

# Deep Metric Learning via Facility Location

Hyun Oh Song<sup>1</sup>, Stefanie Jegelka<sup>2</sup>, Vivek Rathod<sup>1</sup>, and Kevin Murphy<sup>1</sup>

<sup>1</sup>Google Research, <sup>2</sup>MIT

<sup>1</sup>{hyunsong, rathodv, kpmurphy}@google.com, <sup>2</sup>stefje@csail.mit.edu

## Abstract

*Learning image similarity metrics in an end-to-end fashion with deep networks has demonstrated excellent results on tasks such as clustering and retrieval. However, current methods, all focus on a very local view of the data. In this paper, we propose a new metric learning scheme, based on structured prediction, that is aware of the global structure of the embedding space, and which is designed to optimize a clustering quality metric (NMI). We show state of the art performance on standard datasets, such as CUB200-2011 [37], Cars196 [18], and Stanford online products [30] on NMI and R@K evaluation metrics.*

## 1. Introduction

Learning to measure the similarity among arbitrary groups of data is of great practical importance, and can be used for a variety of tasks such as feature based retrieval [30], clustering [10], near duplicate detection [41], verification [3, 4], feature matching [6], domain adaptation [27], video based weakly supervised learning [38], etc. Furthermore, metric learning can be used for challenging extreme classification settings [23, 5, 40], where the number of classes is very large and the number of examples per class becomes scarce. For example, [2] uses this approach to perform product search with 10M images, and [25] shows superhuman performance on face verification with 260M images of 8M distinct identities. In this setting, any direct classification or regression methods become impractical due to the prohibitively large size of the label set.

Currently, the best approaches to metric learning employ state of art neural networks [19, 31, 28, 9], which are trained to produce an embedding of each input vector so that a certain loss, related to distances of the points, is minimized. However, most current methods, such as [25, 2, 30, 29], are very myopic in the sense that the loss is defined in terms of pairs or triplets inside the training mini-batch. These methods don't take the global structure of the embedding space into consideration, which can result in reduced clustering

and retrieval performance.

Furthermore, most of the current methods [25, 2, 30, 29] in deep metric learning require a separate data preparation stage where the training data has to be first prepared in pairs [8, 2], triplets [39, 25], or n-pair tuples [29] format. This procedure has very expensive time and space cost as it often requires duplicating the training data and needs to repeatedly access the disk.

In this paper, we propose a novel learning framework which encourages the network to learn an embedding function that directly optimizes a clustering quality metric (We use the normalized mutual information or NMI metric [21] to measure clustering quality, but other metrics could be used instead.) and doesn't require the training data to be preprocessed in rigid paired format. Our approach uses a structured prediction framework [35, 14] to ensure that the score of the ground truth clustering assignment is higher than the score of any other clustering assignment. Following the evaluation protocol in [30], we report state of the art results on CUB200-2011 [37], Cars196 [18], and Stanford online products [30] datasets for clustering and retrieval tasks.

## 2. Related work

The seminal work in deep metric learning is to train a siamese network with contrastive loss [8, 4] where the task is to minimize the pairwise distance between a pair of examples with the same class labels, and to push the pairwise distance between a pair of examples with different class labels at least greater than some fixed margin.

One downside of this approach is that it focuses on absolute distances, whereas for most tasks, relative distances matter more. For this reason, more recent methods have proposed different loss functions. We give a brief review of these below, and we compare our method to them experimentally in Section 4.

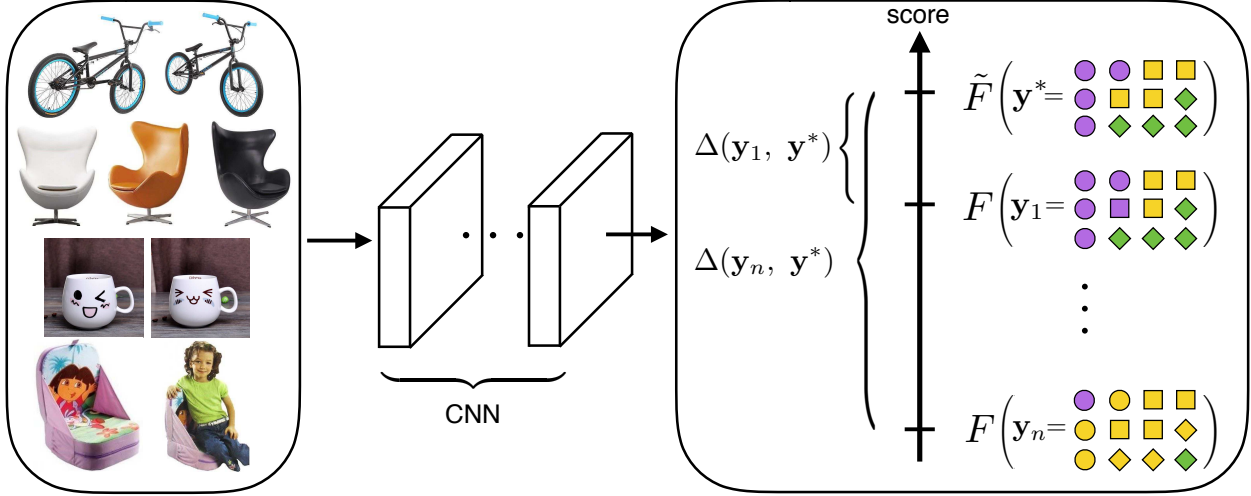


Figure 1. Overview of the proposed framework. The network first computes the embedding vectors for each images in the batch and learns to rank the clustering score  $\tilde{F}$  for the ground truth clustering assignment higher than the clustering score  $F$  for any other assignment at least by the structured margin  $\Delta(\mathbf{y}, \mathbf{y}^*)$ .

