

SADIH: Semantic-Aware DIScrete Hashing

Zheng Zhang,¹ Guo-sen Xie,² Yang Li,¹ Sheng Li,³ Zi Huang¹

¹The University of Queensland, Australia

²Inception Institute of Artificial Intelligence, UAE

³University of Georgia, USA

{darrenzz219, gsxiehm}@gmail.com; y.li9@uq.net.au; sheng.li@uga.edu; huang@itee.uq.edu.au

Abstract

Due to its low storage cost and fast query speed, hashing has been recognized to accomplish similarity search in large-scale multimedia retrieval applications. Particularly, supervised hashing has recently received considerable research attention by leveraging the label information to preserve the pairwise similarities of data points in the Hamming space. However, there still remain two crucial bottlenecks: 1) the learning process of the *full pairwise similarity preservation* is computationally unaffordable and unscalable to deal with big data; 2) the available category information of data are not well-explored to learn discriminative hash functions. To overcome these challenges, we propose a unified Semantic-Aware DIScrete Hashing (SADIH) framework, which aims to directly embed the transformed semantic information into the asymmetric similarity approximation and discriminative hashing function learning. Specifically, a semantic-aware latent embedding is introduced to asymmetrically preserve the full pairwise similarities while skillfully handle the cumbersome $n \times n$ pairwise similarity matrix. Meanwhile, a semantic-aware autoencoder is developed to jointly preserve the data structures in the discriminative latent semantic space and perform data reconstruction. Moreover, an efficient alternating optimization algorithm is proposed to solve the resulting discrete optimization problem. Extensive experimental results on multiple large-scale datasets demonstrate that our SADIH can clearly outperform the state-of-the-art baselines with the additional benefit of lower computational costs.

Introduction

In the big data era, recent years have witnessed the ever-growing volume of multimedia data with high dimensionality. This is made possible by the emergence of large-scale similarity measurement technique with high computational efficiency. Different from the traditional indexing technique (Lew et al. 2006), hashing yields a scalable similarity search mechanism with acceptable accuracies in the fast Hamming space (Wang et al. 2018). Technically, hashing generally compresses the high-dimensional data instances into compact binary codes (typically ≤ 128 -dim), in which the similarity and structural information are preserved from the original data. In this paper, we will mainly focus on the learning-based hashing methods that are formulated by the data-

dependent hash encoding strategy, which has shown better retrieval performance than data-independent (or learning-free) hashing schemes, such as locality-sensitive hashing (LSH) (Gionis et al. 1999) and its variants (Kulis and Grauman 2009; Jiang, Que, and Kulis 2015).

A common problem of learning-based hashing methods is to construct similarity-preserving hash functions, which generate similar binary codes for nearby data items. Many such hashing learning methods have been proposed to enable efficient similarity search and can be broadly grouped into two categories: unsupervised and supervised hashing.

Unsupervised methods typically encode samples as binary codes by exploring data distribution without label or relevances (Zhang et al. 2018b; 2018a). They learn hash codes/functions based on the semantic gap principle (Smeulders et al. 2000), *i.e.*, the difference in structures formed within the high- and low-level descriptors. Representative unsupervised hashing methods include manifold learning based hashing and quantization based hashing. Manifold learning based hashing tries to discover the neighborhood relationship of data points in the learned binary codes, such as spectral hashing (SH) (Weiss, Torralba, and Fergus 2009), multiple feature hashing (MFH) (Song et al. 2011), scalable graph hashing (SGH) (Jiang and Li 2015) and ordinal constraint hashing (OCH) (Liu et al. 2018). Quantization based hashing aims to achieve the minimal quantization error, such as iterative quantization (ITQ) (Gong et al. 2013) and quantization-based hashing (Song et al. 2018). Due to the absence of semantic label information, unsupervised hashing is usually inferior to supervised hashing, which can produce state-of-the-art retrieval results.

Supervised hashing generates discriminative and compact hash codes/functions by leveraging the supervised semantic information from data such as pairwise similarity or relevant feedback. Many supervised hashing methods have been proposed to enable efficient similarity search. Representative methods in this group include semi-supervised hashing (SSH) (Wang, Kumar, and Chang 2012), latent factor hashing (LFH) (Zhang, Zhang, and Li 2014), fast supervised hashing (FastH) (Lin et al. 2014), column sampling based discrete supervised hashing (COSDISH) (Kang, Li, and Zhou 2016), etc. It has been studied to construct encoding functions in a designed kernel space, such as binary reconstructive embedding (BRE) (Kulis and Darrell 2009),

kernel-based supervised hashing (KSH) (Liu et al. 2012), the kernel variant of ITQ (Gong et al. 2013), supervised discrete hashing (SDH) (Shen et al. 2015), SDH with relaxation (SDHR) (Gui et al. 2016; Zhang et al. 2018c) and fast SDH (FSDH) (Gui et al. 2018). The kernel based hashing methods have been shown to achieve promising performance.

To further improve the retrieval performance, many deep hashing models (Erin Liong et al. 2015; Shen et al. 2018; Li, Wang, and Kang 2016; Lai et al. 2015; Lin et al. 2015) have been introduced over the past few years, where the non-linear feature embeddings learned by deep neural networks were typically shown to achieve higher performance than hand-crafted descriptors. As we know, semantic hashing (Salakhutdinov and Hinton 2009) is the pioneering work of using deep machine for hashing. However, these deep hashing models are complicated and need pre-training, which is inefficient in real applications. Moreover, there might be concern about the encoding time of the training and test data.

