# Unsupervised Deep Metric Learning via Orthogonality based Probabilistic Loss

Ujjal Kr Dutta, Mehrtash Harandi, and Chellu Chandra Sekhar

*Abstract*—Metric learning is an important problem in machine learning. It aims to group similar examples together. Existing state-of-the-art metric learning approaches require class labels to learn a metric. As obtaining class labels in all applications is not feasible, we propose an unsupervised approach that learns a metric without making use of class labels. The lack of class labels is compensated by obtaining pseudo-labels of data using a graph-based clustering approach. The pseudo-labels are used to form triplets of examples, which guide the metric learning. We propose a probabilistic loss that minimizes the chances of each triplet violating an angular constraint. A weight function, and an orthogonality constraint in the objective speeds up the convergence and avoids a model collapse. We also provide a stochastic formulation of our method to scale up to large-scale datasets. Our studies demonstrate the competitiveness of our approach against state-of-the-art methods. We also thoroughly study the effect of the different components of our method.

## I. INTRODUCTION

A Key step in artificial intelligence and machine learning algorithms is to find the distance or similarity among examples. Distance Metric Learning (DML) is equivalent to obtaining an embedding where similar examples are grouped together, while moving away dissimilar ones. Recent works in machine learning problems like few-shot object detection [1], zero-shot recognition [2], zero-shot image retrieval and clustering [3], co-saliency detection [4], remote sensing [5] and fine-grained categorization [6] have demonstrated the benefits of learning a distance metric.

Despite their exemplary performance, the problem with existing state-of-the-art DML approaches is that they are *supervised* in nature, *i.e.*, they require class labels (manual annotations) to learn the metric, often in large scale. However, many applications do not have the feasibility of obtaining supervisory signals or class labels. Examples of such applications include medical imaging techniques requiring invasive procedures [7], [8], large-scale 3D point cloud image recognition [9], [10], and image segmentation requiring pixel-level annotations [11], to name a few. Furthermore, obtaining manual annotations also require certain degree of domain expertise, and subjective biases, which often leads to noisy or erroneous labels. Thus, it is essential to be able to learn a metric that could capture the inherent properties of data, without requiring class labels.

To learn a metric, existing supervised DML approaches provide constraints in the form of pairs [12], triplets [13], tuples [14] or batches [15]. These constraints are obtained using the available class labels. Hence, the key challenge that naturally
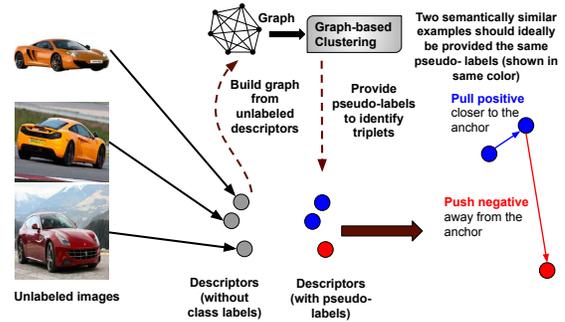
Fig. 1: An illustration of the proposed approach. The input images belong to Cars196 dataset [18]. The figure is best viewed in color.

arises in *unsupervised metric learning* is that of obtaining such constraints. Classically, it has been studied in the context of *manifold learning* [16], [17].

A few recent approaches have attempted this problem from a feature learning perspective [19], [20], [21], [22]. Under the absence of any supervision, some of the approaches [20], [22] have treated each example and its augmentations as a distinct class. The Mining on Manifolds (MOM) approach by Iscen *et al.* [23] follows a paradigm for unsupervised metric learning by making use of a graph. In contrast to other feature learning approaches, rather than treating each instance as a separate class, they actually identify groups of examples that are similar or dissimilar to each other. By performing a random-walk on a graph, they identify manifold and Euclidean neighbors to form triplet constraints, thus indicating a notion of similarity among examples.

The intuitive nature of the MOM approach motivates us to leverage graph-based constraint mining. In contrast to MOM that directly mines hard constraints using manifold and Euclidean neighbors, we simply partition a dataset and obtain pseudo-labels for each example, by using a graph-based clustering approach called Authority Ascent Shift (AAS) [24]. This is because pseudo labels provide us the flexibility to form any type of constraints as desired. With pseudo label information for a mini-batch, one can directly apply existing supervised techniques in a straightforward manner. Moreover, we can also apply online constraint mining strategies (*e.g.*, semi-hard strategy [13]), and recent batch based strategies [25]. Our approach is illustrated in Figure 1. The pseudo-labels are used to identify triplet constraints for metric learning. A triplet essentially consists of an *anchor-positive* pair of semantically similar examples, and a third *negative* example that is dissimilar to the pair.

Following are the major **contributions of our paper**: 1. We

propose a novel probabilistic metric learning objective to ensure that the triplets satisfy an angular property [26]. 2. We further include a weight function, and an orthogonality constraint in our objective. This speeds up the convergence, and avoids a model collapse. 3. We also propose a stochastic formulation of our method to scale up to large datasets, while handling non-linearity in data using deep neural networks. We empirically establish the competitiveness of our method in comparison to existing state-of-the-art approaches, and also carefully analyze the role of different components of our method.

## II. RELATED WORK

Classically, unsupervised DML has been studied as a by-product of *manifold learning* [16], [17] or *diffusion processes* [27], [28]. In the absence of class labels, a natural way to identify the semantic similarity among two examples is to first perform a clustering of the data to obtain *pseudo-labels*, and then learn a metric. The DeepCluster [19] approach follows this approach of jointly learning pseudo-labels obtained from traditional k-means clustering, and using these labels to learn an embedding. However, center-based clustering techniques like k-means are well-known for their drawbacks. For example, sensitivity to initialization, setting the number of clusters apriori, and the lack of effectiveness in higher dimensions.

The Exemplar [21] method randomly samples a set of image patches. It assumes that a few random parameter vectors represent certain elementary transformation operations like translation, scaling, rotation, contrast, and colorization. The goal is to learn these parameters by following a standard supervised setting, where the classes will be obtained from the transformed patches, with the identity of an example denoting the class label.

The NCE [20] and InvariantSpread [22] approaches obtain augmentations of examples, and seek to learn an embedding that pulls augmentations of an example together, while moving away augmentations of different examples. The InvariantSpread method is formulated as an extension of the NCE approach. The Exemplar, NCE, and InvariantSpread methods are *instance-wise* in nature, *i.e.*, they treat each example as a separate class. Ideally, not only we demand augmentations of the same example to be close, we also want to group together two *different* examples, but of the same *semantic similarity*. To address this, the Mining on Manifolds (MOM) [23] approach makes use of a graph-based ranking technique by performing a random-walk. By virtue of the underlying ranking, it identifies manifold and Euclidean neighbors to form triplets. In particular, it identifies *hard* positives and negatives with respect to an anchor example. It then learns a metric by using a standard triplet loss.

## III. BACKGROUND

Let $x_i \in \mathbb{R}^d$ be the descriptor of an example $i$ in a dataset $\mathcal{X}$, which is unlabeled. For $x_i \in \mathbb{R}^d$, let $L^\top x_i \in \mathbb{R}^l$, denote its learned embedding. Here, $L \in \mathbb{R}^{d \times l}$ is the parametric matrix of the squared Mahalanobis-like distance metric $\delta_L^2(x_i, x_j) = (x_i - x_j)^\top L L^\top (x_i - x_j)$, for a pair of examples $x_i, x_j \in \mathbb{R}^d$. Ensuring $l < d$ facilitates dimensionality reduction. The goal of our work is to learn the parametric matrix $L$. As $\mathcal{X}$ is unlabeled, we do not have class labels for learning $L$. To compensate for the lack of class labels, we suggest obtaining pseudo-labels. In our work, we choose the graph-based Authority Ascent Shift (AAS) clustering [24] to obtain the pseudo-labels.

