

DEEP UNSUPERVISED HASHING BY DISTILLED SMOOTH GUIDANCE

Xiao Luo^{1,2*}, Zeyu Ma^{3*}, Daqing Wu^{1,2*}, Huasong Zhong², Chong Chen^{1,2}, Jinwen Ma¹, Minghua Deng¹

¹School of Mathematical Sciences, Peking University, China

²Damo Academy, Alibaba Group, Hangzhou, China

³School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, China

ABSTRACT

Hashing has been widely used in approximate nearest neighbor search for its storage and computational efficiency. Deep supervised hashing methods are not widely used because of the lack of labeled data, especially when the domain is transferred. Meanwhile, unsupervised deep hashing models can hardly achieve satisfactory performance due to the lack of reliable similarity signals. To tackle this problem, we propose a novel deep unsupervised hashing method, namely Distilled Smooth Guidance (DSG), which can learn a distilled dataset consisting of similarity signals as well as smooth confidence signals. To be specific, we obtain the similarity confidence weights based on the initial noisy similarity signals learned from local structures and construct a priority loss function for smooth similarity-preserving learning. Besides, global information based on clustering is utilized to distill the image pairs by removing contradictory similarity signals. Extensive experiments on three widely used benchmark datasets show that the proposed DSG consistently outperforms the state-of-the-art search methods.

Index Terms— Learning to hash, Unsupervised learning, Deep learning

1. INTRODUCTION

Deep learning-based hashing methods can be divided into supervised hashing and unsupervised hashing [1]. At the early stage, many researchers mainly focused on the supervised hashing methods, which utilize semantic labels to greatly improve the performance of image retrieval [2]. However, supervised hashing methods are difficult to be applied in practice when there is not enough labeled data, especially when the domain is transferred. To solve this problem, several deep learning-based unsupervised methods were proposed, including deep binary descriptors (DeepBit) [3], semantic structure-

based unsupervised deep hashing (SSDH) [4] and unsupervised deep hashing by distilling data pairs (DistillHash) [5].

Although deep learning-based unsupervised hashing methods can be applied on unlabelled data, they still have evident limitations. [6] takes clustering information to generate pairwise pseudo-labels. SSDH further studies the deep feature statistics empirically from a pre-trained model and captures the semantic relationships across different data points. Specifically, it selects image pairs with confident pseudo-labels to guide the training of the model. Nevertheless, SSDH discards most image pairs which are hard to be decided whether they are semantically similar or dissimilar, which causes much information loss and thus limits the performance of the model in further image retrieval. In general, these methods only consider either the local information of similarity signals or the global information such as clustering labels but do not consider this task comprehensively.

To tackle the above issues, we propose a novel method, which comprehensively explores correlations among image pairs. First of all, for the correlation among different samples, we adopt the pre-trained deep convolutional neural network (CNN) to generate features for the input images. Then we compute the pairwise cosine distance and construct the similarity pseudo graph. We take all the image pairs into consideration. To alleviate the effect of wrong pseudo labels, different weights are assigned to image pairs according to the confidence of the pseudo labels. Furthermore, for the global robustness, we adopt clustering on the deep image features and then obtain another similarity graph. According to these two similarity graphs, image pairs with different correlation identification are considered contradictory and thus distilled. Finally, we design a deep neural network based on the distilled data pair set and adopt a deep learning framework to perform the deep representation and hash code learning simultaneously. Our main contributions can be summarized as follows:

- We introduce two similarity graphs based on the local (i.e., pairwise cosine similarity) and global information (i.e., clustering) respectively and then obtain the ensemble similarity graph by distilling the contradictory pairs.

The work was done when Xiao Luo and Daqing Wu interned in Damo Academy, Alibaba Group. This work was supported by The National Key Research and Development Program of China (No.2016YFA0502303) and the National Natural Science Foundation of China (No.31871342). * Equal contribution. † Corresponding authors: Chong Chen (cheung.cc@alibaba-inc.com) and Minghua Deng (dengmh@pku.edu.cn)

- We construct a priority loss function for similarity-preserving learning, which prioritizes confident image pairs over fuzzy image pairs based on their pairwise distance to learn deep representations smoothly.
- Experiments on three popular benchmark datasets show that our method DSG outperforms current state-of-the-art unsupervised hashing methods by a large margin.