## 2.1. Triplet learning with semi-hard negative mining

One improvement over contrastive loss is to use triplet loss [39, 26]. This first constructs a set of triplets, where each triplet has an *anchor*, a *positive*, and a *negative* example, where the anchor and the positive have the same class labels and the negative has the different class label. It then tries to move the anchor and positive closer than the distance between the anchor and the negative with some fixed margin. More precisely, it minimizes the following loss:

$$\ell(X, y) = \frac{1}{|\mathcal{T}|} \sum_{(i,j,k) \in \mathcal{T}} [D_{i,j}^2 + \alpha - D_{i,k}^2]_+ \quad (1)$$

where  $\mathcal{T}$  is the set of triples,  $D_{i,j} = \|f(X_i) - f(X_j)\|_2$  is the Euclidean distance in embedding space, the operator  $[\cdot]_+$  denotes the hinge function which takes the positive component of the argument, and  $\alpha$  denotes a fixed margin constant.

In practice, the performance of these methods depends highly on the triplet sampling strategy. FaceNet [25] suggested the following online hard negative mining strategy. The idea is to construct triplets by associating with each positive pair in the minibatch a “semi-hard” negative example. This is an example which is further away from the anchor  $i$  than the positive exemplar  $j$  is, but still hard because the distance is close to the  $i - j$  distance. More precisely, it minimizes

$$\ell(X, \mathbf{y}) = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} [D_{i,j}^2 + \alpha - D_{i,k^*(i,j)}^2]_+$$

where

$$k^*(i, j) = \arg \min_{k: \mathbf{y}[k] \neq \mathbf{y}[i]} D_{i,k}^2 \text{ s.t. } D_{i,k}^2 > D_{i,j}^2$$

and  $\mathcal{P}$  is the set of pairs with the same class label. If there is no such negative example satisfying the constraint, we just pick the furthest negative example in the minibatch, as follows:

$$k^*(i, j) = \arg \max_{k: \mathbf{y}[k] \neq \mathbf{y}[i]} D_{i,k}^2$$

In order to get good results, the FaceNet paper had to use very large minibatches (1800 images), to ensure they picked enough hard negatives. This makes it hard to train the model on a GPU due to the GPU memory constraint. Below we describe some other losses which are easier to minimize using small minibatches.

## 2.2. Lifted structured embedding

Song *et al.* [30] proposed lifted structured embedding where each positive pair compares the distances against all the negative pairs weighted by the margin constraint violation. The idea is to have a differentiable smooth loss which incorporates the online hard negative mining functionality using the *log-sum-exp* formulation.

$$\ell(X, \mathbf{y}) = \frac{1}{2|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \left[ \log \left( \sum_{(i,k) \in \mathcal{N}} \exp \{ \alpha - D_{i,k} \} + \sum_{(j,l) \in \mathcal{N}} \exp \{ \alpha - D_{j,l} \} \right) + D_{i,j} \right]_+^2, \quad (2)$$

where  $\mathcal{N}$  denotes the set of pairs of examples with different class labels.

### 2.3. N-pairs embedding

Recently, Sohn *et al.* [29] proposed N-pairs loss which enforces softmax cross-entropy loss among the pairwise similarity values in the batch.

$$\ell(X, \mathbf{y}) = \frac{-1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \log \frac{\exp\{S_{i,j}\}}{\exp\{S_{i,j}\} + \sum_{k: \mathbf{y}[k] \neq \mathbf{y}[i]} \exp\{S_{i,k}\}} + \frac{\lambda}{m} \sum_i \|f(X_i)\|_2, \quad (3)$$

where  $S_{i,j}$  means the feature dot product between two data points in the embedding space;  $S_{i,j} = f(X_i)^\top f(X_j)$ ,  $m$  is the number of the data, and  $\lambda$  is the regularization constant for the  $\ell_2$  regularizer on the embedding vectors.

### 2.4. Other related work

In addition to the above work on metric learning, there has been some recent work on learning to cluster with deep networks. Hershey *et al.* [10] uses Frobenius norm on the residual between the binary ground truth and the estimated pairwise affinity matrix; they apply this to speech spectrogram signal clustering. However, using the Frobenius norm directly is suboptimal, since it ignores the fact that the affinity matrix is positive definite.

