

# Semi-Supervised Deep Hashing with a Bipartite Graph

Xinyu Yan, Lijun Zhang, Wu-Jun Li

National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China  
 {yanxy, zhanglj}@lamda.nju.edu.cn, liwujun@nju.edu.cn

## Abstract

Recently, deep learning has been successfully applied to the problem of hashing, yielding remarkable performance compared to traditional methods with hand-crafted features. However, most of existing deep hashing methods are designed for the supervised scenario and require a large number of labeled data. In this paper, we propose a novel semi-supervised hashing method for image retrieval, named Deep Hashing with a Bipartite Graph (BGDH), to simultaneously learn embeddings, features and hash codes. More specifically, we construct a bipartite graph to discover the underlying structure of data, based on which an embedding is generated for each instance. Then, we feed raw pixels as well as embeddings to a deep neural network, and concatenate the resulting features to determine the hash code. Compared to existing methods, BGDH is a universal framework that is able to utilize various types of graphs and losses. Furthermore, we propose an inductive variant of BGDH to support out-of-sample extensions. Experimental results on real datasets show that our BGDH outperforms state-of-the-art hashing methods.

## 1 Introduction

With the explosion in the volume of image data, it has been raised as a big challenge about how to index and organize these data efficiently and accurately. Approximate Nearest Neighbor (ANN) search [Indyk and Motwani, 1998] has become a popular way to retrieve content from images with both computational efficiency and search quality. Among existing ANN search methods, hashing is advantageous due to its fast query speed and low memory complexity [Gionis *et al.*, 1999; Gong *et al.*, 2013]. It aims to transform high-dimensional images into a set of short binary codes while maintaining similarity of the original data.

Generally speaking, hashing methods can be divided into two categories: unsupervised and supervised. Unsupervised methods utilize some kinds of distance metrics to learn a hash function from unlabeled data. Methods in this category include data-independent ones like Locally Sensitive Hashing

(LSH) [Gionis *et al.*, 1999] and data-dependent ones like Iterative Quantization (ITQ) [Gong *et al.*, 2013], Spectral Hashing (SH) [Weiss *et al.*, 2009], Anchor Graph Hashing (AGH) [Liu *et al.*, 2011]. On the other hand, to deal with more complicated semantic similarity, supervised hashing methods are proposed to exploit label information to improve the hashing quality. Representative supervised methods include Latent Factor Hashing (LFH) [Zhang *et al.*, 2014], Fast Supervised Hashing (FastH) [Lin *et al.*, 2014], Supervised Discrete Hashing (SDH) [Shen *et al.*, 2015]. However, labeling large-scale image dataset is inefficient and time-consuming. As a result, Semi-Supervised Hashing (SSH) [Wang *et al.*, 2012] has been developed to make use of labeled data as well as the abundant unlabeled data.

In traditional hashing methods, images are represented by hand-crafted features such as GIST [Oliva and Torralba, 2001], and the choice of features requires heavy manual interventions. Motivated from the great success of deep neural networks in image analysis [Krizhevsky *et al.*, 2012], recently some deep hashing methods have been proposed to learn features and hash codes simultaneously [Li *et al.*, 2016; Zhu *et al.*, 2016; Liu *et al.*, 2016]. Although those deep hashing methods yield better performance compared with the traditional methods, they usually need a large number of labeled instances as training data. To address this limitation, a semi-supervised deep hashing, named SSDH, have been developed [Zhang *et al.*, 2016]. SSDH is fundamentally built upon graph-based semi-supervised learning [Zhou *et al.*, 2004] and the loss function contains a graph regularization term which involves both the labeled and unlabeled data. In theory, SSDH needs to construct a nearest neighbor graph of all the data. Unfortunately, this step takes  $O(n^2)$  time, where  $n$  is the number of instances, and thus intractable for large scale data.

In this paper, we propose a novel semi-supervised hashing method, named Deep Hashing with a Bipartite Graph (BGDH), which performs graph embedding, feature learning and hash code learning in a unified framework. First, we construct a bipartite graph to capture the information hidden in the labeled and unlabeled data. The bipartite graph could be a semantic graph that describes relationships between images and concepts, an anchor graph that describes similarities between images and landmarks [Liu *et al.*, 2010], or a traditional nearest neighbor graph. Then, inspired by the recent work on graph embedding [Yang *et al.*, 2016], we learn an embed-

ding for each instance to predict the neighborhood context in the graph. Finally, we feed both raw pixels and embeddings to a deep neural network, and concatenate the corresponding hidden layers when producing binary codes. BGDH is a general learning framework in the sense that any loss function of hashing and any type of graph can be incorporated.

Graph-based methods are usually transductive, because they can only handle instances that are already appeared in the graph. Since embeddings of instances are learnt from the graph, the basic BGDH is also transductive. To address the out-of-sample problem, we further propose an inductive variant of BGDH, in which the embeddings are defined as a parametric function of the raw features. In this way, we can produce hash codes for new instances that have not seen during training. To demonstrate the effectiveness of our approach, we conduct extensive experiments on two large-scale datasets: CIFAR-10 and NUS-WIDE. Experimental results show that BGDH outperforms other methods and achieves the state-of-the-art performance in image retrieval.

Finally, we emphasize that although both BGDH and SSDH are semi-supervised deep hashing methods, the proposed BGDH differs from SSDH in the following two aspects:

- a. While SSDH is built upon graph regularization, our BGDH relies on graph embedding.
- b. SSDH uses graphs to exploit the unlabeled data, in contrast BGDH makes use of bipartite graphs, which can be constructed more efficiently since building an anchor graph only costs  $O(n)$  time.