2. RELATED WORK

Deep Supervised Hashing. Deep supervised hashing methods usually map the data points into Hamming space where the semantic similarities can be preserved by learning a deep neural network [1]. Typically, as the first supervised deep hashing method, CNNH [2] splits the hash learning into two stages based on a convolutional neural network. Hash codes are learned on the first stage, a specific deep network is learned on the second stage to map the input samples to the learned hash codes. Deep Supervised Hashing [7] uses a loss function with product form based on the pairwise similarities and Hamming distances, which can be trained by the end-to-end back propagation algorithm.

Deep Unsupervised Hashing. Unsupervised deep hashing methods aim to turn unsupervised problems into supervised problems by constructing pseudo labels based on deep features. Semantic Structure-based Unsupervised Deep Hashing (SSDH) [4] uses a specific truncated function on the pairwise distances and constructs the similarity matrices. Distill-Hash [5] improves the performance of SSDH by distilling the data pairs for confident similarity signals. MLS³RDUH [8] utilizes the intrinsic manifold structure in feature space to reconstruct the local semantic similarity structure, and achieves the state-of-the-art performance.

3. METHOD

In the problem of deep unsupervised hashing, $\mathcal{X} = \{x_i\}_{i=1}^N$ denotes the training set with N samples without label annotations, it aims to learn a hash function

$$\mathcal{H}: x \rightarrow b \in \{-1, 1\}^L,$$

where x is the input sample and b is a compact L -bit hash code. It is noticed that $x_i \in \mathbb{R}^d$ is the normalized extracted deep feature for the i -th image through the pre-trained neural network by removing the last fully-connected layer. Here we use VGG-F [9] to be consistent with other articles.

3.1. Pseudo-Similarity Graph

In our model, the pseudo-similarity graph is constructed firstly. The pseudo-similarity graph is used to capture pairwise similarity information from a local perspective. Based

on the pre-trained deep feature x_i , the cosine distance between the i -th and the j -th samples can be computed by $d_{ij} = 1 - \frac{x_i \cdot x_j}{\|x_i\|_2 \|x_j\|_2}$. We set a large threshold t , and consider data points with the cosine distance smaller than t as potential similar and data points with the cosine distance larger than t as potential dissimilar. Based on the threshold t , we construct the pseudo-graph S as:

$$S_{ij} = \begin{cases} 1 & d_{ij} \leq t, \\ -1 & d_{ij} > t \end{cases} \quad (1)$$

Where S_{ij} is set to 1 if points x_i and x_j are potential similar, and -1 if points x_i and x_j are potential dissimilar.

3.2. Smooth Weight Matrix

Although pseudo-similarity graph is constructed, the semantic confidence of pseudo-label for each pair is different. In this section, we construct the weight matrix for the pseudo-graph based on the semantic confidence.

By observing the distribution of cosine distance for deep feature pairs, [4] finds that each distance histogram is similar to two half Gaussian distributions, where m_l and σ_l denote the mean and the standard deviation of the first (left half) distribution and m_r and σ_r denote the mean and the standard deviation of the second (right half) distribution. Accordingly, we obtain two distance thresholds $d_l = m_l - \alpha\sigma_l$ and $d_r = m_r + \beta\sigma_r$, where the hyper-parameters α and β control the values of the distance thresholds d_l and d_r respectively and dictate the percentage of similar points and dissimilar points from all data points respectively as well. According to the theory of confidence interval, the pairs with distances smaller than d_l or larger than d_r have confident semantic similarity information. As a result, we set the weights for confident pairs to 1. Since the pairs with distance in the interval $[d_l, d_r]$ probably have no certain semantic information, we obtain the weight as:

$$W_{1,ij} = \begin{cases} \frac{(t-d_{ij})^2}{(t-d_l)^2} & d_l < d_{ij} \leq t, \\ \frac{(d_{ij}-t)^2}{(d_r-t)^2} & t < d_{ij} < d_r, \\ 1 & d_{ij} \leq d_l \text{ or } d_{ij} \geq d_r \end{cases} \quad (2)$$

From the equation, the weight for each image pair with the distance out of the interval is set to 1 and the weight for each image pair with the distance in the interval is smaller if their distance is closer to the threshold in a quadratic form. In this way, all the image pairs are taken into consideration with the guidance of the smooth weight matrix. Accordingly, the confidence of their similarity relationship is guided by the smooth weight W_1 .