Although achieving progress, current supervised similarity-preserving hashing methods are still facing severe challenges. First, to *avoid* using the full $n \times n$ pairwise similarity matrix, these methods employ sampling strategies in the training phase to reduce the large computation and memory overhead. In such a case, they fail to capture the full structures residing on the entire data, which inevitably results in information loss and unsatisfactory performance. Their objectives would be suboptimal for realistic search tasks, and such methods become inappropriate for large-scale retrieval tasks. Second, only preserving the pairwise similarities transformed from labels clearly excludes the category information of data from the training step. In this way, these methods can not transfer the discriminative information from labels into the learned binary codes, resulting in inferior performance. Third, since the discrete optimization introduced by binary constraint leads to an NP-hard mixed integer programming problem, most of them usually solve it by relaxing the binary variables into continuous ones, followed by thresholding or quantization. However, such relaxation strategy can amplify the quantization errors, which may greatly influence the quality of the learned binary codes and degrade the performance.

To address the aforementioned problems, we propose a novel discriminative binary code learning framework, dubbed Semantic-Aware DIScrete Hashing (SADIH), for fast scalable supervised hashing. Specifically, we introduce an asymmetric similarity-preserving strategy that can preserve the discrete constraint and reduce the accumulated quantization error between binary code matrix and the well-designed latent semantic-aware embedding. During the learning step, such trick can skillfully handle the huge $n \times n$ pairwise similarity matrix, and preserve the discriminative category information into the learned binary codes. Meanwhile, we also develop a novel semantic-aware encoder-decoder paradigm to guarantee the high-quality latent embedding. In particular, an encoder projects the visual features of an image into a latent semantic space, and in turn consider the latent semantic-aware embedding as an input to a decoder which reconstructs the original visual representation. As such, our learning framework not only can ef-

fectively preserve the discriminative semantic information into the learned binary codes and hashing functions, but efficiently approximates the full pairwise similarities without information loss. Furthermore, an alternating algorithm is developed to solve the resulting problem, where each sub-problem can be optimized efficiently, yielding satisfactory solutions. To sum up, the main contributions of this work are:

(1) A novel semantic-aware discrete hashing framework is proposed to simultaneously consider the full pairwise similarities ($n \times n$) and the category information into the joint learning objective. SADIH aims to generate discriminative binary codes which can successfully capture the entire pairwise similarities as well as the intrinsic correlations between visual features and semantics from different categories.

(2) We introduce a latent semantic embedding space which can reconcile the structural difference between the visual and semantic spaces, meanwhile preserve the discriminative structures in the learned binary codes.

(3) An asymmetric similarity approximation loss is developed to reduce the accumulated quantization error between the learned binary codes and the latent semantic-aware embeddings. Meanwhile, a supervised semantic-aware autoencoder is constructed to jointly perform the data structural preservation and data reconstruction. The well-designed alternating optimization algorithm with guaranteed convergence is applied to produce the high-quality hash codes.

Semantic-Aware Discrete Hashing

Basic Formulation

This work mainly focuses on supervised hashing to enable efficient semantic similarity search by Hamming ranking of compact hash codes. Suppose we have n d -dimensional data points, denoted as $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$, and their associated semantic labels are $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \{0, 1\}^{c \times n}$, where c is the number of classes. The i -th column of matrix \mathbf{Y} , *i.e.* $\mathbf{y}_i = [0, \dots, 1, \dots, 0]^T \in \mathbb{R}^c$, is the label vector of the i -th sample, and $y_{ji} = 1$ indicates \mathbf{x}_i belongs to the j -th class. Notably, supervised hashing also contains a pairwise similarity matrix $\mathbf{S} \in \{-1, 1\}^{n \times n}$ obtained from semantic correlations such as labels used in this paper. Specifically, $s_{ij} = 1$ means that data items i and j are semantically similar and share at least one label, while $s_{ij} = -1$ indicates items i and j are semantically dissimilar.

The goal of supervised hashing aims to learn l hashing functions to project the data \mathbf{X} into a discriminative Hamming space, and generate a binary code matrix $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n] \in \{-1, 1\}^{l \times n}$. Moreover, the learned binary codes should preserve the semantic similarities indicated in \mathbf{S} . The commonly-used objective function (Liu et al. 2012) quantizes the approximation error between the Hamming affinity and semantic similarity matrix using

$$\min_{\mathbf{B}} \|\mathbf{r}\mathbf{S} - \mathbf{B}^T \mathbf{B}\|_F^2 \text{ s.t. } \mathbf{B} \in \{-1, 1\}^{l \times n}, \quad (1)$$

where $\|\cdot\|_F$ is the Frobenius norm. In this model, the inner product of any two binary codes reflects the opposite of the Hamming distance, and can be used to approximate the corresponding similarity labels. Due to its effectiveness, this model has become a standard formulation for

supervised hashing learning (Zhang, Zhang, and Li 2014; Lin et al. 2014; Kang, Li, and Zhou 2016). However, there are still several deficiencies. **First**, such a symmetric binary affinity form is limited in matching the real-valued ground truth. Importantly, the optimization on symmetric discrete constraint is time-consuming, which makes it hard to adapt for large-scale datasets (Neyshabur et al. 2013). **Second**, owing to its computation and memory prohibition, the full similarity matrix \mathbf{S} is usually avoided using in the training step. An alternative strategy is to sample a small subset for training, which inevitably causes information loss and suboptimal results. **Third**, directly transforming labels into pairwise similarities loses the category information of training data, which can not preserve the discriminative characteristics into the learned binary codes. **Finally**, most methods solve the discrete optimization problem by relaxing the discrete constraint by omitting the sign function. However, the approximate solution is obviously suboptimal and often generates low-quality hashing codes.

Therefore, we present an efficient **Semantic-Aware Discrete Hashing (SADIH)** framework to alleviate the above limitations. In the method, the asymmetric hamming affinity approximation, latent semantic embedding and encoder-decoder paradigm are simultaneously considered to guarantee discriminative binary codes and hashing functions.

