

# Transfer Metric Learning: Algorithms, Applications and Outlooks

Yong Luo, Yonggang Wen, *Senior Member, IEEE*, Ling-Yu Duan, *Member, IEEE*,  
and Dacheng Tao, *Fellow, IEEE*

**Abstract**—Distance metric learning (DML) aims to find an appropriate way to reveal the underlying data relationship. It is critical in many machine learning, pattern recognition and data mining algorithms, and usually require large amount of label information (such as class labels or pair/triplet constraints) to achieve satisfactory performance. However, the label information may be insufficient in real-world applications due to the high-labeling cost, and DML may fail in this case. Transfer metric learning (TML) is able to mitigate this issue for DML in the domain of interest (target domain) by leveraging knowledge/information from other related domains (source domains). Although achieved a certain level of development, TML has limited success in various aspects such as selective transfer, theoretical understanding, handling complex data, big data and extreme cases. In this survey, we present a systematic review of the TML literature. In particular, we group TML into different categories according to different settings and metric transfer strategies, such as direct metric approximation, subspace approximation, distance approximation, and distribution approximation. A summarization and insightful discussion of the various TML approaches and their applications will be presented. Finally, we indicate some challenges and provide possible future directions.

**Index Terms**—Distance metric learning, transfer learning, survey, machine learning, data mining

## 1 INTRODUCTION

IT is critical to evaluate the distances between samples in pattern analysis and machine learning applications. If an appropriate distance metric can be obtained, even the simple  $k$ -nearest neighbor ( $k$ -NN) classifier, or  $k$ -means clustering can perform well [1], [2]. In addition, for large-scale and efficient information retrieval, the results are usually obtained directly according to the distances to the query [3], and a good distance metric is also the key of many other important applications, such as face verification [4] and person re-identification [5].

To learn a reliable distance metric, we usually need large amount of label information, which can be the class labels or target values as used in the typical machine learning approaches (such as classification or regression), and it is more common to utilize some pair or triplet-based constraints [6]. Such constraints are weakly-supervised since the exact label for an individual sample is unknown. However, in real-world applications, the label information is often scarce since manually labeling is labor-intensive and it is exhausted or even impossible to collect abundant side information for a new learning problem.

Transfer learning [7], which aims to mitigate the label deficiency issue in model training, is thus introduced to improve the performance of distance metric learning (DML)

when the label information is insufficient in a target domain. This leads to the so-called transfer metric learning (TML), which has been found to be very useful in many applications. For example, in face verification [8], the main step is to estimate the similarities/distances between face images. The data distributions of the images captured under different scenarios vary due to the varied background, illumination, etc. Therefore, the metric learned in one scenario may be not effective in a new scenario and TML would be helpful. In person re-identification [5], [9], the key is to estimate the similarities/distances between images of persons appeared in different cameras. The data distributions of the images captured using different cameras vary due to the varied camera setting and scenario. In addition, the distribution for the same camera may change over time. Hence, calibration is needed to achieve satisfactory performance and TML is able to reduce such effort. A more general example is image retrieval, where the data distributions of images in different datasets vary [10]. It would also be very useful to utilize expensive or semantic features to help learn a metric for cheap features or the ones that are hard to be interpreted [11], [12].

In the past decade, dozens of works have been proposed in this area and we provide in this survey a comprehensive overview of these methods. In this survey, we aim to make the machine learners quickly grasp the TML research area, and facilitate the chosen of appropriate methods for machine learning practitioners. Besides, there still be many issues to be tackled in TML, and we hope that some new ideas can be inspired from this survey.

The rest of this survey is organized as follows. We first present the background and overview of TML in Section 2, which includes a brief history of TML, the main notations used throughout the paper, and a categorization of

- Y. Luo and Y. Wen are with the School of Computer Science and Engineering, Nanyang Technological University, Singapore.  
E-mail: yluo180@gmail.com, ygwen@ntu.edu.sg
- L. Duan is with the School of Electronics Engineering and Computer Science, Institute of Digital Media, Peking University, China.  
E-mail: lingyu@pku.edu.cn
- D. Tao is with the UBTECH Sydney Artificial Intelligence Centre and the School of Information Technologies in the Faculty of Engineering and Information Technologies at University of Sydney, 6 Cleveland St, Darlington, NSW 2008, Australia.  
E-mail: dacheng.tao@sydney.edu.au

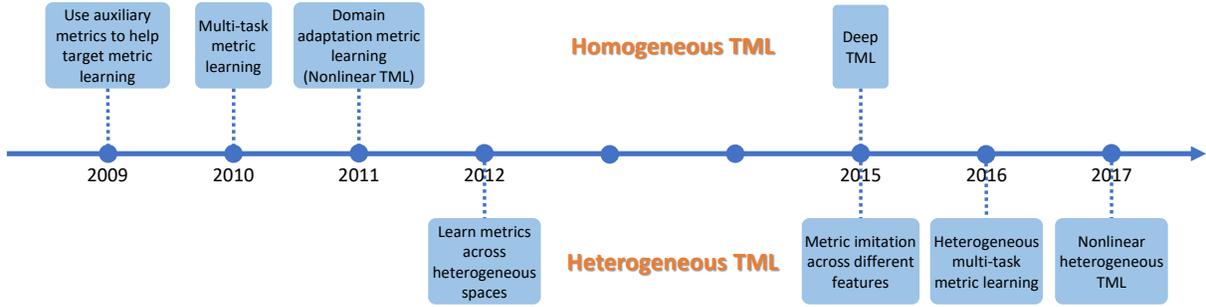


Fig. 1. Evolution of transfer metric learning, which has been studied for almost ten years.

the TML approaches. In the subsequent two sections, we give a detailed description of the approaches in the two main categories, i.e., homogeneous and heterogeneous TML respectively. Section 5 is a summarization of the different applications of TML and finally, we conclude this survey and identify some possible future directions in Section 6.

## 2 BACKGROUND AND OVERVIEW

### 2.1 A brief history of transfer metric learning

Transfer metric learning (TML) is a relatively new research field. The works that explicitly applying transfer learning to improve DML start around the year of 2009. For example, multiple auxiliary (source) datasets are utilized in [14] to help the metric learning on the target set. The main idea is to enforce the target metric to be close to the different source metrics. An adaptive weight is learned to reflect the contribution of each source metric to the target metric. In [15], such contribution is determined by learning a covariance matrix between the different metrics. Instead of directly learning the target metric, the decomposition based method [16] assumes that the target metric can be represented as a linear combination of multiple base metrics, which can be derived from the source metric. Hence, the metric learning is casted as learning combination coefficients, where the parameters to be learned can be much fewer.

We can not only using source metrics to help the target metric learning, but also make the different DML tasks help each other. The latter is often called multi-task metric learning (MTML). One representative work is the multi-task extension [17] of a well-known DML algorithm LMNN [2]. Some other related works including GPMTML [18], MtMCML [5] and CP-mtML [10]. In addition, there are a few domain adaptation metric learning approaches [19], [20]. Most of the above methods can only learn linear metric for the target domain. The domain adaptation metric learning (DAML) approach presented in [19] is able to learn nonlinear target metric based on the kernel method. Recently, neural network is also employed to conduct nonlinear metric transfer [8] by taking the advantage of deep learning technique [21].

The study of heterogeneous TML is a bit later than homogeneous TML and there are much fewer works than those in the homogeneous setting. To the best of our best knowledge, the first work that explicitly designed for heterogeneous TML is the one presented in [22], but it is limited in that only two domains (one source and target domain)

can be handled. There exist a few tensor based approaches [23], [24] for heterogeneous MTML, where the high-order correlations between all domains are exploited. A main disadvantage of these approaches is that the computational complexity is high. Dai et al. [11] proposes an unsupervised heterogeneous TML algorithm, which aims to use some “expensive” (sophisticated, off-the-shelf) features to help learn a metric for relatively “cheap” feature. This is also termed metric imitation. Recently, a general heterogeneous TML framework is proposed in [12], [25]. The framework first extracts some knowledge fragments (linear or nonlinear mappings) from pre-trained source metric, and then using these fragments to help the target domain learn either linear or nonlinear distance metric. The framework is flexible and easy-to-use. An illustration figure for the evolution of TML is shown in Fig. 1.

### 2.2 Notations and definitions

In this survey, we assume there are  $M$  different domains, and the  $m$ 'th domain is associated with a feature space  $\mathcal{X}_m$  and marginal distribution  $P_m(X_m)$ . Without loss of generality, we assume the  $M$ 'th (the last) domain is the target domain, and all the remained ones are source domains. If there is only one source domain, we signify it using the script “ $S$ ”. In distance metric learning (DML), the task is to learn a distance function for any two instances, i.e.,  $d_\phi(\mathbf{x}_i, \mathbf{x}_j)$ , which must satisfy several properties including nonnegativity, identity, symmetry and triangle inequality [6]. Here,  $\phi$  is the parameter of the distance function, and we call it *distance metric* in this survey. For a nonlinear distance metric,  $\phi$  is often given by a nonlinear feature mapping. The linear metric is denoted as  $A$ , which is a positive semi-definite (PSD) matrix and adopted in the popular Mahalanobias metric learning [1].