To overcome this, matrix backpropagation [12] first projects the true and predicted affinity matrix to a metric space where Euclidean distance is appropriate. Then it applies this to normalized cuts for unsupervised image segmentation. However, this approach requires computing the eigenvalue decomposition of the data matrix, which has cubic time complexity in the number of data and is thus not very practical for large problems.

## 3. Methods

One of the key attributes which the recent state of the art deep learning methods in Section 2 have in common is that they are all local metric learning methods. Figure 2 illustrates a case where this can fail. In particular, whenever a positive pair (such as the two purple bold dots connected by the blue edge) is separated by examples from other classes, the attractive gradient signal from the positive pair gets outweighed by the repulsive gradient signal from the negative data points (yellow and green data points connected with the red edges). This failure can lead to groups of examples with the same class labels being separated into partitions in the embedding space that are far apart from each other. This

can lead to degradation in the clustering and nearest neighbor based retrieval performance. For example, suppose we incorrectly created 4 clusters in Figure 2. If we asked for the 12 nearest neighbors of one of purple points, we would retrieve points belonging to other classes.

To overcome this problem, we propose a method that learns to embed points so as to minimize a clustering loss, as we describe below.

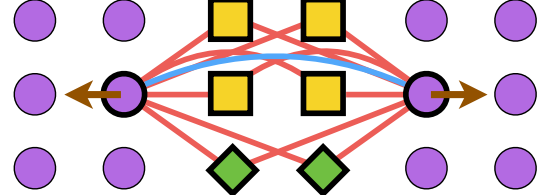


Figure 2. Example failure mode for local metric learning methods. Whenever a positive pair (linked with the blue edge) is separated by negative examples, the gradient signal from the positive pair (attraction) gets outweighed by the negative pairs (repulsion). Illustration shows the failure case for 2D embedding where the purple clusters can't be merged into one cluster.

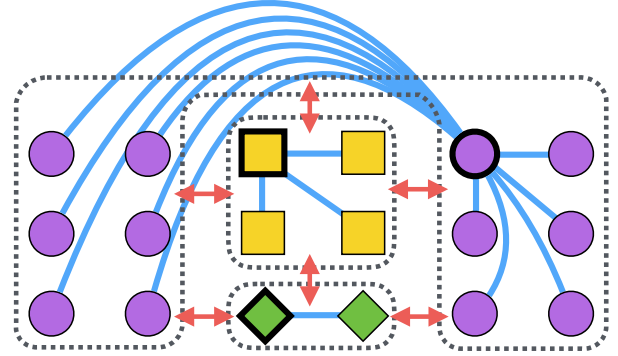


Figure 3. Proposed clustering loss for the same embedding layout in figure 2. Nodes highlighted in bold are the cluster medoids. The proposed method encourages small sum of distances within each cluster, while discouraging different clusters from getting close to each other.

### 3.1. Facility location problem

Suppose we have a set of inputs  $X_i$ , and an embedding function  $f(X_i; \Theta)$  that maps each input to a point in some  $K$  dimensional space. Now suppose we compress this set of points by mapping each example  $i$  to its nearest point from a chosen set of landmarks  $S \subseteq \mathcal{V}$ , where  $\mathcal{V} = \{1, \dots, |X|\}$  is the ground set. We can define the resulting function as follows:

$$F(X, S; \Theta) = - \sum_{i \in |X|} \min_{j \in S} \|f(X_i; \Theta) - f(X_j; \Theta)\|, \quad (4)$$

This is called the facility location function, and has been widely used in data summarization and clustering [20, 34].

The idea is that this function measures the sum of the travel distance for each customer in  $X$  to their respective nearest facility location in  $S$ . In terms of clustering, data points in  $S$  correspond to the cluster medoids, and the cluster assignment is based on the nearest medoid from each data point. Maximizing equation 4 with respect to subset  $S$  is NP-hard, but there is a well established worst case optimality bound of  $O(1 - \frac{1}{e})$  for the greedy solution of the problem via submodularity [17].

Below we show how to use the facility location problem as a subroutine for deep metric learning.

### 3.2. Structured facility location for deep metric learning

The oracle scoring function  $\tilde{F}$  measures the quality of the clustering given the ground truth clustering assignment  $\mathbf{y}^*$  and the embedding parameters  $\Theta$ :

$$\tilde{F}(X, \mathbf{y}^*; \Theta) = \sum_{k=1}^{|\mathcal{Y}|} \max_{j \in \{i: \mathbf{y}^*[i]=k\}} F(X_{\{i: \mathbf{y}^*[i]=k\}}, \{j\}; \Theta), \quad (5)$$

where  $\{i: \mathbf{y}^*[i] = k\}$  denotes the subset of the elements in  $\mathcal{V}$  with the ground truth label equal to  $k$ .

