.72±1.66	51.83±1.96	53.74±0.82	53.11±1.36	52.74±1.28	54.43±1.15	53.46±1.92	53.84±1.07	56.74±1.35
	700	54.19±1.28	53.52±1.74	51.95±1.36	53.84±0.59	52.54±1.83	53.61±0.94	54.38±1.19	54.62±1.02	55.04±1.30	56.19±0.77
Mnist Ovs5	80	97.50±1.82	88.75±4.99	93.14±1.87	92.25±2.20	93.92±3.22	92.74±2.03	95.70±2.17	95.92±2.23	94.37±2.16	98.54±1.08
	160	97.56±1.28	90.75±3.02	91.78±2.43	95.70±1.26	94.34±1.54	95.87±1.85	95.53±1.61	95.29±1.80	94.16±1.85	98.61±0.57
	320	97.61±0.90	92.72±1.76	90.35±2.67	96.06±0.98	95.20±0.96	95.74±1.73	95.22±1.33	95.04±1.39	94.61±1.17	98.73±0.64
Gisette	100	91.58±2.87	86.24±4.76	91.28±2.56	90.48±3.29	89.87±3.62	90.26±2.71	95.08±2.35	94.36±1.88	92.88±2.06	96.12±1.18
	200	90.68±1.79	88.41±3.00	90.92±2.84	90.69±2.73	92.22±2.03	91.23±1.89	94.88±1.39	93.81±1.50	93.14±1.83	95.94±1.72
	300	91.18±1.13	89.42±2.27	92.13±2.14	92.52±1.71	91.83±1.61	92.06±1.36	94.65±1.12	93.91±1.45	93.08±1.26	95.71±1.68
USPS Ovs5	120	97.48±0.12	94.95±2.46	92.87±2.16	96.08±1.36	93.83±2.13	95.84±1.71	94.77±1.62	94.61±1.69	93.92±1.53	98.57±0.94
	160	97.56±0.19	95.17±2.18	93.05±1.94	96.24±1.55	94.21±1.66	95.49±1.33	94.13±1.54	94.51±1.30	93.75±1.28	98.68±0.65
	240	97.37±0.41	95.58±1.06	92.72±2.33	97.57±0.67	94.62±1.25	95.84±2.03	93.92±1.47	93.62±1.16	92.64±1.09	98.39±0.72
USPS Ovs3vs5	180	92.03±1.57	89.22±3.21	90.86±2.87	90.91±1.81	88.61±2.84	89.73±2.43	84.05±2.27	83.34±2.36	82.94±2.36	93.11±1.87
	240	90.90±1.40	89.13±2.05	90.24±2.93	91.98±2.04	89.62±2.17	90.62±1.87	84.68±1.73	84.49±1.92	83.75±1.81	93.23±1.58
	300	90.48±1.24	89.52±1.59	89.85±3.16	92.61±1.23	89.46±1.80	90.84±1.56	83.25±1.61	83.17±1.66	84.23±1.42	93.13±1.55
ImageNet	10000	52.76±1.55	48.11±2.30	47.82±3.06	49.04±2.83	49.74±2.49	49.43±2.59	50.42±1.57	50.33±2.26	48.42±2.01	53.92±1.62
	12000	53.81±1.44	48.15±2.06	48.82±2.47	50.17±2.42	50.63±1.42	49.95±2.51	52.36±2.13	53.09±2.28	51.64±1.93	54.95±2.17
	14000	55.82±2.05	48.93±3.28	48.36±2.66	50.93±2.87	51.62±2.62	50.51±2.18	53.62±1.83	54.96±2.10	51.63±2.14	57.03±1.92
PAMAP2	600	89.07±1.25	86.71±1.22	82.83±2.17	85.82±1.55	83.06±1.39	84.28±1.93	90.42±1.25	90.37±1.52	91.28±0.66	92.75±0.97
	700	88.47±1.58	87.05±1.42	83.23±1.69	85.52±1.74	82.12±1.26	85.68±1.64	88.63±1.33	90.18±0.94	90.57±1.24	92.51±0.92
	800	88.74±1.03	88.05±1.64	82.19±2.26	85.48±1.94	84.95±1.68	86.26±1.33	90.87±0.92	91.73±1.15	91.62±0.81	92.38±1.03