To learn the metric in the  $m$ 'th domain, we assume there is a training set  $\mathcal{D}_m$ , which contains  $N_m$  samples with  $\mathbf{x}_{mi} \in \mathbb{R}^{d_m}$  to be the feature representation for the  $i$ 'th sample. In a fully-supervised scenario, the corresponding label  $y_{mi}$  is also given. However, DML is usually conducted in a weakly-supervised manner, where only some similar/dissimilar constraints on training sample pairs  $(\mathbf{x}_{mi}, \mathbf{x}_{mj})$  are provided. Alternatively, the constraint can be a relative comparison for a training triplet  $(\mathbf{x}_{mi}, \mathbf{x}_{mj}, \mathbf{x}_{mk})$ , e.g.,  $\mathbf{x}_{mi}$  is more similar to  $\mathbf{x}_{mj}$  than to  $\mathbf{x}_{mk}$  [6].

In traditional DML, we are often provided with abundant labeled data (such as samples with similar/dissimilar

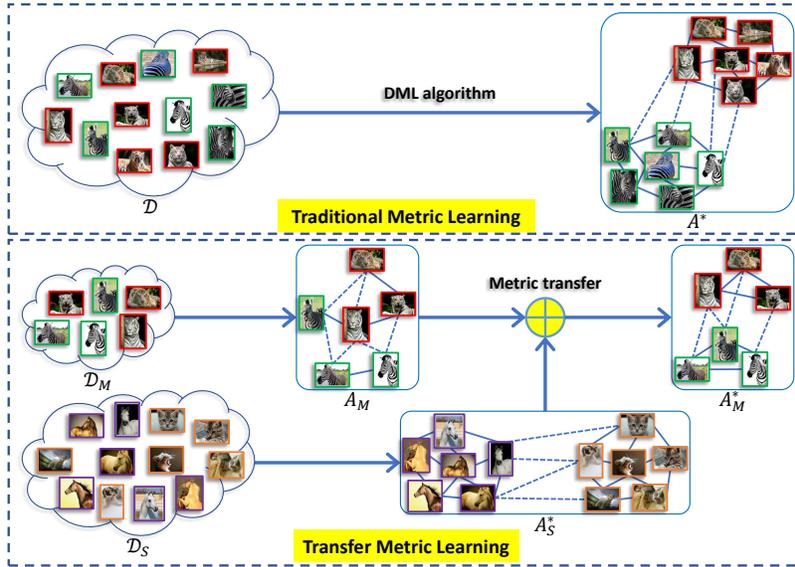


Fig. 2. An illustration of traditional distance metric learning (DML) and transfer metric learning (TML). Given abundant labeled data, DML aims to learn a distance function between samples so that their distance is small if semantically similar and large otherwise. TML improves DML when the labeled data are insufficient in the target domain by utilizing information from related source domains, which have better distance estimations between samples. For example, it may be hard to distinguish “zebra” from “tiger” by observing only a few labeled samples due to the very similar stripe texture. But this task can be much easier if we have enough labeled samples to well distinguish “horse” from “cat”. The sample images are from the NUS-WIDE [13] dataset.

constraints) so that the learned metric  $A^*$  can well separate semantically similar data from dissimilar ones, such as “zebra” and “tiger” shown in Fig. 2. While in real-world applications, the learned target metric  $A_M$  may be not satisfactory since the labeled data are insufficient in the target domain. For example, it may be hard to distinguish “zebra” from “tiger” given only a few labeled samples since the two types of animals have very similar stripe texture. To mitigate the label deficiency issue in the target metric learning, we may utilize the information from other related source domain, where the distance metric  $A_S^*$  is good enough or a good metric can be learned using large amounts of labeled data. For example, if we have enough labeled samples to well distinguish “horse” from “cat”, then it may be very easy for us to recognize “zebra” and “tiger” by observing only a few labeled samples. The source metric cannot be directly used in the target domain due to the different data distributions [14] or representations [22] between the source and target domains. Therefore, (homogeneous or heterogeneous) transfer metric learning (TML) is developed to improve the target metric by transferring knowledge (particularly, the metric information) from the source domain. A summarization and discussion of the various TML methods is given as follows.

### 2.3 A categorization of transfer metric learning techniques

As shown in Fig. 3, we can classify TML into different categories according to various principals. Firstly, TML can be generally grouped as *homogeneous TML* and *heterogeneous TML* according to the feature setting. In the former group, the samples of different domains lie in the same feature space ( $\mathcal{X}_1 = \mathcal{X}_2 = \dots = \mathcal{X}_M$ ), and only the data distributions vary ( $P_1(X_1) \neq P_2(X_2) \neq \dots \neq P_M(X_M)$ ).

Whereas in heterogeneous TML, the feature spaces are different ( $\mathcal{X}_1 \neq \mathcal{X}_2 \neq \dots \neq \mathcal{X}_M$ ) and there may be semantic gap between the source and target domains. For example, in the problem of image matching, we may have only a few labeled images in a new scenario due to the high labeling cost, but there are large amounts of labeled images in some other scenarios. The data distributions of different scenarios vary since there are different backgrounds, illuminations, etc. Besides, the web images are usually associated with text descriptions, and it is useful to utilize the semantic textual features to help learn a better distance metric for visual features [22]. The data representations are quite different for the textual and visual domains.

We can also categorize the different TML approaches as *inductive TML*, *transductive TML*, and *unsupervised TML* according to whether the label information is available in the source or target domains. The relationship of the three learning settings are summarized in Table 1. This is similar to the categorization of transfer learning presented in [7].

Furthermore, we summarize the TML approaches into four different cases according to the utilized transfer strategies. Some early works of TML directly enforce the target metric to be close the source metric, and we thus refer it to as *TML via metric approximation*. Since the main difference between the source and target domains in homogeneous TML is the distribution divergence, some approaches enable metric transfer by minimizing the distribution difference. We refer this case to as *TML via distribution approximation*. There is a large amount of TML approaches that enable knowledge transfer by finding a common subspace for the source and target domains, especially in heterogeneous TML. This context is referred to as *TML via subspace approximation*. Finally, there is a few works that let the distance functions of different domains share some common parts

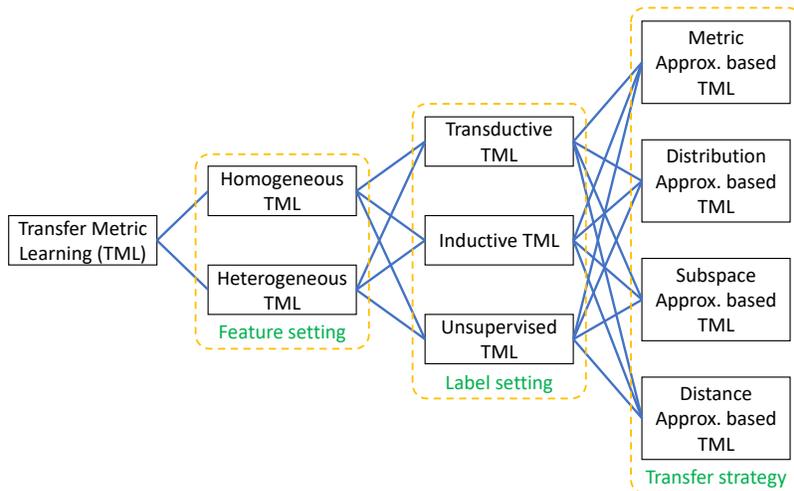


Fig. 3. A categorization of the TML approaches according to different principals. TML can be categorized according to the feature setting, label setting or utilized transfer strategy. Terms on each path from the left to the right make up a certain TML category, e.g., “distribution approximation based transductive homogeneous TML”.

TABLE 1  
Different categorizes of TML according to label setting.

TML categorizes	Source domain label information	Target domain label information
Inductive TML	Available or Unavailable	Available
Transductive TML	Available	Unavailable
Unsupervised TML	Unavailable	Unavailable

TABLE 2  
Different approaches to TML.

TML approaches	Brief description
Metric approximation	Use the target metric to approximate the source metric [14], [15], [16], [17], [18], [26].
Distribution approximation	Conduct metric transfer by minimizing the data distributions of different domains [8], [19], [20], [27], [28].
Subspace approximation	Conduct metric transfer by finding a common subspace for different domains [12], [22], [23], [24], [25], [29].
Distance approximation	Share common parts between distance functions or enforce agreement between distances of corresponding sample pairs in different domains [10], [11], [30].

or enforce the distances of corresponding sample pairs to agree with each other in different domains, and we refer it to as *TML via distance approximation*. The former two cases are usually used in homogeneous TML, and the latter two cases can be adopted for heterogenous TML. Table 2 is a brief description of these cases.

In Table 3, we show which strategies are currently employed for different settings. In homogeneous TML, most of the current algorithms are inductive, and the transductive ones are usually conducted via distribution approximation. There is still no unsupervised method and a possible solution is to extend some unsupervised DML (e.g., [31]) or transfer learning (e.g., [32]) algorithms for unsupervised TML. One challenge is how to ensure the metric learned in the source domain is better since there are no labeled data in both the source and target domains. In the heterogeneous setting [33], since feature dimensions of different domains do not have correspondences, it is inappropriate to conduct TML via direct metric approximation. Most of the current heterogeneous TML approaches first find a common sub-

space for different domains, and then conduct knowledge transfer in the subspace. Unsupervised heterogeneous TML can be easily extended for the transductive heterogeneous setting by further utilizing source labels, and it is possible to adopt the distribution approximation strategy in the heterogeneous setting by first finding a common representation for the different domains.