3.3. Pair Distilling based on clustering

As we know, the obtained pseudo-label is very coarse. In this section, we try to use the clustering method to distill the

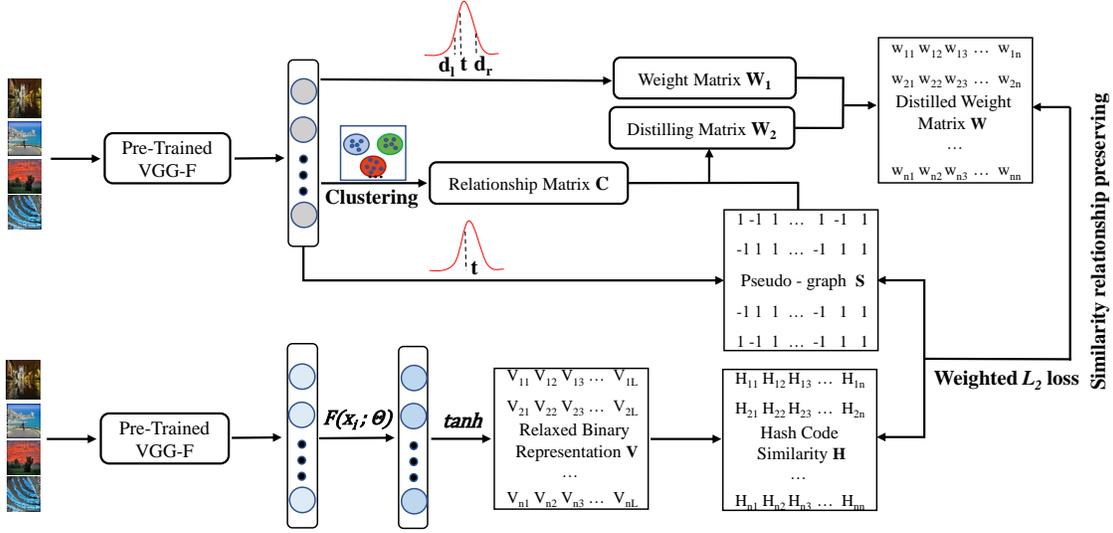


Fig. 1. The framework of our model. First, deep features of the dataset are extracted through the pre-trained VGG-F network. Second, the pseudo-graph S and the smooth weight matrix W_1 are constructed based on the local information of the cosine distance distribution of deep features. Third, the relationship matrix C is constructed based on the global information of the K-means clustering on deep features and the distilling matrix W_2 is obtained by comparing the relationship matrix C with the pseudo-graph S . Lastly, iterations are done aiming to minimize the weighted L_2 loss for preserving similarity relationship of image pairs (i.e., Equation 6) by using the mini-batch stochastic gradient descent method.

pairs by removing contradictory results. To be specific, we first use the extracted features $\{x_i\}_{i=1}^N$ to construct a K-means clustering model. Based on the K-means clustering result, we construct a relationship matrix C :

$$C_{ij} = \begin{cases} 1 & c_i = c_j, \\ -1 & c_i \neq c_j \end{cases} \quad (3)$$

where c_i represents the label of the cluster that the data point x_i belongs to. $c_i \in \{0, 1, \dots, K\}$ and K is the parameter that represents the number of clusters. By comparing the relationship matrix C with the pseudo-graph S , we obtain the distilling matrix W_2 as:

$$W_{2,ij} = \begin{cases} 1 & S_{ij} = C_{ij}, \\ 0 & S_{ij} \neq C_{ij} \end{cases} \quad (4)$$

in which $S_{ij} = C_{ij}$ means that the result of pseudo-graph and clustering for the pair x_i and x_j is consistent.

After setting the weight of contradictory pairs to zero, the final distilled weight matrix is obtained as

$$W = W_1 \odot W_2,$$

in which \odot denotes the element-wise product of vectors.

