

Unsupervised Multi-modal Hashing for Cross-Modal Retrieval

Jun Yu, Xiao-Jun Wu*, Donglin Zhang

Abstract—With the advantage of low storage cost and high efficiency, hashing learning has received much attention in the domain of Big Data. In this paper, we propose a novel unsupervised hashing learning method to cope with this open problem to directly preserve the manifold structure by hashing. To address this problem, both the semantic correlation in textual space and the locally geometric structure in the visual space are explored simultaneously in our framework. Besides, the $\ell_{2,1}$ -norm constraint is imposed on the projection matrices to learn the discriminative hash function for each modality. Extensive experiments are performed to evaluate the proposed method on the three publicly available datasets and the experimental results show that our method can achieve superior performance over the state-of-the-art methods.

Index Terms—Multimodal Hashing, Cross-modal Retrieval, Unsupervised learning, manifold preserving

I. INTRODUCTION

RECENTLY, the explosive growth of multimedia data brings enormous challenge in information retrieval [1], [2], data mining [3], [4], and computer vision [5]. It is necessary to develop methods to support retrieving relevant objects from such massive database. Binary codes learning, a.k.a. hashing, has achieved great success because of its low storage and high efficiency. Among hashing methods [6], [7], [8], [9], [10], [11], [12], [13], Neighborhood Preserving Hashing (NPH) [9], Scalable Deep Hashing (SCADH) [10], Similarity Preserving Linkage Hashing (SPLH) [11], Weakly Supervised Multimodal Hashing (WMH) [12] and Discrete Locally Linear embedding (DLLH) [13] have achieved promising performance. Nevertheless, these methods are assumed in single-modal circumstances and do not directly apply to multi-modal applications.

Cross-modal retrieval is a more interesting scenario because multimodal data are often available in multimedia domains. The major task of cross-modal retrieval is to find the same semantic data from different modal spaces when given query data. Most of the previous works pay attention to supervised and semi-supervised multimodal hashing learning algorithms that focus on learning discriminative features by utilizing available semantic labels. Label Consistent Matrix Factorization Hashing (LCMFH) [14] learns a latent common space where data classified into the same category share a common representation. Multi-view Feature Discrete Hashing

(MFDH) [15] jointly performs classifier learning and subspace learning for cross-modal retrieval. Semantic correlation maximization (SCM) [16] reconstructs the semantic similarity matrix calculated by the label vectors in hamming space to learn the discriminative hash codes. Semantics-Preserving Hashing (SePH) [17] transforms the semantic affinity into a probability distribution and approximates the distribution in Hamming space. Semi-supervised Hashing [18] learns the hash functions by utilizing the label information of partial data. Although the above methods are very efficient to realize cross-modal retrieval, they depend on the labeled data and it is time-consuming and labor-intensive to obtain them in real applications.

Unsupervised cross-modal hashing methods aim to learn the high-quality hash codes which preserve the structural and topological information of data. Cross-View Hashing (CVH) [19] is a pioneering work that extends the traditional unimodal spectral hashing [20] to the multimodal situation. Robust Cross-view Hashing (RCH) [21] learns a common Hamming space in which the binary codes of the paired different modalities are as consistent as possible. Canonical Correlation Analysis (CCA) [5] transforms multiple views into a common latent subspace in which the correlation between two views is maximized. Fusion Similarity Hashing (FSH) [22] embeds the graph-based fusion similarity into a common Hamming space. The main idea of Inter-Media Hashing (IMH) [23] is that the learned binary codes preserve inter-media and intra-media consistency simultaneously. Unsupervised multimodal hashing generally needs to solve two basic problems: how to preserve the geometric structure among data points by hash codes and how to simultaneously select discriminative features for multiple modalities. Although existing unsupervised hashing methods have been developed, but above problems are not well addressed simultaneously. In fact, some tags or texts associated with uploaded images in social media contain the weakly semantic information. In this paper, we proposed a unsupervised multi-modal hashing where both the weakly semantic structure information provided by textual modality and the visually underlying manifold structure are explored simultaneously. Besides, the projection matrices are constrained by $\ell_{2,1}$ -norm to learn the discriminative and compact binary codes. The overview of the proposed method is shown in Fig.1 and the advantages of our method are summarized as follows

- (1) We propose a sparse multi-modal hashing method by which the learned hash codes preserve the semantically and visually structural information.
- (2) Our model jointly performs the multi-modal graph embedding and discriminative features learning, which further

J. Yu, X.-J. Wu (corresponding author) and D. Zhang are with the School of Artificial Intelligence and Computer Science, Jiangnan University, 214122, Wuxi, China. J. Yu, X.-J. Wu and D. Zhang are also with the Jiangsu Provincial Engineering Laboratory of Pattern Recognition and Computational Intelligence, Jiangnan University 214122, Wuxi, China. e-mail: (yujunjason@aliyun.com; wu_xiaojun@jiangnan.edu.cn; dlinzhang@163.com).

improves the performance.

(2) Experimentally, a comparative evaluation of the proposed method on three available datasets with other state-of-the-art hashing methods shows that our method boosts the retrieval performance.

Structurally, the rest of this paper falls into four sections. In section II, we simply introduce the related work in this field. Our model and the optimization algorithm are presented in section III. In section IV, we discuss the experimental results on three available datasets and analyze the sensitivity of some parameters. Finally, the conclusions are drawn in section V.

II. RELATED WORK

In this section, we preliminarily review the related work in the field of cross-modal hashing. Cross-modal hashing algorithms are roughly divided into supervised cross-modal hashing and unsupervised cross-modal hashing which are distinguished by whether the label information is utilized or not.