### 3 HOMOGENEOUS TRANSFER METRIC LEARNING

In homogeneous TML, the utilized features (data representations) are the same, but the data distributions vary for different domains. For example, in sentiment classification as shown in Fig. 4, we would like to determine the sentiment polarity (positive, negative or neutral) for a review of electronics. The performance of a sentiment classifier depends much on the distance estimation between reviews. To obtain reliable distance estimation, we usually need large amounts of labeled reviews to learn a good distance metric. However, we may only have a few labeled electronics reviews due to the high labeling cost and thus the obtained metric is

TABLE 3  
Different transfer strategies used in different TML settings.

TML strategies	Homogeneous TML			Heterogeneous TML		
	Inductive	Transductive	Unsupervised	Inductive	Transductive	Unsupervised
Metric	✓			×	×	×
Distribution		✓				
Subspace	✓			✓		✓
Distance	✓			✓		✓

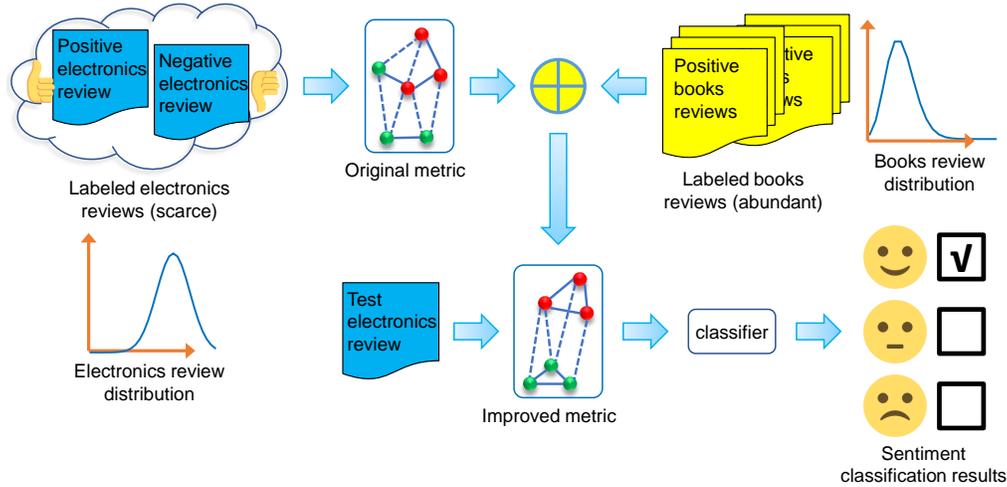


Fig. 4. An example of homogeneous transfer metric learning. In sentiment classification, distance metric learned for target (such as electronics) reviews may be not satisfactory due to the insufficient labeled data. Homogeneous TML improves the metric by using abundant labeled source (such as book) reviews, where the data distribution is different from the target reviews.

not satisfactory. Fortunately, we may have abundant labeled book reviews, which are often easier to collect. Directly applying the metric learned using the labeled book reviews to the sentiment classification of electronics reviews is not appropriate due to the distribution difference between the electronics and book reviews. Transfer metric learning is able to deal with this issue and learn improved distance metric for the target sentiment classification of electronics reviews by using labeled book reviews.

### 3.1 Inductive TML

Under the inductive setting, we are provided with a few labeled data in the target domain. The number of labeled data in the source domain is large enough so that a good distance metric can be obtained, i.e.,  $N_S \gg N_M > 0$ . In inductive transfer learning [7], there may be no labeled source data ( $N_S = 0$ ), but we have not seen such works in homogeneous TML.

#### 3.1.1 TML via metric approximation

An intuitive idea for homogeneous TML is to first use the source domain data  $\{\mathcal{D}_m\}$  to learn the source distance metrics  $\{\phi_m\}$  beforehand, and then enforce the target metric to be close to the pre-trained source metrics. Therefore, the general formulation for learning the target metric  $\phi_M$  is given by

$$\arg \min_{\phi_M} \epsilon(\phi_M) = L(\phi_M; \mathcal{D}_M) + \gamma R(\phi_M; \phi_1, \dots, \phi_{M-1}), \quad (1)$$

where  $L(\phi_M; \mathcal{D}_M)$  is the empirical loss w.r.t. the metric,  $R(\phi_M; \phi_1, \dots, \phi_{M-1})$  is a regularization term that exploits the relationship between the source and target metrics, and  $\gamma \geq 0$  is a trade-off hyper-parameter. Any loss function used in standard DML can be adopted, and the key is how to design an appropriate regularization term. In [14], two different regularization terms are developed. The first one is to minimize the *LogDet* divergence [34] between the source and target Mahalanobias metrics, i.e.,

$$\begin{aligned} R(A_M; A_1, \dots, A_{M-1}) &= \sum_{m=1}^{M-1} \alpha_m D_{LD}(A_M, A_m) \\ &= \sum_{m=1}^{M-1} \alpha_m (\text{tr}(A_m^{-1} A_M) - \log \det(A_M)). \end{aligned} \quad (2)$$

Here,  $\{A_m \succeq 0\}_{m=1}^M$  are constrained to be PSD matrices and  $D_{LD}(\cdot, \cdot)$  indicates the *LogDet* divergence of two matrices. This is more appropriate than the Frobenius norm of matrix difference due to the desirable properties of the *LogDet* divergence, such as scale invariance [34]. The coefficients  $\{\alpha_m\}$  that satisfy  $\alpha_m \geq 0$  and  $\sum_{m=1}^{M-1} \alpha_m = 1$  is learned to reflect the contributions of different source metrics to the target metric. Secondly, to exploit the geometric structure of data distribution, Zha et al. [14] propose a regularization term based on manifold regularization [35]:

$$R(A_M; A_1, \dots, A_{M-1}) = \sum_{m=1}^{M-1} \alpha_m \text{tr} \left( X^U L_m (X^U)^T A_M \right), \quad (3)$$

where  $X^U$  is the feature matrix of unlabeled data, and  $L_m$  is the Laplacian matrix of the data adjacency graph calculated based on the metric  $A_m$ . In [15], the importance of the source metrics to the target metric is exploited by learning a task covariance matrix over the metrics. The matrix can model the correlations between different tasks. This approach allows negative and zero transfer.

Both of the above two approaches incorporate the source metrics into a regularization term to penalize the target metric learning. Different from them, a novel decomposition-based TML method is proposed in [16], which constructs the target metric by using the base metrics derived from the source metrics, that is,

$$A_M = U_M \text{diag}(\theta) U_M^T = \sum_{r=1}^{N_B} \theta_{Mr} \mathbf{u}_{Mr} \mathbf{u}_{Mr}^T = \sum_{r=1}^{N_B} \theta_{Mr} B_{Mr}, \quad (4)$$

where  $\{\mathbf{u}_{Mr}\}$  are eigenvectors of source metrics (which are PSD matrices),  $\{\theta_{Mr} \geq 0\}$  are combination coefficients of different base metrics, and  $N_B$  is the number of bases. This transforms the metric learning into coefficient learning. Hence, the number of parameters to be learned is reduced significantly, and the performance can be improved since the labeled samples in the target domain is scarce. Another advantage of the model is that the PSD constraint of the target metric can be automatically satisfied, and thus the computational cost is low. A semi-supervised extension was presented in [26] by combining it with manifold regularization.

In addition to utilizing the source metrics to help the target metric learning, there exist some multi-task metric learning (MTML) approaches that enable different metrics to help each other in metric learning. A representative work is the large margin multi-task metric learning (mtLMNN) [17], which is a multi-task extension of a well-known DML algorithm, i.e., large margin nearest neighbor (LMNN) [2]. In mtLMNN, all the different metrics are learned simultaneously by assuming that each metric consists of a common metric  $A_0$  and task-specific metric  $\hat{A}_m$ , i.e.,  $A_m = A_0 + \hat{A}_m$ . Based on the same idea, a semi-supervised MTML method is developed in [36], where the unlabeled data is utilized by designing a loss to preserve neighborhood relationship. Then a regularization term is designed to control the amount of information to be shared among all tasks. In [15], a MTML approach is presented by first vectorizing the Mahalanobias metrics and then using a task covariance matrix to exploit the task relationship. Similarly, the metrics are vectorized in [5], but the different metrics are enforced to be close under the graph-based regularization theme [37]. In addition, a general MTML framework is proposed in [18], which enables knowledge transfer by enforcing different metrics  $\{A_m\}$  to be close to a common metric  $A_0$ . The general Bregman matrix divergence [38] is introduced to measure the difference between two metrics. The framework incorporates mtLMNN as a special case and the geometry is preserved in the transfer by adopting a special Bregman divergence, i.e., the von Neumann divergence [38].

### 3.1.2 TML via subspace approximation

Most of the TML approaches via direct metric approximation have a main drawback, i.e., when the feature dimension

is high, the model is prone to overfitting due to the large number of parameters to be learned. This also leads to high computational cost in both training and prediction. To tackle this issue, some low-rank TML methods are proposed. They usually decompose the metric as  $A_m = U_m U_m^T$ , where  $U_m \in \mathbb{R}^{d_m \times r}$  is a low-rank transformation matrix. This leads to a common subspace for different domains, and the knowledge transfer is conducted in the subspace. For example, a low-rank multi-task metric learning framework is proposed in [29], [39], which assumes that each transformation is a product of a common transformation and task-specific one, i.e.,  $U_m = \hat{U}_m U_0$ . As a special case, the large margin component analysis (LMCA) [40] is extended to multi-task LMCA (mtLMCA), which is shown to be superior to mtLMNN.