3.4. Hash code learning

For the purpose of preserving the similarity relationship of data points, similar data points are expected to be mapped into

similar hash codes and dissimilar data points are expected to be mapped into dissimilar hash codes. The similarity output $H_{N \times N}$ of hash codes is formulated as

$$H_{ij} = \frac{1}{L} \mathbf{b}_i \top \mathbf{b}_j, \quad \mathbf{b}_i \in \{-1, +1\}^L \quad (5)$$

To preserve the obtained semantic structures, we minimize the weighted L_2 loss between the hash code similarity and the pseudo-graph. In formulation,

$$\mathcal{L} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N W_{ij} (H_{ij} - S_{ij})^2$$

in which $\mathbf{b}_i = \text{sign}(F(x_i; \Theta))$ and $F(x_i; \Theta)$ denotes the output of the neural network, and Θ is the parameters. In this way, we can integrate the above loss function into the deep architecture. However, it is infeasible to train the neural network with binary outputs by the standard Back Propagation algorithm because of the ill-posed gradient problem. As a result, $\tanh(\cdot)$ is utilized to relieve the binary constraint. Thus we adopt the following objective function:

$$\begin{aligned} \min_{\Theta} \mathcal{L}(\Theta) &= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m W_{ij} (H_{ij} - S_{ij})^2 \\ \text{s.t. } H_{ij} &= \frac{1}{L} \mathbf{v}_i \top \mathbf{v}_j, \quad \mathbf{v}_i = \tanh(F(x_i; \Theta)) \end{aligned} \quad (6)$$

in which \mathbf{v}_i denotes the relaxed binary representation.

For the point q_i not in the training set, its hash code \mathbf{b}_i is obtained by directly forward propagating it through the

Algorithm 1 Training algorithm for our model

Input: Training images $\mathcal{I} = \{I_1, \dots, I_N\}$;

Number of clusters: K ;

Cosine distance threshold: t

Output: Parameters Θ for the neural network;

Hash codes B for training images.

- 1: Get deep features of \mathcal{I} through VGG-F: $\mathcal{X} = \{x_1, \dots, x_N\}$.
 - 2: Construct the pseudo-graph S by Equation 1
 - 3: Construct the weight matrix W_1 by Equation 2
 - 4: Cluster \mathcal{X} into K different groups by K -means. and construct the distilling matrix W_2 by Equation 3
 - 5: Calculate distilled weight matrix W by Equation 4
 - 6: **repeat**
 - 7: Sample N points from \mathcal{X} and then construct a mini-batch.
 - 8: Calculate the outputs by forward-propagating through the network.
 - 9: Update parameters of the VGG-F network by Back propagation by Equation 6
 - 10: **until** convergence
-

learned neural network.

$$b_i = \text{sign}(F(q_i; \Theta)) \quad (7)$$

3.5. Optimization

To optimize the problem, we construct the pseudo-graph S from the pre-trained neural network by using Equation 1. Then the smooth weight matrix W_1 and the distilling matrix W_2 are constructed to get the distilled weight matrix W . Lastly, we minimize Equation 6 by using the standard stochastic gradient descent (SGD) method. The whole learning procedure is summarized in Algorithm 1.

4. EXPERIMENTS

We implement extensive experiments on three datasets to evaluate our DSG by comparing with several state-of-the-art unsupervised hashing methods.

4.1. Datasets and Baselines

CIFAR-10 [10] is a dataset for image classification and retrieval, containing 60K images from 10 different categories. For each class, we randomly select 1,000 images as queries and 500 as training images, resulting in a query set containing 10,000 images and a training set made up of 5,000 images. All images except for the query set are used as the retrieval set. NUS-WIDE [11] contains 269,648 images, each of the images is annotated with multiple labels referring to 81 concepts. The subset containing the 10 most popular concepts is

used here. We randomly select 5,000 images as a test set; the remaining images are used as a retrieval set, and 5000 images are randomly selected from the retrieval set as the training set. FLICKR25K [12] contains 25,000 images collected from the Flickr website. Each image is manually annotated with at least one of the 24 unique labels provided. We randomly select 2,000 images as a test set; the remaining images are used as a retrieval set, from which we randomly select 10,000 images as a training set.

Our method is compared with both traditional hashing methods and state-of-the-art unsupervised deep learning methods. Traditional methods includes SpH [13], DSH [14] and SGH [15]. Deep unsupervised hashing methods includes DeepBits [3], SSDH [4], DistillHash [5], CUDH [16], and MLS³RUDH [8].