Let, the AAS clustering be denoted by a function $c : \mathbb{R}^d \to \mathbb{Z}^+$ such that $c(x_i)$, a positive integer, denotes the *pseudo-label* assigned to $x_i \in \mathbb{R}^d$. Briefly, AAS requires constructing a weighted graph with nodes representing the examples, and edges between the nearest neighbors. Edge weights denote *affinities* between examples. With $\omega$ denoting the stationary probability distribution of a random walker on the graph, the *node relevancy* from node $i$ to node $j$ can be defined as [24]:

$$\psi(i, j) = d_i T_{ij} \exp(-\gamma(\nabla_\omega(i, j))^2). \tag{1}$$

Here, $d_i$ is the out-degree of node $i$, $T_{ij}$ is the transition probability from node $i$ to node $j$, $\exp(.)$ is the exponential function, $\nabla_\omega(i, j) = [\omega(j) - \omega(i)]$ and $\gamma > 0$ is a hyperparameter. The set of *relevant neighbors* of node $i$ can be defined as:

$$\mathcal{N}_\epsilon(i) = \{j \in \mathcal{V} : \psi(i, j) > \epsilon\} \cup \{i\}. \tag{2}$$

Here, $\epsilon > 0$ is a hyperparameter and $\mathcal{V}$ is the vertex set of the graph. *Authority ascent* of a node $i$ can be performed by moving towards the node $j^*$ such that $j^* = \arg\max_{j \in \mathcal{N}_\epsilon(i)} T_{ij} \nabla_\omega(i, j)$. By subsequently performing authority ascent on neighboring nodes, we can associate a *authority mode* [24] to node $i$. Nodes sharing a common authority mode build a tree. Disjoint trees represent the distinct, arbitrary-shaped clusters present in the data [24].

**Motivation to use AAS:** AAS does not require to fix the number of clusters apriori, which is a crucial benefit in the unsupervised setting. It is able to detect arbitrary-shaped clusters in the data, while being robust to noise and outliers. Li *et al*. [29] pointed the necessity to capture *intra-class variances* that may occur in visual data due to minor pose, illumination or viewpoint differences. Such variances can be captured by AAS, as it relies on a notion of geometric similarity.

## IV. PROPOSED METHODOLOGY

Using the clustering function $c(.)$, we can form a set of triplets: $\mathcal{T} = \{(x_i, x_i^+, x_i^-)\}_{i=1}^{|\mathcal{T}|}$, each element of which consists of the following: i) $x_i$, an arbitrary example with a value $c(x_i)$ . ii) $x_i^+$, another arbitrary example with $c(x_i^+) = c(x_i)$. iii) $x_i^-$, such that $c(x_i^-) \neq c(x_i)$. The examples $x_i$, $x_i^+$ and $x_i^-$ are referred to as the *anchor*, *positive* and *negative* respectively (Figure 1). Here, $|\mathcal{T}|$ is the number of triplets formed. Using $\mathcal{T}$, our goal is to learn $L$ in $\delta_L^2(x_i, x_j)$. In particular, we make use of the *semi-hard* strategy [13] of mining triplets using the pseudo-labels.

### A. Probabilistic Metric Learning Objective with Orthogonality

Given a triplet $(x_i, x_i^+, x_i^-)$, we seek to minimize the following *angular constraint* [26] based hinge-loss term:

$$[z_i]_+ = [\delta_L^2(x_i, x_i^+) - 4 \tan^2\alpha \; \delta_L^2(x_i^-, x_{i-avg})]_+. \tag{3}$$

Here, $[z_i]_+ = \max(0, z_i)$. Figure 2 illustrates the angular constraint using a circle centered at $x_{i-avg} = \frac{x_i + x_i^+}{2}$, while

Fig. 2: Illustration of the angular constraint.
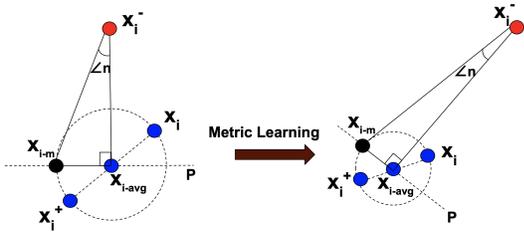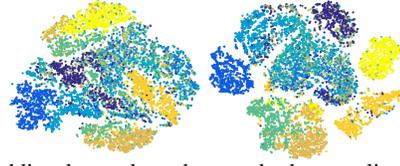


Fig. 3: Embedding learned on the standard test split of the Fashion-MNIST dataset [30] without orthogonality (left-panel), and with orthogonality (right panel) on the metric parameters. While there may be a model collapse without orthogonality, imposing orthogonality leads to well-separated and compact groups.

having the points $\boldsymbol{x}_i$ and $\boldsymbol{x}_i^+$ on its perimeter, and $\boldsymbol{x}_i^-$ outside it. The line $P$ passing through $\boldsymbol{x}_{i-avg}$ intersects the circle at $\boldsymbol{x}_{i-m}$, and is perpendicular to the line segment $\boldsymbol{x}_{i-avg} - \boldsymbol{x}_i^-$. By constraining the $\angle n$ with an upper bound $\angle \alpha > 0°$ (a hyperparameter), one can push the negative $\boldsymbol{x}_i^-$ away from the center $\boldsymbol{x}_{i-avg}$ of the local cluster defined by $\boldsymbol{x}_i$ and $\boldsymbol{x}_i^+$, while dragging the latter two closer. This is achieved by constraining: $\tan(n) \leq \tan(\alpha) \Rightarrow \delta_{\boldsymbol{L}}(\boldsymbol{x}_{i-m}, \boldsymbol{x}_{i-avg})/\delta_{\boldsymbol{L}}(\boldsymbol{x}_i^-, \boldsymbol{x}_{i-avg}) \leq \tan(\alpha) \Rightarrow 0.5\delta_{\boldsymbol{L}}(\boldsymbol{x}_i, \boldsymbol{x}_i^+)/\delta_{\boldsymbol{L}}(\boldsymbol{x}_i^-, \boldsymbol{x}_{i-avg}) \leq \tan(\alpha) \Rightarrow \delta_{\boldsymbol{L}}^2(\boldsymbol{x}_i, \boldsymbol{x}_i^+) \leq 4\tan^2\alpha\ \delta_{\boldsymbol{L}}^2(\boldsymbol{x}_i^-, \boldsymbol{x}_{i-avg})$, which leads to the loss in (3).

As $[z_i]_+$ is non-smooth, we propose to minimize the following smooth version: $m_i = \log(1 + \exp(z_i))$ instead of $[z_i]_+$. This is possible as $\log(\exp(a) + \exp(b)) \geq \max(a, b)$ where $a, b \in \mathbb{R}$. Let, $\sigma(a) = \frac{1}{1+\exp(-a)}, a \in \mathbb{R}$ be the logistic function. We now formulate a novel probabilistic metric learning objective to satisfy the above angular constraint in $m_i$. For this, we define the probability of a triplet $(\boldsymbol{x}_i, \boldsymbol{x}_i^+, \boldsymbol{x}_i^-)$ not violating the angular constraint, as follows:

$$p\{(\boldsymbol{x}_i, \boldsymbol{x}_i^+, \boldsymbol{x}_i^-) \text{ does not violate the angular constraint}\} = p_i = \sigma(-f_i). \tag{4}$$

Here, $f_i$ is a *weighted loss* defined as: $f_i = w_i m_i$. Note that minimizing the loss $f_i$ will maximize $p_i$. $w_i$ is a weight term on the smooth loss $m_i$, which is defined as a function $w_i : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d \to [0, 1]$ that provides a scaling weight for the loss term on a triplet. $w_i$ has a parametric matrix $\boldsymbol{R} \in \mathbb{R}^{d \times l}$, which we discuss later.

To ensure that a triplet satisfies the angular constraint, we need to maximize $p_i$. As each triplet is independent, the joint probability that all the triplets satisfy the angular constraint is expressed as $\prod_i p_i$. Thus, the log likelihood of the triplet set $\mathcal{T}$ satisfying the angular constraint is expressed as: $\log \prod_i p_i = \sum_i \log p_i$. To learn the distance metric, we propose to maximize this log likelihood (*i.e.*, *minimizing* the negative log likelihood), by solving the following optimization problem:

$$\min_{\boldsymbol{R}, \boldsymbol{L}} \mathcal{L}(\boldsymbol{R}, \boldsymbol{L}) = \sum_{i=1}^{|\mathcal{T}|} (-\log p_i). \tag{5}$$

Let, the learned embeddings of the anchor and the positive be denoted as $\boldsymbol{L}^\top \boldsymbol{x}_i$ and $\boldsymbol{L}^\top \boldsymbol{x}_i^+$ respectively. The bilinear similarity between them can be expressed as $(\boldsymbol{L}^\top \boldsymbol{x}_i)^\top \boldsymbol{L}^\top \boldsymbol{x}_i^+ = \boldsymbol{x}_i^\top \boldsymbol{L} \boldsymbol{L}^\top \boldsymbol{x}_i^+$. A lower value of similarity indicates that the pair is too hard to be an anchor-positive one, and should not have been grouped together by the clustering. Hence, we should down-weigh the loss associated with this triplet via a weight term $w_i^+$.