### 3.1.3 TML via distance approximation

Both the models of mtLMNN and mtLMCA are trained based on labeled sample triplets. Different from them, CP-mtML [10] learn the metrics using labeled pairs, which are often easier to collect. Similar to mtLMCA, CP-mtML decomposes the metric as  $A_m = U_m U_m^T$ , but the different projections  $\{U_m\}$  are coupled by assuming that the distance function consists of a common part and task-specific one, i.e.,

$$d_{U_m}^2(\mathbf{x}_i, \mathbf{x}_j) = d_{U_0}^2(\mathbf{x}_i, \mathbf{x}_j) + d_{U_m}^2(\mathbf{x}_i, \mathbf{x}_j). \quad (5)$$

A main advantage of CP-mtML is that the optimization problem can be solved efficiently using stochastic gradient descent (SGD), and hence the model is scalable for high-dimensional features and large amounts of training data. Besides, the learned transformation can be used to derive low-dimensional features, which are desirable in large-scale information retrieval.

## 3.2 Transductive TML

Under the transductive setting, there are no labeled data in the target domain and we only have large amounts of labeled source data, i.e.,  $N_S \gg N_M = 0$ .

### 3.2.1 TML via distribution approximation

In homogeneous TML, the data distributions vary for different domains. Therefore, we can minimize the distribution difference between the source and target domains, so that the source domain samples can be reused in the target metric learning. In [19], a domain adaptation metric learning (DAML) approach is proposed. In DAML, the distance metric is parameterized by a feature mapping  $\phi_M$ . The mapping is learned by first transforming the samples in the source and target domains using the mapping, and then minimizing the distribution difference of the source and target domains in the transformed space. At the same time,  $\phi_M$  is learned to make the transformed samples satisfy the similar/dissimilar constraints in the source domain. The general formulation for learning  $\phi_M$  is given by

$$\arg \min_{\phi_M} \epsilon(\phi_M) = L(\phi_M; \mathcal{D}_S) + \gamma D_{PD}(P_M(X_M), P_S(X_S)), \quad (6)$$

where  $D_{PD}(\cdot, \cdot)$  is a measure of the difference between two probability distributions. Maximum mean discrepancy

(MMD) [41] is adopted as the measure in DAML. The nonlinear mapping  $\phi_M$  is learned in the reproducing kernel Hilbert space (RKHS), and the solution is found using the kernel method. Since the source and target samples in the transformed space follow similar distribution, the mapping learned using the source label information is also discriminative in the target domain. The same idea is adopted in deep TML (DTML) [8], and the main difference is that the nonlinear mapping is assumed to be a multi-layer neural network. The knowledge transfer is conducted at the output layer and each hidden layer, and some weight hyperparameters are set to balance the importance of the losses in different layers. A major limitation of these works is that they only consider the marginal distribution difference. This limitation is overcome in [27], where a novel TML method is developed by simultaneously reducing the marginal and conditional distribution divergences between the source and target domains. The conditional distribution divergence is reduced by first assigning pseudo labels to target domain data using the classifiers trained on source domain data, and then applying the class-wise MMD [42].

Different from these methods, which reduce the distribution difference in a new space, the importance sampling [43] is introduced in [20] to handle DML under covariate shift. The formulation is given as follows,

$$\arg \min_{A_M \geq 0} \epsilon(A_M) = \sum_{i,j} w_{ij} l(A_M; \mathbf{x}_{Si}, \mathbf{x}_{Sj}, y_{Sij}), \quad (7)$$

where  $l(\cdot)$  is some pre-defined loss function over a training pair  $(\mathbf{x}_{Si}, \mathbf{x}_{Sj})$  with  $y_{Sij} = \pm 1$  indicating the two samples are similar or not. The weight  $w_{ij} = \frac{P_M(\mathbf{x}_{Si})P_M(\mathbf{x}_{Sj})}{P_S(\mathbf{x}_{Si})P_S(\mathbf{x}_{Sj})}$  indicates the importance of the pair in the source domain for learning the target metric. Intuitively, if the pair of source samples have large probability to be occurred in the target domain, they should contribute highly in the target metric learning. In particular, for the distance (such as the popular Mahalanobias distance) which is induced by a norm, i.e.,  $d(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i - \mathbf{x}_j)$ , we can calculate the weight as  $w_{ij} = \frac{P_M(\delta_{Sij})}{P_S(\delta_{Sij})}$ , where  $\delta_{Sij} = \mathbf{x}_{Si} - \mathbf{x}_{Sj}$ . In [20], the weights and target metric are learned separately and this may lead to error propagation across them. The issue is tackled by [28], where the weights and target metric are learned simultaneously in a unified framework.

### 3.3 Discussion

TML via metric approximation is straightforward in that divergence between the source and target metrics (parameterized by PSD matrices) are directly minimized. A major difference of the various metric approximation based approaches is that the source and target metrics are enforced to be close in different ways, e.g., by adopting different types of divergence. These approaches are often limited in that the training complexity is high due to the PSD constraint and the distance calculation in the inference stage is not efficient for high-dimensional data. Subspace approximation based TML compensates for these shortcomings by reformulating the metric learning as learning a transformation or mapping. The PSD constraint is automatically satisfied and the learned transformation can be used to derive compressed

representation, which would facilitate efficient distance estimation or sample matching, where the hash technique [44] can be involved. This is critical in many applications, such as information retrieval. The main disadvantage of the subspace approximation based methods is that their optimization problems are often non-convex and hence only local optimum can be obtained. The recent work [10] based on distance approximation also learn a projection instead of the metric but the optimization is more efficient. All of these approaches do not explicitly deal with the distribution difference, which is the main issue that transfer learning would like to tackle. Distribution approximation based methods focus on this point by usually minimizing the MMD measure or utilizing the importance sampling strategy.

### 3.4 Related work

TML is quite related to transfer subspace learning (TSL) [45], [46] or transfer feature learning (TFL) [47]. An early work on TSL is presented in [45] that finds a low-dimensional latent space, where the distribution difference between the source and target domain is minimized. This algorithm is conducted in a transductive manner and not convenient to derive a representation for new samples. This issue is tackled by Si et al. [46], where a generic regularization framework is proposed for TSL based on Bregman divergence [48]. A low-rank TSL (LTSL) framework is proposed in [49], [50], where the subspace is found by reconstructing the projected target data using the projected source data under the low-rank representation [51], [52] theme. The main advantage of the framework is that only relevant source data are utilized to find the subspace and noisy information can be filtered out. That is, it can avoid negative transfer. The framework is further extended in [53] to help recover missing modality in the target domain and improved in [54] by exploiting both low-rank and sparse structures on the reconstruction matrix.

TFL is very similar to TSL and a representative method is presented in [47], where the typical MMD is modified to take both the marginal and class-conditional distributions into consideration. More recent works on TFL are built upon the powerful deep feature learning. For example, considering that the features in deep neural networks are usually general in the first layers and task-specific in higher layers, Long et al. [55] propose the deep adaptation networks (DAN), which freezes the general layers in convolutional neural networks (CNN) [56] and only conduct adaption in the task-specific layers. Besides, multi-kernel MMD (MK-MMD) [57] is employed to improve kernel selection in MMD. In DAN, only the marginal distribution difference between the source and target domains is exploited. This is improved by the joint adaptation networks (JAN) [58], which is able to reduce the joint distribution divergence using a proposed joint MMD (JMMD). The JMMD can involve both the input features and output labels in domain adaptation. The constrained deep TSL [59] method can also exploit the joint distribution and the target domain knowledge is incorporated gradually during a progressive transfer procedure.

All of these TSL or TFL approaches have very close relationships to the subspace and distribution approximation

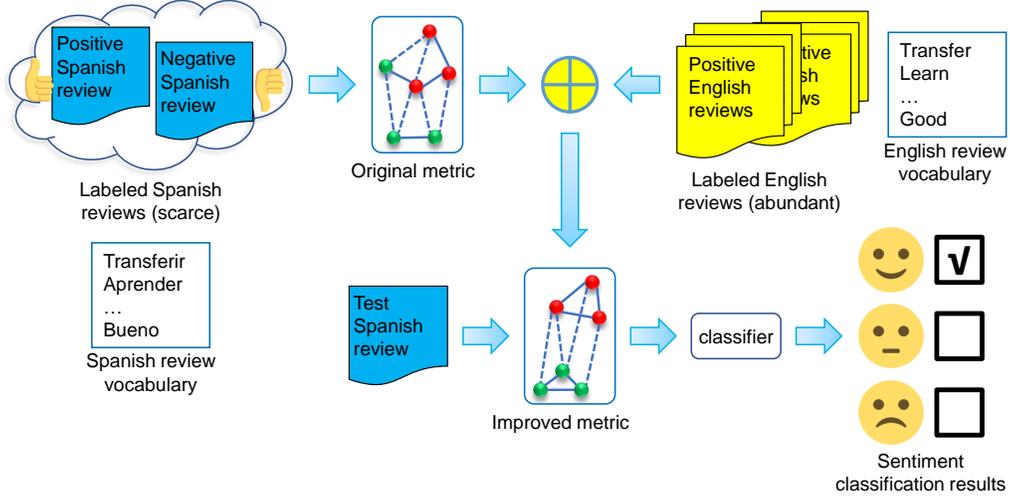


Fig. 5. An example of heterogeneous transfer metric learning. In multi-lingual sentiment classification, distance metric learned for target reviews (such as the ones written in Spanish) may be not satisfactory due to the insufficient labeled data. Heterogeneous TML improves the metric by using abundant labeled source reviews (such as the ones written in English), where the data representation is different from the target reviews (e.g., due to the different vocabularies).

based TML. Although they do not aim to learn metrics, it is not hard to adapt them for TML by adopting some metric learning loss in these models.