4.2. Implementation Details

The framework is implemented by Pytorch V1.4. The mini-batch size is set to 24 and the momentum to 0.9. The learning rate is fixed at 0.001. The initial weights of the first seven layers of the neural network are from the model pre-trained with ImageNet, and the last fully-connected layer is learnt from scratch. The parameter α and β are all set following [4] and the threshold t to 0.1. The number of clusters is 70.

4.3. Evaluation

The ground-truth similarity information for evaluation is constructed from the ground-truth image labels: two data points are considered similar if they share the same label (for CIFAR-10) or share at least one common label (for FLICKR25K and NUSWIDE).

The retrieval quality are evaluated by the following three evaluation metrics: Mean Average Precision (MAP), Precision-Recall curve and Top N precision curve.

MAP is a widely-used criteria to evaluate retrieval accuracy. Given a query and a list of R ranked retrieval results, the average precision (AP) for the given query can be computed. MAP is defined as the average of APs for all queries. For datasets FLICKR25K and NUSWIDE, we set R as 5000 for the experiments. For CIFAR-10, R is set as 50000. Precision-recall curve reveals the precision at different recall levels and is a good indicator of overall performance. In addition, Top N precision curve, which is the precision curve with respect to the top K retrieved instances, also visualizes the performance from a different perspective .

4.4. Overall Performance and Ablation Study

Table 1 shows the MAPs for different methods on three datasets with hash code lengths varying from 16 to 128. It can be seen that the performances of deep learning-based algorithms are overall better than traditional methods. For our

Table 1. MAP for different methods on FLICKR25K, CIFAR-10 and NUS-WIDE datasets.

Method	FLICKR25K				CIFAR-10				NUS-WIDE			
	16bits	32bits	64bits	128bits	16bits	32bits	64bits	128bits	16bits	32bits	64bits	128bits
SH	0.6091	0.6105	0.6033	0.6014	0.1605	0.1583	0.1509	0.1538	0.4350	0.4129	0.4062	0.4100
SpH	0.6119	0.6315	0.6381	0.6451	0.1439	0.1665	0.1783	0.1840	0.4458	0.4537	0.4926	0.5000
SGH	0.6362	0.6283	0.6253	0.6206	0.1795	0.1827	0.1889	0.1904	0.4994	0.4869	0.4851	0.4945
DeepBit	0.5934	0.5933	0.6199	0.6349	0.2204	0.2410	0.2521	0.2530	0.3844	0.4341	0.4461	0.4917
SSDH	0.7240	0.7276	0.7377	0.7343	0.2568	0.2560	0.2587	0.2601	0.6374	0.6768	0.6829	0.6831
DistillHash	0.6964	0.7056	0.7075	0.6995	0.2844	0.2853	0.2867	0.2895	0.6667	0.6752	0.6769	0.6747
CUDH	0.7332	0.7426	0.7549	0.7561	0.2856	0.2903	0.3025	0.3000	0.6996	0.7222	0.7451	0.7418
MLS ³ RUDH	0.7587	0.7754	0.7870	0.7927	0.2876	0.2962	0.3139	0.3117	0.7056	0.7384	0.7629	0.7818
DSG	0.7994	0.8172	0.8197	0.8245	0.3225	0.3241	0.3453	0.3450	0.7795	0.7981	0.8098	0.8187

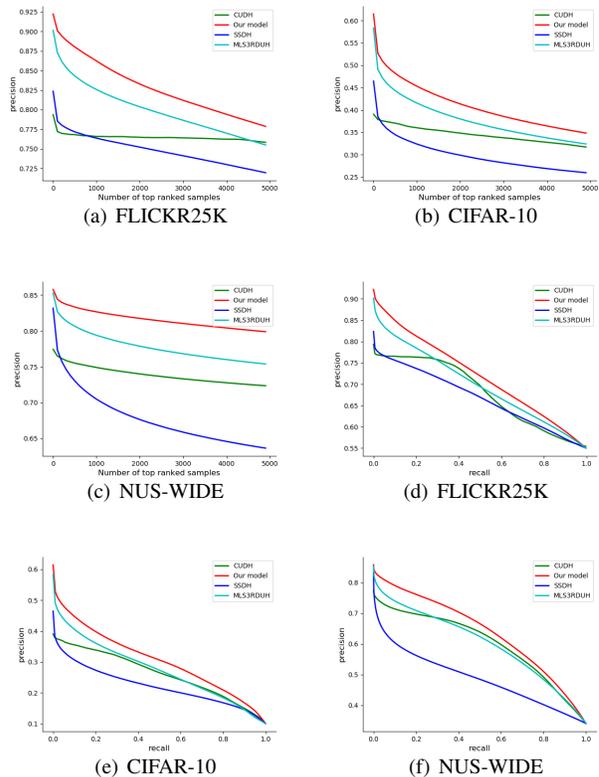
Table 2. Ablation analysis on datasets FLICKR25K, CIFAR-10 and NUS-WIDE.