On the contrary, $\boldsymbol{x}_{i-avg}^\top \boldsymbol{L} \boldsymbol{L}^\top \boldsymbol{x}_i^-$ represents the bilinear similarity of the negative $\boldsymbol{x}_i^-$, w.r.t. the average representation $\boldsymbol{x}_{i-avg}$ of the (anchor, positive) pair in the triplet. A higher similarity indicates that all three of them should have been grouped together by the clustering. Hence, we down-weigh the loss associated with this triplet using a weight term $w_i^-$.

We express the final weight term $w_i$ for defining $f_i$ in (4), as follows: $w_i = (w_i^+ + w_i^-)/2$. However, parameterizing the function $w_i$ in terms of $\boldsymbol{L}$ is too restrictive on $\boldsymbol{L}$. Hence, we use a separate parametric matrix $\boldsymbol{R}$ of the same dimensions, to represent $\boldsymbol{L}$. We finally encode the bilinear similarities using the logistic function, and define $w_i^+$ and $w_i^-$ as follows: $w_i^+ = \sigma(\boldsymbol{x}_i^\top \boldsymbol{R} \boldsymbol{R}^\top \boldsymbol{x}_i^+)$, and $w_i^- = 1 - \sigma(\boldsymbol{x}_{i-avg}^\top \boldsymbol{R} \boldsymbol{R}^\top \boldsymbol{x}_i^-)$. In the expression for $w_i^-$, the subtraction from 1 has been done because a higher value of $\boldsymbol{x}_{i-avg}^\top \boldsymbol{R} \boldsymbol{R}^\top \boldsymbol{x}_i^-$ indicates that we should give a lower weightage.

A common practice in metric learning is to impose a regularizer on the parameters of the metric. In our work, we impose orthogonality on the matrix $\boldsymbol{L}$. This is because without orthogonality, there may be a model collapse, leading to degenerate embeddings. Figure 3 shows that the embeddings obtained for the standard test split of the Fashion-MNIST dataset [30] are relatively better with orthogonality. In the context of metric learning, Xie *et al.* [31] further pointed that orthogonality helps in learning a compact set of projection vectors, helps reducing the adverse affects of class imbalance, and avoids overfitting. The general benefits of orthogonality has also been studied by many recent works [32], [33], [31].

---

**Algorithm 1** stochastic OPML (sOPML)

---

1: **Input:** Unlabeled data $\mathcal{X}$, initial $\Phi$; $c : \mathbb{R}^d \to \mathbb{Z}^+$; $\alpha$, $maxiter > 0$.
2: Perform $c(.)$ on $\mathcal{X}$ to obtain pseudo-labels.
3: Initialize $\boldsymbol{R}_{prev}, \boldsymbol{L}_{prev}$.
4: **while** not converged **do**
5:     **for** MB **in** $dataloader$ **do**       ▷ MB: Mini-Batch
6:         images, pseudo-labels ← MB.     ▷ $|MB|$: batch size
7:         $\{\boldsymbol{x}_j, c(\boldsymbol{x}_j)\}_{j=1}^{|MB|} \leftarrow (\Phi(\text{images}), \text{pseudo-labels})$.
8:         $\mathcal{T} = \{(\boldsymbol{x}_i, \boldsymbol{x}_i^+, \boldsymbol{x}_i^-)\}_{i=1}^{|\mathcal{T}|} \leftarrow$ Semi-Hard-Mine$(\{\boldsymbol{x}_j, c(\boldsymbol{x}_j)\}_{j=1}^{|MB|})$ [13].
9:         $[\boldsymbol{R}, \boldsymbol{L}] \leftarrow$ RCGD$(maxiter, \mathcal{T}, \mathcal{L}(\boldsymbol{R}_{prev}, \boldsymbol{L}_{prev}))$   ▷ Using (5)
10:         $\boldsymbol{R}_{prev} \leftarrow \boldsymbol{R}; \boldsymbol{L}_{prev} \leftarrow \boldsymbol{L}$.
11:         loss ← $\mathcal{L}(\boldsymbol{R}, \boldsymbol{L})$.         ▷ Using (5)
12:         loss.backward() and optimizer.step() to learn $\Phi$.  ▷ BackPropagation
13:     **end for**
14: **end while**
15: **return** $\boldsymbol{L}$

---

Upon adding the *orthogonality constraint* $\boldsymbol{L}^\top \boldsymbol{L} = \mathbf{I}_l \in \mathbb{R}^{l \times l}$, the matrix $\boldsymbol{L}$ lies on a orthogonal Stiefel manifold $\text{St}(l, d)$ [34]. However, the objective of (5) is invariant to the right action of the orthogonal group $\mathcal{O}(l) = \{\boldsymbol{B} \in \mathbb{R}^{l \times l} : \boldsymbol{B}\boldsymbol{B}^\top = \boldsymbol{B}^\top \boldsymbol{B} = \mathbf{I}_l\}$, *i.e.*, for $\boldsymbol{B} \in \mathcal{O}(l)$, we have $\mathcal{L}(\boldsymbol{R}, \boldsymbol{L}) = \mathcal{L}(\boldsymbol{R}, \boldsymbol{L}\boldsymbol{B})$.

Therefore, to jointly learn the parameters $\boldsymbol{R}$ and $\boldsymbol{L}$, we must constrain the optimization problem in (5) on the following Riemannian product manifold: $\mathcal{M}_p \triangleq \mathbb{R}^{d \times l} \times \mathcal{G}(l, d)$. Here, $\mathcal{G}(l, d)$ is the Grassmann manifold [34], which is the quotient of $\mathrm{St}(l, d)$. A detailed description on Riemannian geometry and related optimization can be found in [34]. We perform Riemannian Conjugate Gradient Descent (RCGD) to jointly learn the parameters of our method. We call our method as **Orthogonality based Probabilistic Unsupervised Metric Learning (OPML)**.

### B. Incorporation with deep neural networks for scaling up, and computational complexity

To handle large datasets, while capturing non-linearity using a Convolutional Neural Network (CNN), we suggest a stochastic formulation of OPML (Algorithm 1). Given the unlabeled dataset $\mathcal{X}$, let $\Phi$ denote the parameters of a CNN $z : \mathcal{X} \to \mathbb{R}^d$ that provides the non-linear embedding $\boldsymbol{x}_i \in \mathbb{R}^d$ of the $i^{th}$ raw example in $\mathcal{X}$, for use in (5). One can perform RCGD to learn $\boldsymbol{L}$ and $\boldsymbol{R}$ within the loss layer, without requiring any modifications in the architecture (see Figure 4). As our objective is smooth, and fully differentiable, the gradients can be back-propagated to jointly learn $(\Phi, \boldsymbol{R}, \boldsymbol{L})$.

The computational complexity of AAS [24] depends on the number of power iterations required for computing $\omega$. As pointed by the authors, for real applications with sparse graphs, the complexity is much lower, and grows roughly linearly with the dataset size. [35] showed that it is possible to obtain t-SNE embeddings of datasets with millions of objects in $O(N \log N)$. Using that as a preceding step before performing AAS makes AAS scalable to large datasets, while achieving fairly good clustering results. Alternately, one may randomly sample a sufficiently large partition of the data, and perform AAS on that partition. During training we can first pick a random partition, and sample mini-batches from within the partition. Also, (5) is linear in $|\mathcal{T}|$ and quadratic in terms of $d$, and hence, efficient as well.
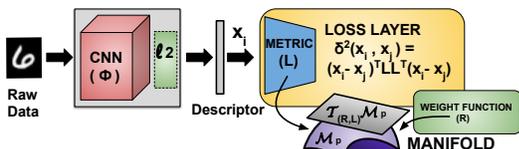


Fig. 4: Architecture required for stochastic OPML (sOPML). The raw image belongs to the MNIST [36] dataset.

## V. Experimental Results

In this section, we study the different components of our proposed method, and make empirical comparisons against state-of-the-art approaches on a number of benchmark datasets.

**Performance Metrics**: All comparisons in our experiments are made with respect to either the clustering or the retrieval performance of an approach on a dataset. The performance on the clustering task is measured in terms of Normalized Mutual Information (NMI), F-measure (F), Precision (P) and Recall (R). NMI is defined as the ratio of mutual information and the average entropy of clusters and entropy of actual ground truth class labels. F-measure is the harmonic mean of precision and recall. For the retrieval task, we use the Recall@K metric that gives us the percentage of test examples that have at least one K nearest neighbor from the same class. A higher value of a metric indicates a better performance for an approach.