Objective Function

Asymmetric Similarity Approximation Loss To fully explore the entire similarities on n available points, we introduce a simple but effective semantic-aware constraint, *i.e.*, $\mathcal{V} = \mathbf{W}^T \mathbf{Y}$ where $\mathbf{W} \in \mathbb{R}^{c \times l}$, to asymmetrically approximate the ground-truth affinity. Meanwhile, the label information are embedded into the latent semantic embedding \mathcal{V} . Particularly, the matrix \mathbf{W}^T can be viewed as the category-level basis matrix of the latent semantic features \mathcal{V} , because $\mathbf{v}_i = \mathbf{W}^T \mathbf{y}_i$, where each item in \mathbf{v}_i contains the category partition factor. Moreover, using the real-valued embeddings can produce more accurate approximation of similarity, and reduce the accumulated quantization error (Dong, Charikar, and Li 2008; Luo, Wu, and Xu 2018). Based on the asymmetric hashing learning (Shrivastava and Li 2014), we replace one of the binary codes \mathbf{B} in (1), and consider its robust model

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{W}} \|\mathbf{S} - \mathcal{V}^T \mathbf{B}\|_{21} \\ \text{s.t. } \mathcal{V} = \mathbf{W}^T \mathbf{Y}, \mathbf{B} \in \{-1, 1\}^{l \times n}, \end{aligned} \quad (2)$$

which $\|\mathbf{A}\|_{21} = \sum_{i=1}^n \|\mathbf{a}^i\|_2$ denotes the l_{21} -norm of matrix \mathbf{A} and \mathbf{a}^i is the i -th row of matrix \mathbf{A} . The l_{21} -norm is robust to noise or outliers based on the rotation-invariance property. It is noteworthy that this simple constraint can enable the model to effectively exploit all the n data points for training (shown in optimization) without any sampling. Moreover, it also can more precisely measure the quantization between the given similarities \mathbf{S} and the learned asymmetric affinity.

Semantic-Aware Autoencoder The discriminative binary codes for training data can be learned based on model (2), but there still remain two concerns. On one hand, the latent

semantic embedding only leverages the label information, while the inherent characteristics embedded in the training data are not well-explored to capture the instance-level features. On the other hand, (2) can not be generalized to the out-of-sample cases for efficient query generation. To this end, we formulate the linear semantic-aware autoencoder scheme, which optimizes against the following objective:

$$\begin{aligned} \min_{\{\mathbf{P}_i, \mathbf{c}_i\}_{i=1}^2, \mathcal{V}, \mathbf{W}} \|\mathbf{X} - (\mathbf{P}_2(\mathbf{P}_1 \mathbf{X} + \mathbf{c}_1 \mathbf{1}^T) + \mathbf{c}_2 \mathbf{1}^T)\|_F^2 \\ + \gamma \mathcal{R}(\mathbf{P}_2, \mathbf{P}_1) \text{ s.t. } \mathcal{V} = \mathbf{W}^T \mathbf{Y}, \mathcal{V} = \mathbf{P}_1 \mathbf{X} + \mathbf{c}_1 \mathbf{1}^T, \end{aligned} \quad (3)$$

where $\mathbf{P}_1 \in \mathbb{R}^{l \times d}$ and $\mathbf{P}_2 \in \mathbb{R}^{d \times l}$ are the encoding and decoding matrices, $\mathcal{R}(\cdot) = \|\cdot\|_F^2$ is the regularization term to avoid overfitting, \mathbf{c}_1 and \mathbf{c}_2 are the biased vectors, and γ is a weighting parameter. It is clear that *this model can make use of semantic attributes in \mathcal{V} as an intermediate level clue to associate low-level visual features with high-level semantic information*. However, when we project the visual d -dim features into the lower l -dim (typically $l \ll d$) semantic space, this may encounter the imbalanced projection problem, *i.e.* the variances of the projected dimensions vary severely (Wang, Kumar, and Chang 2012). To this end, inspired by ITQ (Gong et al. 2013), we may change the coordinates of the whole feature space through an adjustment rotation. For eliminating the bias variables, we reformulate the above problem into a relaxed optimization with an orthogonal transformation:

$$\begin{aligned} \min_{\mathbf{P}_1, \mathbf{P}_2, \mathcal{V}, \mathbf{W}} \|\mathbf{X} - \mathbf{P}_2 \mathcal{V}\|_F^2 + \beta \|\mathcal{V} - \mathbf{P}_1 \mathbf{X}\|_F^2 + \gamma \mathcal{R}(\mathbf{P}_2) \\ \text{s.t. } \mathcal{V} = \mathbf{W}^T \mathbf{Y}, \mathbf{P}_1^T \mathbf{P}_1 = \mathbf{I}. \end{aligned} \quad (4)$$

From Eqn. (4), we can see that the preferred latent attributes satisfy $\mathcal{V} = \mathbf{P}_1 \mathbf{X}$ with minimum reconstruction error $\mathbf{X} = \mathbf{P}_2 \mathcal{V}$. Given the orthogonal transformation, the overall variance can be effectively diffused into all the learned dimensions through the adjustment rotation, and the underlying characteristics hidden in data \mathbf{X} are uncovered and transferred into the semantic embeddings. Importantly, the encoder can be used as a linear hash function for new queries.

Joint Objective Function To preserve the interconnection between the semantic-aware similarity approximation and preferable latent semantic space construction, SADIH combines the asymmetric similarity approximation loss given in (2) and semantic-aware autoencoder scheme given in (4) into one unified learning framework. Such a learning framework can minimize the intractable full pairwise similarity-preserving error, meanwhile interactively enhances the qualities of the learned binary codes and discriminative latent semantic-aware embeddings:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{W}, \mathbf{P}_1, \mathbf{P}_2, \mathcal{V}} \|\mathbf{S} - \mathcal{V}^T \mathbf{B}\|_{21} + \alpha \|\mathbf{X} - \mathbf{P}_2 \mathcal{V}\|_F^2 \\ + \beta \|\mathcal{V} - \mathbf{P}_1 \mathbf{X}\|_F^2 + \gamma \mathcal{R}(\mathbf{P}_2) \\ \text{s.t. } \mathbf{B} \in \{-1, 1\}^{l \times n}, \mathcal{V} = \mathbf{W}^T \mathbf{Y}, \mathbf{P}_1^T \mathbf{P}_1 = \mathbf{I}, \end{aligned} \quad (5)$$

where α is a weighting parameter. Furthermore, it is clear that the first term makes sure the asymmetric correlations between the discriminative binary codes \mathbf{B} and latent semantic representations \mathcal{V} . Therefore, the encoding matrix

P_1 can capture the discriminative information embedded in the latent semantic space, and the hashing codes for out-of-sample \mathbf{x}_t can be directly generated by $\mathbf{b} = \text{sgn}(P_1 \mathbf{x}_t)$, where $\text{sgn}(\cdot)$ denotes the element-wise sign function.