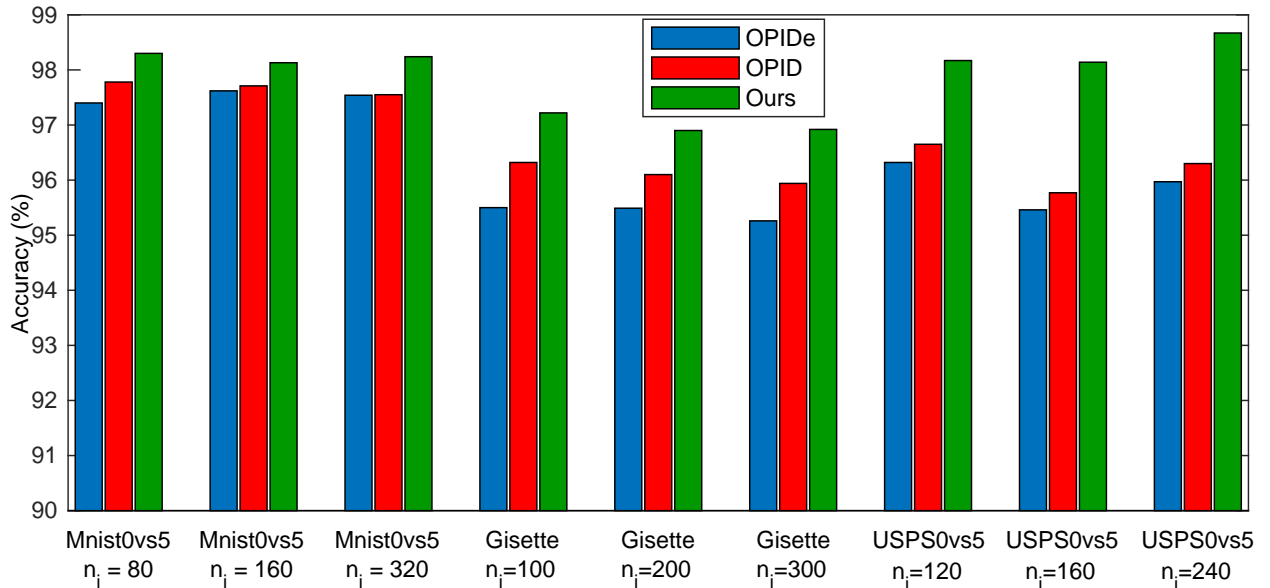


Fig. 7. Comparisons in terms of accuracy (%) on three datasets: mean and standard errors averaged over fifty random runs in multi-shot scenario for Task II.

verifies the superior performance of our EML model.

3) *Ablation Studies*: In this subsection, we conduct extensive variant experiments on Task I and Task II to investigate the efficiency of each component of our EML model in the multi-shot scenario, as introduced in Table VII and Table VIII. We have the following conclusions from the presented results: 1) All designed components in our EML model could cooperate well to achieve the best performance for both Task I and Task II in the multi-shot scenario, which validates the effectiveness and necessity of each module. 2) Two complementary strategies (*i.e.*, T-stage and I-stage) effectively compress the important information from vanished features and inherit the metric performance from the previous stage. They play an indispensable role in addressing both feature and instance evolutions simultaneously under the Wasserstein distance metric. 3) The performance degradation of Ours-woW illustrates the

TABLE VII

ABLATION STUDIES OF OUR PROPOSED EML MODEL IN MULTI-SHOT SCENARIO FOR TASK I.

Dataset	n_i	Ours-woT	Ours-woI	Ours-woW	Ours
Mnist Ovs5	80	94.93±0.34	94.26±0.38	96.04±0.62	98.54±1.08
	160	94.17±0.55	93.41±0.82	95.52±0.39	98.61±0.57
	320	96.13±0.59	95.47±0.85	96.84±1.03	98.73±0.64
Gisette	100	92.45±0.83	91.17±0.76	93.61±0.35	96.12±1.18
	200	92.84±0.72	92.04±0.28	94.12±0.46	95.94±1.72
	300	93.05±0.80	92.36±0.73	93.88±1.14	95.71±1.68
USPS Ovs5	120	95.54±0.75	94.18±0.93	96.33±0.41	98.57±0.94
	160	95.06±0.83	94.27±0.53	96.84±0.65	98.68±0.65
	240	95.36±0.32	94.91±0.77	97.05±0.41	98.39±0.72
USPS Ovs3vs5	180	90.15±0.19	89.35±0.87	90.94±0.51	93.11±1.87
	240	90.62±0.30	89.87±0.64	91.58±0.74	93.23±1.58
	300	90.86±0.81	90.22±0.63	92.08±0.26	93.13±1.55

effectiveness of the smoothed Wasserstein distance to explore the similarity relationships for heterogeneous samples among different stages.