Supervised cross-modal hashing methods learn the discriminative hashing feature via exploiting the available label information. Semantic Correlation Maximization (SCM) [16] utilizes the semantic label to calculate the cosine similarity which is preserved in hamming space. Supervised Matrix Factorization Hashing (SMFH) [24] integrates the graph regularization and matrix factorization into an overall hashing learning framework. Semantics-Preserving Hashing (SePH) [17] transforms the affinity matrix into a probability distribution and approximates it in Hamming space via minimizing their Kullback-Leibler divergence. Generalized Semantic Preserving Hashing (GSePH) [25] preserves the semantic similarity by the unified binary codes. Semi-supervised NMF (CPSNMF) [26] uses a constraint propagation approach to get more supervised information, which can greatly improve the retrieval performance. Cross-Modal Hamming Hashing (CMHH) [27] designs a pairwise focal loss to generate compact and highly concentrated hash codes. In spite that supervised hashing methods have achieved promising performance, they overly depend on massive labeled data. Fortunately, unsupervised cross-modal hashing methods can handle effectively the problem.

Unsupervised cross-modal hashing methods mainly explore the structure, distribution, correlation and geometry among data and make these information be preserved well in hamming space. Canonical Correlation Analysis (CCA) [5] learns a common space where the correlation between different two modalities is maximized. Inter-Media Hashing (IMH) [23] introduces inter-media consistency and intra-media consistency to discover a common Hamming space. Cross View Hashing (CVH) [19] extends the classical unimodal spectral hashing to the multi-modal scenario. Robust Cross-view Hashing(RCH) [21] learns a common Hamming space where the binary codes representing the same semantic content but different modalities should be as consistent as possible. Collective Reconstructive Embeddings (CRE) [28] directly learns the unified binary codes via reconstructive embeddings collectively. Robust and Flexible Discrete Hashing(RFDH) [29] adopts the discrete matrix decomposition to learn the binary codes, which avoids the

large quantization error caused by relaxation. Fusion Similarity Hashing(FSH) [22] constructs an undirected asymmetric graph to model the similarity among objects.

Different from the above approaches, we propose a sparse multi-modal hashing approach that explores the local manifold structure and the weakly semantic correlation to learn the robust hash functions. The $\ell_{2,1}$ -norm regularization is incorporated to select the discriminative and relevant features from multi-modal data simultaneously.

III. UNSUPERVISED MULTIMODAL HASHING

A. Notation and Problem Statement

Suppose that the training set $O = \{o_i\}_{i=1}^n$ contains n instances of image-text pair. $V = [v_1, v_2, \dots, v_n] \in R^{d_1 \times n}$ and $T = [t_1, t_2, \dots, t_n] \in R^{d_2 \times n}$ denote the image modality and text modality respectively. Each instance $o_i = (v_i, t_i)$ consists of an image $v_i \in R^{d_1}$ and a text $t_i \in R^{d_2}$. Without loss of generality, samples in each modality are zero-centered, i.e. $\sum_i v_i = 0$ and $\sum_i t_i = 0$. Given the code length r , all instance O can be represented by the binary codes $B = [b_1; b_2; \dots; b_n] \in R^{n \times r}$ in hamming space. We first calculate the kernel matrices $X^{(m)} = [x_1^{(m)}, x_2^{(m)}, \dots, x_n^{(m)}] \in R^{d \times n}$ ($m = 1, 2$) of the m -th modality by employing the RBF kernel function. Taking the image modality for an example, $x_i^{(1)} = [exp(\|v_i - a_1\|^2/\sigma), \dots, exp(\|v_i - a_d\|^2/\sigma)]^T$, where $\{a_j\}_{j=1}^d$ are d anchor points that are randomly selected from the image modality of the training data. The aims of our method is to learn the mapping functions from the kernel spaces to the common Hamming space, that is, $f : R^d \rightarrow \{1, -1\}^r$ for image-modality and $g : R^d \rightarrow \{1, -1\}^r$ for text-modality.

Notations. Given an example matrix M and its i -th row is M_i , the $\ell_{2,1}$ -norm of M is defined as $\|M\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^m M_{ij}^2}$. $sgn(\cdot)$ signifies the sign function, specifically,

$$sgn(x) = \begin{cases} -1 & x < 0 \\ 1 & x \geq 0 \end{cases} \quad (1)$$

B. Structure Preservation

The local manifold structure in the original space should be preserved in the Hamming space. In the visual space, data point can be well approximated by the linear combination of its k -nearest neighbor points. In multi-modal applications, the text content associated with an image can provide the weakly supervised semantic information. Thus, the semantic correlation of textual space should be considered in the process of the hashing learning.

1) *Visual Model:* We hope that the similar locally manifold structure within visual modality can be projected into the same hash bin, and vice versa. The Discrete Locally Linear Embedding (DLLE) [13] is employed to preserve the local linear structure in the discrete Hamming space. The reconstruction error is written as follows

$$\min_S \frac{1}{2} \sum_{i=1}^n \|x_i^{(1)} - \sum_{j \neq i} S_{ji} x_j^{(1)}\|^2 \quad (2)$$