## 2 Notations and Problem Definitions

In this section, we introduce notations and problem definitions.

### 2.1 Notations

We use script letters like  $\mathcal{X}$  to denote sets, boldface lowercase letters like  $\mathbf{e}$  to denote vectors and boldface uppercase letters like  $\mathbf{E}$  to denote matrices. We denote the element at the  $i$ -th row and  $j$ -th column of  $\mathbf{E}$  by  $E_{ij}$ .  $\mathbf{E}^T$  is the transpose of  $\mathbf{E}$ .  $\|\cdot\|_2$  denotes the Euclidean norm of a vector and  $\|\cdot\|_F$  denotes the Frobenius norm of a matrix.  $\text{sgn}(\cdot)$  is an element-wise sign function and  $\sigma(\cdot)$  is the sigmoid function.  $[\mathbf{u}; \mathbf{v}]$  denotes the concatenation of two vectors  $\mathbf{u}$  and  $\mathbf{v}$ .

### 2.2 Problem Definitions

Given a set of  $n$  instances/images  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$  where  $\mathbf{x}_i$  is the feature vector of the  $i$ -th instance. Without loss of generality, we assume the first  $l$  instances  $\{\mathbf{x}_1, \dots, \mathbf{x}_l\}$  are labeled and the rest are unlabeled. We assume the supervised information is given in term of pairwise labels, though our method can also support other kinds of labels. Specifically, for the first  $l$  instances, we have a set of pairwise labels  $\mathcal{S} = \{s_{ij}\}$  with  $s_{ij} \in \{0, 1\}$ , where  $s_{ij} = 1$  means that  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are similar,  $s_{ij} = 0$  means that  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are dissimilar. In addition, a bipartite graph  $\mathcal{G} = (\mathcal{X}, \mathcal{O}, \mathcal{E})$  between  $n$  instances and  $m$  objects are given, where  $\mathcal{O}$  is a set of objects such as concepts or landmarks, and  $\mathcal{E}$  is the set of edges.

The goal of our semi-supervised hashing method is to learn a mapping function  $\mathcal{H} : \mathbf{x}_i \rightarrow \{-1, 1\}^c$ , which encodes each

point  $\mathbf{x}_i$  into a  $c$ -dimensional binary code  $\mathbf{b}_i = \mathcal{H}(\mathbf{x}_i) = [h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_c(\mathbf{x}_i)]^T \in \{-1, 1\}^c$ . The binary codes  $\mathcal{B} = \{\mathbf{b}_i\}_{i=1}^n$  should preserve the semantic similarity and structure similarity in the Hamming space.

## 3 Semi-Supervised Deep Hashing with a Bipartite Graph

In this section, we first present the details of our semi-supervised Deep Hashing with a Bipartite Graph (BGDH), then introduce an inductive variant and finally discuss the learning procedure.

### 3.1 The Proposed BGDH Framework

The end-to-end deep learning architecture of our BGDH is shown in Figure 1, which includes three main components: graph embedding, feature learning, and hash code learning. Similar to other semi-supervised learning methods [Yang *et al.*, 2016], the loss function of BGDH can be expressed as

$$\mathcal{L}_s + \lambda \mathcal{L}_g \quad (1)$$

where  $\mathcal{L}_s$  is a supervised loss designed to preserve the similarity between pairwise instances, and  $\mathcal{L}_g$  is an unsupervised loss of predicting the graph context. In the following, we first introduce  $\mathcal{L}_g$  which aims to learn embeddings from the bipartite graph  $\mathcal{G}$ , then formulate  $\mathcal{L}_s$  which is used to learn both features and binary codes from hidden layers of deep networks.

#### Graph embedding

We propose to use a bipartite graph  $\mathcal{G} = (\mathcal{X}, \mathcal{O}, \mathcal{E})$  to capture the information hidden in the unlabeled data. It can be constructed in different ways as stated below.

- An anchor graph constructed from the dataset  $\mathcal{X}$ . In this case,  $\mathcal{O}$  contains  $m$  landmarks and the construction of  $\mathcal{G}$  takes  $O(n)$  time [Liu *et al.*, 2010].
- A nearest neighbor graph. In this case,  $\mathcal{O} = \mathcal{X}$  and the construction of  $\mathcal{G}$  takes  $O(n^2)$  time.
- A semantic graph constructed from external data. In this case,  $\mathcal{O}$  may contain  $m$  concepts, styles, or owners.

In the following, we briefly introduce one way to construct an anchor graph. We first randomly sample  $m$  instances from  $\mathcal{X}$  as landmarks, denoted by  $\mathcal{O} = \{\mathbf{o}_1, \dots, \mathbf{o}_m\}$ . Then, we put an edge between  $\mathbf{x}_i$  and  $\mathbf{o}_j$  if  $\mathbf{o}_j$  is among  $k$  nearest landmarks of  $\mathbf{x}_i$ , or if the distance between them is smaller than some threshold  $\epsilon$ . Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$  be the similarity matrix of  $\mathcal{G}$ , where  $A_{ij}$  denotes the weight of the edge between  $\mathbf{x}_i$  and  $\mathbf{o}_j$ . The value of  $A_{ij}$  may be binary, that is,  $A_{ij} = 1$  if there is an edge, otherwise 0. If a real value is preferred,  $A_{ij}$  can be set according to the heat kernel:  $A_{ij} = e^{-\|\mathbf{x}_i - \mathbf{o}_j\|_2^2 / \rho}$ , where  $\rho > 0$  is a parameter.