Optimization

The key of our optimization is to learn discriminative binary codes \mathbf{B} and find a proper latent-embedding space \mathcal{V} with the data structural preservation P_1 and data reconstruction P_2 . However, problem (5) is non-convex to all variables, and involves a discrete constraint, which leads to an NP-hard problem. Thus, we propose an alternating optimization algorithm to obtain a local minima.

We first rewrite problem (5) as an equivalent form:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{W}, P_1, P_2} & \|l\mathbf{S} - (\mathbf{W}^T \mathbf{Y})^T \mathbf{B}\|_{21} + \alpha \|\mathbf{X} - P_2 \mathbf{W}^T \mathbf{Y}\|_F^2 \\ & + \beta \|\mathbf{W}^T \mathbf{Y} - P_1 \mathbf{X}\|_F^2 + \gamma \mathcal{R}(P_2, \mathbf{W}^T \mathbf{Y}) \\ \text{s.t. } & \mathbf{B} \in \{-1, 1\}^{l \times n}, P_1^T P_1 = \mathbf{I}. \end{aligned} \quad (6)$$

Next, we iteratively update each variable with an alternative manner, *i.e.*, updating one when fixing others.

B-Step: When fixing \mathbf{W} , P_1 and P_2 , we can update \mathbf{B} by using two discrete learning strategies. The first learning scheme optimizes \mathbf{B} with a constant-approximation solution, inspired by (Kang, Li, and Zhou 2016). Specifically, we push the l_{21} -norm loss to a more strict l_1 -norm loss, *i.e.*,

$$\min_{\mathbf{B}} \|l\mathbf{S} - (\mathbf{W}^T \mathbf{Y})^T \mathbf{B}\|_1 \text{ s.t. } \mathbf{B} \in \{-1, 1\}^{l \times n}. \quad (7)$$

where $\|\cdot\|_1$ denotes the l_1 -norm. Problem (7) reaches its minimum when $\mathbf{B} = \text{sgn}(\mathbf{W}^T \mathbf{Q})$, where $\mathbf{Q} = l\mathbf{Y}\mathbf{S}$ can be calculated beforehand. We denote this as **SADIH-LI**.

Alternatively, we can directly optimize the l_{21} -norm loss with an equivalent transformation. We first define $\mathbf{D} \in \mathbb{R}^{n \times n}$ as a diagonal matrix, the i -th diagonal element of which is defined as $d_{ii} = 1/2\|\mathbf{u}^i\|$, where \mathbf{u}^i is the i -th row of $(l\mathbf{S} - \mathbf{Y}^T \mathbf{W}\mathbf{B})$. In this way, we need to solve

$$\begin{aligned} \min_{\mathbf{B}} & \text{Tr} \left((l\mathbf{S} - \mathbf{Y}^T \mathbf{W}\mathbf{B})^T \mathbf{D} (l\mathbf{S} - \mathbf{Y}^T \mathbf{W}\mathbf{B}) \right) \\ \text{s.t. } & \mathbf{B} \in \{-1, 1\}^{l \times n}, \end{aligned} \quad (8)$$

which can be equivalently rewritten as

$$\begin{aligned} \min_{\mathbf{B}} & \text{Tr} (\mathbf{R}^T \mathbf{B} \mathbf{D} \mathbf{B}^T \mathbf{R}) - 2 \text{tr} (\mathbf{B}^T \mathbf{M}) \\ \text{s.t. } & \mathbf{B} \in \{-1, 1\}^{l \times n}, \end{aligned} \quad (9)$$

where $\mathbf{R} = \mathbf{W}^T \mathbf{Y}$, $\mathbf{M} = \mathbf{W}^T \mathbf{D} \mathbf{Q}$, and $\text{Tr}(\cdot)$ is the trace norm. For this binary quadratic program problem, we employ the discrete cyclic coordinate descent (DCC) method (Shen et al. 2015) to sequentially learn each row of \mathbf{B} while fixing other rows. Let \mathbf{b}^T be the k -th row of \mathbf{B} , $k = 1, \dots, l$, and $\bar{\mathbf{B}}$ be the matrix of \mathbf{B} excluding \mathbf{b} . Similarly, let \mathbf{r}^T and \mathbf{q}^T be the k -th row of \mathbf{R} and \mathbf{M} , respectively. $\bar{\mathbf{R}}$ and $\bar{\mathbf{M}}$ are the matrix of \mathbf{R} excluding \mathbf{r} and the matrix of \mathbf{M} excluding \mathbf{q} , respectively. Then, we have

$$\begin{aligned} & \text{Tr} (\mathbf{R}^T \mathbf{B} \mathbf{D} \mathbf{B}^T \mathbf{R}) \\ & = \text{Tr} \left((\bar{\mathbf{B}}^T \bar{\mathbf{R}} + \mathbf{b} \mathbf{r}^T) \mathbf{D} (\bar{\mathbf{R}} \bar{\mathbf{B}} + \mathbf{r} \mathbf{b}^T) \right) \\ & = 2 \text{Tr} (\mathbf{b} \mathbf{r}^T \mathbf{D} \bar{\mathbf{R}} \bar{\mathbf{B}}) + \text{Tr} (\mathbf{b} \mathbf{r}^T \mathbf{D} \mathbf{r} \mathbf{b}^T) + \text{const}. \end{aligned} \quad (10)$$

Since $\text{Tr}(\mathbf{b} \mathbf{r}^T \mathbf{D} \mathbf{r} \mathbf{b}^T) = \text{Tr}(\mathbf{r}^T \mathbf{D} \mathbf{r} \mathbf{b} \mathbf{b}^T) = n \sum_i d_{ii} \mathbf{r}^T \mathbf{r}$, this term is a constant *w.r.t.* \mathbf{b} .