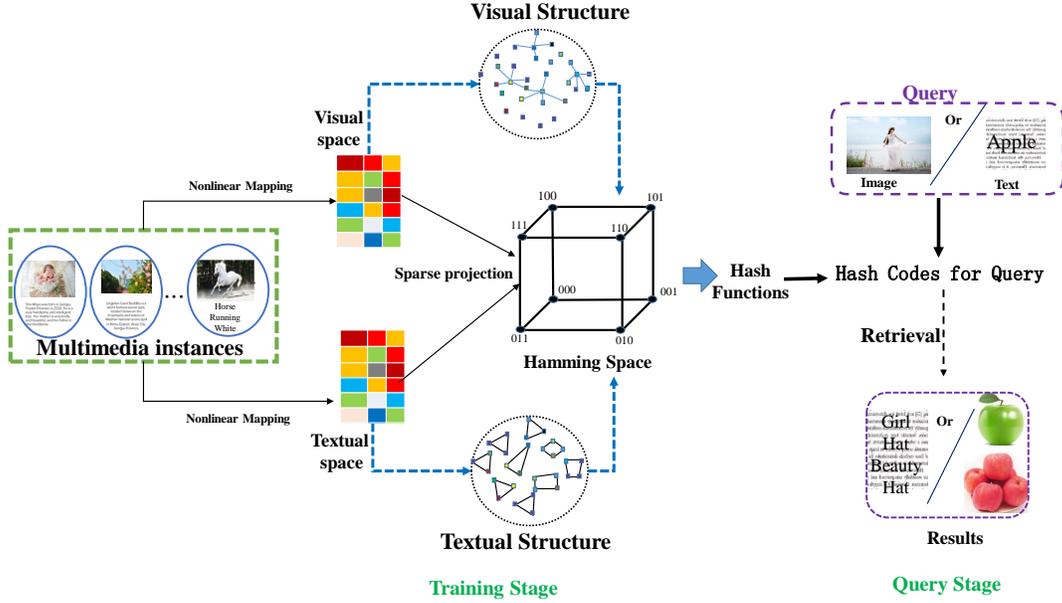


Fig. 1: Illustration of the proposed approach. The proposed framework finds a discrete hamming space where the local geometric structure of visual space and the semantic correlation information provided by textual modality can be preserved simultaneously. In the query phase, we can obtain the hash codes of an arbitrary query according to the learned hash functions, and other modal data with the nearest hamming distance are returned. Best viewed in color.

where $S \in R^{n \times n}$ is an affinity matrix. The optimal solution can be obtained as follows

$$S_i = \frac{G_i^{-1} \mathbf{1}}{\mathbf{1}^T G_i^{-1} \mathbf{1}} \quad (3)$$

where $X_j^{(1)}$ and $X_l^{(1)}$ are K-nearest neighbor points of $X_i^{(1)}$ and G_i is the local Gram matrix defined as $(X_i^{(1)} - X_j^{(1)})^T (X_i^{(1)} - X_l^{(1)})$ for $X_i^{(1)}$. Each point can be approximated by an affine combination of its K nearest neighbors in Hamming space. The reconstruction error as

$$\min_B \|B - SB\|^2 \quad (4)$$

2) *Textual Model*: Texts associated with social images are usually provided by web users. These tagged texts with rich semantic information are beneficial to hash functions learning. The pairwise similarity matrix $Z \in R^{n \times n}$ is calculated using the cosine similarity function. The textual affinity between the o_i and o_j is defined as follows

$$Z_{ij} = \frac{(x_i^{(2)})^T x_j^{(2)}}{\|x_i^{(2)}\|_2 \|x_j^{(2)}\|_2} \quad (5)$$

The higher textual similarity two instances have, the more similar binary codes they have. The above idea can be transformed the following weighted maximization problem

$$\begin{aligned} L(B, X^{(2)}) &= \arg \max_B \sum_{i,j=1}^n Z_{ij} b_i b_j^T \\ &= \arg \max_B \text{Tr}(B^T Z B) \end{aligned} \quad (6)$$

The embedding scheme in Eq. (6) is termed Discrete Locally Projection Preservation (DLPP) in this paper. To achieve

maximal information entropy, each hash bit is expected to be balanced on the training data [20]. More specifically, the number of +1 and that of -1 should be consistent as much as possible for each bit. We integrate Eq. (5) and Eq. (6) into the following Eq. (7) to obtain compact binary codes.

$$\begin{aligned} \min_B \|B - SB\|^2 - \beta \text{Tr}(B^T Z B) + \rho \|\mathbf{1}_n^T B\|_F^2 \\ \text{s.t. } B \in \{-1, 1\}^{n \times r} \end{aligned} \quad (7)$$

C. Hash Functions Learning

The $\ell_{2,1}$ -norm has been proven to be effective to obtain the discriminative features by some recent works [31], [32]. We impose the $\ell_{2,1}$ -norm constraint on the projection matrices to learn the discriminative representation, which leads to the following problem

$$\begin{aligned} \min_{P^{(m)}, B} \sum_{m=1}^2 \alpha^{(m)\gamma} (\|X^{(m)T} P^{(m)} - B\|_F^2 + \lambda_m \|P^{(m)}\|_{2,1}) \\ \text{s.t. } B \in \{-1, 1\}^{n \times r}, \sum_{m=1}^2 \alpha_m = 1, \alpha_m > 0 \end{aligned} \quad (8)$$

where $P^{(m)}$ denotes the projection matrix of the m -th modality and $\alpha^{(m)}$ is the weight factor with the adjustment coefficient γ and λ_m is a penalty parameter.

Then the overall objective function combining Eq. (7) and Eq. (8) is given as follows

$$\begin{aligned} \min_{B, P^{(m)}, \alpha^{(m)}} & \sum_{m=1}^2 \alpha^{(m)\gamma} (\|X^{(m)T} P^{(m)} - B\|_F^2 + \lambda_m \|P^{(m)}\|_{2,1}) \\ & + \eta \|B - SB\|_F^2 - \beta \text{Tr}(B^T ZB) + \rho \|\mathbf{1}_n^T B\|_F^2 \\ \text{s.t. } & B \in \{-1, 1\}^{n \times r}, \sum_{m=1}^2 \alpha_m = 1, \alpha_m > 0 \end{aligned} \quad (9)$$

where β and ρ are two hyper-parameters.

