

Deep Metric Learning with Alternating Projections onto Feasible Sets

Oğul Can* Yeti Z. Gürbüz† A. Aydın Alatan‡

Center for Image Analysis (OGAM), Middle East Technical University, Turkey

*ogul.can@metu.edu.tr, †ygurbuz@metu.edu.tr, ‡alatan@metu.edu.tr

Abstract

During the training of networks for distance metric learning, minimizers of the typical loss functions can be considered as "feasible points" satisfying a set of constraints imposed by the training data. To this end, we reformulate distance metric learning problem as finding a feasible point of a constraint set where the embedding vectors of the training data satisfy desired intra-class and inter-class proximity. The feasible set induced by the constraint set is expressed as the intersection of the relaxed feasible sets which enforce the proximity constraints only for particular samples (a sample from each class) of the training data. Then, the feasible point problem is to be approximately solved by performing alternating projections onto those feasible sets. Such an approach introduces a regularization term and results in minimizing a typical loss function with a systematic batch set construction where these batches are constrained to contain the same sample from each class for a certain number of iterations. Moreover, these particular samples can be considered as the class representatives, allowing efficient utilization of hard class mining during batch construction. The proposed technique is applied with the well-accepted losses and evaluated on Stanford Online Products, CAR196 and CUB200-2011 datasets for image retrieval and clustering. Outperforming state-of-the-art, the proposed approach consistently improves the performance of the integrated loss functions with no additional computational cost and boosts the performance further by hard negative class mining.

1. Introduction

Distance metric learning (DML) is the problem of finding a proper function that satisfies metric axioms and assesses the semantic dissimilarity of the data samples from its domain. This task is generally realized by learning proper representations for the data samples so that the semantically similar ones are embedded to the small vicinity in the representation space as the dissimilar samples are placed relatively apart in the Euclidean sense. The representations are learned through an optimization framework in which

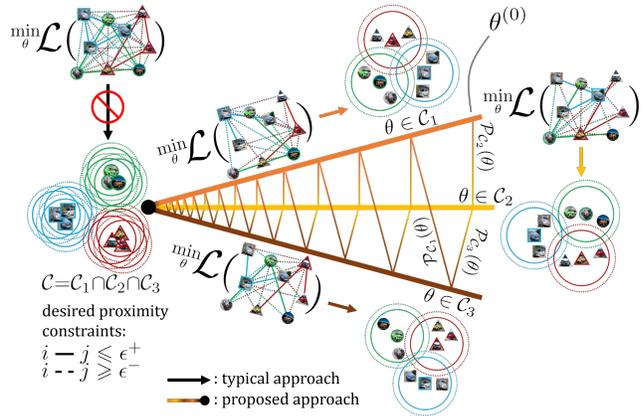


Figure 1: Proposed approach to DML problem. Instead of directly minimizing a loss function composed of the penalty terms enforcing the all proximity constraints (left), we alternatively minimize loss functions of proximity constraints only for particular samples (outlined samples). Each of the three lines (orange, yellow, brown) depicts a subset where the proximity constraints are satisfied for outlined samples as in Eq. (4.2). The solutions are related by projection to obtain a solution at the intersection. (Best viewed once magnified).

the objective function utilizes the loss terms to impose the desired intra-class and inter-class proximity constraints in the representation space [11, 14, 25, 30, 35, 36, 38, 45, 48]. The optimization is performed with mini-batch gradient updates and the procedure is generally guided by providing deliberately selected exemplars [10, 12, 35, 36, 39, 45, 46, 48, 53].

Existing approaches focus on inventing loss functions to enforce all the proximity constraints at once. However, mini-batch gradient update nature of the optimization procedure implies alternatively considering only the subsets of the proximity constraints and eventually, the obtained representations are fail to satisfy the desired proximity constraints holistically due to possible traps in local minima.

To alleviate this problem, we revisit the proximity constraints in the representation space implied by the loss terms for proper DML. To develop a general framework for DML, we focus on finding a feasible point satisfying the proximity

constraints. Such an alternative track is novel in the development of the DML frameworks and addresses the challenge of satisfying the proximity constraints for all data pairs to better match the distance in the representation space with semantic dissimilarity.

In contrast to existing methods, we approach the problem by posing it as a set intersection problem and propose to solve it by performing alternating projections onto the relaxed sets defined by the subsets of the proximity constraints. Our formulation results in relatively easier subproblems to be solved by minimizing the regularized version of the typical loss functions for DML with a systematic batch construction, where the batches are constrained to contain a particular sample from each class for a certain number of iterations. Not only that structure allows efficient utilization of hard negative class mining (HNCM) to guide the optimization without offline processing during batch construction but also the subproblems better fit the mini-batch gradient update procedure and are less prone to be trapped in local minima. Different from the existing methods, the subproblems are related by a regularization term owing to projection of the solutions to the subsets. Our approach to the DML problem is depicted in Fig. 1.

The implications and contributions as the results of our formulation are 1) a general framework that better exploits the proximity constraints to improve the performances of the current and possibly future DML loss functions, 2) idea of re-utilization of the particular class samples in the consecutive mini-batches during the optimization iterations, 3) a solid background through why such a biased batch construction should work, 4) relating the mini-batch updates in general with a regularization term and 5) an efficient class mining method for the batch sampling with $\mathcal{O}(L)$ complexity in contrast with similar methods [10, 12, 39] of $\mathcal{O}(N^2)$.

2. Notations and Definitions

We consider dataset, $\{(x_i, y_i)_{i \in \mathcal{N}} \mid i \in \mathcal{N}\}$, of two-tuples, where $x_i \in \mathcal{X}$ denotes a sample vector from the data space (e.g. images), $y_i \in \mathcal{Y} = \{1, \dots, L\}$ denotes the corresponding label of the sample among L many classes and $\mathcal{N} = \{1, \dots, N\}$ denotes the set of indexes to represent samples from the dataset of size N . Indicator of the two samples indexed by i and j belonging to the same class is denoted as $y_{i,j} \in \{0, 1\}$ where $y_{i,j} = 1$ if $y_i = y_j$. We call j positive/negative sample for i if $y_{i,j} = 1/0$.