The goal of graph embedding is to learn an embedding for each instance that predicts the context in the graph. Given an instance and its context, the objective of graph embedding is usually formulated as minimizing certain loss of predicting the context using the embedding of an instance as input feature [Weston *et al.*, 2012; Mikolov *et al.*, 2013]. The context of an instance can be simply defined as its neighbors in the graph, or generated by sophisticated methods such as random

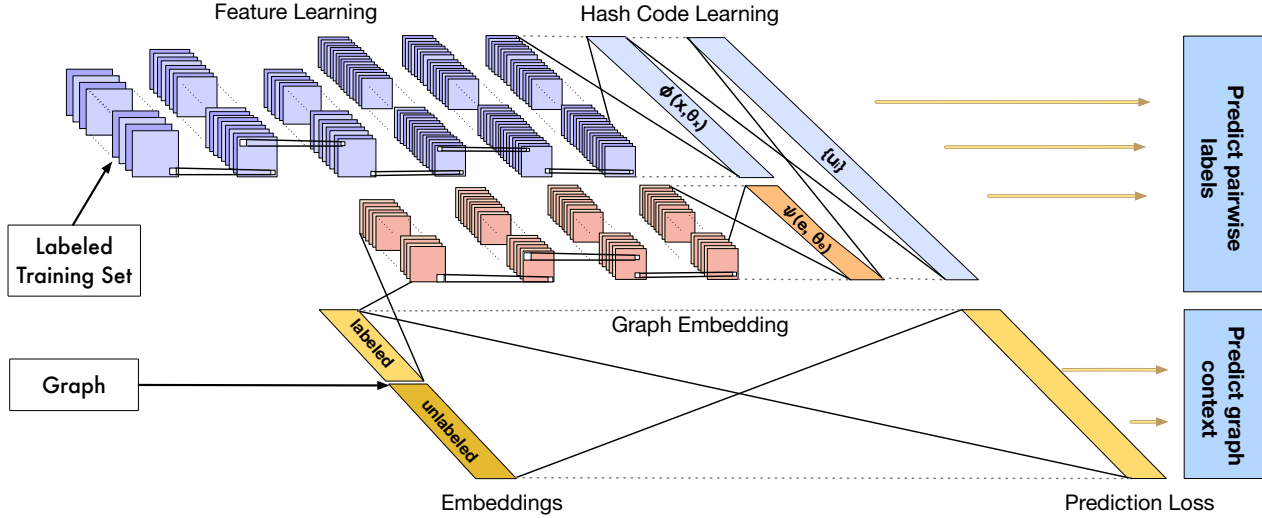


Figure 1: The end-to-end deep learning architecture of BGDH.

walk [Perozzi *et al.*, 2014]. Following previous studies [Yang *et al.*, 2016], we present a simple procedure for context generation in Algorithm 1, where the parameter  $d$  is the length of the random walk and  $r \in (0, 1)$  determines the ratio of positive contexts.

By invoking Algorithm 1  $t$  times, we obtain a set of triples  $\{(i_j, c_j, \gamma_j)\}_{j=1}^t$  where  $\gamma_j = 1$  indicates node  $c_j$  is a context of node  $i_j$ , and  $\gamma_j = -1$  indicates  $c_j$  is not a context. Let  $\mathbf{e}_{i_j}$  be the embedding of node  $i_j$ , and  $\mathbf{w}_{c_j}$  be the parameters for predicting node  $c_j$  as a context. We define the objective function of graph embedding

$$\mathcal{L}_g = \frac{1}{t} \sum_{j=1}^t \ell(\mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}, \gamma_j) \quad (2)$$

where  $\ell(\cdot, \cdot)$  is a loss that measures the discrepancy between  $\mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}$  and  $\gamma_j$ . In machine learning, the following losses are commonly used.

- The square loss

$$\ell(\mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}, \gamma_j) = (\gamma_j - \mathbf{e}_{i_j}^\top \mathbf{w}_{c_j})^2$$

- The logistic loss

$$\ell(\mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}, \gamma_j) = \log \left( 1 + e^{-\gamma_j \mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}} \right)$$

To constrain the solution space, we may further impose sparse constraints or nonnegative constraints [Lee and Seung, 1999].

### Feature learning and hash code learning

We utilize deep neural network model to learn features from the raw pixels and embeddings of labeled instances, and then combine them together to learn the hash codes. BGDH contains a CNN model to learn features from raw image pixels, and the model has seven layers as those of CNN-F [Chatfield *et al.*, 2014] while other networks like AlexNet [Krizhevsky and Hinton, 2009] can be used too. The configuration of the

---

### Algorithm 1 Context generation based on random walk

---

- 1: **Input:** A bipartite graph  $\mathcal{G} = (\mathcal{X}, \mathcal{O}, \mathcal{E})$ , parameters  $r$  and  $d$
  - 2: Uniformly sample an instance  $i$  from  $\mathcal{X}$
  - 3: Uniformly sample a random variable  $u$  from  $[0, 1]$
  - 4: **if**  $u < r$  **then**
  - 5:      $\gamma \leftarrow +1$
  - 6:     Uniformly sample a random walk sequence  $S$  of length  $d$  started from  $i$
  - 7:     Uniformly sample a context  $c$  from  $S$  except  $i$
  - 8: **else**
  - 9:      $\gamma \leftarrow -1$
  - 10:    Uniformly sample context  $c$  from  $\mathcal{X}$
  - 11: **end if**
  - 12: **return**  $(i, c, \gamma)$
- 

network is presented in Table 1, and a detailed explanation can be found in [Li *et al.*, 2016]. The output of the last feature learning layer (full7) of labeled instance  $\mathbf{x}_i$  is represented by  $\phi(\mathbf{x}_i; \theta_x)$ , where  $\theta_x$  denotes all the parameters in the seven layers of feature learning part. In contrast with existing supervised hashing methods, we also learn features from embeddings of labeled instances. The output associated with embedding  $\mathbf{e}_i$  is denoted by  $\psi(\mathbf{e}_i; \theta_e)$ , where  $\theta_e$  contains all the parameters in hidden layers. In this paper, we only add one fully-connected layer for embeddings, of which the size is determined by the dimension of embeddings.