TABLE VIII
ABLATION STUDIES OF OUR PROPOSED EML MODEL IN MULTI-SHOT
SCENARIO FOR TASK II.

Dataset	n_i	Ours-woT	Ours-woI	Ours-woW	Ours
Mnist Ovs5	80	94.58±1.48	93.26±1.71	95.93±1.22	98.30±1.18
	160	95.71±1.29	93.62±1.48	96.04±0.98	98.13±1.16
	320	95.25±0.93	94.54±0.89	95.87±1.13	98.24±1.34
Gisette	100	94.32±0.88	93.68±1.09	95.02±0.95	97.22±1.37
	200	93.10±1.27	92.29±1.48	94.21±1.07	96.90±1.18
	300	93.74±1.18	92.16±1.26	94.23±0.89	96.92±1.70
USPS Ovs5	120	96.15±1.36	95.62±1.26	96.87±0.88	98.17±1.18
	160	94.77±1.38	94.46±1.52	95.45±1.13	98.14±0.97
	240	95.32±1.36	94.37±1.65	96.14±0.76	98.67±0.62
USPS Ovs3vs5	180	90.38±1.47	90.56±1.51	91.06±1.05	93.94±1.50
	240	90.43±1.26	89.85±1.34	91.88±1.17	93.34±1.38
	300	90.35±1.18	90.83±1.43	91.46±0.84	94.58±1.25

VI. CONCLUSION

In this paper, an online Evolving Metric Learning (EML) model is proposed for both instance and feature evolutions, which is successfully applied to one-shot and multi-shot scenarios. Our proposed EML model contains two essential stages, *i.e.*, Transforming stage (T-stage) and Inheriting stage (I-stage). To be specific, for the T-stage, we utilize the survived features to characterize the effective information extracted from vanished and survived features by exploiting a common discriminative metric space. In the I-stage, we inherit the metric performance of survived features from T-stage, and extend it into the new augmented features. Furthermore, we apply the smoothed Wasserstein distance to T-stage and I-stage to better explore the similarity relations of heterogeneous streaming data among different evolution stages. Extensive experiments show the superior performance of our proposed EML model on several representative datasets. In the future, we will consider lifelong machine learning for both instance and feature evolutions, which continually learns a sequence of new streaming evolution tasks without the catastrophic forgetting for the previous learned evolution tasks.

REFERENCES

- [1] H. V. Nguyen and L. Bai, "Cosine similarity metric learning for face verification," in *Computer Vision – ACCV 2010*, R. Kimmel, R. Klette, and A. Sugimoto, Eds., 2011, pp. 709–720.
- [2] J. Xu, L. Luo, C. Deng, and H. Huang, "Bilevel distance metric learning for robust image recognition," in *Proceedings of the 32Nd International Conference on Neural Information Processing Systems*, 2018, pp. 4202–4211.
- [3] Z. Boukrouvalas, "Distance metric learning for medical image registration," 2011.
- [4] G. Chechik, V. Sharma, U. Shalit, and S. Bengio, "Large scale online learning of image similarity through ranking," *J. Mach. Learn. Res.*, vol. 11, pp. 1109–1135, Mar. 2010.
- [5] P. Jain, B. Kulis, I. S. Dhillon, and K. Grauman, "Online metric learning and fast similarity search," in *Advances in Neural Information Processing Systems 21*, 2009, pp. 761–768.
- [6] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *J. Mach. Learn. Res.*, vol. 10, pp. 207–244, Jun. 2009.
- [7] 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.
- [8] R. Jin, S. Wang, and Y. Zhou, "Regularized distance metric learning: theory and algorithm," in *Advances in Neural Information Processing Systems 22*, 2009, pp. 862–870.
- [9] B. Shaw, B. Huang, and T. Jebara, "Learning a distance metric from a network," in *Advances in Neural Information Processing Systems 24*, 2011, pp. 1899–1907.