Similarly, $\text{tr} (\mathbf{B}^T \mathbf{M}) = \begin{bmatrix} \mathbf{b}^T \mathbf{q} & \mathbf{b}^T \bar{\mathbf{M}} \\ \bar{\mathbf{B}}^T \mathbf{q} & \bar{\mathbf{B}}^T \bar{\mathbf{M}} \end{bmatrix}$, and then

$$\text{tr} (\mathbf{B}^T \mathbf{M}) = \mathbf{q}^T \mathbf{b} + \text{const}. \quad (11)$$

Therefore, Eqn. (9) can be reformulated as

$$\begin{aligned} \min_{\mathbf{b}} & \text{Tr} \left((\mathbf{r}^T \mathbf{D} \bar{\mathbf{R}} \bar{\mathbf{B}} - \mathbf{q}^T) \mathbf{b} \right) \\ \text{s.t. } & \mathbf{b} \in \{-1, 1\}^l, \end{aligned} \quad (12)$$

which has a closed-form solution:

$$\mathbf{b} = \text{sgn} (\mathbf{q} - \bar{\mathbf{B}}^T \bar{\mathbf{R}} \mathbf{D} \mathbf{r}). \quad (13)$$

We can see that each bit \mathbf{b} is calculated based on the pre-learned $(l-1)$ bits $\bar{\mathbf{B}}$. We iteratively update each bit until it converges to a set of optimal codes \mathbf{B} .

W-Step: When fixing \mathbf{B} , P_1 and P_2 , the objective function (6) *w.r.t.* \mathbf{W} is degenerated to

$$\begin{aligned} \min_{\mathbf{W}} & \text{Tr} \left((l\mathbf{S} - \mathbf{Y}^T \mathbf{W}\mathbf{B})^T \mathbf{D} (l\mathbf{S} - \mathbf{Y}^T \mathbf{W}\mathbf{B}) \right) \\ & + \alpha \|\mathbf{X} - P_2 \mathbf{W}^T \mathbf{Y}\|_F^2 + \beta \|\mathbf{W}^T \mathbf{Y} - P_1 \mathbf{X}\|_F^2 + \gamma \|\mathbf{W}^T \mathbf{Y}\|_F^2. \end{aligned} \quad (14)$$

The closed-form solution of \mathbf{W} is given by setting the derivation of (14) to zero, *i.e.*,

$$\begin{aligned} \mathbf{W} & = \mathbf{L} (\mathbf{Q} \mathbf{D} \mathbf{B}^T + \mathbf{Y} \mathbf{X}^T (\alpha P_2 + \beta P_1)) \\ & (\mathbf{B} \mathbf{D} \mathbf{B}^T + \alpha P_2 P_2^T + (\beta + \gamma) \mathbf{I})^{-1}, \end{aligned} \quad (15)$$

where $\mathbf{L} = (\mathbf{Y} \mathbf{Y}^T)^{-1}$ can be calculated beforehand.

F-Step: When fixing \mathbf{B} , \mathbf{W} and P_2 , problem (6) *w.r.t.* P_1 becomes

$$\min_{P_1} \|\mathbf{W}^T \mathbf{Y} - P_1 \mathbf{X}\|_F^2 \text{ s.t. } P_1^T P_1 = \mathbf{I}, \quad (16)$$

which can be solved by the following lemma.

Lemma 1. $P_1 = \mathbf{U} \mathbf{V}^T$ is the optimal solution to the problem in Eqn. (16), where \mathbf{U} and \mathbf{V} are the left and right singular matrices of the compact Singular Value Decomposition (SVD) on $(\mathbf{X} \mathbf{Y}^T \mathbf{W})$.

P-Step: Similarly, when fixing other variables, problem (6) *w.r.t.* P_2 can be re-written as:

$$\min_{P_2} \alpha \|\mathbf{X} - P_2 \mathbf{W}^T \mathbf{Y}\|_F^2 + \gamma \|P_2\|_F^2. \quad (17)$$

The minimal P_2 can be obtained by setting the partial derivative of Eqn. (17) to zero, and we have

$$P_2 = \mathbf{T} \mathbf{X} \mathbf{Y}^T \mathbf{W}, \quad (18)$$

where $\mathbf{T} = (\alpha \mathbf{X} \mathbf{X}^T + \gamma \mathbf{I})^{-1}$ can be computed beforehand.

The proposed optimization iteratively updates four variables until satisfying the convergence criteria. The convergence of the proposed optimization algorithm is guaranteed by the following theorem.

Theorem 1. *The alternating optimization steps of our method will monotonously decrease the value of the objective function until it converges to a local optima.*