We concatenate  $\phi(\mathbf{x}_i; \theta_x)$  and  $\psi(\mathbf{e}_i; \theta_e)$  as a new feature for instance  $i$ , then send it to a hash code learning layer as:

$$\mathbf{u}_i = \mathbf{M}^T [\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] + \mathbf{v} \quad (3)$$

where  $\mathbf{M} \in \mathbb{R}^{(4096+d) \times c}$  denotes a weight matrix,  $d$  is the dimension of embedding, and  $\mathbf{v} \in \mathbb{R}^{c \times 1}$  is a bias vector. Note that any supervised loss function of hashing can be used in our framework to learn parameters  $\mathbf{M}$  and  $\mathbf{v}$ . In this paper, we

Table 1: Configuration of the feature learning network

Layer	Configuration
conv1	filter $64 \times 11 \times 11$ , stride $4 \times 4$ , pad 0, LRN, pool $2 \times 2$
conv2	filter $256 \times 5 \times 5$ , stride $1 \times 1$ , pad 2, LRN, pool $2 \times 2$
conv3	filter $256 \times 3 \times 3$ , stride $1 \times 1$ , pad 1
conv4	filter $256 \times 3 \times 3$ , stride $1 \times 1$ , pad 1
conv5	filter $256 \times 3 \times 3$ , stride $1 \times 1$ , pad 1, pool $2 \times 2$
full6	4096
full7	4096

choose the loss function of deep pairwise-supervised hashing (DPSH) [Li *et al.*, 2016], and  $\mathcal{L}_s$  is given by

$$\mathcal{L}_s = - \sum_{s_{ij} \in \mathcal{S}} (s_{ij} \Theta_{ij} - \log(1 + e^{\Theta_{ij}})) + \eta \sum_{i=1}^l \|\mathbf{b}_i - \mathbf{u}_i\|_2^2 \quad (4)$$

where  $\Theta_{ij} = \frac{1}{2} \mathbf{u}_i^T \mathbf{u}_j$  and  $\eta > 0$  is a regularization parameter. By substituting Eq. (3) into Eq. (4), we obtain the final loss function of the supervised part:

$$\begin{aligned} \mathcal{L}_s = & - \sum_{s_{ij} \in \mathcal{S}} (s_{ij} \Theta_{ij} - \log(1 + e^{\Theta_{ij}})) \\ & + \eta \sum_{i=1}^l \|\mathbf{b}_i - (\mathbf{M}^T [\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] + \mathbf{v})\|_2^2. \end{aligned} \quad (5)$$

### The objective of BGDH

We combine the supervised part and unsupervised part to form the transductive version of BGDH. From (2) and (5), the loss function of BGDH is

$$\begin{aligned} \mathcal{L}_s + \lambda \mathcal{L}_g = & - \sum_{s_{ij} \in \mathcal{S}} (s_{ij} \Theta_{ij} - \log(1 + e^{\Theta_{ij}})) \\ & + \eta \sum_{i=1}^l \|\mathbf{b}_i - (\mathbf{M}^T [\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] + \mathbf{v})\|_2^2 \\ & + \frac{\lambda}{t} \sum_{j=1}^t \ell(\mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}, \gamma_j) \end{aligned} \quad (6)$$

where  $\lambda > 0$  is a constant weighting factor. The first two terms are the loss of predicting pairwise labels and the third one is the loss of predicting context. As a result, our BGDH can simultaneously learn embeddings, features, and hash codes. During the training phase, semantic similarity can affect graph embeddings, at the same time structure of data also influence the prediction of pairwise labels.

### 3.2 An Inductive Variant

Note that the basic BGDH is transductive, because the embeddings of instances are learnt from the graph. Since the hash code of an instance depends on both the raw pixels and its embedding, we need to design a way to infer the embedding of a unseen instance. To this end, we insert hidden layers to connect the raw pixels and embedding [Yang *et al.*, 2016], and in this way, the embedding  $\mathbf{e}_i$  becomes a parameterized function of  $\mathbf{x}_i$ , denoted by  $\xi(\mathbf{x}_i; \vartheta_x)$ . The loss function of

inductive hashing model can be written as:

$$\begin{aligned} \mathcal{L}_s + \lambda \mathcal{L}_g = & - \sum_{s_{ij} \in \mathcal{S}} (s_{ij} \Theta_{ij} - \log(1 + e^{\Theta_{ij}})) \\ & + \eta \sum_{i=1}^l \|\mathbf{b}_i - (\mathbf{M}^T [\phi(\mathbf{x}_i; \theta_x); \psi(\xi(\mathbf{x}_i; \vartheta_x); \theta_e)] + \mathbf{v})\|_2^2 \\ & + \frac{\lambda}{t} \sum_{j=1}^t \ell(\mathbf{e}_{i_j}^\top \mathbf{w}_{c_j}, \gamma_j) \end{aligned} \quad (7)$$

We can predict the hash code of any point  $\mathbf{x}_q \notin \mathcal{X}$  as:

$$\mathbf{b}_q = \text{sgn}(\mathbf{M}^T [\phi(\mathbf{x}_q; \theta_x); \psi(\xi(\mathbf{x}_q; \vartheta_x); \theta_e)] + \mathbf{v}). \quad (8)$$

### 3.3 Learning

In the transductive version of BGDH, the optimization variables include  $\mathbf{M}$ ,  $\mathbf{v}$ ,  $\{\mathbf{b}_i\}$ ,  $\theta_x$ ,  $\theta_e$ ,  $\{\mathbf{e}_i\}$  and  $\{\mathbf{w}_c\}$ . We adopt stochastic gradient descent (SGD) [Bottou, 2010] to train our model.