The parametric distance between x_i and x_j is defined as:

$$d_{i,j}^f(\theta) \triangleq \|f(x_i; \theta) - f(x_j; \theta)\|_2 \quad (2.1)$$

which is the Euclidean distance equipped with a D -dimensional vector valued parametric function, $f : \mathcal{X} \xrightarrow{f} \mathbb{R}^D$, with parameters θ . The projection of θ onto a set \mathcal{S} is:

$$\mathcal{P}_{\mathcal{S}}(\theta) \triangleq \arg \min_{\vartheta \in \mathcal{S}} \frac{1}{2} \|\theta - \vartheta\|_2^2. \quad (2.2)$$

For a set defined by an inequality $\mathcal{S} = \{\theta \mid g(\theta) \leq 0\}$, we denote its indicator as $\iota_{\mathcal{S}}(\theta)$ which is:

$$\iota_{\mathcal{S}}(\theta) \triangleq \lim_{\lambda \rightarrow 0} \left[\frac{1}{\lambda} g(\theta) \right]_+, \quad (2.3)$$

where $[z]_+ = \max\{0, z\}$ and $g(\cdot)$ is an arbitrary function defining \mathcal{S} . We are to approximate $\iota_{\mathcal{S}}(\theta)$ for small λ as:

$$\iota_{\mathcal{S}}(\theta) \approx \frac{1}{\lambda} [g(\theta)]_+ \triangleq \hat{\iota}_{\mathcal{S}}(\theta). \quad (2.4)$$

3. Review of the Related Works

We restrict ourselves to the distance metric learning problem which is posed as learning the parameters, θ , of an embedding function, $f(\cdot; \theta)$, so that the parametric distance, $d_{i,j}^f(\theta)$, between the data samples reflects their semantic dissimilarity. f as a linear mapping is considered in earlier approaches [29, 47, 49] that later inspire most of state-of-the-art frameworks [11, 14, 25, 30, 35, 36, 38, 45, 48] in which f is a nonlinear mapping realized by deep neural networks. The general framework for learning the parameters is to minimize an overall loss function of the proximity constraints:

$$\mathcal{L}^{(\cdot)}(\theta; \mathcal{T}) = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \ell_t^{(\cdot)}(\theta), \quad (3.1)$$

where \mathcal{T} is the set of index tuples (e.g. pair (i, j) , triplet (i, j^+, j^-) , etc.) and $\ell_t^{(\cdot)}(\theta)$ is a loss term penalizing the ranking violations among the samples indexed by t . We omit dataset, \mathcal{D} , dependency of, \mathcal{L} , for clarity.

Learning a proper linear mapping for the parametric distance is initially formulated as a convex optimization problem in [49]. To prevent null mappings, a constraint that enforces mapping of the samples from different classes to be at least separated by some margin is added to the formulation. Moving that constraint to the objective via hinge loss [11] results in the well-known contrastive loss:

$$\ell_{i,j}^{ctrv}(\theta) = y_{i,j} d_{i,j}^f(\theta)^2 + (1 - y_{i,j}) [\varepsilon - d_{i,j}^f(\theta)]_+^2, \quad (3.2)$$

for sample pairs, (i, j) . Contrastive loss ignores intra-class variations of the classes. Triplet loss introduced in [47] and popularized in deep metric learning frameworks [12, 35] alleviates this problem by constraining the distance to any positive sample to be at least some margin smaller than the distance to any negative sample for each sample:

$$\ell_{i,j^+,j^-}^{triplet}(\theta) = [d_{i,j^+}^f(\theta)^2 - d_{i,j^-}^f(\theta)^2 + \varepsilon]_+, \quad (3.3)$$

where j^+ and j^- are a positive and a negative sample for the sample i , respectively. Minimizing triplet loss entails deliberately selection of the triplets to have nonzero loss terms. Thus, either large batch size or mining for exemplars violating the triplet constraint is required [35]. Such an effort makes the computation of the triplet loss is less attractive

than of the contrastive loss. Margin-based loss is introduced in [48] to provide the flexibility in the distribution of the classes in the embedding space without using triplets as the exemplars. It expresses the margin constraint of the triplet loss as separate loss terms of the distances between positive and negative sample pairs by relaxing the constraint of the contrastive loss on the positive pairs:

$$\begin{aligned} \ell_{i,j}^{margin}(\theta) &= y_{i,j} [d_{i,j}^f(\theta) - \varepsilon + \delta]_+ \\ &+ (1 - y_{i,j}) [\varepsilon + \delta - d_{i,j}^f(\theta)]_+, \end{aligned} \quad (3.4)$$

where δ controls the separation margin and ε is a trainable parameter for the boundary between positive and negative pairs. The contrastive, triplet and margin-based loss terms are the simplest forms of the pairwise distance ranking based losses. Proceeding approaches utilize smoothed versions of these losses by replacing hinge loss with log-sum-exp [46, 51] or soft-max [36, 52] expression. In a different aspect, angular loss [44] that constrains the local geometry of the samples in the embedding space is proposed to better exploit the relation among triplets. In those losses, only 2 or 3 data samples contribute to the loss terms. Ranking among more samples are considered in quadruplet [6, 15, 20], histogram [4, 42] and soft-batch-mining [25, 36, 46, 52] based losses. Soft-batch-mining exploits log-sum-exp expression as the approximation of max operator to select samples from the batch. Thus, ranking among multiple samples are considered.

The aforementioned approaches consider minimization of an overall loss function to enforce proximity constraints. However, the optimization is performed upon mini-batches due to the vast amount of constraints. Therefore, intra-class variations are prone to be missed, leading to poor generalization [16, 22, 32]. Utilizing informative tuples with non-trivial settings in the gradient updates is considered to address the side effects of the mini-batch gradient updates to improve the representations [10, 12, 35, 36, 39, 48, 53]. Yet, global mining [10, 12, 39] for such informative tuples brings additional computational burden and limits scalability to the large datasets, on the other hand, approximate mining [35, 36, 45, 46, 48, 53] is limited to cover enough settings to capture variations in the data. Generation of synthetic hard negative samples through adversarial [9, 55] and variational [22, 56] models are addressed in recent approaches.

To enhance the diversity of the semantic information embedded to the vector representations, ensemble techniques are also combined with deep metric learning framework [17, 26, 34, 50]. The general idea is simply concatenating the vectors from multiple embedding functions whose parameters are learned by considering different local features of the samples. Hence, better embedding space can be obtained by integrating the vectors that are specialized to different aspects of the samples. Modeling of the class variations through clustering track are also addressed in [8, 30, 31]. Rather than modeling, eliminating intra-class variances is

considered in [7, 16, 22, 32]. Though the general framework is to disentangle the intra-class variance upon global representation, differently in [16], variations in the local features are addressed and higher-order moments are considered to regularize the local features so that their aggregation is of less variance.

Dealing with the relaxed proximity constraints for better representation learning is not only the motivation of us but also the motivation of the recent approaches [5, 7, 23, 30]. The underlying idea is to better solve the relatively simpler problems. Existing approaches either are limited to the solution of the relaxed problem [7, 23, 30] or fail to effectively exploit the relations among the relaxed problems to improve the solutions of the global problem [5, 34]. Contrarily to existing approaches, our motivation is to improve generalization of the learned representations by satisfying proximity constraints more holistically through effectively combining the solutions of the simpler relaxed problems. To this end, we propose a feasible point formulation to inherently relate the relaxed problems and introduce a regularizer to be exploited not only in the proposed but also in a typical DML framework.