#### 4 HETEROGENEOUS TRANSFER METRIC LEARNING

In heterogeneous TML, the different domains have different features (data representations), and sometimes have semantic gap, such as the textual and visual domains. A typical example is the multi-lingual sentiment classification as shown in Fig. 5, where we would like to determine the sentiment polarity for a review written in Spanish. The labeled Spanish reviews may be scarce but it is much easier to collect abundant labeled reviews written in English. Directly applying the metric learned using the labeled English reviews to the sentiment classification of Spanish reviews is infeasible since the representations of Spanish and English reviews are different due to the varied vocabularies. This issue can be tackled by heterogeneous TML, which improves the distance metric for the target sentiment classification of Spanish reviews using labeled English reviews.

##### 4.1 Inductive heterogeneous TML

Different from the inductive homogenous setting, the number of labeled data in the source domain can be zero under the inductive heterogeneous setting. This is because the source feature may have much stronger representation power than that in the target domain, and thus no labeled data are required to obtain a good distance function in the source domain.

###### 4.1.1 Heterogeneous TML via subspace approximation

To our best knowledge, heterogeneous TML under the inductive setting is only studied in recent years. For example, a heterogeneous multi-task metric learning (HMTML) method is proposed in [23]. HMTML assumes that the

similar/dissimilar constraints are limited in multiple heterogeneous domains, but there are large amounts of unlabeled data that have representations in all domains, i.e.,  $\mathcal{D}^U = \{(\mathbf{x}_{1n}^U, \dots, \mathbf{x}_{Mn}^U)\}_{n=1}^{N^U}$ . To build a connection between different domains, the linear metrics  $\{A_m\}$  are decomposed as  $\{A_m = U_m U_m^T\}$ , and then the different representations of an unlabeled data are transformed into a common subspace using  $\{U_m\}$ . The general formulation is given by

$$\arg \min_{\{A_m \succeq 0\}} \epsilon(\{A_m\}) = \sum_{m=1}^M L(A_m; \mathcal{D}_m) + \gamma R(U_1, \dots, U_M; \mathcal{D}^U). \quad (8)$$

Since the different representations corresponding to the same (unlabeled) sample, the transformed representations should be close to each other in the subspace. By minimizing the divergence of transformed representations (or equivalently maximizing their correlations), each transformation is learned by using the information from all domains. This results in an improved transformation, and thus better metric than learning them separately. In [23], a tensor-based regularization term is designed to exploit the high-order correlations between different domains. A variant of the model is presented in [24], which uses the class labels to build domain connection.

In [25], a general heterogeneous TML approach is proposed based on the knowledge fragments transfer [60] strategy. The optimization problem is given by

$$\arg \min_{\phi_M} \epsilon(\phi_M) = L(\phi_M; \mathcal{D}_M) + \gamma R(\{\phi_{M_c}(\cdot)\}, \{\varphi_{S_c}(\cdot)\}; \mathcal{D}^U), \quad (9)$$

where  $\phi_{M_c}(\cdot)$  is the  $c$ 'th coordinate of the mapping  $\phi_M$ , and  $\varphi_{S_c}(\cdot)$  is the  $c$ 'th fragment of the knowledge in the source domain. The source knowledge fragments are represented by some mapping functions, which are learned by applying existing DML algorithms in the source domain beforehand. Then the target metric (which also consists of multiple mapping functions) is enforced to agree with the source fragments on the unlabeled corresponding data. This helps

learn an improved metric in the target domain since the pre-trained source distance function is assumed to be superior than the target distance function without knowledge transfer. Intuitively, the target subspace is enforced to approach a better source subspace. An improvement of the model is presented in [12], where the locality of the geometric structure of the data distribution is preserved via manifold regularization [35].

#### 4.1.2 Heterogeneous TML via distance approximation

We can not only enforce the subspace representations of corresponding sample in different domains to be close, but also let the distances of corresponding sample pairs to agree with each other in different domains. For example, an online heterogeneous TML approach is proposed in [30], which also assumes that there are abundant unlabeled corresponding data, but the target labeled sample pairs are provided in a sequential manner (one by one). Given a new labeled training pair, the target metric is updated as:

$$\begin{aligned} A_M^{k+1} &= \arg \min_{A_M \succeq 0} \epsilon(A_M) \\ &= L(A_M) + \gamma_A D_{LD}(A_M, A_M^k) + \gamma_I R(d_{A_M}, d_{A_S}; \mathcal{D}^U), \end{aligned} \quad (10)$$

where  $L(A_M)$  is the empirical loss w.r.t. the current labeled pair,  $D_{LD}(\cdot, \cdot)$  is the *LogDet* divergence [34], and  $R(d_{A_M}, d_{A_S}; \mathcal{D}^U)$  is a regularization term that enforces agreements between the source and target distances (of corresponding pairs). Here,  $A_M^k$  is the target metric obtained previously and initialized as an identity matrix. The source metric  $A_S$  can be an identity matrix if the source feature is much more powerful than the target feature. By pre-calculating  $A_S$  and formulating the term  $R(\cdot)$  under the manifold regularization theme [35], an online algorithm is developed to update the target metric  $A_M$  efficiently.

## 4.2 Unsupervised heterogeneous TML

There exist a few unsupervised heterogeneous TML approaches that utilize unlabeled corresponding data for metric transfer and no label information is provided in either the source or target domains ( $N_S = N_M = 0$ ). Under this unsupervised paradigm, the utilized source feature should be more expressive or interpretable than the target feature, so that the estimated distances in the source domain can be better than those in the target domain.

#### 4.2.1 Heterogeneous TML via subspace approximation

An early work is done in [22], where the main idea is to maximize the similarity of any unlabeled corresponding pairs in a common subspace, i.e.,

$$\arg \min_{A_M \succeq 0} \epsilon(A_M) = \sum_{n=1}^{N^U} l(\varphi(\theta)), \quad (11)$$

where  $\varphi(\theta) = \frac{1}{1 + \exp(-\theta)}$  with  $\theta = (\mathbf{x}_{Mn}^U)^T G \mathbf{x}_{Sn}^U$  and  $G = U_M^T U_S$ . Here,  $l(\cdot)$  is chosen to be the negative logistic loss and the proximal gradient method is adopted for optimization. A main disadvantage of this TML approach is that the computational complexity is high since the costly singular value decomposition (SVD) is involved in each iteration of the optimization.

#### 4.2.2 Heterogeneous TML via distance approximation

Instead of directly maximizing the likelihood between unlabeled sample pairs, Dai et al. [11] propose to use the target samples to approximate the source manifold. The method is inspired by locally linear embedding (LLE) [61], and metric transfer is conducted by enforcing embeddings of target samples to preserve local properties in the source domain. The optimization problem is given by

$$\arg \min_{U_M} \epsilon(U_M) = \sum_{i=1}^{N^U} \left\| U_M^T \mathbf{x}_{Mi}^U - \sum_{j=1}^{N^U} w_{Sij} (U_M^T \mathbf{x}_{Mj}^U) \right\|, \quad (12)$$

where  $w_{Sij}$  is the weight in the adjacency graph calculated using the source domain feature. This enables the distances (between samples) in the source and target domains to agree with each other on the manifold. The optimization is much more efficient than [22] since only a generalized eigenvector problem is needed to be solved.

## 4.3 Discussion

It is nature to conduct heterogeneous TML via subspace approximation since the representations of different domains vary and finding a common representation can facilitate the knowledge transfer. Similar to that in the homogeneous setting, the main drawback is that the optimization problem is usually non-convex. Although this drawback can be remedied by directly learning a PSD matrix, such as using the distance approximation strategy, it is nontrivial to perform efficient distance inference for high-dimensional data and extend the algorithm to learn nonlinear metric. Due to the strong ability and rapid development of deep learning, it may be more promising to learn transformation or mapping than PSD matrix in TML, based on either subspace or distance approximation.

## 4.4 Related work

Some early heterogeneous transfer learning approaches are not specially designed for DML, but the learned feature transformation or mapping for each domain can be used to derive a metric. For example, in the work of heterogeneous domain adaptation via manifold alignment (DAMA) [62], the class labels are utilized to align different domains. A mapping function is learned for each domain and all functions are learned together. After being projected into a common subspace, the samples should be close to each other if they belong to the same class and separated otherwise. This is conducted for all samples from either the same domain or different domains. The label information of all different domains can be utilized to learn the shared subspace, and thus better embeddings (representations) can be learned for different domains than learning them separately. In [63], a multi-task discriminant analysis (MTDA) approach is proposed to deal with heterogeneous feature spaces in different domains. MTDA assumes the linear transformation of the  $m$ 'th domain is given by  $U_m = W_m H$ , which consists of a task-specific part  $W_m$  and a common part  $H$  for all tasks. Then all the transformations are learned in a single optimization problem, which is similar to that of the well-known

TABLE 4  
A summarization of the different applications in which TML utilized.

Homogeneous TML	Computer vision	Handwritten letter/digit classification, face recognition/verification, image retrieval.
	Speech recognition	English alphabet recognition, vowel classification.
	Other applications	Social network, customer behavior analysis.
Heterogeneous TML	Computer vision	Face recognition/verification, scene categorization, image clustering/retrieval/super-resolution.
	Text analysis	Multilingual text categorization, sentiment classification, email spam detection.

linear discriminant analysis (LDA) [64]. In [65], a multi-task nonnegative matrix factorization (MTNMF) approach is proposed to learn the different mappings for all domains by simultaneously factorizing their data representation and feature-class correlation matrices. The factorized class representation matrix is assumed to be shared by all tasks. This leads to a common subspace for different domains.