**Datasets**: The following datasets have been used:

- **MNIST** [36]: It consists of 70000 gray-scale images of handwritten digits, of $28 \times 28$ pixels. The standard split consists of 60000 training images and 10000 test images.
- **JHMDB** [37]: It is an action recognition dataset that consists of 960 video sequences belonging to 21 actions.
- **HMDB** [38]: It is another action recognition dataset that consists of 6766 videos belonging to 51 classes.
- **Animals with Attributes 2 (AwA2)** [39]: It consists of 37322 images belonging to 50 different classes of animals.
- **Fashion-MNIST** [30]: It consists of images from 10 categories of fashion products. The standard split is as follows: 60000 training images, and 10000 test images.
- **CUB-200** [40]: It consists of images of 200 species of birds with first 100 species for training (5864 examples) and remaining for testing (5924 examples).
- **Cars-196** [18]: This dataset consists of images of cars belonging to 196 models. Usually, the first 98 models containing 8054 images are used for training. The remaining 98 models containing 8131 images are used for testing.
- **Stanford Online Products (SOP)** [15]: This dataset consists of images of online products belonging to 22634 product classes, with a total of 120053 images. Usually, the first 11318 classes containing 59551 images are used for training. The remaining 11316 classes containing 60502 images are used for testing.



Fig. 5: Comparison of AAS (left panel) and k-means (right panel) clustering on a synthetic noisy dataset (Best viewed when zoomed).

### A. Study of the graph-based AAS clustering

Our proposed approach makes use of the AAS clustering to obtain pseudo-labels for mining triplets. The quality of the learned metric depends on the triplets formed, which in turn are dependent on how well the underlying clustering models the data. Therefore, it is important to study the impact of changing the various hyperparameters of this approach. It is a well known fact that k-means clustering forms spherical clusters due to its center-based model. When data has arbitrary-shaped clusters with uniformly distributed noise, the AAS method is able to capture the data distribution in a better way than k-means.

To validate this claim, we generate a 2D synthetic dataset as shown in Figure 5. It consists of 6 classes with 900 examples

in total, where 3 classes are of crescent shape, and 3 are Gaussians. There are 300 additional points uniformly distributed to represent noise in the data. As observed in Figure 5, k-means clustering does not represent the underlying data distribution. On the other hand, the AAS clustering method models the representation fairly well, in addition to filtering out the noise.

We make use of the available ground truth class labels and obtained cluster labels to empirically evaluate the clustering performance in terms of the NMI and F metrics. The AAS clustering approach obtained an NMI of 80.9 and F-measure of 82.9, whereas, the k-means method obtained NMI and F-measure values of 50.0 and 42.7 respectively.

We also conducted experiments on the standard test split of the benchmark MNIST [36] digit dataset containing 10000 examples. Each image is represented by the concatenation of the raw pixels as a vector, followed by a t-SNE embedding [35]. We observed that the k-means clustering achieves NMI and F-measure of 73.4 and 68.7 respectively, whereas the AAS clustering obtains NMI and F-measure of 77.9 and 71.0 respectively. The performances in the synthetic and MNIST datasets indicate that the AAS method outperforms the k-means clustering approach in certain scenarios.

The AAS clustering approach has the following major hyperparameters: i) $k\_graph$: the number of nearest neighbors in the underlying kNN graph, ii) the scale parameter $\gamma$ in (1), and iii) the node relevancy threshold $\epsilon$ in (2). The AAS method defines *cluster authority* as the sum of authority scores of all the elements present in a cluster. A cluster with an unusually low cluster authority (*e.g.*, one that contributes less than $\theta_{min}\%$ of the total cluster authority) can be filtered as a *noise cluster*. One may choose to filter out such (less informative/ noisy) clusters by setting an appropriate value of $\theta_{min}$.

For the above experiment in the MNIST dataset, we set the following values for the AAS approach: $\theta_{min} = 5$, $k\_graph = 500$, $\gamma = 10^2$, and $\epsilon = 0.9$. For the AAS method, we create a kNN graph with an edge weight defined as: $W_{ij} = \exp(-2c^2 \frac{\delta_{ij}^2}{\delta_{max}^2})$, where $\delta_{ij}^2$ is the squared Euclidean distance between the examples represented by nodes $i$ and $j$, and $\delta_{max}^2$ is the maximum distance among all the pairs. The scaling constant $c > 0$ has been set to 1. For the k-means method, we set the number of clusters equal to the actual number of classes. We now use the same test split of the MNIST dataset and perform further studies to highlight the importance of each of the hyperparameters in AAS.

TABLE I: Impact of number of neighbors in the graph used in AAS.

| Dataset | MNIST | | | | |
|---|---|---|---|---|---|
| $k\_graph$ | NMI | F | P | R | $N\_clusters$ |
| 50 | 58.5 | 18.5 | **75.6** | 10.5 | 214 |
| 100 | 58.1 | 26.4 | 60.5 | 16.9 | 65 |
| 500 | **68.0** | **56.2** | 62.1 | 51.4 | 13 |
| 1000 | 63.0 | 56.9 | 46.6 | **73.1** | 7 |

*1) Impact of number of neighbors in the graph:* We study the impact of changing the number of nearest neighbors $k\_graph$ in the kNN graph used in AAS. We fix $\gamma = 10^2$, $\epsilon = 0.5$, and do not eliminate clusters (*i.e.*, $\theta_{min} = 0$). The clustering performance of AAS on the MNIST dataset has been shown in Table I, in terms of NMI, F, P and R values. Best performance for each metric is shown in bold against each column. The

number of clusters detected is also reported.

We observe that a lower value of $k\_graph$ leads to a better precision, while a higher value of $k\_graph$ leads to a better recall. An intermediate value of $k\_graph$ (around 500) leads to a balance in the precision and recall metrics. Noticeably, the number of clusters obtained for $k\_graph = 500$ is close to the actual number of classes, *i.e.*, 10. This is indicated by the underline. We also set $k\_graph = 5000$, for which all examples get merged to two clusters in MNIST, which is not meaningful.

*2) Impact of the node relevancy threshold $\epsilon$:* We now fix $k\_graph = 500$, $\gamma = 10^2$, and study the impact of $\epsilon$ (shown in Table II). Intuitively, a higher value of $\epsilon$ ensures that examples are not readily merged into the same cluster, unless they have sufficient relevance for that cluster. Therefore, we observed that a larger value of $\epsilon$ leads to higher values of performance metrics. However, increasing $\epsilon$ beyond 1 could lead to a sharp rise in the number of clusters (degeneracy because of many examples forming clusters of their own). Therefore, $\epsilon > 1$ should be avoided. We observed that for $\epsilon \to 1$, performing *noise elimination* is a good practice because some of the smaller clusters could be noise or outliers. For the general case of AAS without noise elimination, an intermediate value of $\epsilon$ works well (*e.g.*, $\epsilon = 0.65$).

TABLE II: Impact of node relevancy threshold in AAS.

| Dataset | MNIST | | | | |
|---|---|---|---|---|---|
| $\epsilon$ | NMI | F | P | R | $N\_clusters$ |
| 0.1 | 66.2 | 55.9 | 62.1 | 50.7 | 13 |
| 0.2 | 66.8 | 56.5 | 63.0 | 51.2 | 13 |
| 0.5 | 68.0 | 56.2 | 62.1 | 51.4 | 13 |
| 0.6 | 68.6 | 56.5 | 62.2 | **51.7** | 13 |
| 0.7 | 68.9 | 57.0 | 64.4 | 51.0 | 15 |
| 0.8 | 69.6 | 57.1 | 65.4 | 50.7 | 21 |
| 0.9 | **72.7** | **59.8** | **81.2** | 47.3 | 53 |

*3) Impact of the scale parameter $\gamma$:* We now perform our studies on AAS with respect to the scale parameter $\gamma$. For this, we fix $k\_graph = 500$, $\epsilon = 0.9$ and also eliminate noise by setting $\theta_{min} = 5$. The obtained performances of AAS are reported in Table III. We observed that for $\gamma = 10^2$, the performance is fairly accurate and the number of clusters (shown with underline) is closer to the actual number of classes present in the data. It should be noted that the values of performance metrics reported in Table (III) are higher than those reported in Table (I) and Table (II). This is attributed to the noise removal being performed.