4. Proposed Approach

We revisit the early DML formulations [47, 49] which involve pairwise proximity constraints in the embedding space for the negative pairs. Yet, differently, we omit the positive pair distances to be minimized from the objective function and introduce it as a constraint for the constraint set of θ . Given a dataset of N samples, $\{(x_i, y_i)\}_{i \in \mathcal{N}}$, we consider the constraint set $\mathcal{C} = \mathcal{C}^+ \cap \mathcal{C}^-$ for θ :

$$\mathcal{C}^+ = \{\theta \mid d_{i,j}^f(\theta) \leq \varepsilon^+, \forall (i, j) \in \mathcal{N}^2, y_{i,j} = 1\}, \quad (4.1a)$$

$$\mathcal{C}^- = \{\theta \mid d_{i,j}^f(\theta) \geq \varepsilon^-, \forall (i, j) \in \mathcal{N}^2, y_{i,j} = 0\}. \quad (4.1b)$$

Once $\varepsilon^+ < \varepsilon^-$, the constraints \mathcal{C}^+ and \mathcal{C}^- together enforce $f(\cdot; \theta)$ to map the samples from the same class to some neighborhood such that no sample from any other class can be mapped to that neighborhood.

Proposition 4.1. *Any $\theta \in \mathcal{C}$ for some ε^+ and ε^- is a global minimizer of the loss function, $\mathcal{L}^{(\cdot)}(\theta; \mathcal{T})$, defined in Eqn. (3.1) for the loss terms $\ell_{i,j}^{ctrv}(\theta)$, $\ell_{i,j^+,j^-}^{triplet}(\theta)$ and $\ell_{i,j}^{margin}(\theta)$ defined in Eqns. (3.2)-(3.4), respectively.*

Proof. The loss terms defined in Eqns. (3.2)-(3.4) satisfy $\ell_t^{(\cdot)}(\theta) \geq 0, \forall \theta$ which implies $\min \mathcal{L}^{(\cdot)}(\theta; \mathcal{T}) = 0$. Thus, it is enough to show that the proper selection of ε^+ and ε^- yields $\ell_t^{(\cdot)}(\theta) = 0, \forall t \in \mathcal{T}, \forall \theta \in \mathcal{C}$. If $\varepsilon^- = \varepsilon$, any $\theta \in \mathcal{C}$ results in $\ell_{i,j}^{ctrv}(\theta) \rightarrow 0$ as $\varepsilon^+ \rightarrow 0, \forall (i, j) \in \mathcal{T}$. Similarly, choosing $\varepsilon^+ = \varepsilon - \delta$ and $\varepsilon^- = \varepsilon + \delta$ makes $\ell_{i,j}^{margin}(\theta) = 0, \forall (i, j) \in \mathcal{T}, \forall \theta \in \mathcal{C}$. Finally, for any ε^+ , choosing $\varepsilon^- = ((\varepsilon^+)^2 + \varepsilon)^{1/2}$ results in $\ell_{i,j^+,j^-}^{triplet}(\theta) = 0, \forall (i, j^+, j^-) \in \mathcal{T}, \forall \theta \in \mathcal{C}$. \square

Proposition 4.1 can be extended to the other pairwise distance based losses and suggests that finding a feasible point of \mathcal{C} is equivalent to solving the minimization of those loss functions. This is actually the restatement of the motivation of the existing approaches [11, 21, 23, 35, 36, 38, 42, 44, 48] in which the loss functions are developed to impose those constraints in the first place. Therefore, these methods can be considered as implicitly finding a feasible point of the constraint set via developing a loss function to be minimized. We address the problem differently by directly focusing on finding a feasible point and develop the loss function accordingly. To formulate our approach, we consider the relaxed set $\mathcal{C}_k = \mathcal{C}_k^+ \cap \mathcal{C}_k^-$:

$$\begin{aligned}\mathcal{C}_k^+ &= \{\theta \mid d_{k_l, j}^f(\theta) \leq \varepsilon^+, \forall (l, j), l \in \mathcal{Y}, j \in \mathcal{N}, y_{k_l, j} = 1\}, \\ \mathcal{C}_k^- &= \{\theta \mid d_{k_l, j}^f(\theta) \geq \varepsilon^-, \forall (l, j), l \in \mathcal{Y}, j \in \mathcal{N}, y_{k_l, j} = 0\},\end{aligned}\quad (4.2)$$

where $k_l \in \{i \in \mathcal{N} \mid y_i = l\}$ denotes the index of a sample from class l . Note that the set $\{k_l\}_{l=1}^L$ contains a sample index from each class. For all samples, \mathcal{C}_k enforces proximity constraints only relative to the particular class samples indexed by $\{k_l\}_{l=1}^L$. For each k , we consider distinct samples from each class such that $\{k_l\}_{l=1}^L \cap \{k'_l\}_{l=1}^L = \emptyset$ for $k' \neq k$. Then, \mathcal{C} can be expressed as $\mathcal{C} = \bigcap_{k=1}^K \mathcal{C}_k$, where K is the total number of sets, which can be considered as the maximum number of samples for a class. Then, the feasible point problem can be reformulated as finding a point in the intersection of the sets. If the sets, $\{\mathcal{C}_k\}_k$, were closed and convex, the problem would be solvable by alternating projection methods [2, 3]. However, it is not uncommon to perform alternating projection methods to non-convex set intersection problems [27, 37]. Hence, we propose to solve the problem approximately by performing alternating projections onto the feasible sets $\{\mathcal{C}_k\}_k$. Hence, the problem becomes:

$$\theta^* = \lim_{k \rightarrow \infty} \theta^{(k)}, \text{ where } \theta^{(k)} = \mathcal{P}_{\mathcal{C}_k}(\theta^{(k-1)}), \quad (4.3)$$

with $\mathcal{C}_{k+K} \triangleq \mathcal{C}_k$ and $\theta^{(0)}$ is arbitrary. A problem instance corresponding to a projection becomes:

$$\theta^{(k)} = \mathcal{P}_{\mathcal{C}_k}(\theta^{(k-1)}) = \arg \min_{\theta \in \mathcal{C}_k} \frac{1}{2} \|\theta^{(k-1)} - \theta\|_2^2, \quad (4.4)$$

which can be written as an unconstrained problem in terms of the set indicator functions in (2.3) as:

$$\begin{aligned}\theta^{(k)} &= \arg \min_{\theta} \frac{1}{2} \|\theta^{(k-1)} - \theta\|_2^2 \\ &+ \sum_{(l, j) \mid y_{k_l, j} = 1} \iota_{\mathcal{S}_{k_l, j}^+}(\theta) + \sum_{(l, j) \mid y_{k_l, j} = 0} \iota_{\mathcal{S}_{k_l, j}^-}(\theta),\end{aligned}\quad (4.5)$$

where $\mathcal{S}_{k_l, j}^+ = \{\theta \mid d_{k_l, j}^f(\theta) \leq \varepsilon^+\}$, $\mathcal{S}_{k_l, j}^- = \{\theta \mid d_{k_l, j}^f(\theta) \geq \varepsilon^-\}$ and $(l, j) \in \mathcal{Y} \times \mathcal{N}$. If the set indicator functions are to be

approximated for small λ as $\hat{\iota}_{\mathcal{S}_{k_l, j}^\mp}(\theta) = \frac{1}{\lambda} [\mp (d_{k_l, j}^f(\theta) - \varepsilon^\mp)]_+$, as in Eqn. (2.4), the problem becomes after scaling with λ :

$$\begin{aligned}\theta^{(k)} &= \arg \min_{\theta} \sum_{k_l, j} (y_{k_l, j} [d_{k_l, j}^f(\theta) - \varepsilon^+]_+ \\ &+ (1 - y_{k_l, j}) [\varepsilon^- - d_{k_l, j}^f(\theta)]_+) + \frac{\lambda}{2} \|\theta^{(k-1)} - \theta\|_2^2\end{aligned}\quad (4.6)$$

where $(l, j) \in \mathcal{Y} \times \mathcal{N}$. The resultant minimization problem for a projection step is very similar to a typical DML formulation in Eqn. (3.1) with a margin-based loss term [48] defined in Eqn. (3.4). The two significant differences are the utilization of the particular class samples, $\{k_l\}_{l \in \mathcal{Y}}$, for the pairwise distance losses and the regularization term relating the alternating subproblems.

4.1. Solving for Parameters

To obtain the solution, θ^* , defined in Eqn. (4.3), one should cycle through the sets, $\{\mathcal{C}_k\}_k$, and perform projections until the convergence. The nature of the problem is non-convex. Therefore, exploiting diverse combinations of sets might improve the solution. We propose to perform projections by randomly selecting the class samples for the sets. This approach results in different feasible sets $\{\mathcal{C}_k\}_k$ for each cycle so that the procedure does not stick to the specific sets. Performing a projection involves a minimization problem. In this perspective, either convergence can be monitored to pass the next projection or the projection operator can be approximated by M iterations of training. The latter approach gives the flexibility to control the fitting of the parameters to the subproblems. We therefore propose to use M -step approximations of the projection operators as $\theta^{(k)} \approx \mathcal{P}_{\mathcal{C}_k}^{(M)}(\theta^{(k-1)})$ in our framework.

Proposed learning procedure for θ necessitates utilization of the particular class samples for the loss computation. To provide scalability, those particular class samples can also be sampled during batch construction. In this manner, another implication of the proposed framework becomes imposing constraint on the batch construction for the minimization. If we disregard the resultant loss formulation in Eq. (4.6) and consider only batch construction method of our framework, we can formulate minimization of any pairwise distance ranking based loss function within the proposed batch construction method as:

$$\theta^{(k)} = \arg \min_{\theta} \mathcal{L}^{(\cdot)}(\theta; \mathcal{T}_k) + \frac{\lambda}{2} \|\theta^{(k-1)} - \theta\|_2^2, \quad (4.7)$$

where (\cdot) can be any proper loss and \mathcal{T}_k is the tuple set of the sample indices defining \mathcal{C}_k . The proposed DML framework in its most general form is summarized in Algorithm 1. Alternating projections only introduces a constrained batch construction step to the standard optimization procedure.

4.2. Implications

Robustness. The parameters are obtained through solving then relating the subproblems rather than solving an

Algorithm 1 PROFS DML

randomly initialize parameters $\theta^{(0)}$, $k=0$
repeat
 sample $\mathcal{R}=\{k_l \mid y_{k_l}=l\}_{l=1}^L \sim \mathcal{N}$ class representatives
 // *i.e.* an example for each class
 repeat M times // $\theta^{(k+1)} \approx \mathcal{P}_{\mathcal{C}_k}^{(M)}(\theta^{(k)})$
 sample $\mathcal{R}'=\{i_r\}_r \sim \mathcal{R}$ a subset of the class reps.
 // possibly via hard class mining
 sample $\mathcal{E}=\{j_n\}_n \sim \mathcal{N}$ a batch of examples
 construct exemplar batch \mathcal{B} from \mathcal{R}' and \mathcal{E}
 // *e.g.* pairs (i_r, j_n)
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} (\frac{\lambda}{2} \|\theta^{(k)} - \theta\|_2^2 + \mathcal{L}^{(\cdot)}(\theta; \mathcal{B}))$
 $k \leftarrow k+1$, $\theta^{(k)} \leftarrow \theta$
until convergence

overall problem at once. The subproblems are relatively easier problems owing to the relaxed proximity constraints to be considered. Thus, better solutions to the subproblem specific proximity constraints can be obtained. Nevertheless the solutions are expected to be localized, the regularization term entangles the solutions of the alternating subproblems for a better holistic solution. This nature makes our approach more robust to the mini-batch updates.

Class representatives. The subproblems can be seen as learning representations so that the particular class samples become the class representatives, akin to learning class representative vectors for a softmax classifier. Therefore, choosing the relaxed feasible sets as in Eqn. (4.2) implicitly integrates classification framework to the DML problem in our formulation. In this point of view, our formulation is aligned with the studies [23, 29, 30, 40, 54] supporting the superiority of the softmax incorporation into DML.

Hard negative class mining. The nature of the proposed DML framework allows efficient HNCM batch construction. As θ is projected onto the feasible sets, the set-specific class samples can be considered as the class representatives. Thus, using those representatives, approximate global mining can be efficiently performed. To this end, we store the embedding representations of the class representatives and update those representations as we sample the corresponding representatives for the batch construction. In this way, we can perform online HNCM with $\mathcal{O}(L)$ complexity.

Relation to linear metric learning with convex optimization. Convex optimization formulation of linear metric learning in [49] involves similar alternating projections onto feasible sets to perform projected gradient ascent. That approach only considers the entire constraint set induced by the positive pairs, whereas we consider both negative and positive pair distances jointly in the relaxed feasible sets.