First, we sample a batch of labeled instances of which set  $\mathcal{I}_1$  contains indexes. A gradient step is then taken to optimize the supervised loss  $\mathcal{L}_s$ . For all  $i \in \mathcal{I}_1$ ,  $\mathbf{b}_i$  can be directly optimized as follow:

$$\mathbf{b}_i = \text{sgn}(\mathbf{u}_i) = \text{sgn}(\mathbf{M}^T [\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] + \mathbf{v}) \quad (9)$$

For other parameters in  $\mathcal{L}_s$ , i.e.,  $\mathbf{M}$ ,  $\mathbf{v}$ ,  $\theta_x$ ,  $\theta_e$ , and  $\{\mathbf{e}_i : i \in \mathcal{I}_1\}$ , we use back propagation (BP) to optimize them. Derivatives of  $\mathcal{L}_s$  w. r. t.  $\mathbf{u}_i$  are presented as follows:

$$\begin{aligned} \frac{\partial \mathcal{L}_s}{\partial \mathbf{u}_i} = & \frac{1}{2} \sum_{j: s_{ij} \in \mathcal{S}} (a_{ij} - s_{ij}) \mathbf{u}_j + \frac{1}{2} \sum_{j: s_{ji} \in \mathcal{S}} (a_{ji} - s_{ji}) \mathbf{u}_j \\ & + 2\eta(\mathbf{u}_i - \mathbf{b}_i) \end{aligned} \quad (10)$$

where  $a_{ij} = \sigma(\frac{1}{2} \mathbf{u}_i^T \mathbf{u}_j)$ . We can then update other parameters according to

$$\frac{\partial \mathcal{L}_s}{\partial \mathbf{M}} = [\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] \left( \frac{\partial \mathcal{L}_s}{\partial \mathbf{u}_i} \right)^T, \quad (11)$$

$$\frac{\partial \mathcal{L}_s}{\partial \mathbf{v}} = \frac{\partial \mathcal{L}_s}{\partial \mathbf{u}_i}, \quad (12)$$

$$\frac{\partial \mathcal{L}_s}{\partial \phi(\mathbf{x}_i; \theta_x)} = \frac{\partial \mathcal{L}_s}{\partial \psi(\mathbf{e}_i; \theta_e)} = \mathbf{M} \frac{\partial \mathcal{L}_s}{\partial \mathbf{u}_i}. \quad (13)$$

Then, we perform a gradient step to optimize the unsupervised graph embedding loss  $\mathcal{L}_g$  calculated by sampling triples generated in Algorithm 1. In this case parameters  $\{\mathbf{e}_i, \mathbf{w}_c : (i, c) \in \mathcal{I}_2\}$  will be updated where  $\mathcal{I}_2$  denotes the set containing indexes of sampled instances and contexts. The above procedures are repeated for  $T_1$  and  $T_2$  iterations respectively.

The whole learning algorithm of BGDH is summarized in Algorithm 2. Notice that before training jointly, we first train unsupervised part for a number of iterations to learn the initialization embeddings  $\{\mathbf{e}_i\}$ . For the inductive variant, we will update parameters  $\vartheta_x$  instead of embeddings  $\{\mathbf{e}_i\}$ .

---

**Algorithm 2** Learning algorithm for BGDH

---

- 1: **Input:** A bipartite graph  $\mathcal{G} = (\mathcal{X}, \mathcal{O}, \mathcal{E})$ , images  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$ , pairwise labels  $\mathcal{S} = \{s_{ij}\}$ , parameters  $\eta, \lambda$ , batch iterations  $T_1, T_2$  and sizes  $N_1, N_2$
  - 2: **Output:** Binary codes  $\mathcal{B} = \{\mathbf{b}_i\}_{i=1}^n$
  - 3: **Initialization:** Initialize  $\theta_x$  with the pre-trained CNN-F model on ImageNet; Initialize each entry of  $\mathbf{M}, \mathbf{v}$ , and  $\theta_e$  by randomly sampling from a Gaussian distribution with mean 0 and variance 0.01.
  - 4: **REPEAT**
  - 5: **for**  $t \leftarrow 1$  **to**  $T_1$  **do**
  - 6:     Randomly sample  $N_1$  labeled images, and let  $\mathcal{I}_1$  be the set containing indexes of sampled instances
  - 7:     Calculate  $\phi(\mathbf{x}_i; \theta_x)$  and  $\psi(\mathbf{e}_i; \theta_e)$  for all  $i \in \mathcal{I}_1$  by forward propagation
  - 8:     Compute  $\mathbf{u}_i = \mathbf{M}^T[\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] + \mathbf{v}$  and the binary code of  $\mathbf{x}_i$  with  $\mathbf{b}_i = \text{sgn}(\mathbf{u}_i)$  for all  $i \in \mathcal{I}_1$
  - 9:     Compute the derivative of  $\mathcal{L}_s$  w. r. t.  $\{\mathbf{u}_i : i \in \mathcal{I}_1\}$
  - 10:    Update parameters  $\mathbf{M}, \mathbf{v}, \theta_x, \theta_e$ , and  $\{\mathbf{e}_i : i \in \mathcal{I}_1\}$  by back propagation
  - 11: **end for**
  - 12: **for**  $t \leftarrow 1$  **to**  $T_2$  **do**
  - 13:     Randomly generate a batch of triples by invoking Algorithm 1  $N_2$  times, and let  $\mathcal{I}_2$  be the set containing indexes of sampled instances and contexts
  - 14:     Compute the derivative of  $\mathcal{L}_g$  w. r. t.  $\{\mathbf{e}_i, \mathbf{w}_c : (i, c) \in \mathcal{I}_2\}$
  - 15:     Update parameters  $\{\mathbf{e}_i, \mathbf{w}_c : (i, c) \in \mathcal{I}_2\}$
  - 16: **end for**
  - 17: **UNTIL** stopping
  - 18: Calculate  $\mathbf{b}_i = \text{sgn}(\mathbf{M}^T[\phi(\mathbf{x}_i; \theta_x); \psi(\mathbf{e}_i; \theta_e)] + \mathbf{v})$  for  $i = 1, \dots, n$
- 