TABLE III: Impact of scale parameter $\gamma$ on AAS.

| Dataset | MNIST | | | | |
|---|---|---|---|---|---|
| $\gamma$ | NMI | F | P | R | $N\_clusters$ |
| 10 | 76.7 | 69.4 | 79.8 | 61.3 | 12 |
| $10^2$ | **77.9** | 71.0 | **80.7** | 63.3 | 11 |
| $10^3$ | 77.7 | **77.7** | 65.6 | **95.3** | 4 |

*4) Effect of eliminating noisy clusters:* As part of our last experiments on AAS, we study the impact of eliminating *noisy clusters*. We set $k\_graph = 500$, $\epsilon = 0.9$ and $\gamma = 10$. Now we vary the values of $\theta_{min}$ and report our observations in Table IV. Here, $\theta_{min} = 0$ indicates that we are not eliminating any noisy cluster. We observed that for the same values of fixed hyperparameters, simply eliminating noisy clusters leads to a better clustering performance by AAS. Typically, each of

the individual clusters contribute to roughly around $10 - 20\%$ of the total authority score. Thus setting $\theta_{min} > 20$ might falsely eliminate all the clusters as noise. Usually, the range $\theta_{min} \in [1, 5]$ leads to meaningful clustering results.

TABLE IV: Impact of noise removal on AAS.

| Dataset | MNIST | | | | |
|---|---|---|---|---|---|
| $\theta_{min}$ | NMI | F | P | R | $N\_clusters$ |
| 0 | 72.9 | 60.5 | 79.7 | 48.7 | 22 |
| 1 | 73.1 | 60.7 | 79.7 | 49.0 | 18 |
| 2 | 73.5 | 61.7 | 79.7 | 50.3 | 17 |
| 5 | **76.7** | **69.4** | **79.8** | **61.3** | 12 |

*5) Conclusions on AAS:* As already discussed, for real applications with sparse graphs, the complexity of AAS could be much lower. For all our further experiments, we apply the t-SNE approach [35] to the representations before performing AAS. This is not only computationally simpler and lets us scale to large datasets, but also provides us a good visualization of how the high-dimensional data is clustered. It should be noted that the tSNE embeddings are used only for performing the clustering. The metric is then learned on the original data. It took 1.67 seconds to compute the affinity matrix for the embeddings of MNIST dataset, which consists of $10^4$ examples. For the actual clustering, it took merely 22.823 seconds on an Intel Core i7-8700 CPU (@ 3.20GHz x 12) with 32 GB RAM. However, due to the large number of hyperparameters, it is important to clearly identify their impact. We make the following conclusions regarding that:

1) In a real application, it is impossible to tune so many hyperparameters in the unsupervised setting. If we have a related dataset, or a separate validation set whose ground truths are available, one may perform clustering using AAS and compute the performance metrics to determine the hyperparameters.

   However, our claim in this work is to learn a metric in a completely *unsupervised* manner. Hence, we do not *tune* the hyperparameters for the rest of our experiments (as such, our approach *could* lead to better performance, if further tuned adequately). We merely take guidance from our observations on the MNIST dataset to fix the hyperparameters of AAS apriori.

2) $k\_graph$, a crucial hyperparameter, when set to 500 led to a better performance on the MNIST dataset, where each class has 1000 examples. This indicates that essentially $k\_graph$ should have a value comparable to (around $50\%$) that of the average number of examples per class.

   But, for real-world datasets with high overlap (*e.g.*, CUB-200 [40], Cars-196 [18], JHMDB [37] etc) where we have fewer examples per class, it is advisable to set $k\_graph$ to a relatively smaller value. For the rest of our experiments, we follow the MOM [23] approach, which is also a graph-based approach, and set $k\_graph = 50$.

3) For the other hyperparameters, combinations of $\gamma \in \{10^2, 10^3\}$ along with either ($\epsilon \approx 0.65, \theta_{min} = 0$) or ($\epsilon \rightarrow 0.9, \theta_{min} \in \{1, 5\}$) works reasonably well. Having fixed the hyperparameters of AAS, our approach becomes fairly simple to learn a metric, as now it leaves us with only the hyperparameter $\alpha$ in (3).

## B. Comparison of OPML against state-of-the-art traditional DML methods

In this subsection, we compare the standalone version of our OPML method against the following state-of-the-art traditional DML methods:

- **JDRML** [41]: This method performs Joint-Dimensionality Reduction and Metric Learning (JDRML) to learn a better metric.
- **LRGMML** [42]: This method provides a Low-Rank (LR) extension to the state-of-the-art Geometric Mean Metric Learning (GMML) method that learns a metric as a point on the geodesic between two scatter matrices.
- **AML** [43]: The Adversarial Metric Learning (AML) approach generates synthetic adversarial pairs to provide difficult constraints for learning a metric.
- **MDMLCLDD** [31]: The Mahalanobis DML with Convex Log-Determinant Divergence (MDMLCLDD) approach provides a convex relaxation to impose orthogonality on the parameters of the distance metric.

For all the above methods, we learn a metric using the class labels, as all of them are *supervised* in nature. Using the learned metric we project the test data to the embedding space, where we perform clustering and retrieval on the test embeddings. However, for our method we do not make use of class labels during training. For the comparisons, we make use of the JHMDB, HMDB, AwA2 and Fashion-MNIST datasets, as discussed earlier. We now provide the experimental protocol.

TABLE V: Comparison against state-of-the-art traditional supervised DML methods. Bold denotes the best value for a metric.

| Dataset | | JHMDB | | | | | |
|---|---|---|---|---|---|---|---|
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| JDRML | Yes | 70.7 | 63.3 | 84.7 | 90.9 | **95.3** | 97.5 |
| LRGMML | Yes | 69.1 | 61.3 | 86.2 | 90.5 | 94.3 | 97.5 |
| AML | Yes | 69.5 | 61.4 | **87.1** | **91.2** | 94.4 | 97.0 |
| MDML | Yes | 69.7 | 62.3 | 87.0 | 91.1 | 94.6 | 97.3 |
| **Ours** | No | **73.4** | **67.6** | 86.6 | 89.8 | 94.4 | **97.8** |
| Dataset | | HMDB | | | | | |
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| JDRML | Yes | 48.8 | 28.8 | 72.8 | 81.6 | 88.7 | 93.4 |
| LRGMML | Yes | 49.5 | 29.6 | 72.6 | 81.7 | 89.1 | 94.0 |
| AML | Yes | **53.2** | **31.5** | **73.3** | **82.5** | **89.5** | **94.8** |
| MDML | Yes | 51.2 | 30.3 | 73.2 | 82.4 | 89.4 | 94.5 |
| **Ours** | No | 52.2 | 31.1 | 72.9 | 81.9 | 89.2 | 94.5 |
| Dataset | | AwA2 | | | | | |
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| JDRML | Yes | 83.2 | 78.7 | 93.9 | 96.9 | 98.8 | 99.4 |
| LRGMML | Yes | 82.3 | 77.6 | 94.8 | 97.5 | 98.9 | 99.3 |
| AML | Yes | 83.3 | 78.1 | 94.7 | 97.6 | **99.1** | 99.5 |
| MDML | Yes | **84.3** | **78.8** | 94.6 | **97.7** | **99.1** | 99.6 |
| **Ours** | No | 83.4 | **78.8** | **94.9** | 97.6 | 99.0 | **99.7** |
| Dataset | | Fashion-MNIST | | | | | |
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| JDRML | Yes | 59.8 | 47.7 | 77.4 | **87.2** | 91.7 | 95.6 |
| LRGMML | Yes | 59.4 | 49.5 | **78.6** | 86.4 | **92.2** | **96.0** |
| AML | Yes | **64.6** | **52.5** | 77.8 | 86.5 | 91.5 | 95.2 |
| MDML | Yes | 63.9 | 50.2 | 77.7 | 86.3 | 91.6 | 95.3 |
| **Ours** | No | 63.9 | 49.7 | 78.0 | 86.5 | 91.6 | 95.7 |

TABLE VI: Sensitivity of OPML towards $\alpha$ on the JHMDB and HMDB datasets, with respect to NMI and F values on the test data.

| Dataset | $\alpha$ | 35° | 40° | 45° | 50° | 55° | 60° |
|---|---|---|---|---|---|---|---|
| JHMDB | NMI | 69.5 | 70.5 | 73.4 | 71.7 | 72.4 | 68.7 |
| | F | 61.4 | 62.2 | 67.6 | 64.8 | 65.3 | 59.0 |
| HMDB | NMI | 51.2 | 52.3 | 52.2 | 52.5 | 52.0 | 51.6 |
| | F | 30.8 | 30.9 | 31.1 | 31.0 | 31.0 | 30.4 |

