

# Hashing with Binary Autoencoders

Miguel Á. Carreira-Perpiñán  
EECS, University of California, Merced  
<http://eecs.ucmerced.edu>

Ramin Raziperchikolaei  
EECS, University of California, Merced  
<http://eecs.ucmerced.edu>

## Abstract

*An attractive approach for fast search in image databases is binary hashing, where each high-dimensional, real-valued image is mapped onto a low-dimensional, binary vector and the search is done in this binary space. Finding the optimal hash function is difficult because it involves binary constraints, and most approaches approximate the optimization by relaxing the constraints and then binarizing the result. Here, we focus on the binary autoencoder model, which seeks to reconstruct an image from the binary code produced by the hash function. We show that the optimization can be simplified with the method of auxiliary coordinates. This reformulates the optimization as alternating two easier steps: one that learns the encoder and decoder separately, and one that optimizes the code for each image. Image retrieval experiments show the resulting hash function outperforms or is competitive with state-of-the-art methods for binary hashing.*

## 1. Introduction

We consider the problem of binary hashing, where given a high-dimensional vector  $\mathbf{x} \in \mathbb{R}^D$ , we want to map it to an  $L$ -bit vector  $\mathbf{z} = \mathbf{h}(\mathbf{x}) \in \{0, 1\}^L$  using a hash function  $\mathbf{h}$ , while preserving the neighbors of  $\mathbf{x}$  in the binary space. Binary hashing has emerged in recent years as an effective technique for fast search on image (and other) databases. While the search in the original space would cost  $\mathcal{O}(ND)$  in both time and space, using floating point operations, the search in the binary space costs  $\mathcal{O}(NL)$  where  $L \ll D$  and the constant factor is much smaller. This is because the hardware can compute binary operations very efficiently and the entire dataset ( $NL$  bits) can fit in the main memory of a workstation. And while the search in the binary space will produce some false positives and negatives, one can retrieve a larger set of neighbors and then verify these with the ground-truth distance, while still being efficient.

Many different hashing approaches have been proposed in the last few years. They formulate an objective function of the hash function  $\mathbf{h}$  or of the binary codes that tries to

capture some notion of neighborhood preservation. Most of these approaches have two things in common:  $\mathbf{h}$  typically performs dimensionality reduction ( $L < D$ ) and, as noted, it outputs binary codes ( $\mathbf{h}: \mathbb{R}^D \rightarrow \{0, 1\}^L$ ). The latter implies a step function or binarization applied to a real-valued function of the input  $\mathbf{x}$ . Optimizing this is difficult. In practice, most approaches follow a two-step procedure: first they learn a real hash function ignoring the binary constraints and then the output of the resulting hash function is binarized (e.g. by thresholding or with an optimal rotation). For example, one can run a continuous dimensionality reduction algorithm (by optimizing its objective function) such as PCA and then apply a step function. This procedure can be seen as a “filter” approach [14] and is suboptimal: in the example, the thresholded PCA projection is not necessarily the best thresholded linear projection (i.e., the one that minimizes the objective function under all thresholded linear projections). To obtain the latter, we must optimize the objective jointly over linear mappings and thresholds, respecting the binary constraints while learning  $\mathbf{h}$ ; this is a “wrapper” approach [14]. In other words, optimizing real codes and then projecting them onto the binary space is not the same as optimizing the codes in the binary space.

In this paper we show this joint optimization, respecting the binary constraints during training, can actually be done reasonably efficiently. The idea is to use the recently proposed *method of auxiliary coordinates (MAC)* [5, 6]. This is a general strategy to transform an original problem involving a nested function into separate problems without nesting, each of which can be solved more easily. In our case, this allows us to reduce drastically the complexity due to the binary constraints. We focus on *binary autoencoders*, i.e., where the code layer is binary. Section 2 describes the binary autoencoder model and objective function. Section 3 derives a training algorithm using MAC and explains how, with carefully implemented steps, the optimization in the binary space can be carried out efficiently. Our hypothesis is that constraining the optimization to the binary space results in better hash functions and we test this in experiments (section 4). These show that our resulting hash function is consistently competitive with the state-of-the-art.