We would like the clustering score of the oracle clustering assignment to be greater than the score for the maximally violating clustering assignment. Hence we define the following structured loss function:

$$\ell(X, \mathbf{y}^*) = \left[ \underbrace{\max_{\substack{S \subset \mathcal{V} \\ |S|=|\mathcal{Y}|}} \left\{ F(X, S; \Theta) + \gamma \Delta(g(S), \mathbf{y}^*) \right\}}_{(*)} - \tilde{F}(X, \mathbf{y}^*; \Theta) \right]_+ \quad (6)$$

We will define the structured margin  $\Delta(\mathbf{y}, \mathbf{y}^*)$  below. The function  $y = g(S)$  maps the set of indices  $S$  to a set of cluster labels by assigning each data point to its nearest facility in  $S$ :

$$g(S)[i] = \arg \min_j \|f(X_i; \Theta) - f(X_{\{j\}}; \Theta)\| \quad (7)$$

Intuitively, the loss function in Equation 6 encourages the network to learn an embedding function  $f(\cdot; \Theta)$  such that the oracle clustering score  $\tilde{F}$  is greater than the clustering score  $F$  for any other cluster assignments  $g(S)$  at least by the structured margin  $\Delta(\mathbf{y}, \mathbf{y}^*)$ . Figure 1 gives the pictorial illustration of the overall framework.

The structured margin term  $\Delta(\mathbf{y}, \mathbf{y}^*)$  measures the quality of the clustering. The margin term outputs 0 if the clustering quality of  $\mathbf{y}$  with respect to the ground truth clustering assignment  $\mathbf{y}^*$  is perfect (up to a permutation) and 1 if the quality is the worst. We use the following margin term

$$\Delta(\mathbf{y}, \mathbf{y}^*) = 1 - \text{NMI}(\mathbf{y}, \mathbf{y}^*) \quad (8)$$

where NMI is the normalized mutual information (NMI) [21]. This measures the label agreement between the two clustering assignments ignoring the permutations. It is defined by the ratio of mutual information and the square root of the product of entropies for each assignments:

$$\text{NMI}(\mathbf{y}_1, \mathbf{y}_2) = \frac{MI(\mathbf{y}_1, \mathbf{y}_2)}{\sqrt{H(\mathbf{y}_1)H(\mathbf{y}_2)}} \quad (9)$$

The marginal and joint probability mass used for computing the entropy and the mutual information can be estimated as follows:

$$\begin{aligned} P(i) &= \frac{1}{m} \sum_j \mathbb{I}[\mathbf{y}[j] == i] \\ P(i, j) &= \frac{1}{m} \sum_{k, l} \mathbb{I}[\mathbf{y}_1[k] == i] \cdot \mathbb{I}[\mathbf{y}_2[l] == j], \end{aligned} \quad (10)$$

where  $m$  denotes the number of data (also equal to  $|X|$ ).

Figure 3 illustrates the advantages of the proposed algorithm. Since the algorithm is aware of the global landscape of the embedding space, it can overcome the bad local optima in figure 2. The clustering loss encourages small intra cluster (outlined by the dotted lines in figure 3) sum of distances with respect to each cluster medoids (three data points outlined in bold) while discouraging different clusters from getting close to each other via the NMI metric in the structured margin term.

### 3.3. Backpropagation subgradients

We fit our model using stochastic gradient descent. The key step is to compute the derivative of the loss, which is given by the following expression:

$$\begin{aligned} \partial \ell(X, \mathbf{y}^*) &= \mathbb{I}[\ell(X, \mathbf{y}^*) > 0] \left( \nabla_{\Theta} F(X, S_{\text{PAM}}; \Theta) \right. \\ &\quad \left. - \nabla_{\Theta} \tilde{F}(X, \mathbf{y}^*; \Theta) \right) \end{aligned} \quad (11)$$

Here  $S_{\text{PAM}}$  is the solution to the subproblem marked (\*) in Equation 6; we discuss how to compute this in Section 3.4.

The first gradient term is as follows:

$$\nabla_{\Theta} F(X, S; \Theta) = - \sum_{i \in |X|} \left[ \frac{f(X_i; \Theta) - f(X_{j^*(i)}; \Theta)}{\|f(X_i; \Theta) - f(X_{j^*(i)}; \Theta)\|} \cdot \nabla_{\Theta} (f(X_i; \Theta) - f(X_{j^*(i)}; \Theta)) \right] \quad (12)$$

where  $j^*(i)$  denotes the index of the closest facility location in the set  $S_{\text{PAM}}$ . The gradient for the oracle scoring function can be derived by computing

$$\nabla_{\Theta} \tilde{F}(X, \mathbf{y}_i^*; \Theta) = \sum_k \nabla_{\Theta} F(X_{\{i: \mathbf{y}^*[i]=k\}}, \{j^*(k)\}; \Theta) \quad (13)$$

Equation 11 is the formula for the exact subgradient and we find an approximate maximizer  $S_{\text{PAM}}$  in the equation (section 3.4) so we have an approximate subgradient. However, this approximation works well in practice and have been used for structured prediction setting [20, 34].

### 3.4. Loss augmented inference

We solve the optimization problem (\*) in Equation 6 in two steps. First, we use the greedy Algorithm 1 to select an initial good set of facilities. In each step, it chooses the element  $i^*$  with the best marginal benefit. The running time of the algorithm is  $O(|\mathcal{Y}|^3 \cdot |\mathcal{V}|)$ , where  $|\mathcal{Y}|$  denotes the number of clusters in the batch and  $\mathcal{V} = \{1, \dots, |X|\}$ . This time is linear in the size of the minibatch, and hence does not add much overhead on top of the gradient computation. Yet, if need be, we can speed up this part via a stochastic version of the greedy algorithm [22].