All of these approaches have very close relationships to the subspace approximation based heterogeneous TML, but they mainly utilize the fully-supervised class labels to learn feature mappings for different domains. As we mentioned previously, it is common to utilize the weakly-supervised pair/triplet constraints in DML and it is not hard to adapt these approaches for heterogeneous TML by adopting some metric learning loss w.r.t. pair/triplet constraints in these models.

## 5 APPLICATIONS

In general, for any applications where DML is appropriate, TML is a good candidate when the label information is scarce or hard to collect. In Table 4, we summarize the different applications that TML utilized in.

### 5.1 Homogeneous TML

#### 5.1.1 Computer vision

Similar to DML [66], most of the TML approaches are applied in computer vision. For example, effectiveness of many homogeneous TML methods are verified in the common image classification application, which includes handwritten letter/digit classification [15], [16], [18], [39], face recognition [14], [19], natural scene categorization and object recognition [16], [19], [27].

DML is particular suitable and crucial for some applications, such as face verification [8], person re-identification [5] and image retrieval [10]. This is because in these applications, results can be directly inferred from the distances between samples. Face verification aims to decide whether two face images belong to the same person or not. In [8], TML is applied for face verification across different datasets, where the distributions vary. The goal of person re-identification is to decide whether the people appear in multiple cameras are the same person or not, where the cameras often do not have overlapping views. The data distributions of the images captured by different cameras vary due to the varying illumination, background, etc. Besides, distribution may change over time for the same camera. Hence, TML can be very useful in person re-identification [5], [8], [67]. An efficient

MTML approach is proposed in [10] to make use of auxiliary datasets for face retrieval, where the tasks vary for different datasets. Stochastic gradient descent (SGD) is adopted for optimization and the algorithm is scalable to large amounts of training data and high dimensional features.

#### 5.1.2 Speech recognition

Different groups of speakers have different ways in uttering an English alphabet. In [17], [18], [36], alphabet recognition in each group is regarded as a task, and MTML is employed to learn the metrics of different groups together. Similarly, since men and women have different pronunciation styles, vowel classification is performed for two different groups according to the gender, and MTML is adopted to learn their metrics simultaneously by making use of all available labeled data [18].

#### 5.1.3 Other applications

In [68], MTML is used for predictions in social networks. For example, citation prediction is to predict the referencing between articles given their contents. The citation patterns of different areas (such as computer science and engineering) are different but related, and thus MTML is adopted to learn the prediction models of multiple areas simultaneously. Social circle prediction is to assign a person to appropriate social circles given his/her profile. Different types of social circles (such as family members and colleges) are different but related with each other, and hence MTML is applied to improve the performance. In [17], [18], [29], MTML is applied to customer information prediction in insurance company. There are multiple variables that can be used to predict the interest of a person in buying a certain insurance policy. Each variable is a discrete value and can be predicted using other variables. The predictions of different variables can be conducted together since they are correlated with each other.

### 5.2 Heterogeneous TML

#### 5.2.1 Computer vision

Similar to homogeneous TML, heterogeneous TML is also mainly applied to the computer vision community, such as image classification including face recognition [63], natural scene categorization [12], [24], [25] and object recognition [12], [23], [25], image clustering [11], image retrieval [11], [12], [69], and face verification [12]. In these applications, either the feature dimensions vary or different types of features are extracted for the source and target domains. In particular, expensive features (has strong representation

power but high computational cost, such as CNN [70]) can be used to guide learning an improved metric for relatively cheap features (such as LBP [71]), and interpretable text feature can help the metric learning of visual feature, which is often to interpret [22], [65].

In [11], heterogeneous TML is adopted to improve image super-resolution, which is to generate a high-resolution (HR) image for its low-resolution (LR) counterpart. The method is based on JOR [72], which is an example-based super-resolution approach. JOR needs to find the nearest neighbors for the LR images, and a metric is learned in [11] to replace the Euclidean metric in the  $k$ -NN search by leveraging information from the HR domain.

### 5.2.2 Text analysis

In the text analysis area, heterogenous TML is mainly applied by using labeled documents written in one language (such as English) to help analysis of the documents in another language (such as Spanish). The utilized vocabularies vary for different languages, and thus the data representations are heterogeneous for different domains. Some typical examples including text categorization [23], [62], sentiment classification [65] and document retrieval [62]. In [65], heterogenous MTML is applied to email spam detection since the vocabularies for different persons' email vary.

## 6 CONCLUSION AND DISCUSSION

### 6.1 Summary

In this survey, we provide a comprehensive and structured overview of the transfer metric learning (TML) methods and their applications. We generally group TML as homogeneous and heterogeneous TML according to the feature setting. Similar to [7], the TML approaches can also be classified into inductive, transductive and unsupervised TML according to the label setting. According to the transfer strategy, we further categorize the TML approaches into four contexts, i.e., TML via metric approximation, TML via distribution approximation, TML via subspace approximation and TML via distance approximation.

Homogeneous TML has been studied extensively under the inductive setting and various transfer strategies can be adopted. In the transductive setting, TML is mainly conducted by distribution approximation, and there are still no unsupervised methods for homogeneous TML. Unsupervised TML can be carried out under the heterogeneous setting. This is because if more powerful feature is utilized in the source domain, then the distance estimation can be better than that in the target domain [11]. Since the data representations vary for different domains in heterogeneous TML, most of these approaches find a common subspace for knowledge transfer.

### 6.2 Challenges and future directions

We finally identify some challenges in TML and speculate several possible future directions.

#### 6.2.1 Selective transfer in TML

Current transfer learning and TML algorithms usually assume that the source tasks or domain samples are positively related with the target ones. However, this assumption may not hold in real-world applications [15], [73]. The TML algorithm presented in [15] can leverage negatively correlated task by learning task correlation matrix. In [74], the relations of 26 popular visual learning tasks are learned using a large image dataset, where each image has annotations in all tasks. This leads to a task taxonomy map, which can be used to guide the chosen of appropriate supervision policies in transfer learning. Different from these approaches, which consider selective transfer [75] at the task-level, a heterogeneous transfer learning method based on the attention mechanism is proposed in [73], which can avoid negative transfer at the instance-level. The low-rank TML model presented in [50] can also avoid negative transfer to some extent by filtering noisy information in the source domain.

Task correlations have been exploited for metric approximation based TML [15], and the attention scheme can be used for subspace approximation based TML following [73]. It is still unclear how to conduct selective transfer in distribution and distance approximation based TML. Adopting the attention scheme may be a certain choice, but this scheme cannot make use of the negative transfer. Therefore, a promising future direction may be to conduct selective transfer at the hypothesis space-level so that both the positive and negative transfer can be effectively utilized.

#### 6.2.2 More theoretical understanding of TML

There is a theoretical study in [12], which shows that generalization ability of the target metric can be improved by directly enforcing the source feature mappings to agree with the target mappings. But there is still lack of general analysis scheme (such as [76], [77], [78]) and theoretical results for TML. In particular, more theoretical studies should be conducted to understand when and how could the source domain knowledge help the target metric learning.

#### 6.2.3 TML for handling complex data

Most of current TML approaches only learn linear metrics (such as the Mahalanobis metric). However, there may be nonlinear structure in the data, e.g., most of the visual feature representations. Linear metric may fail to capture such structure and hence it is desirable to learn nonlinear metric for the target domain in TML. There have been several works on nonlinear homogeneous TML based on neural networks [8], [55], [58]. But all of them are mainly designed for continuous real-valued data and learn real-valued metrics. More studies can be conducted for histogram data or learning binary target metrics. The histogram data is popular in visual analytic-based applications, and binary metric is efficient in distance calculation. As far as we are concerned, there is only one nonlinear TML work under the heterogeneous setting [25] (with an extension presented in [12]), where gradient boosting regression tree (GBRT) [79], [80] is adopted to learn nonlinear metric in the target domain. Some other nonlinear learning techniques can be investigated, and also binary metrics can be learned to accelerate prediction. In addition, when the structure of

the target data distribution is very complex, it could be a good choice to learn Riemannian metric [81] or multiple local metrics [82] to approximate the geodesic distance in the target domain.

#### 6.2.4 TML for handling changing and big data

In TML, all the training data in the source and target domains are usually assumed to be provided at once and a fixed target metric is learned. However, in real-world applications, the data are usually comes in a sequential order and the data distribution may change overtime. For example, tremendous amounts of data are uploaded on the web everyday, and for a robot, the environment changes overtime and feedbacks are provided continuously. Therefore, it is desirable to develop some TML algorithms to make the metric adapt to different changes. Some quite related topics including online learning [83], [84] and lifelong learning [85]. There is a recent try in [30], where an online heterogeneous TML is developed. However, this approach needs abundant unlabeled corresponding data in the source and target domains for knowledge transfer. Hence, the approach may be not efficient when vast amounts of unlabeled data are needed to achieve satisfactory accuracy.

Although the number of training data in the target domain is often assumed to be small, the continuously changing data is “big” in a long term. In addition, when the feature dimension is high, computational costs of the distances between vast amounts of samples based on a learned Mahalanobias metric is intolerable. A typical example is information retrieval. Therefore, it is desirable to learn some target metric that is efficient in distance calculation, e.g., learn hamming distance metric [86], [87] or feature hashing [44], [88], [89], [90] in the target domain.