The above Eq. (9) is a non-convex problem. We solve the optimization problem by updating each variable with the other variables fixed alternatively.

Update B with other variables fixed. The subproblem is to minimize the following

$$\begin{aligned} \min_B & \sum_{m=1}^2 \alpha^{(m)\gamma} \|X^{(m)T} P^{(m)} - B\|_F^2 + \eta \|B - SB\|_F^2 \\ & - \beta \text{Tr}(B^T ZB) + \rho \|\mathbf{1}_n^T B\|_F^2 \\ \text{s.t. } & B \in \{-1, 1\}^{n \times r} \end{aligned} \quad (10)$$

The Eq.(10) is an NP-hard problem since B is constrained to be discrete value. We relax it to be continuous value H . Thus the optimization problem can be transformed to

$$\begin{aligned} \min_{H, B} & -2\text{Tr}(R^T H) + \text{Tr}(H^T H) + \eta \|CH\|_F^2 - \beta \text{Tr}(H^T ZH) \\ & + \rho \|\mathbf{1}_n^T H\|_F^2 + \xi \|H - B\|_F^2 \end{aligned} \quad (11)$$

where $R = \sum_{m=1}^2 \alpha^{(m)\gamma} X^{(m)T} P^{(m)}$ and $C = S - I$. Then we can get

$$H = (\eta C^T C - \beta Z + \rho \mathbf{1}_n \mathbf{1}_n^T + (\xi + 1)I)^{-1} (R + \xi B) \quad (12)$$

The problem with respect to B can be presented as

$$\begin{aligned} \max_B & \text{tr}(HB^T) \\ \text{s.t. } & B \in \{-1, 1\}^{n \times r} \end{aligned} \quad (13)$$

The solution of B can be directly obtained as

$$B = \text{sgn}(H) \quad (14)$$

Update $P^{(m)}$ with other variables fixed. Keeping terms relating to $P^{(m)}$, the objective function Eq. (9) can be rewritten as follows

$$\min_{P^{(m)}} \|X^{(m)T} P^{(m)} - B\|^2 + \lambda_m \|P^{(m)}\|_{2,1} \quad (15)$$

Setting the derivative of Eq. (15) with respect to $P^{(m)}$ to zero, we can obtain

$$P^{(m)} = (X^{(m)} X^{(m)T} + \lambda_{(m)} D^{(m)})^{-1} X^{(m)} B \quad (16)$$

where $D^{(m)}$ is a diagonal matrix with the i -th diagonal element $D_{ii}^{(m)} = \frac{1}{2\|P_i^{(m)}\|_2 + \epsilon}$, and $P_i^{(m)}$ signifies the i -th row of $P^{(m)}$.

Update weight $\alpha^{(m)}$ with other variables fixed. By dropping terms irrelating to $\alpha^{(m)}$, we get

$$\begin{aligned} \min_{\alpha^{(m)}} & \sum_{m=1}^2 \alpha^{(m)\gamma} C^{(m)} \\ \text{s.t. } & \sum_{m=1}^2 \alpha_m = 1, \alpha_m > 0 \end{aligned} \quad (17)$$

where $C^{(m)} = \|X^{(m)T} P^{(m)} - B\|_F^2 + \lambda_m \|P^{(m)}\|_{2,1}$. We employ the Lagrange multiplier to transform Eq. (17) into the following

$$\min_{\alpha^{(m)}} \sum_{m=1}^2 \alpha^{(m)\gamma} C^{(m)} + \xi (1 - \sum_{m=1}^2 \alpha^{(m)}) \quad (18)$$

Setting the derivative of Eq. (18) with respect to $\alpha^{(m)}$ to zero, we obtain

$$\alpha^{(m)} = \frac{(\gamma C^{(m)})^{1/(1-\gamma)}}{\sum_{m=1}^M (\gamma C^{(m)})^{1/(1-\gamma)}} \quad (19)$$

After acquiring the projection matrix $P^{(m)}$, the binary codes b of query x is computed according to the rule $b = \text{sgn}(x^T P^{(m)})$. The overall optimization procedure is summarized in Algorithm 1. This iteration process is repeated until it converges. As shown in Fig. 4, our algorithm converges quickly on the WiKi, PASCAL-VOC and UCI Handwritten Digit.

Algorithm 1 Unsupervised Multi-modal Hashing

Input: $X^{(m)} \in R^{d \times n}$, ($m = 1, 2$); hash codes length r .

Output: $P^{(m)}$, B , $\alpha^{(m)}$.

Initialize B , $P^{(m)}$, $\alpha^{(m)}$, $\lambda_{(m)}$, ρ , β and η .

Compute similarity matrix S according to (3)

Compute similarity matrix Z according to (5)

- 1: **repeat**
 - 2: Update B according to (12) and (14)
 - 3: Compute $D_{ii}^{(1)}$ by $D_{ii}^{(1)} = \frac{1}{2\|P_i^{(1)}\|_2 + \epsilon}$.
 - 4: Compute $D_{ii}^{(2)}$ by $D_{ii}^{(2)} = \frac{1}{2\|P_i^{(2)}\|_2 + \epsilon}$.
 - 5: Update $P^{(1)}$ using Eq.(16)
 - 6: Update $P^{(2)}$ using Eq.(16)
 - 7: Update $\alpha^{(1)}$ according to Eq.(19)
 - 8: Update $\alpha^{(2)}$ according to Eq.(19)
 - 9: **until**
-

IV. EXPERIMENTS

A. Datasets

Wiki [33] contains 2,866 multimedia documents harvested from Wikipedia. Every document consists of an image and a text description, and each document is classified into one of 10 categories. Each image is represented by a 128-dimensional SIFT histogram vector. A 10-dimensional feature vector generated by latent Dirichlet allocation is used to represent each text. We take 2173 pairs from the dataset to form the training set and database, the resting 973 as a query set.