*1) Experimental protocol:* For the JHMDB dataset, each frame of a video in this dataset is represented by: i) a RGB vector, and ii) a flow vector. These vectors were provided as
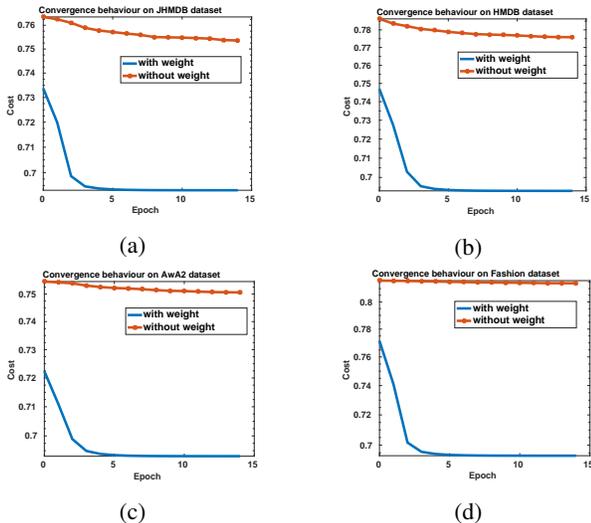
Fig. 6: Convergence behaviour of our method OPML, and role of the weight function, as observed on: a) JHMDB, b) HMDB, c) AwA2, and d) Fashion-MNIST datasets respectively.

part of the work done in Cherian *et al*. [44]. We take average of the *RGB vectors* for all the frames, and represent it as a mean *RGB vector*. We obtain another mean vector obtained by averaging the *flow vectors* for each of the frames. The *mean RGB* and *mean flow* vectors are added together to obtain the final encoding of a video. For the JHMDB dataset, classes 1-8 are for training and 14-21 are for testing.

For the HMDB dataset, we used the split-1 provided as part of the work done in Cherian *et al*. [44]. Similar to JHMDB, we obtain two mean vectors using the representations of all the frames, one for RGB, and the other for flow. These two mean vectors are added to encode a video. For the HMDB dataset, classes 1-21 are for training and 32-51 are for testing. Due to the disjoint nature of training and testing classes in both the datasets, it corresponds to the *zero-shot learning* scenario.

We now design experiments on the AwA2 and Fashion-MNIST datasets to highlight the robustness to noise. For the AwA2 dataset, the features provided in Xian *et al*. [39] have been used in our experiments. The ten largest classes have been picked for our experiments. Based on descending order of the number of examples present in the selected classes, following are the class ids: 7, 23, 40, 49, 31, 27, 38, 29, 1, and 19. From each class, examples 1-200 are used for training, and examples 401-600 are used for testing. To tune the compared supervised baselines, examples 201-400 are used as validation data. To act as noise for the AwA2 dataset, images from the first 10 classes of the CUB dataset [40] have been added. The features for CUB are the same as in AwA2, and are provided by Xian *et al*. [39].

For the Fashion-MNIST dataset, the first 200 examples from each class of the standard training split has been chosen to constitute our training data. To tune the compared supervised baselines, the next 200 examples from each class are used as validation data. The first 200 examples from each class of the standard testing split has been chosen to constitute our test data. To act as noise for the Fashion dataset, we perform the following: i) The first 150 digit images from each class of

the original MNIST dataset is picked. ii) To each pixel of the chosen digit images, we add random noise. iii) This collection of noisy digit images has been added along with the Fashion dataset. The raw pixels have been concatenated to form the feature vectors.

In all datasets, the examples have been $l2$-normalized.

*2) Optimization:* We performed Riemannian Conjugate Gradient Descent (RCGD) using the Manopt toolbox [45]. We kept all the default Manopt parameters (*e.g*., learning rate), except the maximum number of epochs of RCGD, which we set as 30. Due to the product manifold based optimization, our method converges quite fast (Figure 6). This is in accordance with studies in Riemannian optimization [46], [47].

*3) Empirical Observations:* In Table V, we report the results of the compared approaches. For our method, we arbitrarily set $\alpha = 45°$, and embedding size of 128. We also indicate the nature of an approach, *i.e*., whether it makes use of class labels or not. Our method performs competitive, despite not making use of class labels. This shows promise of the probabilistic objective (5) of our method to provide a novel way to learn a metric, while being unsupervised in nature.

### C. Ablation studies studying the role of the pseudo-label mining approach, weight function, and sensitivity to hyperparameter $\alpha$ in the probabilistic objective

Having demonstrated the promise of our proposed orthogonality-based probabilistic objective, we conducted a few experiments to observe the effect of replacing the AAS method with k-means for obtaining pseudo-labels, while setting $k$ to the actual number of classes. By replacing AAS with k-means on the HMDB dataset, we observed a value drop of 6.8 for NMI, 1.2 in R@1, 0.4 in R@2, and 0.5 in R@4 and R@8. Similarly we observed a value drop of 2.4 in NMI, 0.5 in R@1 and 1.6 in R@8, by replacing AAS with k-means on the JHMDB dataset. Interestingly, our method gained a value of 0.5 in NMI by replacing AAS with k-means on MNIST-Fashion, as well as gain of 2.2 in NMI and 0.4 in R@1 upon replacing AAS with k-means on the AwA2 dataset.

This implies that our probabilistic objective could still be used in conjunction with any alternative pseudo-label mining approach. However, exploration of such alternatives is beyond the scope of this paper, as our focus is on the metric learning objective. In our paper, we pick the AAS clustering approach for its robustness, intuitive nature, and conceptual superiority over the naive k-means approach, as shown earlier for certain scenarios. The major drawback of the k-means is its dependency on prior knowledge about the number of classes in data, which is infeasible in the purely unsupervised setting.

We also report the convergence behaviour of our method OPML, both with and without the weight function, in Figure 6. The weight function leads to lower values of the objective and a faster convergence. In Table VI, we report the performance of our method on the JHMDB and HMDB datasets, against varying values of $\alpha$ as present in (3). We observed that the performances are fairly stable, in the range $35° - 60°$.

TABLE VII: Comparisons on FGVC datasets. The best value for each metric is shown in bold, while the second best is shown in underline.

| Dataset | | CUB 200 [40] | | | | | Cars 196 [18] | | | | | SOP [15] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | Labels | NMI | R@1 | R@2 | R@4 | R@8 | NMI | R@1 | R@2 | R@4 | R@8 | NMI | R@1 | R@10 | R@100 |
| Initial (Random) | No | 34.6 | 31.5 | 42.4 | 55.0 | 67.0 | 23.4 | 22.7 | 31.7 | 42.7 | 54.8 | 79.8 | 25.4 | 35.6 | 48.5 |
| DeepCluster [19] | No | 53.0 | 42.9 | 54.1 | 65.6 | 76.2 | 38.5 | 32.6 | 43.8 | 57.0 | 69.5 | 82.8 | 34.6 | 52.6 | 66.8 |
| MOM [23] | No | 55.0 | 45.3 | 57.8 | 68.6 | 78.4 | **38.6** | 35.5 | 48.2 | 60.6 | 72.4 | 84.4 | 43.3 | 57.2 | 73.2 |
| Exemplar [21] | No | 45.0 | 38.2 | 50.3 | 62.8 | 75.0 | 35.4 | 36.5 | 48.1 | 59.2 | 71.0 | <u>85.0</u> | 45.0 | 60.3 | 75.2 |
| InvariantSpread [22] | No | <u>55.4</u> | <u>46.2</u> | <u>59.0</u> | <u>70.1</u> | <u>80.2</u> | 35.8 | <u>41.3</u> | <u>52.3</u> | <u>63.6</u> | <u>74.9</u> | **86.0** | **48.9** | **64.0** | **78.0** |
| **Ours** | No | **55.6** | **47.1** | **59.7** | **72.1** | **82.8** | <u>38.1</u> | **45.0** | **56.2** | **66.7** | **76.6** | 84.2 | <u>45.5</u> | <u>61.6</u> | <u>77.1</u> |

## D. Studies on Fine-Grained Visual Categorization (FGVC)

We now compare the proposed stochastic version of our method (sOPML) against the following unsupervised deep approaches for feature learning: DeepCluster [19], MOM [23], Exemplar [21] and InvariantSpread [22]. As discussed in Section II, these approaches are representative of the various paradigms of unsupervised metric learning. DeepCluster makes use of pseudo-labels using k-means clustering, MOM is a graph-based method utilizing random walk to identify triplets, Exemplar and InvariantSpread are instance-wise methods. We compare the approaches on the CUB-200, Cars-196 and SOP datasets as discussed earlier. These datasets are fine-grained in nature with huge inter-class similarities and intra-class variances, and are standard benchmarks in evaluating deep metric learning approaches, due to the challenging nature of images present in them. The train-test splits mentioned earlier follows standard deep metric learning protocol [15].