## 4 Experiments

In this section, we present experimental results. All the experiments are performed on a NVIDIA K80 GPU server with MatConvNet [Vedaldi and Lenc, 2014].

### 4.1 Datasets and Setting

We conduct experiments on two widely used benchmark datasets: CIFAR-10 and NUS-WIDE. The *CIFAR-10* dataset<sup>1</sup> consists of 60,000 images from 10 classes (6000 images per class). It is a single-label dataset. The *NUS-WIDE* dataset<sup>2</sup> contains 269,648 images from Flickr. Following the settings in [Lai *et al.*, 2015], we use 19,5834 images belonging to the 21 most frequent classes and each class consists at least 5000 images. Additionally, two images are defined as a ground-truth neighbor when they share at least one common label.

In our experiments, BGDH-T denotes transductive model while BGDH-I denotes inductive one. We compare our method with several state-of-the-art methods including unsupervised methods: ITQ [Gong *et al.*, 2013], LSH [Gionis *et al.*, 1999], IsoH [Kong and Li, 2012] and SpH [Heo *et al.*, 2012]; supervised methods: LFH [Zhang *et al.*, 2014], FastH

[Lin *et al.*, 2014], SDH [Shen *et al.*, 2015] and COSDISH [Kang *et al.*, 2016]; supervised deep hashing methods: DPSH [Li *et al.*, 2016], DHN [Zhu *et al.*, 2016] and DSH [Liu *et al.*, 2016]; semi-supervised deep methods SSDH [Zhang *et al.*, 2016].

For hashing methods using hand-crafted features, we represent each image in CIFAR-10 by a 512-dimensional GIST vector. In NUS-WIDE, an image is represented by a 1134-dimensional feature vector, including 500-D bag-of-words features, 64-D color histogram, 144-D color correlogram, 73-D edge direction histogram, 128-D wavelet texture, 225-D block-wise color moments. For deep methods, we use the raw image pixels as input. The network architectures of these deep models are different from each other, for fair comparison, we adopt the same deep network architecture CNN-F with the same initialization parameters pre-trained on ImageNet [Deng *et al.*, 2009] for all deep hashing methods. The bipartite graph of BGDH is constructed based on hand-crafted features with heat kernel, where the hyper-parameter  $\rho$  is set as 1 for CIFAR-10 and 10 for NUS-WIDE. The hyper-parameter  $\eta$  in BGDH is set as 10 for CIFAR-10 and 100 for NUS-WIDE similar to DPSH [Li *et al.*, 2016]. We simply set  $T_1 = 10, T_2 = 5, \lambda = 0.1$  in all the experiments.

### 4.2 Experimental Results

To evaluate the retrieval accuracy, we use the Mean Average Precision (MAP) as the evaluation metric. Following [Lai *et al.*, 2015; Xia *et al.*, 2014], 1000 images (100 images per class) are randomly selected as the query set in CIFAR-10. We use all the rest of images as training set for the unsupervised methods and for building a neighbour graph  $\mathcal{G}$  in BGDH. We randomly select 5000 (500 images per class) images as the training set for supervised methods.

In NUS-WIDE, we randomly sample 2100 query images from 21 most frequent labels following the strategy in [Lai *et al.*, 2015; Xia *et al.*, 2014]. We use all the rest of images as training set for the unsupervised methods. To reduce the time complexity, we only use 5000 landmarks sampled randomly to build a bipartite graph for BGDH. For supervised methods, we randomly select 10500 images from dataset. The MAP values are calculated within the top 5000 returned neighbors.

We present the MAP results in Table 2 where BGDH, SSDH, DSH, DHN, DPSH are deep methods. Except for SSDH which uses triplet labels, all deep methods are trained with pairwise labels. To be fair, parameters of these methods are set according to the suggestions of authors. The results show that in most cases our proposed BGDH method substantially outperforms other baselines including unsupervised methods, supervised methods and semi-supervised methods.