This algorithm is motivated by the fact that the first term,  $F(X, S; \Theta)$ , is a monotone submodular function in  $S$ . We observed that throughout the learning process, this term is large compared to the second, margin term. Hence, in this case, our function is still close to submodular. For approximately submodular functions, the greedy algorithm can still be guaranteed to work well [7].

Yet, since  $A(S)$  is not entirely submodular, we refine the greedy solution with a local search, Algorithm 2. This algorithm performs pairwise exchanges of current medoids  $S[k]$  with alternative points  $j$  in the same cluster. The running time of the algorithm is  $O(T|\mathcal{Y}|^3 \cdot |\mathcal{V}|)$ , where  $T$  is the maximum number of iterations. In practice, it converges quickly, so we run the algorithm for  $T = 5$  iterations only.

Algorithm 2 is similar to the partition around medoids (PAM) [15] algorithm for k-medoids clustering, which independently reasons about each cluster during the medoid swapping step. Algorithm 2 differs from PAM by the structured margin term, which involves all clusters simultaneously.

The following lemma states that the algorithm can only improve over the greedy solution:

---

#### Algorithm 1: Loss augmented inference for (\*)

---

**Input** :  $X \in \mathbb{R}^{m \times d}$ ,  $\mathbf{y}^* \in |\mathcal{Y}|^m$ ,  $\gamma$   
**Output** :  $S \subseteq \mathcal{V}$   
**Initialize**:  $S = \{\emptyset\}$   
**Define** :  $A(S) := F(X, S; \Theta) + \gamma \Delta(g(S), \mathbf{y}^*)$

```

1 while  $|S| < |\mathcal{V}|$  do
2    $i^* = \arg \max_{i \in \mathcal{V} \setminus S} A(S \cup \{i\}) - A(S)$ 
3    $S := S \cup \{i^*\}$ 
4 end
5 return S
```

---

**Lemma 1.** *Algorithm 2 monotonically increases the objective function  $A(S) = F(X, S; \Theta) + \gamma \Delta(g(S), \mathbf{y}^*)$ .*

*Proof.* In any step  $t$  and for any  $k$ , let  $c = S[k]$  be the  $k$ th medoid in  $S$ . The algorithm finds the point  $j$  in the  $k$ th cluster such that  $A((S \setminus \{c\}) \cup \{j\})$  is maximized. Let  $j^*$  be a maximizing argument. Since  $j = c$  is a valid choice, we have that  $A((S \setminus \{c\}) \cup \{j^*\}) \geq A((S \setminus \{c\}) \cup \{c\}) = A(S)$ , and hence the value of  $A(S)$  can only increase.  $\square$

In fact, with a small modification and  $T$  large enough, the algorithm is guaranteed to find a local optimum, i.e., a set  $S$  such that  $A(S) \geq A(S')$  for all  $S'$  with  $|S \Delta S'| = 1$  (Hamming distance one). Note that the overall problem is NP-hard, so a guarantee of global optimality is impossible.

**Lemma 2.** *If the exchange point  $j$  is chosen from  $X$  and  $T$  is large enough that the algorithm terminates because it makes no more changes, then Algorithm 2 is guaranteed to find a local optimum.*

### 3.5. Implementation details

We used Tensorflow [1] package for our implementation. For the embedding vector, we  $\ell_2$  normalize the embedding vectors before computing the loss for our method. The model slightly underperformed when we omitted the embedding normalization. We also tried solving the loss augmented inference using Algorithm 2 with random initialization, but it didn't work as well as initializing the algorithm with the greedy solution from Algorithm 1.

For the network architectures, we used the Inception [32] network with batch normalization [11] pretrained on ILSVRC 2012-CLS [24] and finetuned the network on our datasets. All the input images are first resized to square size ( $256 \times 256$ ) and cropped at  $227 \times 227$ . For the data augmentation, we used random crop with random horizontal mirroring for training and a single center crop for testing. In Npairs embedding [29], they take multiple random crops and average the embedding vectors from the cropped



---

**Algorithm 2:** Loss augmented refinement for (\*)

---

**Input** :  $X \in \mathbb{R}^{m \times d}$ ,  $\mathbf{y}^* \in |\mathcal{Y}|^m$ ,  $S_{\text{init}}$ ,  $\gamma$ ,  $T$   
**Output** :  $S$   
**Initialize:**  $S = S_{\text{init}}$ ,  $t = 0$

```
1 for  $t < T$  do
  // Perform cluster assignment
2    $\mathbf{y}_{\text{PAM}} = g(S)$ 
  // Update each medoids per cluster
3   for  $k < |\mathcal{Y}|$  do
    // Swap the current medoid in
    // cluster  $k$  if it increases the
    // score.
4     $S[k] = \arg \max_{j \in \{i: \mathbf{y}_{\text{PAM}}[i]=k\}} F(X_{\{i: \mathbf{y}_{\text{PAM}}[i]=k\}}, \{j\}; \Theta)$ 
5     $+ \gamma \Delta(g(S \setminus \{S[k]\} \cup \{j\}), \mathbf{y}^*)$ 
6   end
7 end
8 return  $S$ 
```

---

images during testing. However, in our implementation of [29], we take a single center crop during testing for fair comparison with other methods.