PASCAL-VOC [34] consists of 9,963 image-tag pairs.

Each image is represented by a 512-dimensional GIST feature vector and each text is represented as a 399-dimensional word frequency count. Each sample are classified into one of 20 categories. We select 5,649 pairs with only one object in our experiment. 2,808 pairs are taken out as a training set and database, the remaining samples as the query data.

UCI Handwritten Digit is comprised of handwritten numerals(0 - 9) collected from Dutch utility maps. Each of the character shapes is regarded as a class and each class consists of 200 samples. Following [35], we select 76 Fourier coefficients and 64 Karhunen-Love coefficients of the character shapes as the feature of two different modalities respectively. 1,500 samples are treated as the training set and database, the resting 500 as the query set.

B. Experimental Setting

To verify the effectiveness of our method, some comparative experiments are conducted on two cross-modal retrieval tasks: Image (Modality 1) query text (Modality 2) database and Text (Modality 2) query image (Modality 1) database which are termed as 'Task1' and 'Task2' respectively. As our method is a unsupervised hashing method, for a fair comparison, we compare our method with six state-of-the-art unsupervised learning models. Specifically, the baselines include CVH [19], CCA [5], IMH [23], RCH [21], FSH [22] and CRE [28]. Since the source code of RCH and CRE is not available, we implemented it by ourselves. The codes of other baselines are kindly provided by the authors. The value of λ_1 , λ_2 , β and ρ are tuned in the candidate range $\{1e^{-5}, 1e^{-4}, 1e^{-3}, 1e^{-2}, 1e^{-1}, 1\}$. γ is set to 0.5 empirically and the best results are reported in this paper. Our experiments are implemented on MATLAB 2016b and Windows 10 (64-Bit) platform based on desktop machine with 12 GB memory and 4-core 3.6GHz CPU, and the model of the CPU is Intel(R) CORE(TM) i7-7700.

TABLE I: The mAP results on WiKi

Tasks	Methods	The length of hash code			
		16	32	64	128
Task 1	CVH	0.1499	0.1408	0.1372	0.1323
	CCA	0.1699	0.1519	0.1495	0.1472
	IMH	0.2022	0.2127	0.2164	0.2171
	RCH	0.2102	0.2234	0.2397	0.2497
	FSH	0.2346	0.2491	0.2531	0.2573
	CRE	0.2301	0.2446	0.2599	0.2620
	UMH	0.2511	0.2505	0.2578	0.2611
Task 2	CVH	0.1315	0.1171	0.1080	0.1093
	CCA	0.1587	0.1392	0.1272	0.1211
	IMH	0.1648	0.1703	0.1737	0.1720
	RCH	0.2171	0.2497	0.2825	0.2973
	FSH	0.2149	0.2241	0.2332	0.2368
	CRE	0.2442	0.2695	0.2846	0.2897
	UMH	0.4984	0.5057	0.5224	0.5298
Average	CVH	0.1407	0.1290	0.1226	0.1208
	CCA	0.1643	0.1456	0.1384	0.1341
	IMH	0.1835	0.1915	0.1951	0.1946
	RCH	0.2137	0.2365	0.2611	0.2735
	FSH	0.2248	0.2366	0.2431	0.2470
	CRE	0.2372	0.2571	0.2723	0.2759
	UMH	0.3747	0.3781	0.3901	0.3955

TABLE II: The mAP results on PASCAL-VOC

Tasks	Methods	The length of hash code			
		16	32	64	128
Task 1	CVH	0.1484	0.1187	0.1651	0.1411
	CCA	0.1245	0.1267	0.1230	0.1218
	IMH	0.2087	0.2016	0.1873	0.1718
	RCH	0.2633	0.3013	0.3209	0.3330
	FSH	0.2890	0.3173	0.3340	0.3496
	CRE	0.2758	0.3046	0.3216	0.3270
	UMH	0.3225	0.3368	0.3741	0.3701
Task 2	CVH	0.0931	0.0945	0.0978	0.0918
	CCA	0.1283	0.1362	0.1465	0.1553
	IMH	0.1631	0.1558	0.1537	0.1464
	RCH	0.2145	0.2656	0.3275	0.3983
	FSH	0.2617	0.3030	0.3216	0.3428
	CRE	0.2395	0.2713	0.2941	0.2981
	UMH	0.4760	0.5472	0.5825	0.5701
Average	CVH	0.1208	0.1066	0.1315	0.1165
	CCA	0.1264	0.1315	0.1347	0.1386
	IMH	0.1859	0.1787	0.1705	0.1591
	RCH	0.2389	0.2834	0.3242	0.3657
	FSH	0.2753	0.3102	0.3278	0.3462
	CRE	0.2577	0.2880	0.3079	0.3126
	UMH	0.3993	0.4420	0.4783	0.4701