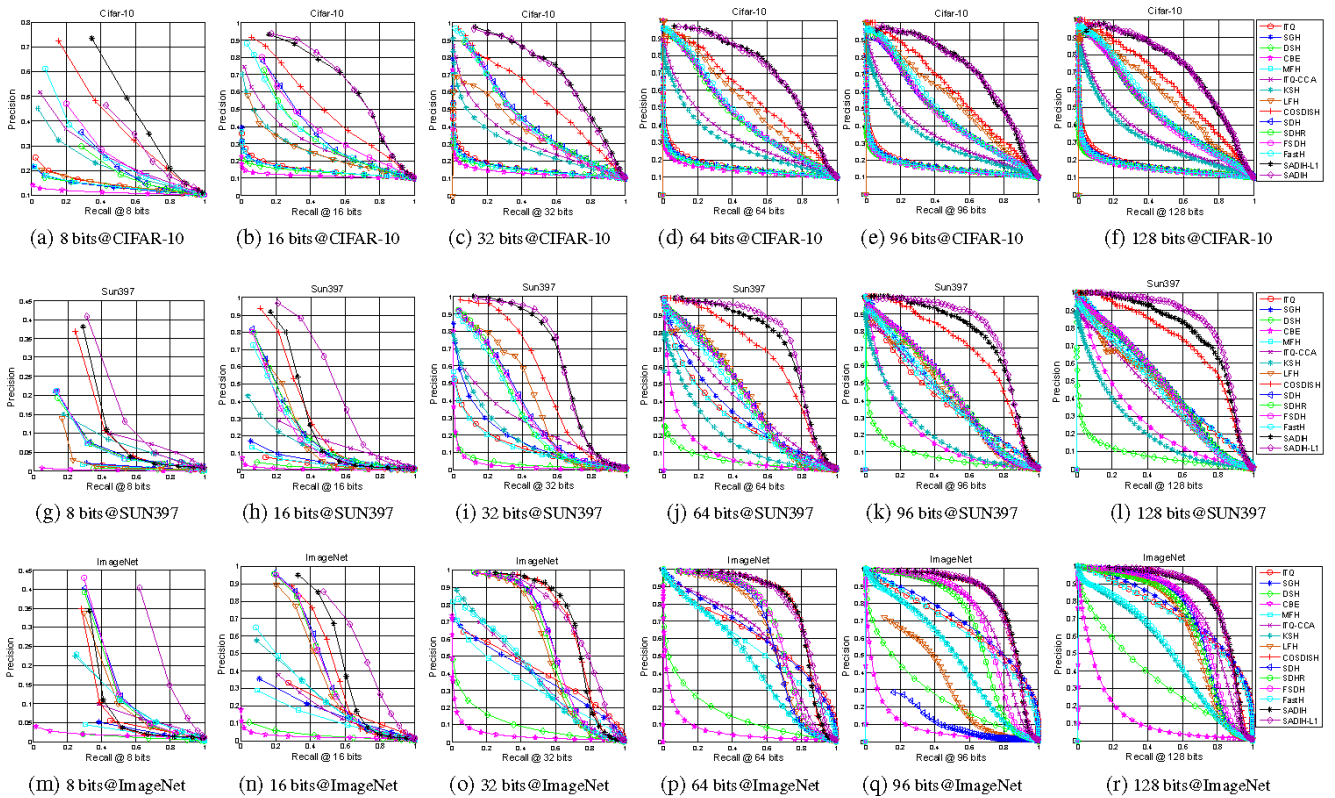


Figure 1: Precision-Recall curves of different methods with different code lengths on CIFAR-10, SUN397 and ImageNet.

Similarly, the 4,096-diml deep CNN features from VGG19 *fc7* were used for evaluation.

Experimental Settings: Following the previous experimental configurations used in (Shen et al. 2015; Kang, Li, and Zhou 2016), we randomly split the CIFAR-10 dataset into a training set (59K images) and a test query set (1,000 images), which has 100 images per category. In SUN-397, we randomly sample 100 images from each of the first 18 largest scene categories to formulate 1,800 query images. For ImageNet, we randomly select 50 images for each category from the validation image dataset to construct the 5,000 query image dataset. The semantic similarities on these datasets are measured whether two images share the same semantic label. For all datasets, we conduct feature normalization to make each dimension have zero-centered mean and equal variance.

Baseline Methods and Implementation Details: In experiments, we compare our SADIH and SADIH-L1 with 14 hashing methods including six unsupervised hashing methods, (*i.e.*, ITQ, SGH, DSH (Jin et al. 2014), CBE, and OCH) and eight supervised hashing methods (*i.e.*, ITQ-CCA (Gong et al. 2013), KSH, LFH, COSDISH, SDH, SDHR, FSDH, FastH). All the compared algorithms were performed five times with different random initializations, and the averaged experimental results were reported. To make fair comparison, all the compared methods were reimplemented using the released source codes given by the corresponding authors. Specifically, we searched the best parameters

carefully for each algorithm by five-folds cross-validation, or directly employed the default parameters suggested by the original papers. For graph based method such as KSH and OCH, using the full semantic information for training is impossible due to the heavy computation complexity, and 5,000 samples were selected from the training data for model construction. For our SSAH, the parameter γ was empirically set to 0.001. For the parameters α and β , we should tune it by cross-validation from the candidate set $\{0.01, 0.1, 1.0, 5, 10\}$. The maximum iteration number t was set to 5, which could assure the best performance.

Evaluation Measures: We adopted three frequently-used performance metrics (Manning, Raghavan, and others 2008) to evaluate different methods, *i.e.* mean average precision (MAP), the precision-recall curves and mean precision rate curves of top 1000 returned samples. Moreover, we also compared the computation time to show efficiency.

Quantitative Results: Table 1 explicitly illustrates the retrieval results of different algorithms with MAP and Precision@top100 as well as the computation time on CIFAR-10, SUN397 and ImageNet datasets. 1) Generally, supervised methods can achieve higher accuracies than unsupervised methods, while OCH and ITQ can lead to competitive performances. Importantly, the unsupervised methods are unsuitable to deal with large-scale image searching with large number of classes, such as SUN397 and ImageNet, while they can handle datasets with fewer categories such as CIFAR-10. 2) Our SADIH and SADIH-L1 in most cases can