Relation to tuplet and proxy-based losses. Minimization of tuplet losses [36, 52] can be considered as performing

our method by using $M=1$ step approximation of the projection operator. Similarly, the minimization of proxy-based losses [7, 23, 29, 30] can be considered as a single projection operation. Therefore, tuplet and proxy-based loss formulations are to be the two extreme cases of our framework.

Regularization. Our formulation suggests a regularization term to entangle the solutions of the different subproblems. The gradient update of a typical DML framework exploiting deliberately constructed mini-batches during the loss minimization can be considered as performing projections with $M=1$ step approximation. The proposed regularization term can be introduced to the loss to improve generalization by entangling the updates of the different batches.

5. Experimental Work

We examine the effectiveness of the proposed DML framework through evaluation on the three widely-used benchmark datasets for the image retrieval and clustering tasks. We perform ablation study on the effect of M -step approximation of the projection operation and the regularization term. Throughout the section, we use PROFS to refer our framework.

5.1. Benchmark Datasets and Evaluation Metrics

We obtain results by utilizing three public benchmark datasets. The conventional protocol of splitting training and testing sets for a zero-shot setting [25] is followed for all datasets. Hence, no image is in the intersection of the training and the test sets. Stanford Online Products (SOP) [25] dataset has 22,634 classes with 120,053 product images. The first 11,318 classes (59,551 images) are split for training and the other 11,316 (60,502 images) classes are used for testing. Cars196 [19] dataset contains 196 classes of cars with 16,185 images. The first 98 classes (8,054 images) are used for training and remaining 98 classes (8,131 images) are reserved for testing. CUB-200-2011 [43] dataset consists of 200 species of birds with 11,788 images. The first 100 species (5,864 images) are split for training, the rest of 100 species (5,924 images) are used for testing.

We follow the standard metric learning experimental protocol defined in [25] to evaluate the performance of the deep metric learning approaches for the retrieval and clustering tasks. We utilize normalized mutual information (NMI) and F_1 score to measure the quality of the clustering task which is performed by conventional k -means clustering algorithm. In order to evaluate clustering performance, normalized mutual information (NMI) and F_1 score are utilized. NMI computes the label agreement between predicted and groundtruth clustering assignments neglecting the permutations while F_1 measures harmonic mean of the precision and recall. Furthermore, we exploit binary Recall@K metric to evaluate the performance of the retrieval task. Recall@K for a test query is 1 if at least one sample from the same class of the query

is in the K nearest neighborhood of the query. The average of the Recall@K for the test queries gives the Recall@K performance on the dataset. We refer [25] for the detailed information related to these evaluation metrics.

5.2. Training Setup

We use Tensorflow [1] and PyTorch [28] deep learning libraries throughout the experiments. Tensorflow is used for trainings on GoogLeNet V1 (Inception V1) [41] and PyTorch is used for trainings on GoogLeNet V2 (with batch normalization) [24] and ResNet-50 [13]. After the images are normalized and scaled to 256×256 , we perform 224×224 random crop and data augmentation by horizontal mirroring as pre-processing. For the embedding function, $f(\cdot; \theta)$, we utilize architectures until the output of the global average pooling layer with the parameters pretrained on ImageNet ILSVRC dataset [33]. After the pooling layer, we add a linear transformation layer (fully connected layer) to obtain the representation vectors of size 512. We fix the embedding size of the samples at 512 throughout experiments, since it is shown in [25] that the embedding size does not have a key role on comparing performances of the deep metric learning loss functions. The parameters of the linear transformation layer are randomly initialized and are learned by using 10 times larger learning rate than the pretrained parameters for the sake of fast convergence. For the hyper-parameters, our framework introduces 2 additional parameters: λ for regularization and M to approximate projection operation. We set the regularization term as the result of the projection based formulation, λ , to a small reasonable value, 10^{-3} . The number of projections steps, M , is determined according to the findings of the ablation study which is presented in subsection 5.4. For the other hyper-parameters coming from adaptation of the baseline framework (e.g. margins, number of positive samples etc.), we follow the settings in the corresponding baseline work. For the optimization procedure, we select the base learning rate as 10^{-4} for SOP dataset whereas we utilize 10^{-5} learning rate to train CUB-200-2011 and Cars196, since they tend to meet over-fitting problem due to the limited dataset size. We exploit Adam [18] optimizer for mini-batch gradient descent with a mini-batch size is 128 and default moment parameters, $\beta_1=.9$ and $\beta_2=.99$. Finally, since the convergence rate of each method is different, we train all the approaches for 100 epochs and post the performance at their best epoch as in [48] instead of following the conventional procedure [25] reporting performance of DML approaches after a certain number of training iterations.

5.3. Baseline Methods and PROFS Framework Adaptation

We apply proposed PROFS framework with and without HNCM on the contrastive [11], triplet [35], lifted structured [25], N -pair [36], angular [44], margin-based (Mar-

gin) [48], multi-similarity (MS) [45] and SoftTriple [30] loss functions in order to directly compare with the state-of-the-art methods. The comparison with the contrastive, triplet and Margin losses is important to examine the effectiveness of our original formulation, while the comparison with the other losses is to show that the formulation can be extended to the other loss functions by exploiting the proposed batch construction and regularization. Furthermore, we compare proxy-based [23,30] loss functions with PROFS owing to its relation to our formulation.

To evaluate the approaches in the same basis, we retrain all the aforementioned loss functions excluding SoftTriple by exploiting the same GoogLeNet V1 architecture with the default hyper-parameters used in the original works except for the mini-batch and the embedding sizes as explained in the subsection 5.2 for a fair comparison. Moreover, we integrate our framework to Margin and SoftTriple losses in their own architectures (ResNet-50 and GoogLeNet V2, respectively) in order to compare our framework with state-of-the-art.

To adapt the lifted structured and MS loss to PROFS framework, the loss is slightly modified by ignoring the pairwise terms between the non-representative samples. For the adaptation of the sampling strategies, we trained Margin loss with the distance weighted sampling in its own architecture. However, we were unable to acquire good results in GoogLeNet V1, since the distance weighted sampling is very sensitive to its parameters and we could not determine well-performing hyper-parameters. On the other hand, due to its similarity with the contrastive loss, we exploit the same hard mining strategy inspired from [53] as in the contrastive loss for the mini-batch sampling method of the Margin loss. It should be noted that we sample one negative pair for each positive pair as in [48] for the contrastive, triplet and Margin loss functions. Such an approach provides balance to the number of positive and negatives pairs. In the conventional hard mining, the number of hard negative pairs should match the number of positive pairs for the contrastive, triplet and Margin loss functions. Thus, in PROFS framework, each class representative should have the same number of its corresponding positive pairs and hard negative pairs to obtain an exemplar set consistent with hard mining without violating the batch construction constraint of PROFS. No adaptation is performed for the N -pair and angular loss, since their formulation is consistent with PROFS. Finally, for SoftTriple loss, the adaptation is not straightforward, since this framework is inherently relaxed formulation of DML problem. On the other hand, in that framework, the loss terms are determined by assigning samples to the trainable cluster centers. At different iterations, the assignment may differ. To this end, we consider each iteration as a subproblem and integrate our framework by exploiting the regularization term to entangle the updates of the iterations. For HNCM, we utilize cluster

centers as representatives in SoftTriple.