- [10] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *Proceedings of the 24th International Conference on Machine Learning*. ACM, 2007, pp. 209–216.
- [11] J. Hu, J. Lu, and Y. Tan, "Deep metric learning for visual tracking," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 26, no. 11, pp. 2056–2068, 2016.
- [12] Y. Duan, J. Lu, J. Feng, and J. Zhou, "Deep localized metric learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 10, pp. 2644–2656, 2018.
- [13] L. Wang, B. Sun, J. P. Robinson, T. Jing, and Y. Fu, "Ev-action: Electromyography-vision multi-modal action dataset," *arXiv preprint arXiv:1904.12602*, 2019.
- [14] C. K. Ho, A. Robinson, D. R. Miller, and M. J. Davis, "Overview of sensors and needs for environmental monitoring," *Sensors*, vol. 5, no. 1, pp. 4–37, 2005.
- [15] S. Shalev-Shwartz, Y. Singer, and A. Y. Ng, "Online and batch learning of pseudo-metrics," in *Proceedings of the Twenty-first International Conference on Machine Learning*. ACM, 2004, p. 94.
- [16] 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, Oct 2019.
- [17] Q. Qian, R. Jin, J. Yi, L. Zhang, and S. Zhu, "Efficient distance metric learning by adaptive sampling and mini-batch stochastic gradient descent (sgd)," *Machine Learning*, vol. 99, no. 3, pp. 353–372, Jun 2015.
- [18] 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*, 2014.
- [19] 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, March 2014.
- [20] Z. Ding, M. Shao, W. Hwang, S. Suh, J.-J. Han, C. Choi, and Y. Fu, "Robust discriminative metric learning for image representation," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 11, pp. 3173–3183, Nov. 2019.
- [21] J. Yu, X. Yang, F. Gao, and D. Tao, "Deep multimodal distance metric learning using click constraints for image ranking," *IEEE Transactions on Cybernetics*, vol. 47, no. 12, pp. 4014–4024, 2017.
- [22] J. Yu, Y. Rui, Y. Y. Tang, and D. Tao, "High-order distance-based multi-view stochastic learning in image classification," *IEEE Transactions on Cybernetics*, vol. 44, no. 12, pp. 2431–2442, 2014.
- [23] B.-J. Hou, L. Zhang, and Z.-H. Zhou, "Learning with feature evolvable streams," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1416–1426.
- [24] C. Hou and Z.-H. Zhou, "One-pass learning with incremental and decremental features," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, pp. 2776–2792, 2018.
- [25] H.-J. Ye, D.-C. Zhan, Y. Jiang, and Z.-H. Zhou, "Rectify heterogeneous models with semantic mapping," in *ICML*, 2018.
- [26] Q. Zhang, P. Zhang, G. Long, W. Ding, C. Zhang, and X. Wu, "Towards mining trapezoidal data streams," *2015 IEEE International Conference on Data Mining*, pp. 1111–1116, 2015.
- [27] Q. Zhang, P. Zhang, G. Long, W. Ding, C. Zhang, and X. Wu, "Online learning from trapezoidal data streams," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 10, pp. 2709–2723, Oct 2016.
- [28] C. Hu, Y. Chen, X. Peng, H. Yu, C. Gao, and L. Hu, "A novel feature incremental learning method for sensor-based activity recognition," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 6, pp. 1038–1050, June 2019.
- [29] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online passive-aggressive algorithms," *J. Mach. Learn. Res.*, vol. 7, pp. 551–585, Dec. 2006.
- [30] J. Xu, L. Luo, C. Deng, and H. Huang, "Multi-level metric learning via smoothed wasserstein distance," in *IJCAI*, 2018, pp. 2919–2925.
- [31] R. Sandler and M. Lindenbaum, "Nonnegative matrix factorization with earth mover's distance metric for image analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, pp. 1590–1602, Aug 2011.
- [32] M. Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," in *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, 2013, pp. 2292–2300.
- [33] A. Rolet, M. Cuturi, and G. Peyré, "Fast dictionary learning with a smoothed wasserstein loss," in *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, vol. 51, pp. 09–11 May 2016, pp. 630–638.