TABLE III: The mAP results on UCI Handwritten Digit

Tasks	Methods	The length of hash code			
		16	32	64	128
Task 1	CVH	0.3421	0.2496	0.1907	0.1759
	CCA	0.3155	0.2360	0.1841	0.2082
	IMH	0.2947	0.2375	0.1892	0.1737
	RCH	0.6181	0.6636	0.6991	0.7056
	FSH	0.6323	0.6776	0.7027	0.7139
	CRE	0.6636	0.7425	0.7516	0.7643
	UMH	0.7496	0.7944	0.8149	0.8043
Task 2	CVH	0.3215	0.2471	0.1939	0.1695
	CCA	0.3160	0.2398	0.1855	0.1102
	IMH	0.2943	0.2315	0.1789	0.1514
	RCH	0.5810	0.6336	0.6768	0.6979
	FSH	0.6460	0.6745	0.7069	0.7149
	CRE	0.6448	0.7357	0.7547	0.7671
	UMH	0.7327	0.7997	0.8333	0.8417
Average	CVH	0.3318	0.2483	0.1923	0.1727
	CCA	0.3157	0.2379	0.1848	0.1592
	IMH	0.2945	0.2345	0.1840	0.1626
	RCH	0.5996	0.6486	0.6880	0.7017
	FSH	0.6392	0.6761	0.7048	0.7144
	CRE	0.6542	0.7391	0.7532	0.7657
	UMH	0.7411	0.7970	0.8241	0.8230

C. Evaluation metric

The Mean Average Precision (mAP) is used to evaluate the performance of our method and comparison methods. Specifically, the Average Precision (AP) for a query q is defined as follows

$$AP(q) = \frac{1}{l_q} \sum_{m=1}^R P_q(m) \delta_q(m) \quad (20)$$

where $P_q(m)$ denotes the accuracy of top m retrieval results; $\delta_q(m) = 1$ if the m -th position is true neighbor of the query q , and otherwise $\delta_q(m) = 0$; l_q is the correct statistics of top R retrieval results. (R is set to the number of entire database). The mAP is defined as the mean of the average precisions of all queries

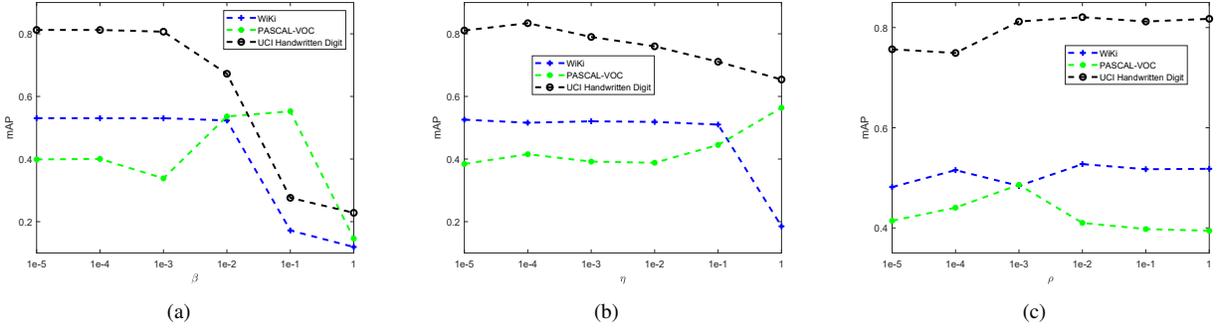


Fig. 2: The mAP variation with respect to β , η and ρ .

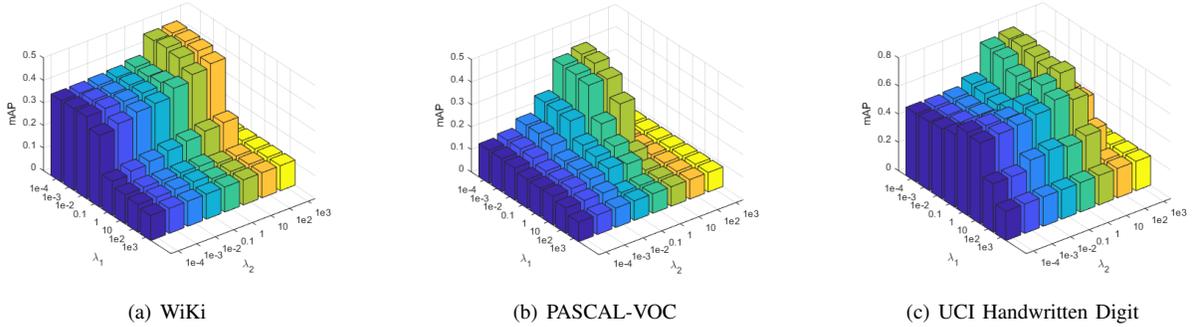


Fig. 3: The mAP variation with respect to different combination of λ_1 and λ_2 on Wiki(a), PASCAL-VOC (b), and UCI Handwritten Digit (c).

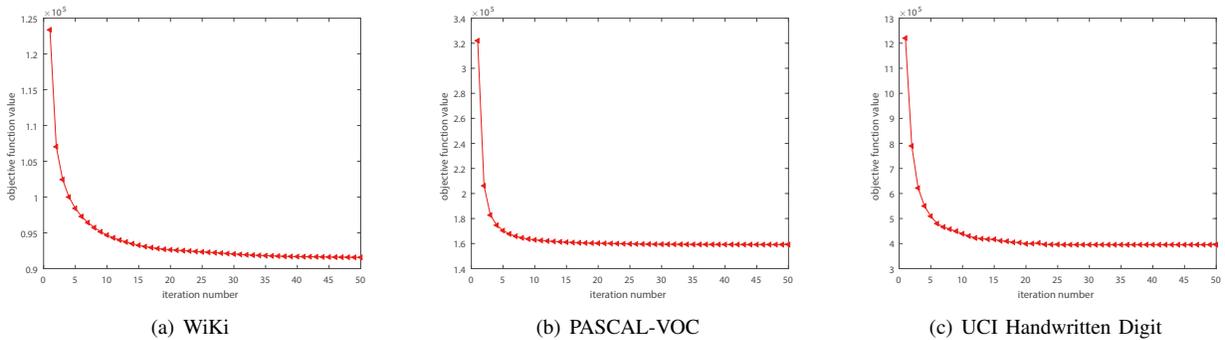


Fig. 4: The convergence curve of algorithm 1 on Wiki(a), PASCAL-VOC (b), and UCI Handwritten Digit (c).