5.4. Ablation Study

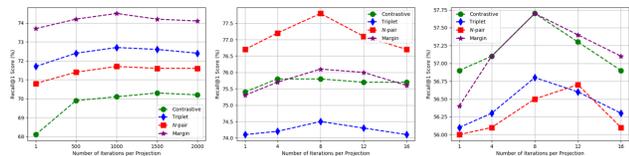


Figure 2: The retrieval performance of the PROFS on the Stanford Online Products [25] (left), the CARS196 [19] (middle) and the CUB-200-2011 [43] (right) for different number of iterations used to compute projection.

We present to ablation studies on the approximation of the projection operator and the effectiveness of the regularization term. For the computational complexity, the proposed PROFS without HNCM framework does not bring any additional computational cost. On the other hand, HNCM only introduces scan of class representatives during batch construction and has 0.008% (CUB&CARS) and 0.9% (SOP) increase in overall computation time. Moreover, we have not observed any significant increase or decrease in the convergence rate owing to re-utilization of the class samples.

5.4.1 Number of Projection Iterations

We perform ablation study to determine approximately how many training iterations, M , required to assure an acceptable approximation of the projection operation. Due to the non-convex nature of the problem, enforcing convergence for each projection might lead to ill-configured parameters that the proceeding projections cannot recover.

To examine effect of M , we apply PROFS framework without hard negative class mining to the contrastive, triplet, N -pair and margin loss functions on SOP, Cars196 and CUB-200-2011. We use 2 positive images, $I=2$, per class in mini-batches for each loss function to have consistency among the loss terms. The retrieval performance curves for varying M values are plotted Fig. 2. The retrieval performance of each loss function improves on all three datasets as M increases up to certain values. Exceeding that certain number of iterations leads to possible over-fitting to the subproblems and the performance drops accordingly. Once the performance curves of the different datasets are compared, a heuristic relation between M and the number of classes in a dataset can be deduced. Hence, instead of fine-tuning M parameter for each problem, we derive M as $M = \lceil \rho / p(y_i) \rceil$ where $\lceil \cdot \rceil$ is the ceiling function, $p(y_i)$ denotes the probability of observing a class i in a mini-batch of size B , and ρ is how many times a class representation is to be used during the minimization. We write $p(y_i) = B / I L$ for a batch containing I samples from

each class it contains among L many classes. Considering the plots in Fig. 2, we set $\rho=6$ throughout the experiments.

5.4.2 Effectiveness of the Regularization Term

We analyze the effect of $\frac{\lambda}{2} \|\theta^{(curr.)} - \theta\|_2^2$ term by sweeping the λ parameter for training settings with Margin loss in GoogLeNet V1. Theoretically λ should be as small as possible (Eqn. (4.6)), thus large values result in performance drop owing to over-regularization. On the other hand, $\lambda=0$ implies no regularization for the projection and leads to over-fitting to the subproblems. According to results in Fig. 3, the performances are similar for the values within $[10^{-2}, 10^{-4}]$. Though we pick 10^{-3} throughout the experiments, one can choose to adaptively reduce λ as the parameters converge. In this way, the relative significance of the loss function for the constraints can be preserved as the parameters converge.



Figure 3: Effect of the $\|\theta^{(curr.)} - \theta\|_2^2$ term

5.5. Quantitative Results

The quantitative results of the proposed PROFS framework trained in GoogLeNet V1 with and without HNCM for the clustering and retrieval tasks on SOP, CARS196 and CUB-200-2011 datasets are provided in Table 1 together with the baseline methods indicated in 5.3 for comparison. It can be observed that PROFS framework consistently outperforms the associated baseline methods. Compared with the original loss functions, the proposed PROFS framework boosts their performance on each dataset for the clustering and retrieval tasks by up to 3.2%, 11.7% and 9.5% points on NMI, F_1 and $R@1$ metrics, respectively. Additionally, the proposed PROFS framework with the contrastive and triplet loss functions produce competitive results in comparison with superior Margin loss. It supports that considering relaxed feasible sets iteratively is beneficial over solving the entire problem at once even with the basic loss functions. This result is important to support the motivation of our formulation. Furthermore, performance improvements on the loss functions which does not directly fit in our formulation in Eq. (4.6) show that the batch construction and regularization implication of the proposed formulation can be generalized to the pairwise distance based loss functions.

The performances of the proposed framework integrated to Margin and SoftTriple architectures are provided in Table 2. In most settings, the proposed framework improves state-of-the-art. Especially, improving SoftTriple is of great importance, since that performance increase comes from

Table 1: Comparison with the existing methods for the clustering and the retrieval tasks on SOP [25], CARS196 [19] and CUB-200-2011 [43] datasets. Red: the overall best. Blue: the overall second best. Bold: the loss term specific best.