**Related work** The most basic hashing approaches are data-independent, such as Locality-Sensitive Hashing (LSH) [1], which is based on random projections and thresholding, and kernelized LSH [17]. Generally, they are outperformed by data-dependent methods, which learn a specific hash function for a given dataset in an unsupervised or supervised way. We focus here on unsupervised, data-dependent approaches. These are typically based on defining an objective function (usually based on dimensionality reduction) either of the hash function or the binary codes, and optimizing it. However, this is usually achieved by relaxing the binary codes to a continuous space and thresholding the resulting continuous solution. For example, spectral hashing [28] is essentially a version of Laplacian eigenmaps where the binary constraints are relaxed and approximate eigenfunctions are computed that are then thresholded to provide binary codes. Variations of this include using AnchorGraphs [20] to define the eigenfunctions, or obtaining the hash function directly as a binary classifier using the codes from spectral hashing as labels [31]. Other approaches optimize instead a nonlinear embedding objective that depends on a continuous, parametric hash function, which is then thresholded to define a binary hash function [27, 26]; or an objective that depends on a thresholded hash function, but the threshold is relaxed during the optimization [23]. Some recent work has tried to respect the binary nature of the codes or thresholds by using alternating optimization directly on an objective function, over one entry or one row of the weight matrix in the hash function [16, 21], or over a subset of the binary codes [19, 18]. Since the objective function involves a large number of terms and all the binary codes or weights are coupled, the optimization is very slow. Also, in [19, 18] the hash function is learned after the codes have been fixed, which is suboptimal.

The closest model to our binary autoencoder is Iterative Quantization (ITQ) [10], a fast and competitive hashing method. ITQ first obtains continuous low-dimensional codes by applying PCA to the data and then seeks a rotation that makes the codes as close as possible to binary. The latter is based on the optimal discretization algorithm of [30], which finds a rotation of the continuous eigenvectors of a graph Laplacian that makes them as close as possible to a discrete solution, as a postprocessing for spectral clustering. The ITQ objective function is

$$\min_{\mathbf{B}, \mathbf{R}} \|\mathbf{B} - \mathbf{V}\mathbf{R}\|^2 \text{ s.t. } \mathbf{B} \in \{-1, +1\}^{N \times L}, \mathbf{R} \text{ orthogonal} \quad (1)$$

where  $\mathbf{V}$  of  $N \times L$  are the continuous codes obtained by PCA. This is an NP-complete problem, and a local minimum is found using alternating optimization over  $\mathbf{B}$ , with solution  $\mathbf{B} = \text{sgn}(\mathbf{V}\mathbf{R})$  elementwise, and over  $\mathbf{R}$ , which is a Procrustes alignment problem with a closed-form solution based on a SVD. The final hash function is  $\mathbf{h}(\mathbf{x}) = \text{sgn}(\mathbf{W}\mathbf{x})$ , which has the form of a thresholded lin-

ear projection. Hence, ITQ is a postprocessing of the PCA codes, and it can be seen as a suboptimal approach to optimizing a binary autoencoder, where the binary constraints are relaxed during the optimization (resulting in PCA), and then one “projects” the continuous codes back to the binary space. Semantic hashing [26] also uses an autoencoder objective with a deep encoder (consisting of stacked RBMs), but again its optimization trains it as a continuous problem and then a threshold is applied to the encoder.

## 2. Our hashing models: binary autoencoder and binary factor analysis

We consider a well-known model for continuous dimensionality reduction, the (continuous) autoencoder, defined in a broad sense as the composition of an encoder  $\mathbf{h}(\mathbf{x})$  which maps a real vector  $\mathbf{x} \in \mathbb{R}^D$  onto a real code vector  $\mathbf{z} \in \mathbb{R}^L$  (with  $L < D$ ), and a decoder  $\mathbf{f}(\mathbf{z})$  which maps  $\mathbf{z}$  back to  $\mathbb{R}^D$  in an effort to reconstruct  $\mathbf{x}$ . Although our ideas apply more generally to other encoders, decoders and objective functions, in this paper we mostly focus on the least-squares error with a linear encoder and decoder. As is well known, the optimal solution is PCA.