- [34] L. Breiman, "Stacked regressions," *Machine Learning*, vol. 24, pp. 49–64, Jul 1996.
- [35] Z.-H. Zhou, *Ensemble Methods: Foundations and Algorithms*. Chapman & Hall/CRC, 2012.
- [36] J.-F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.
- [37] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter, "Pegasos: primal estimated sub-gradient solver for svm," *Mathematical Programming*, vol. 127, pp. 3–30, Mar 2011.
- [38] Y. Zhu, W. Gao, and Z.-H. Zhou, "One-pass multi-view learning," in *ACML*, 2015.
- [39] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang, "Domain adaptation via transfer component analysis," *IEEE Transactions on Neural Networks*, vol. 22, no. 2, pp. 199–210, Feb 2011.
- [40] S. Chen, L. Luo, J. Yang, C. Gong, J. Li, and H. Huang, "Curvilinear distance metric learning," in *Advances in Neural Information Processing Systems 32*, 2019.
- [41] A. Reiss and D. Stricker, "Introducing a new benchmarked dataset for activity monitoring," in *2012 16th International Symposium on Wearable Computers*, 2012, pp. 108–109.
- [42] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *arXiv preprint arXiv:1512.03385*, 2015.



Jiahua Dong Jiahua Dong is currently a Ph. D candidate in State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences. He received the B.S. degree from Jilin University in 2017. His current research interests include computer vision, machine learning, transfer learning, domain adaptation and medical image processing.



Yang Cong Yang Cong (S'09-M'11-SM'15) is a full professor of Chinese Academy of Sciences. He received the he B.Sc. de. degree from Northeast University in 2004, and the Ph.D. degree from State Key Laboratory of Robotics, Chinese Academy of Sciences in 2009. He was a Research Fellow of National University of Singapore (NUS) and Nanyang Technological University (NTU) from 2009 to 2011, respectively; and a visiting scholar of University of Rochester. He has served on the editorial board of the *Journal of Multimedia*. His current research interests include image processing, compute vision, machine learning, multimedia, medical imaging, data mining and robot navigation. He has authored over 70 technical papers. He is also a senior member of IEEE.



Gan Sun Gan Sun (S'19) is an Assistant Professor in State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences. He received the B.S. degree from Shandong Agricultural University in 2013, the Ph.D. degree from State Key Laboratory of Robotics, Chinese Academy of Sciences in 2020, and has been visiting Northeastern University from April 2018 to May 2019, Massachusetts Institute of Technology from June 2019 to November 2019. His current research interests include lifelong machine learning, multi-task learning, medical data analysis, deep learning and 3D computer vision.



Tao Zhang Tao Zhang is currently working toward the Ph.D. degree in pattern recognition and intelligent systems at the State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences, Shenyang, China. His research interests include pattern recognition, image processing, tactile sensing and robotics.



Xu Tang Xu Tang is currently a reserach associate in State Key Laboratory of Robotics, Shenyang Institute of Automation. He received MESc degree from Harbin Institute of Technology in 2017. His current research interests include computer vision and machine learning.



Xiaowei Xu Xiaowei Xu is a professor of Information Science at University of Arkansas at Little Rock (UALR), received a a B.Sc. de. degree in Mathematics from Nankai University in 1983 and a Ph.D. degree in Computer Science from University of Munich in 1998. He holds an adjunct professor position in the Department of Mathematics and Statistics at University of Arkansas at Fayetteville. Before his appointment in UALR, he was a senior research scientist in Siemens. He was a visiting professor in Microsoft Research Asia and Chinese University of Hong Kong. His research spans data mining, machine learning, bioinformatics, data management and high performance computing. He has published over 70 papers in peer reviewed journals and conference proceedings. His groundbreaking work on density-based clustering algorithm DBSCAN has been widely used in many textbooks; and received over 10203 citations based on Google scholar. Dr. Xu is a recipient of 2014 ACM KDD Test of Time Award that "recognizes outstanding papers from past KDD Conferences beyond the last decade that have had an important impact on the data mining research community."