#### 6.2.5 TML for handling extreme cases

One-shot learning [91] and zero-shot learning [92] are two extreme cases of transfer learning. In these cases, the number of labeled data in the target domain is very small (such as only one) and even zero. The main goal is to recognize rare or unseen classes [93], where some additional knowledge (such as descriptions of the relations between existing and unseen classes) may be provided. This is more like human learning, and much useful in practice. They are quite related to the concepts of domain generalization [94], [95], [96].

DML has been found to be useful in learning unknown classifiers [97] (with an extension in [98]), but it does not aim to learn a metric in the target domain. In [99], an unbiased metric is learned across different domains, but no specific information about the target domain is leveraged. Although some existing TML algorithms allow no labeled data in the target domain [8], [11], they need large amounts of unlabeled target data, which can be regarded as additional knowledge. If we do not have unlabeled data, is it possible to utilize other semantic information to help the target metric learning? There exists a try in [100], where the ColorChecker Chart is utilized as additional information for person re-identification under the one-shot setting. But such information is not easy to access and not general for different applications. Hence, more common and easily

accessible knowledge should be identified and explored for general TML under the one/zero-shot setting.

## REFERENCES

- [1] E. P. Xing, M. I. Jordan, S. Russell, and A. Ng, “Distance metric learning with application to clustering with side-information,” in *Advances in Neural Information Processing Systems*, 2002, pp. 505–512.
- [2] K. Q. Weinberger, J. Blitzer, and L. K. Saul, “Distance metric learning for large margin nearest neighbor classification,” in *Advances in Neural Information Processing Systems*, 2005, pp. 1473–1480.
- [3] P. Jain, B. Kulis, I. S. Dhillon, and K. Grauman, “Online metric learning and fast similarity search,” in *Advances in Neural Information Processing Systems*, 2008, pp. 761–768.
- [4] S. Chopra, R. Hadsell, and Y. LeCun, “Learning a similarity metric discriminatively, with application to face verification,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 539–546.
- [5] L. Ma, X. Yang, and D. Tao, “Person re-identification over camera networks using multi-task distance metric learning,” *IEEE Transactions on Image Processing*, vol. 23, no. 8, pp. 3656–3670, 2014.
- [6] A. Bellet, A. Habrard, and M. Sebban, “A survey on metric learning for feature vectors and structured data,” *arXiv preprint arXiv:1306.6709v4*, 2014.
- [7] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [8] J. Hu, J. Lu, and Y.-P. Tan, “Deep transfer metric learning,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 325–333.
- [9] C. Su, F. Yang, S. Zhang, Q. Tian, L. S. Davis, and W. Gao, “Multi-task learning with low rank attribute embedding for multi-camera person re-identification,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1167–1181, 2018.
- [10] B. Bhattarai, G. Sharma, and F. Jurie, “Cp-mtml: Coupled projection multi-task metric learning for large scale face retrieval,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4226–4235.
- [11] D. Dai, T. Kroeger, R. Timofte, and L. Van Gool, “Metric imitation by manifold transfer for efficient vision applications,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3527–3536.
- [12] Y. Luo, Y. Wen, T. Liu, and D. Tao, “Transferring knowledge fragments for learning distance metric from a heterogeneous domain,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [13] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, “NUS-WIDE: a real-world web image database from National University of Singapore,” in *ACM international Conference on Image and Video Retrieval*, 2009, pp. 48:1–48:9.
- [14] Z.-J. Zha, T. Mei, M. Wang, Z. Wang, and X.-S. Hua, “Robust distance metric learning with auxiliary knowledge,” in *International Joint Conference on Artificial Intelligence*, 2009, pp. 1327–1332.
- [15] Y. Zhang and D.-Y. Yeung, “Transfer metric learning by learning task relationships,” in *ACM SIGKDD international conference on Knowledge Discovery and Data mining*, 2010, pp. 1199–1208.
- [16] Y. Luo, T. Liu, D. Tao, and C. Xu, “Decomposition-based transfer distance metric learning for image classification,” *IEEE Transactions on Image Processing*, vol. 23, no. 9, pp. 3789–3801, 2014.
- [17] S. Parameswaran and K. Q. Weinberger, “Large margin multi-task metric learning,” in *Advances in Neural Information Processing Systems*, 2010, pp. 1867–1875.
- [18] P. Yang, K. Huang, and C.-L. Liu, “Geometry preserving multi-task metric learning,” *Machine Learning*, vol. 92, no. 1, pp. 133–175, 2013.
- [19] B. Geng, D. Tao, and C. Xu, “DAML: Domain adaptation metric learning,” *IEEE Transactions on Image Processing*, vol. 20, no. 10, pp. 2980–2989, 2011.
- [20] B. Cao, X. Ni, J.-T. Sun, G. Wang, and Q. Yang, “Distance metric learning under covariate shift,” in *International Joint Conference on Artificial Intelligence*, 2011, pp. 1204–1210.
- [21] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

- [22] G.-J. Qi, C. C. Aggarwal, and T. S. Huang, "Transfer learning of distance metrics by cross-domain metric sampling across heterogeneous spaces," in *SIAM International Conference on Data Mining*, 2012, pp. 528–539.
- [23] Y. Luo, Y. Wen, and D. Tao, "On combining side information and unlabeled data for heterogeneous multi-task metric learning," in *International Joint Conference on Artificial Intelligence*, 2016, pp. 1809–1815.
- [24] —, "Heterogeneous multitask metric learning across multiple domains," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 9, pp. 4051–4064, 2018.
- [25] Y. Luo, Y. Wen, T. Liu, and D. Tao, "General heterogeneous transfer distance metric learning via knowledge fragments transfer," in *International Joint Conference on Artificial Intelligence*, 2017, pp. 2450–2456.
- [26] H. Shi, Y. Luo, C. Xu, Y. Wen, and C. M. I. Center, "Manifold regularized transfer distance metric learning," in *BMVC*, 2015, pp. 158.1–158.11.
- [27] Z. Ding and Y. Fu, "Robust transfer metric learning for image classification," *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 660–670, 2017.
- [28] Y. Xu, S. J. Pan, H. Xiong, Q. Wu, R. Luo, H. Min, and H. Song, "A unified framework for metric transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 6, pp. 1158–1171, 2017.
- [29] P. Yang, K. Huang, and C.-L. Liu, "A multi-task framework for metric learning with common subspace," *Neural Computing and Applications*, vol. 22, no. 7-8, pp. 1337–1347, 2013.
- [30] Y. Luo, T. Liu, Y. Wen, and D. Tao, "Online heterogeneous transfer metric learning," in *International Joint Conference on Artificial Intelligence*, 2018, pp. 2525–2531.
- [31] R. G. Cinbis, J. Verbeek, and C. Schmid, "Unsupervised metric learning for face identification in tv video," in *International Conference on Computer Vision*, 2011, pp. 1559–1566.
- [32] H. Chang, J. Han, C. Zhong, A. M. Snijders, and J.-H. Mao, "Unsupervised transfer learning via multi-scale convolutional sparse coding for biomedical applications," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1182–1194, 2018.
- [33] C. Peng, X. Gao, N. Wang, and J. Li, "Graphical representation for heterogeneous face recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 2, pp. 301–312, 2017.
- [34] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *International Conference on Machine Learning*, 2007, pp. 209–216.
- [35] M. Belkin, P. Niyogi, and V. Sindhwani, "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples," *Journal of Machine Learning Research*, vol. 7, no. 11, pp. 2399–2434, 2006.
- [36] Y. Li and D. Tao, "Online semi-supervised multi-task distance metric learning," in *International Conference on Data Mining Workshops*, 2016, pp. 474–479.
- [37] T. Evgeniou, C. A. Micchelli, and M. Pontil, "Learning multiple tasks with kernel methods," *Journal of Machine Learning Research*, vol. 6, pp. 615–637, 2005.
- [38] I. S. Dhillon and J. A. Tropp, "Matrix nearness problems with bregman divergences," *SIAM Journal on Matrix Analysis and Applications*, vol. 29, no. 4, pp. 1120–1146, 2008.
- [39] P. Yang, K. Huang, and C.-L. Liu, "Multi-task low-rank metric learning based on common subspace," in *International Conference on Neural Information Processing*, 2011, pp. 151–159.
- [40] L. Torresani and K.-c. Lee, "Large margin component analysis," in *Advances in Neural Information Processing Systems*, 2007, pp. 1385–1392.
- [41] K. M. Borgwardt, A. Gretton, M. J. Rasch, H.-P. Kriegel, B. Schölkopf, and A. J. Smola, "Integrating structured biological data by kernel maximum mean discrepancy," *Bioinformatics*, vol. 22, no. 14, pp. e49–e57, 2006.
- [42] M. Long, J. Wang, G. Ding, S. J. Pan, and S. Y. Philip, "Adaptation regularization: A general framework for transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 5, pp. 1076–1089, 2014.
- [43] H. Shimodaira, "Improving predictive inference under covariate shift by weighting the log-likelihood function," *Journal of Statistical Planning and Inference*, vol. 90, no. 2, pp. 227–244, 2000.
- [44] J. Wang, T. Zhang, N. Sebe, H. T. Shen *et al.*, "A survey on learning to hash," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 769–790, 2018.
- [45] S. J. Pan, J. T. Kwok, and Q. Yang, "Transfer learning via dimensionality reduction," in *AAAI Conference on Artificial Intelligence*, 2008, pp. 677–682.
- [46] S. Si, D. Tao, and B. Geng, "Bregman divergence-based regularization for transfer subspace learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 7, pp. 929–942, 2010.
- [47] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu, "Transfer feature learning with joint distribution adaptation," in *International Conference on Computer Vision*, 2013, pp. 2200–2207.
- [48] B. A. Frigiyik, S. Srivastava, and M. R. Gupta, "Functional bregman divergence and bayesian estimation of distributions," *IEEE Transactions on Information Theory*, vol. 54, no. 11, pp. 5130–5139, 2008.
- [49] M. Shao, C. Castillo, Z. Gu, and Y. Fu, "Low-rank transfer subspace learning," in *International Conference on Data Mining*, 2012, pp. 1104–1109.
- [50] M. Shao, D. Kit, and Y. Fu, "Generalized transfer subspace learning through low-rank constraint," *International Journal of Computer Vision*, vol. 109, no. 1-2, pp. 74–93, 2014.
- [51] G. Liu, Z. Lin, and Y. Yu, "Robust subspace segmentation by low-rank representation," in *International Conference on Machine Learning*, 2010, pp. 663–670.
- [52] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust recovery of subspace structures by low-rank representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 171–184, 2013.
- [53] Z. Ding, M. Shao, and Y. Fu, "Latent low-rank transfer subspace learning for missing modality recognition," in *AAAI Conference on Artificial Intelligence*, 2014, pp. 1192–1198.
- [54] Y. Xu, X. Fang, J. Wu, X. Li, and D. Zhang, "Discriminative transfer subspace learning via low-rank and sparse representation," *IEEE Transactions on Image Processing*, vol. 25, no. 2, pp. 850–863, 2016.
- [55] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," in *International Conference on Machine Learning*, 2015, pp. 97–105.
- [56] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [57] A. Gretton, D. Sejdinovic, H. Strathmann, S. Balakrishnan, M. Pontil, K. Fukumizu, and B. K. Sriperumbudur, "Optimal kernel choice for large-scale two-sample tests," in *Advances in Neural Information Processing Systems*, 2012, pp. 1205–1213.
- [58] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *International Conference on Machine Learning*, 2017, pp. 2208–2217.
- [59] Y. Wu and Q. Ji, "Constrained deep transfer feature learning and its applications," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5101–5109.
- [60] V. Vapnik and R. Izmailov, "Learning using privileged information: Similarity control and knowledge transfer," *Journal of Machine Learning Research*, vol. 16, pp. 2023–2049, 2015.
- [61] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *science*, vol. 290, no. 5500, pp. 2323–2326, 2000.
- [62] C. Wang and S. Mahadevan, "Heterogeneous domain adaptation using manifold alignment," in *International Joint Conference on Artificial Intelligence*, 2011, pp. 1541–1546.
- [63] Y. Zhang and D.-Y. Yeung, "Multi-task learning in heterogeneous feature spaces," in *AAAI Conference on Artificial Intelligence*, 2011, pp. 574–579.
- [64] K. Fukunaga, *Introduction to statistical pattern recognition*. Academic Press, 1990.
- [65] X. Jin, F. Zhuang, S. J. Pan, C. Du, P. Luo, and Q. He, "Heterogeneous multi-task semantic feature learning for classification," in *International Conference on Information and Knowledge Management*, 2015, pp. 1847–1850.
- [66] B. Kulis, "Metric learning: A survey," *Foundations and Trends in Machine Learning*, vol. 5, no. 4, pp. 287–364, 2012.
- [67] W. Li, R. Zhao, and X. Wang, "Human reidentification with transferred metric learning," in *Asian Conference on Computer Vision*, 2012, pp. 31–44.