The experimental ablation study reported in [30] suggested that the embedding size doesn't play a crucial role during training and testing phase so we decided to fix the embedding size at  $d = 64$  throughout the experiment (In [30], the authors report the recall@K results with  $d = 512$  and provided the results for  $d = 64$  to us for fair comparison). We used RMSprop [33] optimizer with the batch size  $m$  set to 128. For the margin multiplier constant  $\gamma$ , we gradually decrease it using exponential decay with the decay rate set to 0.94.

As briefly mentioned in section 1, the proposed method does not require the data to be prepared in any rigid paired format (pairs, triplets, n-pair tuples, etc). Instead we simply sample  $m$  (batch size) examples and labels at random. That said, the clustering loss becomes trivial if a batch of data all have the same class labels (perfect clustering merging everything into one cluster) or if the data all have different class labels (perfect clustering where each data point forms their own clusters). In this regard, we guarded against those pathological cases by ensuring the number of unique classes ( $C$ ) in the batch is within a reasonable range. We tried three different settings  $\frac{C}{m} = \{0.25, 0.50, 0.75\}$  and the choice of the ratio did not lead to significant changes in the experimental results. For the CUB-200-2011 [37] and Cars196 [18], we set  $\frac{C}{m} = 0.25$ . For the Stanford Online Products [30] dataset,  $\frac{C}{m} = 0.75$  was the only possible choice because the dataset is extremely fine-grained.

## 4. Experimental results

Following the experimental protocol in [30, 29], we evaluate the clustering and  $k$  nearest neighbor retrieval [13] results on data from previously unseen classes on the CUB-200-2011 [37], Cars196 [18], and Stanford Online Products [30] datasets. We compare our method with three current state of the art methods in deep metric learning: (1) triplet learning with semi-hard negative mining strategy [25], (2) lifted structured embedding [30], (3) N-pairs metric loss [29]. To be comparable with prior work, we  $\ell_2$  normalize the embedding for the triplet (as prescribed by [25]) and our method, but not for the lifted structured loss and the N-pairs loss (as in the implementation sections in [30, 29]).

We used the same train/test split as in [30] for all the datasets. The CUB200-2011 dataset [37] has 11,788 images of 200 bird species; we used the first 100 birds species for training and the remaining 100 species for testing. The Cars196 dataset [18] has 16,185 images of 196 car models. We used the first 98 classes of cars for training and the rest for testing. The Stanford online products dataset [30] has 120,053 images of 22,634 products sold online on eBay.com. We used the first 11,318 product categories for training and the remaining 11,316 categories for testing.

### 4.1. Quantitative results

The training procedure for all the methods converged at  $10k$  iterations for the CUB200-2011 [37] and at  $20k$  iterations for the Cars196 [18] and the Stanford online products [30] datasets.

Tables 1, 2, and 3 shows the results of the quantitative comparison between our method and other deep metric learning methods. We report the NMI score, to measure the quality of the clustering, as well as  $k$  nearest neighbor performance with the Recall@K metric. The tables show that our proposed method has the state of the art performance on both the NMI and R@K metrics outperforming all the previous methods.

	NMI	R@1	R@2	R@4	R@8
Triplet semihard [25]	55.38	42.59	55.03	66.44	77.23
Lifted struct [30]	56.50	43.57	56.55	68.59	79.63
Npairs [29]	57.24	45.37	58.41	69.51	79.49
Clustering (Ours)	<b>59.23</b>	<b>48.18</b>	<b>61.44</b>	<b>71.83</b>	<b>81.92</b>

Table 1. Clustering and recall performance on CUB-200-2011 [37] @10k iterations.

### 4.2. Qualitative results

Figure 4, 5, and 6 visualizes the t-SNE [36] plots on the embedding vectors from our method on CUB200-

	NMI	R@1	R@2	R@4	R@8
Triplet semihard [25]	53.35	51.54	63.78	73.52	82.41
Lifted struct [30]	56.88	52.98	65.70	76.01	84.27
Npairs [29]	57.79	53.90	66.76	77.75	86.35
Clustering (Ours)	<b>59.04</b>	<b>58.11</b>	<b>70.64</b>	<b>80.27</b>	<b>87.81</b>

Table 2. Clustering and recall performance on Cars196 [18] @20k iterations.

	NMI	R@1	R@10	R@100
Triplet semihard [25]	89.46	66.67	82.39	91.85
Lifted struct [30]	88.65	62.46	80.81	91.93
Npairs [29]	89.37	66.41	83.24	93.00
Clustering (Ours)	<b>89.48</b>	<b>67.02</b>	<b>83.65</b>	<b>93.23</b>

Table 3. Clustering and recall performance on Products [30] @20k iterations.