In real applications, the number of labeled instances is usually far less than that of unlabeled instances. To further verify the effectiveness of leveraging both pairwise labels and unlabeled data, we reduce the number of labeled instances while other settings remain the same. For supervised learning, in CIFAR-10, we randomly select 2500 (250 images per class) images and in NUS-WIDE, we randomly select 5000 images. The MAP results with this experiment setting are listed in Table 3. Note that the results of unlisted unsupervised methods are the same as those in Table 2. We observe that our BGDH

<sup>1</sup><https://www.cs.toronto.edu/~kriz/cifar.html>

<sup>2</sup><http://lms.comp.nus.edu.sg/research/NUS-WIDE.htm>

Table 2: Accuracy in terms of MAP. The best MAPs for each category are shown in boldface. Training size for supervised method is 5000 for CIFAR-10 and 10500 for NUS-WIDE.

Method	CIFAR-10 (MAP)				NUS-WIDE (MAP)			
	12-bits	24-bits	32-bits	48-bits	12-bits	24-bits	32-bits	48-bits
BGDH-T	<b>0.805</b>	<b>0.824</b>	<b>0.826</b>	<b>0.833</b>	<b>0.803</b>	<b>0.818</b>	<b>0.822</b>	<b>0.828</b>
BGDH-I	0.803	0.818	0.822	0.829	0.801	0.815	0.816	0.825
SSDH	0.801	0.813	0.812	0.814	0.773	0.779	0.778	0.778
DSH	0.604	0.746	0.781	0.810	0.751	0.765	0.767	0.773
DHN	0.692	0.703	0.726	0.735	0.751	0.785	0.792	0.799
DPSH	0.684	0.734	0.750	0.767	0.788	0.809	0.817	0.823
COSDISH	0.522	0.590	0.599	0.615	0.691	0.749	0.745	0.765
SDH	0.525	0.671	0.686	0.696	0.752	0.745	0.744	0.730
FastH	0.291	0.351	0.367	0.390	0.622	0.660	0.670	0.687
LFH	0.335	0.433	0.509	0.515	0.749	0.751	0.775	0.780
ITQ	0.163	0.170	0.173	0.176	0.447	0.465	0.468	0.473
LSH	0.152	0.167	0.170	0.200	0.367	0.394	0.413	0.416
IsoH	0.158	0.162	0.166	0.169	0.436	0.454	0.461	0.465
SpH	0.141	0.153	0.154	0.158	0.399	0.437	0.454	0.465

Table 3: Accuracy in terms of MAP. The best MAPs for each category are shown in boldface. Training size for supervised method is 2500 for CIFAR-10 and 5000 for NUS-WIDE.

Method	CIFAR-10 (MAP)				NUS-WIDE (MAP)			
	12-bits	24-bits	32-bits	48-bits	12-bits	24-bits	32-bits	48-bits
BGDH-T	<b>0.755</b>	<b>0.791</b>	<b>0.800</b>	<b>0.812</b>	<b>0.772</b>	<b>0.798</b>	<b>0.806</b>	<b>0.816</b>
BGDH-I	0.746	0.776	0.787	0.796	0.768	0.794	0.801	0.811
SSDH	0.581	0.589	0.595	0.596	0.743	0.745	0.746	0.749
DSH	0.617	0.707	0.737	0.761	0.749	0.769	0.771	0.786
DHN	0.591	0.646	0.640	0.662	0.741	0.763	0.766	0.773
DPSH	0.576	0.634	0.642	0.668	0.762	0.789	0.791	0.803
COSDISH	0.312	0.348	0.373	0.398	0.648	0.678	0.699	0.713
SDH	0.327	0.357	0.374	0.377	0.574	0.597	0.591	0.595
FastH	0.267	0.298	0.320	0.341	0.604	0.634	0.650	0.667
LFH	0.244	0.288	0.311	0.391	0.611	0.644	0.653	0.669

outperforms all the other methods. Specifically, compared to the best baseline in Table 2, we conclude that when labeled data are insufficient, BGDH is able to leverage unlabeled data to deliver a good result.

### 4.3 Parameter Selection

In BGDH, there is a hyper-parameter  $\lambda$  which controls the tradeoff between supervised loss and unsupervised loss. Figure 2 displays the impacts of  $\lambda$  on the performance of BGDH with the experiment settings being the same as those in Table 3. As can be seen, there is a wide range of  $\lambda$  that BGDH performs well. Thus, to a large extent, BGDH is insensitive to  $\lambda$  and the parameter selection is not a crucial problem in our algorithm. Additionally, by comparing the MAP of  $\lambda = 0$  and  $\lambda = 1$ , we verify the importance of graph embedding.

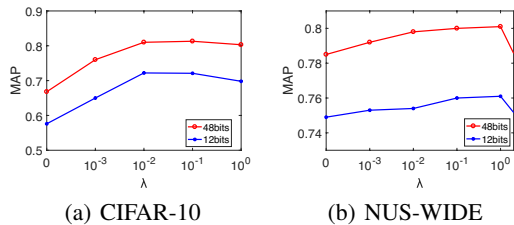


Figure 2: Hyper-parameter Sensitivity

## 5 Conclusion

In this paper, we propose a novel semi-supervised hashing method, named Deep Hashing with a Bipartite Graph (BGDH). To the best of our knowledge, BGDH is the first method that performs graph embedding, feature learning, and hash code learning simultaneously. BGDH constructs a bipartite graph to discover the underlying structure of data, and is much more efficient than methods based on neighborhood graph. Experimental results demonstrate that BGDH outperforms state-of-the-art methods in image retrieval.

## Acknowledgements

This work was partially supported by the NSFC (61603177, 61472182), JiangsuSF (BK20160658), and the Collaborative Innovation Center of Novel Software Technology and Industrialization of Nanjing University.