- [68] C. Fang and D. N. Rockmore, "Multi-task metric learning on network data," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 2015, pp. 317–329.
- [69] Q. Fu, Y. Luo, Y. Wen, D. Tao, Y. Li, and L. Duan, "Towards intelligent product retrieval for tv-to-online (t2o) application: A transfer metric learning approach," *IEEE Transactions on Multimedia*, 2018.
- [70] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," *arXiv preprint arXiv:1405.3531*, 2014.
- [71] T. Ojala, M. Pietikainen, and T. Maenpaa, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, 2002.
- [72] D. Dai, R. Timofte, and L. Van Gool, "Jointly optimized regressors for image super-resolution," in *Computer Graphics Forum*, vol. 34, no. 2, 2015, pp. 95–104.
- [73] S. Moon and J. G. Carbonell, "Completely heterogeneous transfer learning with attention - what and what not to transfer," in *International Joint Conference on Artificial Intelligence*, 2017, pp. 2508–2514.
- [74] A. R. Zamir, A. Sax, W. Shen, L. Guibas, J. Malik, and S. Savarese, "Taskonomy: Disentangling task transfer learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3712–3722.
- [75] W.-S. Chu, F. De la Torre, and J. F. Cohn, "Selective transfer machine for personalized facial expression analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 3, pp. 529–545, 2017.
- [76] M. Gong, K. Zhang, T. Liu, D. Tao, C. Glymour, and B. Schölkopf, "Domain adaptation with conditional transferable components," in *International Conference on Machine Learning*, 2016, pp. 2839–2848.
- [77] T. Liu, D. Tao, M. Song, and S. J. Maybank, "Algorithm-dependent generalization bounds for multi-task learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 2, pp. 227–241, 2017.
- [78] T. Liu, Q. Yang, and D. Tao, "Understanding how feature structure transfers in transfer learning," in *International Joint Conference on Artificial Intelligence*, 2017, pp. 2365–2371.
- [79] J. H. Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of Statistics*, pp. 1189–1232, 2001.
- [80] D. Kedem, S. Tyree, F. Sha, G. R. Lanckriet, and K. Q. Weinberger, "Non-linear metric learning," in *Advances in Neural Information Processing Systems*, 2012, pp. 2573–2581.
- [81] Z. Huang, R. Wang, L. Van Gool, X. Chen *et al.*, "Cross Euclidean-to-Riemannian metric learning with application to face recognition from video," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2827–2840, 2018.
- [82] Y.-K. Noh, B.-T. Zhang, and D. D. Lee, "Generative local metric learning for nearest neighbor classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 1, pp. 106–118, 2018.
- [83] T. Anderson, *The theory and practice of online learning*. Athabasca University Press, 2008.
- [84] B. Wang, G. Wang, K. L. Chan, and L. Wang, "Tracklet association by online target-specific metric learning and coherent dynamics estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 3, pp. 589–602, 2017.
- [85] S. Thrun and L. Pratt, *Learning to learn*. Springer Science & Business Media, 1998.
- [86] M. Norouzi, D. J. Fleet, and R. R. Salakhutdinov, "Hamming distance metric learning," in *Advances in Neural Information Processing Systems*, 2012, pp. 1061–1069.
- [87] Z. Wang, L.-Y. Duan, J. Lin, X. Wang, T. Huang, and W. Gao, "Hamming compatible quantization for hashing," in *International Joint Conference on Artificial Intelligence*, 2015, pp. 2298–2304.
- [88] Z. Wang, L.-Y. Duan, T. Huang, and W. Gao, "Affinity preserving quantization for hashing: A vector quantization approach to learning compact binary codes," in *AAAI Conference on Artificial Intelligence*, 2016, pp. 1102–1108.
- [89] L.-Y. Duan, J. Lin, Z. Wang, T. Huang, and W. Gao, "Weighted component hashing of binary aggregated descriptors for fast visual search," *IEEE Transactions on multimedia*, vol. 17, no. 6, pp. 828–842, 2015.
- [90] L.-Y. Duan, Y. Wu, Y. Huang, Z. Wang, J. Yuan, and W. Gao, "Minimizing reconstruction bias hashing via joint projection learning and quantization," *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 3127–3141, 2018.
- [91] L. Fei-Fei, R. Fergus, and P. Perona, "One-shot learning of object categories," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 4, pp. 594–611, 2006.
- [92] M. Palatucci, D. Pomerleau, G. E. Hinton, and T. M. Mitchell, "Zero-shot learning with semantic output codes," in *Advances in Neural Information Processing Systems*, 2009, pp. 1410–1418.
- [93] G.-J. Qi, W. Liu, C. Aggarwal, and T. Huang, "Joint intermodal and intramodal label transfers for extremely rare or unseen classes," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 7, pp. 1360–1373, 2017.
- [94] K. Muandet, D. Balduzzi, and B. Schölkopf, "Domain generalization via invariant feature representation," in *International Conference on Machine Learning*, 2013, pp. 10–18.
- [95] M. Ghifary, D. Balduzzi, W. B. Kleijn, and M. Zhang, "Scatter component analysis: A unified framework for domain adaptation and domain generalization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 7, pp. 1414–1430, 2017.
- [96] W. Li, Z. Xu, D. Xu, D. Dai, and L. Van Gool, "Domain generalization and adaptation using low rank exemplar SVMs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1114–1127, 2018.
- [97] T. Mensink, J. Verbeek, F. Perronnin, and G. Csurka, "Metric learning for large scale image classification: Generalizing to new classes at near-zero cost," in *European Conference on Computer Vision*, 2012, pp. 488–501.
- [98] —, "Distance-based image classification: Generalizing to new classes at near-zero cost," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 11, pp. 2624–2637, 2013.
- [99] C. Fang, Y. Xu, and D. N. Rockmore, "Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias," in *International Conference on Computer Vision*, 2013, pp. 1657–1664.
- [100] S. Bak and P. Carr, "One-shot metric learning for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.