2011 [37], Cars196 [18], and Stanford online products [30] datasets respectively. The plots are best viewed on a monitor when zoomed in. We can see that our embedding does a great job on grouping similar objects/products despite the significant variations in view point, pose, and configuration.



Figure 4. Barnes-Hut t-SNE visualization [36] of our embedding on the CUB-200-2011 [37] dataset. Best viewed on a monitor when zoomed in.



Figure 5. Barnes-Hut t-SNE visualization [36] of our embedding on the Cars196 [18] dataset. Best viewed on a monitor when zoomed in.

## 5. Conclusion

We described a novel learning scheme for optimizing the deep metric embedding with the learnable clustering function and the clustering metric (NMI) in a end-to-end fashion within a principled structured prediction framework.

Our experiments on CUB200-2011 [37], Cars196 [18], and Stanford online products [30] datasets show state of the art performance both on the clustering and retrieval tasks.

The proposed clustering loss has the added benefit that it doesn't require rigid and time consuming data preparation (i.e. no need for preparing the data in pairs [8], triplets [39, 25], or n-pair tuples [29] format). This characteristic of the proposed method opens up a rich class of possibilities for advanced data sampling schemes.

In the future, we plan to explore sampling based gradient averaging scheme where we ask the algorithm to cluster several random subsets of the data within the training batch and then average the loss gradient from multiple sampled subsets in similar spirit to Bag of Little Bootstraps (BLB) [16].

## References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster,



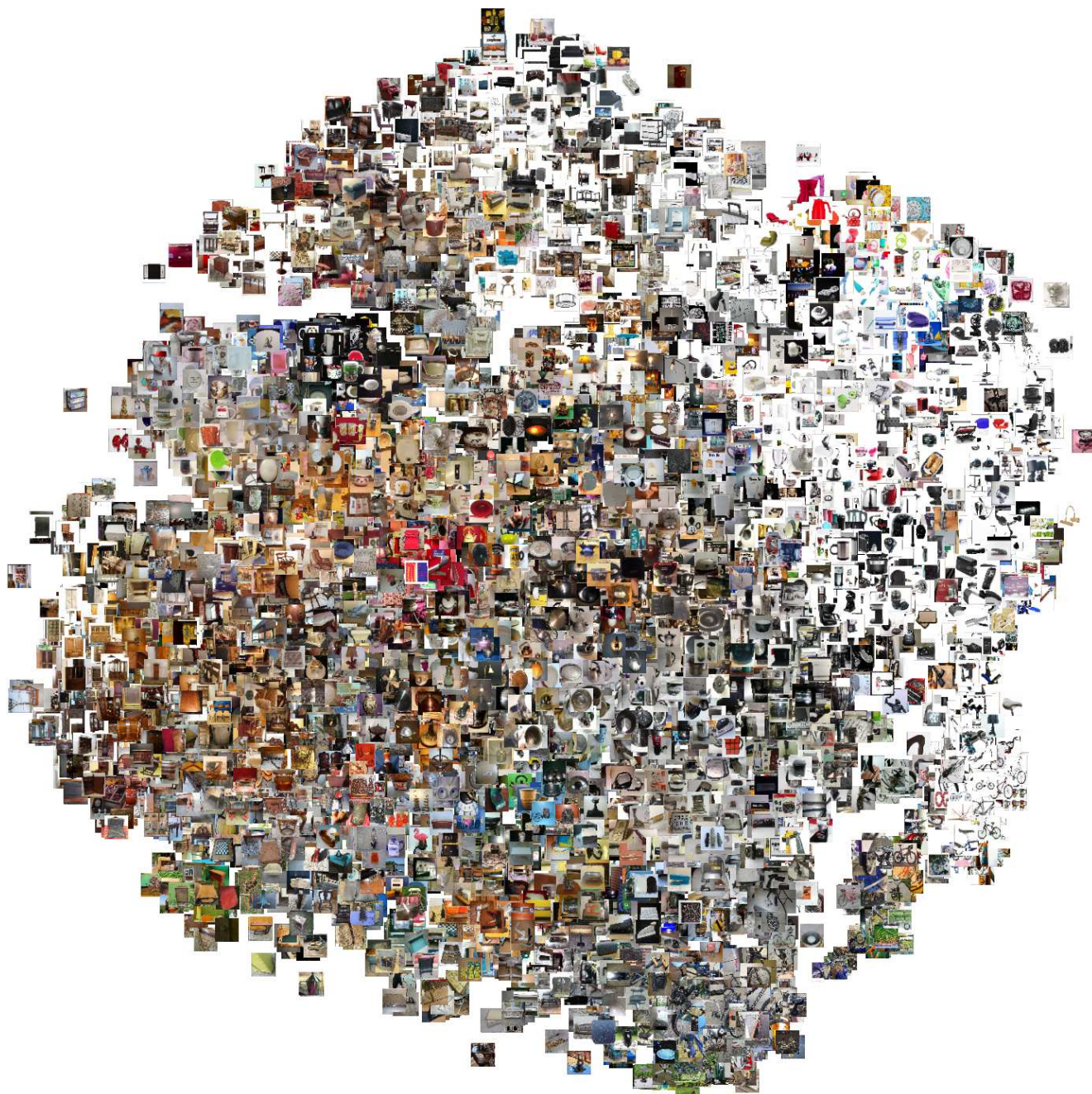


Figure 6. Barnes-Hut t-SNE visualization [36] of our embedding on the Stanford online products dataset [30]. Best viewed on a monitor when zoomed in.