D. Retrieval Performance Evaluation

The mAP scores on Wiki, PASCAL-VOC, and UCI Handwritten Digit are shown in TABLE I, II, and III respectively. We can observe the following points: (1) The performance of our method is superior to the baselines. Among baselines, RCH is a method with $\ell_{2,1}$ -norm constraint imposing on projection matrices, but RCH does not take the manifold structure within each modality into account. FSH constructs an undirected asymmetric graph to model the similarity among samples. CRE utilizes domain-specific method to model different modalities and the intra-modal similarity is preserved in the process of learning unified binary codes. However, FSH and CRE do not explore the discriminative features in their

frameworks. Compared with the above methods, our method boosts retrieval performance. The significant improvement of the proposed method can be attributed to the combination of the $\ell_{2,1}$ -norm regularization and Structure Preserving. (2) Our method outperforms all comparison methods in terms of the average performance for two retrieval tasks on three datasets. With the increasing of hash code length, the retrieval performance on the Task 1 and Task 2 is further improved. The reason for the better performance is that the discriminative information will be more sufficient with the longer hash codes. (4) The results on Task 2 are consistently higher than that on Task 1. This may be because the text modality itself is a weak supervision information which can benefit to improve retrieval

TABLE IV: Experiment results (mAP@64bit) on ablation study

	WiKi	Pascal VOC	UCI Handwritten Digit
$\beta = 0$	0.3801	0.3531	0.8014
$\eta = 0$	0.3850	0.3056	0.7998
$\rho = 0$	0.1102	0.4542	0.2082
Ours	0.3901	0.4783	0.8241

performance.

Ablation study Some ablation experiments are conducted to investigate the influence of different terms in Eq. (9). Three hyper-parameters (β, η, ρ) steer one of the terms respectively. $\beta = 0$ means our method ignores the inner structure within text modality. $\eta = 0$ indicates our model does not consider the visually geometric information. $\rho = 0$ implies the balanced bits term is ruled out from our model. The comparison results are shown in Table IV. From Table IV, the importance of the three terms is dissimilar for different datasets. It is apparent that each term of the objective function collaboratively contributes to the retrieval performance.

E. Parameter Sensitivity Analysis

In our model, $\rho, \beta, \eta, \lambda_1$ and λ_2 are set manually. In this subsection, we explore the influence of different parameters setting on retrieval performance. The empirical analysis is performed for each parameter by varying its value in the candidate range. To discuss the above parameters conveniently, the hash code length is fixed at 64 bit in our experiments. In Fig.2, we plot the performance variation curves with respect to β, η and ρ . On WiKi, PASCAL-VOC and UCI Handwritten Digit, our method can achieve the highest mAP score when β is set to $1e^{-5}, 1e^{-1}$ and $1e^{-6}$ respectively and η is set to $1e^{-1}, 1$ and $1e^{-4}$ respectively. ρ represents the importance of the balanced term in Eq. (9). In Fig.2, we can find that the ρ should not be too large. λ_1 and λ_2 are two penalty parameters controlling the sparse constraint items of two modalities respectively. The mAP scores as a function of λ_1 and λ_2 is plotted in Fig.3, which shows that the optimal combination falls a fixed small range on three datasets.

V. CONCLUSION

In this paper, we propose a unsupervised multi-modal hashing method for cross-modal retrieval. Our model explores the underlying neighborhood structure of the visual space and the semantic correlation provided by textual modality to learn the compact unified hash codes. The sparse constraint is imposed on our model to learn discriminative hash functions for multi-modal data. Encouraging experimental results demonstrate that the effectiveness of the proposed framework on cross-modal retrieval tasks. In the future, we plan to extend the proposed method into the deep learning networks.

ACKNOWLEDGMENT

THE PAPER IS SUPPORTED BY THE NATIONAL NATURAL SCIENCE FOUNDATION OF CHINA (GRANT NO.61672265, U1836218), AND THE 111 PROJECT OF MINISTRY OF EDUCATION OF CHINA (GRANT NO. B12018).