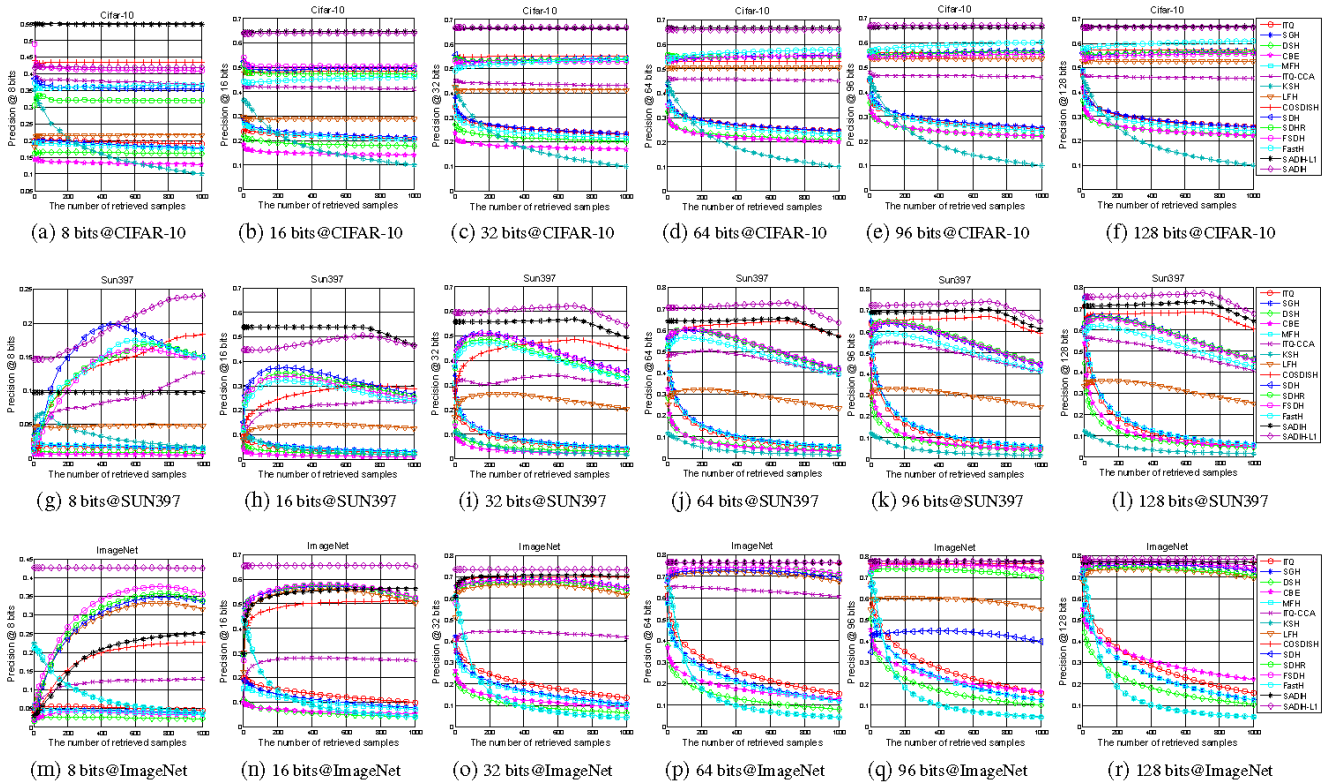


Figure 2: Precision curves of up to 1000 retrieved samples of different methods with different code lengths on CIFAR-10, SUN397 and ImageNet. (Better to view in color)

achieve the highest MAP scores with different code lengths, which demonstrate the efficacy of the proposed framework. Our SADIH and SADIH-L1 always produce superior performance to SDH, FSDH, FastH and SDHR, since our methods simultaneously consider the category information and pairwise similarities, however, other mentioned methods neglect the similarity preservation. Compared to KSH, LFH and COSDISH, our methods can make use of the *full* similarities and discriminative category-level information for learning effective binary codes, yielding superior performance. 3) SADIH-L1 has a tendency to outperform SADIH, since l_1 -norm regularization can generate more accurate approximation measurement. Increasing the coding lengths can improve the retrieval accuracies, but needs more training time. 4) Supervised methods seem to consume longer computation time compared to unsupervised ones. However, our SADIH-L1 is the fastest one in supervised methods, and SADIH can provide good balance between performance and time.

We further show the precision-recall curves of the compared methods with varying code lengths in Fig. 1. It can be observed that our methods are consistently better than all the competing methods, which indicates that our methods can retrieve more similar samples for a query at any fixed code length. Moreover, the precision variations *w.r.t.* different number of retrieved samples are illustrated in Fig. 2. We can observe that our methods are always superior to other methods, and their precisions are relatively stable with vary-

ing number of returned samples.

Conclusion

In this paper, we proposed a novel joint discriminative hashing framework, dubbed semantic-aware DIScrete Hashing (SADIH), which could efficiently guarantee the full semantic similarity preservation and discriminative semantic space construction. SADIH leveraged the asymmetric similarity approximation loss to preserve the full $n \times n$ similarities of the complete dataset. Meanwhile, the supervised semantic-aware autoencoder was designed to construct the discriminative semantic embedding space with full data variation preservation and good data reconstruction. The resulting problem was efficiently solved by the proposed discrete optimization algorithm. Extensive experimental results demonstrated the superiority of our methods on different large-scale datasets in terms of different evaluation protocols.

Acknowledgment This work is partially supported by ARC FT130101530.