Method	FLICKR25K				CIFAR-10				NUS-WIDE			
	16bits	32bits	64bits	128bits	16bits	32bits	64bits	128bits	16bits	32bits	64bits	128bits
DSG-v2	0.6896	0.6570	0.6228	0.6178	0.2356	0.2099	0.1704	0.1567	0.6967	0.6736	0.6257	0.6062
DSG-v1	0.7584	0.7432	0.7355	0.7240	0.2612	0.2496	0.2441	0.1988	0.7632	0.7718	0.7719	0.7703
DSG	0.7994	0.8172	0.8197	0.8245	0.3225	0.3241	0.3453	0.3450	0.7795	0.7981	0.8098	0.8187

proposed method, we find that DSG has a significant improvement over the state-of-the-art deep learning-based methods in all cases, which implies the superiority of our model. By comparing with MLS³RUDH, which has the best performance of MAP among the deep hashing methods, DSG improves the MAP by 3.18%, 3.33% and 3.69% for 128 bit length on datasets FLICKR25K, CIFAR-10 and NUS-WIDE respectively. We also find that the larger the bit length, the greater the improvement of DSG over other methods, which implies that DSG can generate more independent hash bits.

We also compare the performance of the DSG full model, the DSG model without the weight matrix W_1 , which is denoted as DSG-v1 and the DSG model without the weight matrix W_1 and the distilling matrix W_2 , which is denoted as DSG-v2. The results are shown in Table 2. It is easy to find that DSG always achieves the highest MAPs, which implies that both the confidence information provided by the weight matrix W_1 and the global consistency that is ensured by the distilling matrix W_2 are necessary in our model. We can also find that DSG-v1 always outperforms DSG-v2, which proves that the clustering alone can help improve performance by removing contradictory similarity relationships between the pseudo-graph S and the relationship matrix C .

For a more comprehensive comparison, we draw precision-recall curves and Top N precision curves for our method DSG and three state-of-art methods CUDH, SSDH and MLS³RUDH with the hash code length of 128. Figure 2 (a), (b) and (c) show the Top N precision curves of CUDH, DSG, SSDH and MLS³RUDH on datasets FLICKR25K, CIFAR-10 and NUS-WIDE. It can be seen that DSG always has the highest precision among these four models and MLS³RUDH always has the second-highest precision. Since the precision curves are based on the ranks of Hamming dis-

**Fig. 2.** (a), (b) and (c) are the Top N precision curves with code length 128 on FLICKR25K, CIFAR-10 and NUS-WIDE. (d) and (e) and (f) are the precision-recall curves with code length 128 on FLICKR25K, CIFAR-10 and NUS-WIDE.

tance, DSG is able to achieve the highest recall if we directly use Hamming distance for retrieval. It's known that hash codes can also be used for coarse filtering in the form of hash table lookup; we also plot the precision-recall curves for these four models on the same datasets, which are shown in Figure 2 (d), (e) and (f). It can be clearly seen that the curves of DSG are always on top of the other three models' curves, which implies that the hash codes obtained by DSG are also more suitable for the hash table lookup search strategy, which further demonstrates the superiority of our method.

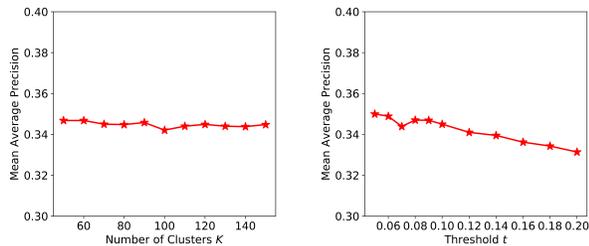


Fig. 3. MAP w.r.t different numbers of clusters and threshold t values with code length 128 on CIFAR-10.