REFERENCES

- [1] L. Lin, G. Wang, W. Zuo, X. Feng, and L. Zhang, Cross-domain visual matching via generalized similarity measure and feature learning, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1089-1102, 2017.
- [2] Y. Guo, G. Ding, L. Liu, J. Han, and L. Shao, Learning to hash with optimized anchor embedding for scalable retrieval, *IEEE Transactions on Image Processing*, vol. 26, no. 3, pp. 1344-1354, 2017.
- [3] Xiao W, Shi G, Li B, et al. Fast hash-based inter-block matching for screen content coding[J]. *IEEE Transactions on Circuits and Systems for Video Technology*, 2016, 28(5): 1169-1182.
- [4] C. Wu, J. Zhu, D. Cai, C. Chen, and J. Bu, Semi-supervised nonlinear hashing using bootstrap sequential projection learning, *IEEE Transactions on Knowledge and Data Engineering*, vol. 25, no. 6, pp. 1380-1393, 2013.
- [5] Gong Y, Lazebnik S, Gordo A, et al. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval[J]. *IEEE transactions on pattern analysis and machine intelligence*, 2012, 35(12): 2916-2929.
- [6] Wang J, Kumar S, Chang S F. Semi-supervised hashing for scalable image retrieval[C]//2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. IEEE, 2010: 3424-3431.
- [7] Norouzi, M. and D.J. Fleet. Minimal loss hashing for compact binary codes. in *International Conference on International Conference on Machine Learning* 2011.
- [8] Yorozu T, Hirano M, Oka K, et al. Electron spectroscopy studies on magneto-optical media and plastic substrate interface[J]. *IEEE translation journal on magnetics in Japan*, 1987, 2(8): 740-741.
- [9] Li S, Chen Z, Lu J, et al. Neighborhood Preserving Hashing for Scalable Video Retrieval[C]//Proceedings of the IEEE International Conference on Computer Vision. 2019: 8212-8221.
- [10] Cui H, Zhu L, Li J, et al. Scalable deep hashing for large-scale social image retrieval[J]. *IEEE Transactions on Image Processing*, 2019, 29: 1271-1284.
- [11] Lin M, Ji R, Chen S, et al. Similarity-Preserving Linkage Hashing for Online Image Retrieval[J]. *IEEE Transactions on Image Processing*, 2020, 29: 5289-5300.
- [12] Tang J, Li Z. Weakly supervised multimodal hashing for scalable social image retrieval[J]. *IEEE Transactions on Circuits and Systems for Video Technology*, 2017, 28(10): 2730-2741.
- [13] Ji R, Liu H, Cao L, et al. Toward optimal manifold hashing via discrete locally linear embedding[J]. *IEEE Transactions on Image Processing*, 2017, 26(11): 5411-5420.
- [14] Wang D, Gao X B, Wang X, et al. Label consistent matrix factorization hashing for large-scale cross-modal similarity search[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [15] Yu J, Wu X J, Kittler J. Learning discriminative hashing codes for cross-modal retrieval based on multi-view features[J]. *Pattern Analysis and Applications*, 2020.
- [16] Zhang D, Li W J. Large-scale supervised multimodal hashing with semantic correlation maximization[C]//Twenty-Eighth AAAI Conference on Artificial Intelligence. 2014.
- [17] Z. Lin, G. Ding, M. Hu, and J. Wang, Semantics-preserving hashing for cross-view retrieval, in *Proceedings of the 28th International Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3864-3872.
- [18] J. Yu, X. Wu and J. Kittler, "Semi-supervised Hashing for Semi-Paired Cross-View Retrieval," 2018 24th International Conference on Pattern Recognition (ICPR), Beijing, 2018, pp. 958-963.
- [19] S. Kumar and R. Udupa, Learning hash functions for cross-view similarity search, in *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*, vol. 22, no. 1. AAAI Press, 2011, pp. 1360-1367.
- [20] Y. Weiss, A. Torralba, and R. Fergus, Spectral hashing, in *Advances in Neural Information Processing Systems. NIPS*, 2009, pp. 1753-1760.
- [21] Shen, X., et al., Robust Cross-view Hashing for Multimedia Retrieval. *IEEE Signal Processing Letters*, 2016. 23(6): p. 893-897.
- [22] Liu, H., et al. Cross-modality binary code learning via fusion similarity hashing. in *Proceedings of CVPR*. 2017.
- [23] Song, J., et al. Inter-media hashing for large-scale retrieval from heterogeneous data sources. in *ACM SIGMOD International Conference on Management of Data*. 2013.
- [24] Liu H, Ji R, Wu Y, et al. Supervised matrix factorization for cross-modality hashing[C]//Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence. 2016: 1767-1773.
- [25] Mandal, Devraj, Kunal N. Chaudhury, and Soma Biswas. "Generalized semantic preserving hashing for n-label cross-modal retrieval." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017.

- [26] Wang D, Gao X, Wang X. Semi-supervised nonnegative matrix factorization via constraint propagation[J]. *IEEE transactions on cybernetics*, 2015, 46(1): 233-244.
- [27] Cao Y, Liu B, Long M, et al. Cross-Modal Hamming Hashing[C]//*Proceedings of the European Conference on Computer Vision (ECCV)*. 2018: 202-218.
- [28] Hu M , Yang Y , Shen F , et al. Collective Reconstructive Embeddings for Cross-Modal Hashing[J]. *IEEE Transactions on Image Processing*, 2019, 28(6):2770-2784.
- [29] Wang D, Wang Q, Gao X. Robust and flexible discrete hashing for cross-modal similarity search[J]. *IEEE Transactions on Circuits and Systems for Video Technology*, 2017, 28(10): 2703-2715.
- [30] Yu J, Wu X J, Kittler J. Discriminative Supervised Hashing for Cross-Modal Similarity Search[J]. *Image and Vision Computing*, 2019, 89: 50-56.
- [31] Wang K , He R , Wang L , et al. Joint Feature Selection and Subspace Learning for Cross-modal Retrieval[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016:1-14.
- [32] Chen Z , Wu X J , Yin H F , et al. Noise-Robust Dictionary Learning with Slack Block-Diagonal Structure for Face Recognition[J]. *Pattern Recognition*, 2019, 100:107118.
- [33] Rasiwasia, N., et al. A new approach to cross-modal multimedia retrieval. in *International Conference on Multimedia*. 2010.
- [34] Hwang, S.J. and K. Grauman, Reading between the lines: Object localization using implicit cues from image tags. *IEEE transactions on pattern analysis and machine intelligence*, 2012. 34(6): p. 1145-1158.
- [35] He R, Zhang M, Wang L, et al. Cross-modal subspace learning via pairwise constraints[J]. *IEEE Transactions on Image Processing*, 2015, 24(12): 5543-5556.