- J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org. 5
- [2] S. Bell and K. Bala. Learning visual similarity for product design with convolutional neural networks. In *SIGGRAPH*, 2015. 1
- [3] J. Bromley, I. Guyon, Y. Lecun, E. Sckinger, and R. Shah. Signature verification using a "siamese" time delay neural network. In *NIPS*, 1994. 1
- [4] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, volume 1, June 2005. 1
- [5] A. Choromanska, A. Agarwal, and J. Langford. Extreme multi class classification. In *NIPS*, 2013. 1
- [6] C. B. Choy, J. Gwak, S. Savarese, and M. Chandraker. Universal correspondence network. In *NIPS*, 2016. 1



- [7] A. Das and D. Kempe. Submodular meets spectral: Greedy algorithms for subset selection, sparse approximation and dictionary selection. In *ICML*, 2011. 5
- [8] R. Hadsell, S. Chopra, and Y. Lecun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, 2006. 1, 7
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. 1
- [10] J. R. Hershey, Z. Chen, J. L. Roux, and S. Watanabe. Deep clustering: Discriminative embeddings for segmentation and separation. In *ICASSP*, 2016. 1, 3
- [11] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5
- [12] C. Ionescu, O. Vantzos, and C. Sminchisescu. Training deep networks with structured layers by matrix backpropagation. In *ICCV*, 2015. 3
- [13] H. Jegou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. In *PAMI*, 2011. 6
- [14] T. Joachims, T. Finley, and C.-N. Yu. Cutting-plane training of structural svms. *JMLR*, 2009. 1
- [15] L. Kaufman and P. Rousseeuw. Clustering by means of medoids. In *Statistical Data Analysis Based on the L1-Norm and Related Methods*, 1987. 5
- [16] A. Kleiner, A. Talwalkar, P. Sarkar, and M. I. Jordan. The big data bootstrap. In *ICML*, 2012. 7
- [17] A. Krause and D. Golovin. Submodular function maximization. *Tractability: Practical Approaches to Hard Problems*, 3(19):8, 2012. 4
- [18] J. Krause, M. Stark, J. Deng, and F.-F. Li. 3d object representations for fine-grained categorization. *ICCV 3dRR-13*, 2013. 1, 6, 7
- [19] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1
- [20] H. Lin and J. Bilmes. Learning mixtures of submodular shells with application to document summarization. In *UAI*, 2012. 3, 5
- [21] C. D. Manning, P. Raghavan, and H. Schütze. *Introduction to Information Retrieval*. Cambridge university press, 2008. 1, 4
- [22] B. Mirzasoleiman, A. Badanidiyuru, A. Karbasi, J. Vondrák, and A. Krause. Lazier than lazy greedy. In *Proc. Conf. on Artificial Intelligence (AAAI)*, 2015. 5
- [23] Y. Prabhu and M. Varma. Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning. In *SIGKDD*, 2014. 1
- [24] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 5
- [25] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 1, 2, 6, 7
- [26] M. Schultz and T. Joachims. Learning a distance metric from relative comparisons. In *NIPS*, 2004. 2
- [27] O. Sener, H. O. Song, A. Saxena, and S. Savarese. Learning transferrable representations for unsupervised domain adaptation. In *NIPS*, 2016. 1
- [28] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 1
- [29] K. Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NIPS*, 2016. 1, 3, 5, 6, 7
- [30] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, 2016. 1, 2, 6, 7, 8
- [31] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 1
- [32] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 5
- [33] T. Tieleman and G. Hinton. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSE: Neural Networks for Machine Learning, 2012. 6
- [34] S. Tschiatschek, R. Iyer, H. Wei, and J. Bilmes. Learning mixtures of submodular functions for image collection summarization. In *NIPS*, 2014. 3, 5
- [35] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. Support vector machine learning for interdependent and structured output spaces. In *ICML*, 2004. 1
- [36] L. van der maaten. Accelerating t-sne using tree-based algorithms. In *JMLR*, 2014. 6, 7, 8
- [37] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 1, 6, 7
- [38] X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, 2015. 1
- [39] K. Q. Weinberger, J. Blitzer, and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. In *NIPS*, 2006. 1, 2, 7
- [40] I. E. Yen, X. Huang, K. Zhong, P. Ravikumar, and I. S. Dhillon. Pd-sparse: A primal and dual sparse approach to extreme multiclass and multilabel classification. In *ICML*, 2013. 1
- [41] S. Zheng, Y. Song, T. Leung, and I. Goodfellow. Improving the robustness of deep neural networks via stability training. In *CVPR*, 2016. 1