4.5. Parameter Sensitivity

We further study the number of clusters K and threshold t . Figure 3 shows the MAPs of different K values ranging [50, 150] on the dataset CIFAR-10 with code length 128 and the MAPs of different threshold t values ranging [0.05, 0.20] on the dataset CIFAR-10 with code length 128 as well. We find the MAP with the K value of 100 is slightly lower than other K values'. In general, DSG's performance with different K values ranging [50, 150] is relatively stable, indicating that the model performance is not sensitive to different K values ranging [50, 150]. Furthermore, we show that the MAP decreases as the threshold t is over 0.1 with K fixed to 70. In addition, the performance of DSG with different t ranging [0.05, 0.1] is relatively stable, which indicates that the suitable interval for t value is [0.05, 0.1]. Accordingly, we set K as 70 and t as 0.1 in our other experiments as default.

5. CONCLUSION

In this paper, we proposed Distilled Smooth Guidance (DSG) for deep unsupervised hashing. DSG not only considers local similarity signals but also considers the confidence of similarity signals from local structure for smoothness. What's more, global information is also explored and the image pairs are distilled by removing contradictory image pairs from two views for the purpose of accuracy. Numeric experiments demonstrates that DSG outperforms the existing state-of-the-art methods.

6. REFERENCES

- [1] X. Luo, C. Chen, H. Zhong, H. Zhang, M. Deng, J. Huang, and X. Hua, "A survey on deep hashing methods," *arXiv preprint arXiv:2003.03369*, 2020.
- [2] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan, "Supervised hashing for image retrieval via image representation learning," in *Twenty-eighth AAAI conference on artificial intelligence*, 2014.
- [3] K. Lin, J. Lu, C.-S. Chen, and J. Zhou, "Learning compact binary descriptors with unsupervised deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1183–1192.
- [4] E. Yang, C. Deng, T. Liu, W. Liu, and D. Tao, "Semantic structure-based unsupervised deep hashing," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018, pp. 1064–1070.
- [5] E. Yang, T. Liu, C. Deng, W. Liu, and D. Tao, "Distillhash: Unsupervised deep hashing by distilling data pairs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2946–2955.
- [6] Q. Hu, J. Wu, J. Cheng, L. Wu, and H. Lu, "Pseudo label based unsupervised deep discriminative hashing for image retrieval," in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 1584–1590.
- [7] H. Liu, R. Wang, S. Shan, and X. Chen, "Deep supervised hashing for fast image retrieval," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2064–2072.
- [8] R.-C. Tu, X.-L. Mao, and W. Wei, "Mls3rduh: Deep unsupervised hashing via manifold based local semantic similarity structure reconstructing," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, C. Bessiere, Ed., 7 2020, pp. 3466–3472.
- [9] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [10] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [11] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, "Nus-wide: a real-world web image database from national university of singapore," in *Proceedings of the ACM international conference on image and video retrieval*. ACM, 2009, p. 48.
- [12] M. J. Huiskes and M. S. Lew, "The mir flickr retrieval evaluation," in *Proceedings of the 1st ACM international conference on Multimedia information retrieval*, 2008, pp. 39–43.
- [13] J.-P. Heo, Y. Lee, J. He, S.-F. Chang, and S.-E. Yoon, "Spherical hashing," in *2012 IEEE Conference on Com-*

puter Vision and Pattern Recognition. IEEE, 2012, pp. 2957–2964.

- [14] Z. Jin, C. Li, Y. Lin, and D. Cai, “Density sensitive hashing,” *IEEE transactions on cybernetics*, vol. 44, no. 8, pp. 1362–1371, 2013.
- [15] B. Dai, R. Guo, S. Kumar, N. He, and L. Song, “Stochastic generative hashing,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR.org, 2017, pp. 913–922.
- [16] Y. Gu, S. Wang, H. Zhang, Y. Yao, and L. Liu, “Clustering-driven unsupervised deep hashing for image retrieval,” *Neurocomputing*, vol. 368, pp. 114–123, 2019.