## References

[Bottou, 2010] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of 19th International Conference on Computational Statistics*, pages 177–186. Springer, 2010.

[Chatfield *et al.*, 2014] Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Return of the devil in the de-

- tails: Delving deep into convolutional nets. In *Proceedings of British Machine Vision Conference*, 2014.
- [Deng *et al.*, 2009] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [Gionis *et al.*, 1999] Aristides Gionis, Piotr Indyk, Rajeev Motwani, et al. Similarity search in high dimensions via hashing. In *Proceedings of 25th International Conference on Very Large Data Bases*, volume 99, pages 518–529, 1999.
- [Gong *et al.*, 2013] Yunchao Gong, Svetlana Lazebnik, Albert Gordo, and Florent Perronnin. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2916–2929, 2013.
- [Heo *et al.*, 2012] Jae-Pil Heo, Youngwoon Lee, Junfeng He, Shih-Fu Chang, and Sung-Eui Yoon. Spherical hashing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2957–2964, 2012.
- [Indyk and Motwani, 1998] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613, 1998.
- [Kang *et al.*, 2016] Wang-Cheng Kang, Wu-Jun Li, and Zhi-Hua Zhou. Column sampling based discrete supervised hashing. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 1230–1236, 2016.
- [Kong and Li, 2012] Weihao Kong and Wu-Jun Li. Isotropic hashing. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2012.
- [Krizhevsky and Hinton, 2009] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- [Krizhevsky *et al.*, 2012] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.
- [Lai *et al.*, 2015] Hanjiang Lai, Yan Pan, Ye Liu, and Shuicheng Yan. Simultaneous feature learning and hash coding with deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3270–3278, 2015.
- [Lee and Seung, 1999] Daniel D. Lee and H. Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- [Li *et al.*, 2016] Wu-Jun Li, Sheng Wang, and Wang-Cheng Kang. Feature learning based deep supervised hashing with pairwise labels. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pages 1711–1717, 2016.
- [Lin *et al.*, 2014] Guosheng Lin, Chunhua Shen, Qinfeng Shi, Anton van den Hengel, and David Suter. Fast supervised hashing with decision trees for high-dimensional data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1963–1970, 2014.
- [Liu *et al.*, 2010] Wei Liu, Junfeng He, and Shih-Fu Chang. Large graph construction for scalable semi-supervised learning. In *Proceedings of the 27th International Conference on Machine Learning*, pages 679–686, 2010.
- [Liu *et al.*, 2011] Wei Liu, Jun Wang, Sanjiv Kumar, and Shih-Fu Chang. Hashing with graphs. In *Proceedings of the 28th International Conference on Machine Learning*, pages 1–8, 2011.
- [Liu *et al.*, 2016] Haomiao Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Deep supervised hashing for fast image retrieval. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2064–2072, 2016.
- [Mikolov *et al.*, 2013] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119, 2013.
- [Oliva and Torralba, 2001] Aude Oliva and Antonio Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International Journal of Computer Vision*, 42(3):145–175, 2001.
- [Perozzi *et al.*, 2014] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: online learning of social representations. In *Proceedings of the 20th ACM International Conference on Knowledge Discovery and Data Mining*, pages 701–710, 2014.
- [Shen *et al.*, 2015] Fumin Shen, Chunhua Shen, Wei Liu, and Heng Tao Shen. Supervised discrete hashing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 37–45, 2015.
- [Vedaldi and Lenc, 2014] Andrea Vedaldi and Karel Lenc. Matconvnet - convolutional neural networks for MATLAB. *CoRR*, abs/1412.4564, 2014.
- [Wang *et al.*, 2012] Jun Wang, Sanjiv Kumar, and Shih-Fu Chang. Semi-supervised hashing for large-scale search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(12):2393–2406, 2012.
- [Weiss *et al.*, 2009] Yair Weiss, Antonio Torralba, and Rob Fergus. Spectral hashing. In *Advances in Neural Information Processing Systems*, pages 1753–1760, 2009.
- [Weston *et al.*, 2012] Jason Weston, Frédéric Ratle, Hossein Mobahi, and Ronan Collobert. Deep learning via semi-supervised embedding. In *Neural Networks: Tricks of the Trade*, pages 639–655. Springer, 2012.
- [Xia *et al.*, 2014] Rongkai Xia, Yan Pan, Hanjiang Lai, Cong Liu, and Shuicheng Yan. Supervised hashing for image retrieval via image representation learning. In *Proceedings of AAAI Conference on Artificial Intelligence*, pages 2156–2162, 2014.
- [Yang *et al.*, 2016] Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 40–48, 2016.
- [Zhang *et al.*, 2014] Peichao Zhang, Wei Zhang, Wu-Jun Li, and Minyi Guo. Supervised hashing with latent factor models. In *Proceedings of the 37th International ACM Conference on Research and Development in Information Retrieval*, pages 173–182, 2014.
- [Zhang *et al.*, 2016] Jian Zhang, Yuxin Peng, and Junchao Zhang. Ssdh: semi-supervised deep hashing for large scale image retrieval. *arXiv preprint arXiv:1607.08477*, 2016.
- [Zhou *et al.*, 2004] Dengyong Zhou, Olivier Bousquet, Thomas Navin Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. In *Advances in Neural Information Processing Systems*, pages 321–328, 2004.
- [Zhu *et al.*, 2016] Han Zhu, Mingsheng Long, Jianmin Wang, and Yue Cao. Deep hashing network for efficient similarity retrieval. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 2415–2421, 2016.