We used GoogLeNet [48] pretrained on ImageNet [49], as the backbone CNN, using the MatConvNet [50] tool, and follow the standard protocol for training (on a Tesla V100-PCIE-16GB). Following the graph-based label mining approach MOM [23], the initial features for AAS are formed by the Regional Maximum Activation of Convolutions (R-MAC) [51] right before the average pool layer, and aggregated over three input scales $(512, 512/\sqrt{2}, 256)$. Similar to the training data, we also used the same scales to obtain the final R-MAC embeddings of the test examples using the learned model. We used mini-batch size of 120 and set $maxiter = 10$ in Algorithm 1, *i.e.*, 10 RCGD updates are done for each mini-batch in the loss layer. All other parameters are kept the same as default. The embedding size has been kept as 128, and we set $\alpha = 45°$. The test embeddings are used for evaluating the clustering and retrieval performances of the compared approaches.
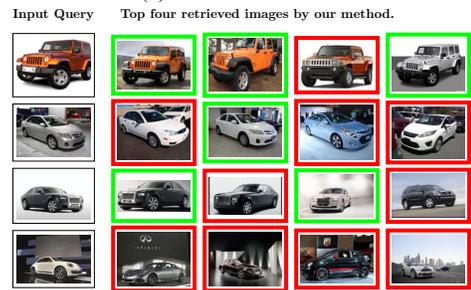
The results on the FGVC datasets are reported in Table VII. Our method achieves state-of-the-art performance on the CUB-200 and Cars-196 datasets. On the SOP dataset, where each class has extremely few examples, the InvariantSpread approach by virtue of its augmentation and spread-out properties performs the best. In Figure 7a-7b, we present a few retrieval results of our method on the Cars-196 dataset. Additionally, Table X compares our method against a few supervised deep baselines.

## E. Studies comparing against MOM, with respect to orthogonality, and role of the probabilistic objective

We performed experiments to study the role of the orthogonality constraint in our approach. A look at Table IX demonstrates that the main benefit of orthogonality comes into the picture when considering the clustering performance, as



(a) Successful cases.



(b) Failure cases.

Fig. 7: (a-b) Retrieval results of our method on the Cars-196 dataset (Best viewed in color). For a query image, a retrieved image is shown in green if it is correct (from same class), and in red if it is otherwise.

also expected from Figure 3. In that case, it is essential to avoid a model collapse (leading to many examples having nearby embeddings). For the retrieval task, it suffices to have just one correct example from the same class within the top K retrieved examples to boost the Recall@K value. This implies that one may not observe a profound impact of orthogonality for retrieval. However, including orthogonality does lead to further improvements in general. Table VIII shows the benefit of orthogonality in the FGVC task. On Cars-196, we observed significant gains in retrieval performance as well.

Using Table IX, we would also like to highlight the gains obtained by our method over MOM, just by using our proposed probabilistic loss in (5), without using orthogonality (Ours-wo_ortho). When compared against MOM, using our loss alone our method achieves: i) 4.3 better R@1 on JHMDB, ii) 3.4, 4.6, 3.7, 2.1 better R@1, R@2, R@4 and R@8 respectively, on HMDB, iii) 6.4, 4.2, 2.1, 0.4 better R@1, R@2, R@4 and R@8 respectively, on AwA2, and iv) 3.5, 2.3, 2.4, 0.5 better R@1, R@2, R@4 and R@8 respectively, on Fashion-MNIST. Similarly, just by using our loss, we outperform MOM consistently in the FGVC datasets as well (Table VIII).

## VI. CONCLUSION AND FUTURE WORK

In this paper, we propose an approach to learn a discriminative distance metric without using manually annotated class labels. The motivation for our paper is the infeasibility of obtaining class labels in many crucial machine learning

TABLE VIII: Ablation study comparing different variants of our method against the MOM approach.

| Dataset | | CUB 200 [40] | | | | | Cars 196 [18] | | | | | SOP [15] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | Labels | NMI | R@1 | R@2 | R@4 | R@8 | NMI | R@1 | R@2 | R@4 | R@8 | NMI | R@1 | R@10 | R@100 |
| MOM [23] | No | 55.0 | 45.3 | 57.8 | 68.6 | 78.4 | 38.6 | 35.5 | 48.2 | 60.6 | 72.4 | 84.4 | 43.3 | 57.2 | 73.2 |
| Ours-wo_ortho | No | 55.5 | 46.7 | 59.2 | 71.7 | 82.6 | 36.8 | 38.8 | 49.4 | 60.1 | 71.1 | 82.5 | 43.8 | 59.4 | 74.5 |
| **Ours** | No | **55.6** | **47.1** | **59.7** | **72.1** | **82.8** | 38.1 | **45.0** | **56.2** | **66.7** | **76.6** | 84.2 | **45.5** | **61.6** | **77.1** |

TABLE IX: Comparison against the graph-based MOM method, and ablation studies with respect to orthogonality.

| Dataset | | JHMDB | | | | | |
|---|---|---|---|---|---|---|---|
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| MOM | No | 65.2 | 57.1 | 80.9 | 88.2 | **94.5** | 96.5 |
| Ours-wo_ortho | No | 66.6 | 53.6 | 85.2 | **90.3** | 94.1 | 96.2 |
| **Ours** | No | **73.4** | **67.6** | **86.6** | 89.8 | 94.4 | **97.8** |
| Dataset | | HMDB | | | | | |
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| MOM | No | 41.4 | 24.4 | 67.1 | 75.9 | 84.2 | 91.7 |
| Ours-wo_ortho | No | 47.8 | 28.7 | 70.5 | 80.5 | 87.9 | 93.8 |
| **Ours** | No | **52.2** | **31.1** | **72.9** | **81.9** | **89.2** | **94.5** |
| Dataset | | AwA2 | | | | | |
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| MOM | No | 68.1 | 62.2 | 87.1 | 93.3 | 96.6 | 98.7 |
| Ours-wo_ortho | No | 73.6 | 65.4 | 93.5 | 97.5 | 98.7 | 99.1 |
| **Ours** | No | **83.4** | **78.8** | **94.9** | **97.6** | **99.0** | **99.7** |
| Dataset | | Fashion-MNIST | | | | | |
| Method | Labels | NMI | F | R@1 | R@2 | R@4 | R@8 |
| MOM | No | 53.7 | 42.2 | 73.1 | 82.3 | 89.3 | 94.8 |
| Ours-wo_ortho | No | 54.8 | 42.5 | 76.6 | 84.6 | **91.7** | 95.3 |
| **Ours** | No | **63.9** | **49.7** | **78.0** | **86.5** | 91.6 | **95.7** |

TABLE X: Comparison against supervised deep metric learning methods on the FGVC datasets.

| Dataset | Method | Labels | R@1 | R@2 | R@4 | R@8 |
|---|---|---|---|---|---|---|
| **CUB 200** | Triplet-SH [13] | Yes | 40.6 | 52.3 | 64.2 | 75.0 |
| | Angular [26] | Yes | 53.6 | 65.0 | 75.3 | 83.7 |
| | Multi-Sim [52] | Yes | 65.7 | 77.0 | 86.3 | 91.2 |
| | Ours | No | 47.1 | 59.7 | 72.1 | 82.8 |
| Dataset | Method | Labels | R@1 | R@2 | R@4 | R@8 |
| **Cars 196** | Triplet-SH [13] | Yes | 53.2 | 65.4 | 74.3 | 83.6 |
| | Angular [26] | Yes | 71.3 | 80.7 | 87.0 | 91.8 |
| | Multi-Sim [52] | Yes | 84.1 | 90.4 | 94.0 | 96.5 |
| | Ours | No | 45.0 | 56.2 | 66.7 | 76.6 |
| Dataset | Method | Labels | R@1 | R@10 | R@100 |
| **SOP** | Triplet-SH [13] | Yes | 57.8 | 75.3 | 88.1 |
| | Angular [26] | Yes | 67.9 | 83.2 | 92.2 |
| | Multi-Sim [52] | Yes | 78.2 | 90.5 | 96.0 |
| | Ours | No | 45.5 | 61.6 | 77.1 |

applications. In particular, we proposed a novel orthogonality-based probabilistic loss for metric learning, that inherently seeks to preserve an angular constraint on a triplet of examples obtained using pseudo-labels. In our paper, we employ a graph-based clustering method to obtain pseudo-labels due to its intuitive and conceptual prowess. However, the same can be replaced with an alternative pseudo-label mining method. The orthogonality component in our method avoids a model collapse of the embeddings obtained, and in general leads to a better metric. In future, our method could be extended along with a complementary augmentation-based self-supervised approach to obtain the initial representations.