Method	Stanford Online Products						CARS196						CUB-200-2011					
	NMI	F ₁	R@1	R@10	R@100	R@1000	NMI	F ₁	R@1	R@2	R@4	R@8	NMI	F ₁	R@1	R@2	R@4	R@8
Contrastive-Hard	89.7	34.5	67.9	83.8	93.2	97.9	66.0	36.6	75.8	84.5	90.1	94.1	63.4	31.8	56.7	68.4	78.8	86.3
C-PROFS	90.3	36.4	70.1	85.4	94.1	98.2	67.8	39.0	76.0	84.8	90.1	94.6	64.1	32.0	57.7	69.0	78.9	86.9
C-PROFS-HNCM	91.0	41.2	74.5	87.9	94.9	98.3	66.9	38.1	77.0	85.3	90.9	94.4	64.6	33.1	57.9	69.0	79.4	87.0
Triplet-Semi	87.1	23.4	57.0	75.0	88.2	96.4	62.4	30.0	71.2	80.7	87.6	92.5	60.9	27.8	55.5	67.8	78.2	86.5
Triplet-Hard	88.1	30.6	65.4	81.4	91.7	97.4	62.5	30.6	71.4	81.1	87.5	92.7	63.1	30.7	56.8	68.7	78.6	86.5
T-PROFS	90.7	37.7	72.4	86.3	94.1	98.1	64.7	33.4	74.5	82.8	89.1	93.5	63.5	31.3	57.9	68.9	78.7	86.7
T-PROFS-HNCM	91.3	42.3	74.9	87.5	94.4	98.1	64.9	33.7	74.5	83.5	89.4	93.7	64.5	32.3	58.1	69.0	78.9	86.8
Lifted	88.9	31.4	66.7	83.2	91.7	97.4	60.1	27.7	67.5	77.3	84.9	90.7	60.6	26.9	53.5	65.3	75.4	84.8
L-PROFS	89.5	33.8	68.1	84.0	92.2	97.6	61.2	27.9	68.4	78.1	85.6	91.1	61.4	28.1	54.5	66.1	76.2	85.1
L-PROFS-HNCM	89.9	35.1	69.3	85.1	93.6	98.1	61.5	30.0	70.7	79.6	86.2	91.5	62.0	29.5	54.6	66.1	76.7	85.5
N-pair	89.9	35.7	70.8	86.0	94.0	98.1	67.4	38.2	76.7	84.8	91.0	95.0	64.6	33.0	56.0	68.9	79.3	87.4
N-PROFS	90.3	37.3	71.7	86.6	94.0	98.2	68.1	38.3	77.8	85.9	91.6	95.2	64.9	33.6	56.5	69.0	79.3	87.7
N-PROFS-HNCM	90.5	37.5	71.9	86.6	94.1	98.2	68.6	39.9	77.6	86.2	91.7	95.2	65.2	33.7	56.7	68.8	79.7	87.5
Angular	90.0	36.1	72.5	86.6	93.6	97.6	66.0	35.9	77.4	85.3	91.0	94.7	61.9	29.4	54.5	66.4	76.8	84.9
A-PROFS	90.1	37.2	73.0	86.8	93.7	97.7	66.3	36.8	77.5	85.6	91.2	94.7	63.0	31.7	54.6	66.7	76.9	85.8
A-PROFS-HNCM	90.4	38.4	73.7	86.9	93.8	97.7	66.3	37.9	77.9	85.6	91.4	94.8	64.6	33.4	55.8	68.1	78.8	87.3
Margin-Hard	90.6	38.1	73.9	87.7	94.8	98.2	64.2	34.6	75.1	83.7	89.5	93.8	64.0	30.9	55.3	67.2	77.9	87.5
M-PROFS	91.3	42.7	74.5	88.0	95.0	98.2	64.6	35.1	76.0	84.3	89.5	93.8	64.3	32.7	57.7	69.5	79.5	87.5
M-PROFS-HNCM	91.4	43.2	76.3	88.8	95.0	98.3	66.8	37.3	77.0	85.1	90.8	94.6	64.8	32.4	58.5	69.6	79.7	87.6
MS	90.5	38.2	71.3	86.4	94.2	98.1	65.3	35.5	76.2	84.2	89.9	94.0	64.7	33.8	56.4	69.1	79.3	87.5
MS-PROFS	90.8	39.5	72.6	87.2	94.3	98.2	66.6	37.1	76.6	84.9	90.6	94.4	65.0	34.3	57.5	69.4	79.6	87.5
MS-PROFS-HNCM	91.0	42.9	74.6	87.7	94.6	98.2	68.4	39.2	77.8	86.0	91.8	95.3	65.2	34.4	58.1	69.3	79.6	87.7

Table 2: Comparison with state-of-the-art methods on SOP [25], CARS196 [19] and CUB-200-2011 [43] datasets. Red: the overall best. Blue: the overall second best. Bold: the loss term specific best.

Method	Stanford Online Products				CARS196				CUB-200-2011			
	R@1	R@10	R@100	R@1000	R@1	R@2	R@4	R@8	R@1	R@2	R@4	R@8
Proxy-NCA [23]	73.7	-	-	-	73.2	82.4	86.4	88.7	49.2	61.9	67.9	72.4
Clustering [38]	67.0	83.7	93.2	-	58.1	70.6	80.3	87.8	48.2	61.4	71.8	59.2
HDC [53]	69.5	84.4	92.8	97.7	73.7	83.2	89.5	93.8	53.6	65.7	77.0	85.6
HTL [10]	74.8	88.3	94.8	98.4	81.4	88.0	92.7	95.7	57.1	68.8	78.7	86.5
MS [45]	78.2	90.5	96.0	98.7	84.1	90.4	94.0	96.5	65.7	77.0	86.3	91.2
TML [52]	78.0	91.2	96.7	99.0	86.3	92.3	95.4	97.3	62.5	73.9	83.0	89.4
Margin [48]	72.7	86.2	93.8	98.0	79.6	86.5	91.9	95.1	63.6	74.4	83.1	90.0
M-PROFS	76.5	89.0	95.2	98.5	81.1	88.1	92.7	95.8	64.9	75.8	84.2	90.4
M-PROFS-HNCM	76.9	89.5	95.3	98.5	81.3	88.0	93.0	95.8	64.1	75.0	84.2	90.3
SoftTriple [30]	78.3	90.3	95.9	-	84.5	90.7	94.5	96.9	65.4	76.4	84.5	90.4
SoftTriple-PROFS	78.6	91.4	96.0	99.1	86.1	91.9	94.7	97.4	66.0	76.8	85.0	90.7
SoftTriple-PROFS-HNCM	78.7	91.7	96.8	99.2	86.3	92.5	95.0	97.5	65.7	76.2	84.5	90.6

the regularization purely. This result supports the effectiveness of the regularization term to be used to improve generalization. Lastly, utilizing HNCM is more efficient on SOP dataset, since it has almost 100 times larger number of classes than CARS196 and CUB-200-2011 datasets.

6. Conclusion

We have presented a novel DML formulation based on alternating projections onto the feasible sets which impose relaxed proximity constraints. The resultant framework in-