References

- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*, 248–255. IEEE.
- Dong, W.; Charikar, M.; and Li, K. 2008. Asymmetric dis-

- tance estimation with sketches for similarity search in high-dimensional spaces. In *SIGIR*, 123–130. ACM.
- Erin Liong, V.; Lu, J.; Wang, G.; Moulin, P.; and Zhou, J. 2015. Deep hashing for compact binary codes learning. In *CVPR*, 2475–2483.
- Gionis, A.; Indyk, P.; Motwani, R.; et al. 1999. Similarity search in high dimensions via hashing. In *VLDB*, volume 99, 518–529.
- Gong, Y.; Lazebnik, S.; Gordo, A.; and Perronnin, F. 2013. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE TPAMI* 35(12):2916–2929.
- Gui, J.; Liu, T.; Sun, Z.; Tao, D.; and Tan, T. 2016. Supervised discrete hashing with relaxation. *TNNLS*.
- Gui, J.; Liu, T.; Sun, Z.; Tao, D.; and Tan, T. 2018. Fast supervised discrete hashing. *TPAMI* 40(2):490–496.
- Jiang, Q.-Y., and Li, W.-J. 2015. Scalable graph hashing with feature transformation. In *IJCAI*, 2248–2254.
- Jiang, K.; Que, Q.; and Kulis, B. 2015. Revisiting kernelized locality-sensitive hashing for improved large-scale image retrieval. In *CVPR*, 4933–4941.
- Jin, Z.; Li, C.; Lin, Y.; and Cai, D. 2014. Density sensitive hashing. *IEEE Trans. Cybern.* 44(8):1362–1371.
- Kang, W.-C.; Li, W.-J.; and Zhou, Z.-H. 2016. Column sampling based discrete supervised hashing. In *AAAI*, 1230–1236.
- Krizhevsky, A., and Hinton, G. 2009. Learning multiple layers of features from tiny images. Technical report, Citeseer.
- Kulis, B., and Darrell, T. 2009. Learning to hash with binary reconstructive embeddings. In *NIPS*, 1042–1050.
- Kulis, B., and Grauman, K. 2009. Kernelized locality-sensitive hashing for scalable image search. In *ICCV*, 2130–2137. IEEE.
- Lai, H.; Pan, Y.; Liu, Y.; and Yan, S. 2015. Simultaneous feature learning and hash coding with deep neural networks. In *CVPR*, 3270–3278.
- Lew, M.; Sebe, N.; Djeraba, C.; and Jain, R. 2006. Content-based multimedia information retrieval: State of the art and challenges. *ACM TOMM* 2(1):1–19.
- Li, W.-J.; Wang, S.; and Kang, W.-C. 2016. Feature learning based deep supervised hashing with pairwise labels. In *IJCAI*, 1711–1717.
- Lin, G.; Shen, C.; Shi, Q.; Van den Hengel, A.; and Suter, D. 2014. Fast supervised hashing with decision trees for high-dimensional data. In *CVPR*, 1963–1970.
- Lin, K.; Yang, H.-F.; Hsiao, J.-H.; and Chen, C.-S. 2015. Deep learning of binary hash codes for fast image retrieval. In *CVPRW*, 27–35.
- Liu, W.; Wang, J.; Ji, R.; Jiang, Y.; and Chang, S. 2012. Supervised hashing with kernels. In *CVPR*, 2074–2081. IEEE.
- Liu, H.; Ji, R.; Wang, J.; and Shen, C. 2018. Ordinal constraint binary coding for approximate nearest neighbor search. *TPAMI* (99):1–1.
- Luo, X.; Wu, Y.; and Xu, X. 2018. Scalable supervised discrete hashing for large-scale search. In *WWW*, 1603–1612.
- Manning, C. D.; Raghavan, P.; et al. 2008. *Introduction to information retrieval*. Cambridge University Press.
- Neyshabur, B.; Srebro, N.; Salakhutdinov, R. R.; Makarychev, Y.; and Yadollahpour, P. 2013. The power of asymmetry in binary hashing. In *NIPS*, 2823–2831.
- Salakhutdinov, R., and Hinton, G. 2009. Semantic hashing. *IJAR* 50(7):969–978.
- Shen, F.; Shen, C.; Liu, W.; and Tao Shen, H. 2015. Supervised discrete hashing. In *CVPR*, 37–45.
- Shen, F.; Xu, Y.; Liu, L.; Yang, Y.; Huang, Z.; and Shen, H. T. 2018. Unsupervised deep hashing with similarity-adaptive and discrete optimization. *TPAMI*.
- Shrivastava, A., and Li, P. 2014. Asymmetric lsh (alsh) for sublinear time maximum inner product search (mips). In *NIPS*, 2321–2329.
- Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Smeulders, A. W.; Worring, M.; Santini, S.; Gupta, A.; and Jain, R. 2000. Content-based image retrieval at the end of the early years. *TPAMI* (12):1349–1380.
- Song, J.; Yang, Y.; Huang, Z.; Shen, H. T.; and Hong, R. 2011. Multiple feature hashing for real-time large scale near-duplicate video retrieval. In *ACMM*, 423–432. ACM.
- Song, J.; Gao, L.; Liu, L.; Zhu, X.; and Sebe, N. 2018. Quantization-based hashing: a general framework for scalable image and video retrieval. *PR* 75:175–187.
- Wang, J.; Zhang, T.; Sebe, N.; Shen, H. T.; et al. 2018. A survey on learning to hash. *IEEE TPAMI* 40(4):769–790.
- Wang, J.; Kumar, S.; and Chang, S. 2012. Semi-supervised hashing for large-scale search. *TPAMI* (12):2393–2406.
- Weiss, Y.; Torralba, A.; and Fergus, R. 2009. Spectral hashing. In *NIPS*, 1753–1760.
- Xiao, J.; Hays, J.; Ehinger, K. A.; Oliva, A.; and Torralba, A. 2010. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 3485–3492. IEEE.
- Zhang, Z.; Liu, L.; Qin, J.; Zhu, F.; Shen, F.; Xu, Y.; Shao, L.; and Shen, H. 2018a. Highly-economized multi-view binary compression for scalable image clustering. In *ECCV*, 717–732.
- Zhang, Z.; Liu, L.; Shen, F.; Shen, H. T.; and Shao, L. 2018b. Binary multi-view clustering. *IEEE TPAMI* 99:1–1.
- Zhang, Z.; Shao, L.; Xu, Y.; Liu, L.; and Yang, J. 2018c. Marginal representation learning with graph structure self-adaptation. *IEEE TNNLS* 29:4645–4659.
- Zhang, P.; Zhang, W.; and Li, W.-J. 2014. Supervised hashing with latent factor models. In *SIGIR*, 173–182. ACM.