## REFERENCES

[1] Leonid Karlinsky, Joseph Shtok, Sivan Harary, Eli Schwartz, Amit Aides, Rogerio Feris, Raja Giryes, and Alex M. Bronstein. Repmet: Representative-based metric learning for classification and few-shot object detection. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1

[2] Huajie Jiang, Ruiping Wang, Shiguang Shan, and Xilin Chen. Adaptive metric learning for zero-shot recognition. *IEEE Signal Processing Letters*, 2019. 1

[3] Binghui Chen and Weihong Deng. Energy confused adversarial metric learning for zero-shot image retrieval and clustering. In *Proc. of Association for the Advancement of Artificial Intelligence (AAAI)*, 2019. 1

[4] Junwei Han, Gong Cheng, Zhenpeng Li, and Dingwen Zhang. A unified metric learning-based framework for co-saliency detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2473–2483, 2017. 1

[5] Gong Cheng, Ceyuan Yang, Xiwen Yao, Lei Guo, and Junwei Han. When deep learning meets metric learning: Remote sensing image scene classification via learning discriminative cnns. *IEEE transactions on geoscience and remote sensing*, 56(5):2811–2821, 2018. 1

[6] Junwei Han, Xiwen Yao, Gong Cheng, Xiaoxu Feng, and Dong Xu. P-cnn: Part-based convolutional neural networks for fine-grained visual categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019. 1

[7] Christian S Perone, Pedro Ballester, Rodrigo C Barros, and Julien Cohen-Adad. Unsupervised domain adaptation for medical imaging segmentation with self-ensembling. *NeuroImage*, 194:1–11, 2019. 1

[8] Liu Yang, Rong Jin, Lily Mummert, Rahul Sukthankar, Adam Goode, Bin Zheng, Steven CH Hoi, and Mahadev Satyanarayanan. A boosting framework for visuality-preserving distance metric learning and its application to medical image retrieval. *IEEE transactions on pattern analysis and machine intelligence*, 32(1):30–44, 2008. 1

[9] Loic Landrieu and Mohamed Boussaha. Point cloud oversegmentation with graph-structured deep metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1

[10] Lin Zhao and Wenbing Tao. Jsnet: Joint instance and semantic segmentation of 3d point clouds. In *AAAI*, pages 12951–12958, 2020. 1

[11] Rui Qian, Yunchao Wei, Honghui Shi, Jiachen Li, Jiaying Liu, and Thomas Huang. Weakly supervised scene parsing with point-based distance metric learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8843–8850, 2019. 1

[12] S Chopra, R Hadsell, and Y LeCun. Learning a similarity metric discriminatively, with application to face verification. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 539–546. IEEE, 2005. 1

[13] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 815–823, 2015. 1, 2, 3, 9

[14] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Proc. of Neural Information Processing Systems (NeurIPS)*, pages 1857–1865, 2016. 1

[15] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4004–4012, 2016. 1, 4, 8, 9

[16] Xiaofei He and Partha Niyogi. Locality preserving projections. In *Proc. of Neural Information Processing Systems (NeurIPS)*, pages 153–160, 2003. 1, 2

[17] Xiaofei He, Deng Cai, Shuicheng Yan, and Hong-Jiang Zhang. Neighborhood preserving embedding. In *Proc. of IEEE International Conference on Computer Vision (ICCV)*, pages 1208–1213, 2005. 1, 2

[18] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proc. of IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 554–561, 2013. 1, 4, 6, 8, 9

[19] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proc. of European Conference on Computer Vision (ECCV)*, pages 132–149, 2018. 1, 2, 8

[20] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3733–3742, 2018. 1, 2

[21] A Dosovitskiy, P Fischer, JT Springenberg, M Riedmiller, and T Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 38(9):1734–1747, 2016. 1, 2, 8

[22] Mang Ye, Xu Zhang, Pong C Yuen, and Shih-Fu Chang. Unsupervised embedding learning via invariant and spreading instance feature. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6210–6219, 2019. 1, 2, 8

[23] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Mining on manifolds: Metric learning without labels. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2, 6, 8, 9

[24] Minsu Cho and Kyoung Mu Lee. Mode-seeking on graphs via random walks. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 606–613. IEEE, 2012. 1, 2, 4

[25] Xun Wang, Haozhi Zhang, Weilin Huang, and Matthew R. Scott. Cross-batch memory for embedding learning. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1

[26] Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. Deep metric learning with angular loss. In *Proc. of IEEE International Conference on Computer Vision (ICCV)*, 2017. 2, 9

[27] Michael Donoser and Horst Bischof. Diffusion processes for retrieval revisited. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1320–1327, 2013. 2

[28] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, Teddy Furon, and Ondrej Chum. Efficient diffusion on region manifolds: Recovering small objects with compact cnn representations. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2077–2086, 2017. 2

[29] Dong Li, Wei-Chih Hung, Jia-Bin Huang, Shengjin Wang, Narendra Ahuja, and Ming-Hsuan Yang. Unsupervised visual representation learning by graph-based consistent constraints. In *Proc. of European Conference on Computer Vision (ECCV)*, pages 678–694. Springer, 2016. 2

[30] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. 3, 4

[31] Pengtao Xie, Wei Wu, Yichen Zhu, and Eric P Xing. Orthogonality-promoting distance metric learning: convex relaxation and theoretical analysis. In *Proc. of International Conference on Machine Learning (ICML)*, 2018. 3, 6

[32] Soumava Kumar Roy, Zakaria Mhammedi, and Mehrtash Harandi. Geometry aware constrained optimization techniques for deep learning. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4460–4469, 2018. 3

[33] Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regularizations in training deep networks? In *Proc. of Neural Information Processing Systems (NeurIPS)*, pages 4261–4271, 2018. 3

[34] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009. 3, 4

[35] Laurens Van Der Maaten. Accelerating t-sne using tree-based algorithms. *The Journal of Machine Learning Research*, 15(1):3221–3245, 2014. 4, 5, 6

[36] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 4, 5

[37] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action recognition. In *Proc. of IEEE International Conference on Computer Vision (ICCV)*, pages 3192–3199, December 2013. 4, 6

[38] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human motion recognition. In *Proc. of IEEE International Conference on Computer Vision (ICCV)*, 2011. 4

[39] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning-a comprehensive evaluation of the good, the bad and the ugly. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018. 4, 7

[40] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010. 4, 6, 7, 8, 9

[41] Mehrtash Harandi, Mathieu Salzmann, and Richard Hartley. Joint dimensionality reduction and metric learning: A geometric take. In *Proc. of International Conference on Machine Learning (ICML)*, 2017. 6

[42] Mukul Bhutani, Pratik Jawanpuria, Hiroyuki Kasai, and Bamdev Mishra. Low-rank geometric mean metric learning. *arXiv preprint arXiv:1806.05454*, 2018. 6

[43] Shuo Chen, Chen Gong, Jian Yang, Xiang Li, Yang Wei, and Jun Li. Adversarial metric learning. In *Proc. of International Joint Conference on Artificial Intelligence (IJCAI)*, 2018. 6

[44] Anoop Cherian, Suvrit Sra, Stephen Gould, and Richard Hartley. Non-linear temporal subspace representations for activity recognition. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2197–2206, 2018. 7

[45] Nicolas Boumal, Bamdev Mishra, P-A Absil, and Rodolphe Sepulchre. Manopt, a Matlab toolbox for optimization on manifolds. *The Journal of Machine Learning Research*, 15(1):1455–1459, 2014. 7

[46] Hongyi Zhang, Sashank J Reddi, and Suvrit Sra. Riemannian svrg: Fast stochastic optimization on riemannian manifolds. In *Proc. of Neural Information Processing Systems (NeurIPS)*, pages 4592–4600, 2016. 7

[47] Silvere Bonnabel. Stochastic gradient descent on riemannian manifolds. *IEEE Transactions on Automatic Control*, 58(9):2217–2229, 2013. 7

[48] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015. 8

[49] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision (IJCV)*, 115(3):211–252, 2015. 8

[50] Andrea Vedaldi and Karel Lenc. Matconvnet: Convolutional neural networks for matlab. In *Proceedings of the 23rd ACM international conference on Multimedia*, pages 689–692. ACM, 2015. 8

[51] Giorgos Tolias, Ronan Sicre, and Hervé Jégou. Particular object retrieval with integral max-pooling of cnn activations. In *Proc. of International Conference on Learning Representations (ICLR)*, 2016. 8

[52] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5022–5030, 2019. 9