troduces a simple, yet effective, batch construction scheme and a regularizer to improve the generalization. Notably, the proposed framework is applicable with the pairwise distance based state-of-the-art DML loss functions without introducing any additional computational cost. Extensive evaluations on the benchmark datasets show that the performances of the several state-of-the-art loss functions are improved by the proposed framework.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- [2] Heinz H Bauschke and Adrian S Lewis. Dykstras algorithm with bregman projections: A convergence proof. *Optimization*, 48(4):409–427, 2000.
- [3] Lev M Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 7(3):200–217, 1967.
- [4] Fatih Cakir, Kun He, Xide Xia, Brian Kulis, and Stan Sclaroff. Deep metric learning to rank. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [5] Guangyi Chen, Tianren Zhang, Jiwen Lu, and Jie Zhou. Deep meta metric learning. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [6] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 403–412, 2017.
- [7] Thanh-Toan Do, Toan Tran, Ian Reid, Vijay Kumar, Tuan Hoang, and Gustavo Carneiro. A theoretically sound upper bound on the triplet loss for improving the efficiency of deep distance metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [8] Yueqi Duan, Jiwen Lu, Jianjiang Feng, and Jie Zhou. Deep localized metric learning. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2644–2656, 2017.
- [9] Yueqi Duan, Wenzhao Zheng, Xudong Lin, Jiwen Lu, and Jie Zhou. Deep adversarial metric learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2780–2789, 2018.
- [10] Weifeng Ge. Deep metric learning with hierarchical triplet loss. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–285, 2018.
- [11] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *null*, pages 1735–1742. IEEE, 2006.
- [12] Ben Harwood, BG Kumar, Gustavo Carneiro, Ian Reid, Tom Drummond, et al. Smart mining for deep metric learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2821–2829, 2017.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- [14] Junlin Hu, Jiwen Lu, and Yap-Peng Tan. Discriminative deep metric learning for face verification in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1875–1882, 2014.
- [15] Chen Huang, Chen Change Loy, and Xiaoou Tang. Local similarity-aware deep feature embedding. In *Advances in neural information processing systems*, pages 1262–1270, 2016.
- [16] Pierre Jacob, David Picard, Aymeric Histace, and Edouard Klein. Metric learning with horde: High-order regularizer for deep embeddings. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [17] Wonsik Kim, Bhavya Goyal, Kunal Chawla, Jungmin Lee, and Keunjoo Kwon. Attention-based ensemble for deep metric learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 736–751, 2018.
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] Andreas Krause and Daniel Golovin. Submodular function maximization., 2014.
- [20] Marc T Law, Nicolas Thome, and Matthieu Cord. Quadruplet-wise image similarity learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 249–256, 2013.
- [21] Marc T Law, Raquel Urtasun, and Richard S Zemel. Deep spectral clustering learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1985–1994. JMLR. org, 2017.
- [22] Xudong Lin, Yueqi Duan, Qiyuan Dong, Jiwen Lu, and Jie Zhou. Deep variational metric learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 689–704, 2018.
- [23] Yair Movshovitz-Attias, Alexander Toshev, Thomas K Leung, Sergey Ioffe, and Saurabh Singh. No fuss distance metric learning using proxies. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 360–368, 2017.
- [24] Batch Normalization. Accelerating deep network training by reducing internal covariate shift. *CoRR*.–2015.–Vol. abs/1502.03167.–URL: <http://arxiv.org/abs/1502.03167>, 2015.
- [25] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4004–4012, 2016.
- [26] Michael Opitz, Georg Waltner, Horst Possegger, and Horst Bischof. Bier-boosting independent embeddings robustly. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5189–5198, 2017.
- [27] CH Pang. Nonconvex set intersection problems: From projection methods to the newton method for super-regular sets. *arXiv preprint arXiv:1506.08246*, 2015.
- [28] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [29] Michaël Perrot and Amaury Habrard. Regressive virtual metric learning. In *Advances in Neural Information Processing Systems*, pages 1810–1818, 2015.
- [30] Qi Qian, Lei Shang, Baigui Sun, Juhua Hu, Hao Li, and Rong Jin. Softtriple loss: Deep metric learning without triplet sam-

- pling. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [31] Oren Rippel, Manohar Paluri, Piotr Dollar, and Lubomir Bourdev. Metric learning with adaptive density discrimination. *International Conference on Learning Representations (ICLR)*, 2016.
- [32] Karsten Roth, Biagio Brattoli, and Bjorn Ommer. Mic: Mining interclass characteristics for improved metric learning. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [33] 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*, 115(3):211–252, 2015.
- [34] Artsiom Sanakoyeu, Vadim Tschernezki, Uta Buchler, and Bjorn Ommer. Divide and conquer the embedding space for metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [35] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [36] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in Neural Information Processing Systems*, pages 1857–1865, 2016.
- [37] Justin Solomon, Fernando De Goes, Gabriel Peyré, Marco Cuturi, Adrian Butscher, Andy Nguyen, Tao Du, and Leonidas Guibas. Convolutional wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Transactions on Graphics (TOG)*, 34(4):66, 2015.
- [38] Hyun Oh Song, Stefanie Jegelka, Vivek Rathod, and Kevin Murphy. Deep metric learning via facility location. In *Computer Vision and Pattern Recognition (CVPR)*, volume 8, 2017.
- [39] Yumin Suh, Bohyung Han, Wonsik Kim, and Kyoung Mu Lee. Stochastic class-based hard example mining for deep metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [40] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Advances in neural information processing systems*, pages 1988–1996, 2014.
- [41] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [42] Evgeniya Ustinova and Victor Lempitsky. Learning deep embeddings with histogram loss. In *Advances in Neural Information Processing Systems*, pages 4170–4178, 2016.
- [43] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [44] Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. Deep metric learning with angular loss. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2612–2620. IEEE, 2017.
- [45] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R. Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [46] Xinshao Wang, Yang Hua, Elyor Kodirov, Guosheng Hu, Romain Garnier, and Neil M. Robertson. Ranked list loss for deep metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [47] Kilian Q Weinberger, John Blitzer, and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. In *Advances in neural information processing systems*, pages 1473–1480, 2006.
- [48] Chao-Yuan Wu, R Manmatha, Alexander J Smola, and Philipp Krahenbuhl. Sampling matters in deep embedding learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2840–2848, 2017.
- [49] Eric P Xing, Michael I Jordan, Stuart J Russell, and Andrew Y Ng. Distance metric learning with application to clustering with side-information. In *Advances in neural information processing systems*, pages 521–528, 2003.
- [50] Hong Xuan, Richard Souvenir, and Robert Pless. Deep randomized ensembles for metric learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 723–734, 2018.
- [51] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Deep metric learning for person re-identification. In *2014 22nd International Conference on Pattern Recognition*, pages 34–39. IEEE, 2014.
- [52] Baosheng Yu and Dacheng Tao. Deep metric learning with triplet margin loss. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [53] Yuhui Yuan, Kuiyuan Yang, and Chao Zhang. Hard-aware deeply cascaded embedding. In *Proceedings of the IEEE international conference on computer vision*, pages 814–823, 2017.
- [54] Xiaofan Zhang, Feng Zhou, Yuanqing Lin, and Shaoting Zhang. Embedding label structures for fine-grained feature representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1114–1123, 2016.
- [55] Yiru Zhao, Zhongming Jin, Guo-jun Qi, Hongtao Lu, and Xian-sheng Hua. An adversarial approach to hard triplet generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 501–517, 2018.
- [56] Wenzhao Zheng, Zhaodong Chen, Jiwen Lu, and Jie Zhou. Hardness-aware deep metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.