For hashing, the encoder maps continuous inputs onto binary code vectors with  $L$  bits,  $\mathbf{z} \in \{0, 1\}^L$ . Let us write  $\mathbf{h}(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x})$  ( $\mathbf{W}$  includes a bias by having an extra dimension  $x_0 = 1$  for each  $\mathbf{x}$ ) where  $\mathbf{W} \in \mathbb{R}^{L \times (D+1)}$  and  $\sigma(t)$  is a step function applied elementwise, i.e.,  $\sigma(t) = 1$  if  $t \geq 0$  and  $\sigma(t) = 0$  otherwise (we can fix the threshold at 0 because the bias acts as a threshold for each bit). Our *desired hash function* will be  $\mathbf{h}$ , and it should minimize the following problem, given a dataset of high-dimensional patterns  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ :

$$E_{\text{BA}}(\mathbf{h}, \mathbf{f}) = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{f}(\mathbf{h}(\mathbf{x}_n))\|^2 \quad (2)$$

which is the usual least-squares error but where the code layer is binary. Optimizing this nonsmooth function is difficult and NP-complete. Where the gradients do exist wrt  $\mathbf{W}$  they are zero nearly everywhere. We call this a *binary autoencoder (BA)*.

We will also consider a related model (see later):

$$E_{\text{BFA}}(\mathbf{Z}, \mathbf{f}) = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{f}(\mathbf{z}_n)\|^2 \text{ s.t. } \mathbf{z}_n \in \{0, 1\}^L \quad n = 1, \dots, N \quad (3)$$

where  $\mathbf{f}$  is linear and we optimize over the decoder  $\mathbf{f}$  and the binary codes  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_N)$  of each input pattern. Without the binary constraint, i.e.,  $\mathbf{Z} \in \mathbb{R}^{L \times N}$ , this model dates back to the 50s and is sometimes called least-squares factor analysis [29], and its solution is PCA. With the binary constraints, the problem is NP-complete because it includes as particular case solving a linear system over  $\{0, 1\}^L$ , which

is an integer LP feasibility problem. We call this model (least-squares) *binary factor analysis (BFA)*. We believe this model has not been studied before, at least in hashing. A hash function  $\mathbf{h}$  can be obtained from BFA by fitting a binary classifier of the inputs to each of the  $L$  code bits. It is a filter approach, while the BA is the optimal (wrapper) approach, since it optimizes (2) jointly over  $\mathbf{f}$  and  $\mathbf{h}$ .

Let us compare BFA (without bias term in  $\mathbf{f}$ ) with ITQ in eq. (1). BFA takes the form  $\min_{\mathbf{Z}, \mathbf{A}} \|\mathbf{X} - \mathbf{AZ}\|^2$  s.t.  $\mathbf{Z} \in \{0, 1\}^{NL}$ . The main difference is that in BFA the binary variables do not appear in separable form and so the step over them cannot be solved easily.

### 3. Optimization of BA and BFA using the method of auxiliary coordinates (MAC)

We use the recently proposed *method of auxiliary coordinates (MAC)* [5, 6]. The idea is to break nested functional relationships judiciously by introducing variables as equality constraints. These are then solved by optimizing a penalized function using alternating optimization over the original parameters and the coordinates, which results in a coordination-minimization (CM) algorithm. Recall eq. (2), this is our nested problem, where the model is  $\mathbf{y} = \mathbf{f}(\mathbf{h}(\mathbf{x}))$ . We introduce as auxiliary coordinates the outputs of  $\mathbf{h}$ , i.e., the codes for each of the  $N$  input patterns, and obtain the following equality-constrained problem:

$$\min_{\mathbf{h}, \mathbf{f}, \mathbf{Z}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{f}(\mathbf{z}_n)\|^2 \text{ s.t. } \mathbf{z}_n = \mathbf{h}(\mathbf{x}_n) \in \{0, 1\}^L \quad n = 1, \dots, N. \quad (4)$$

Note the codes are binary. We now apply the quadratic-penalty method (it is also possible to apply the augmented Lagrangian method instead; [22]) and minimize the following objective function while progressively increasing  $\mu$ , so the constraints are eventually satisfied:

$$E_Q(\mathbf{h}, \mathbf{f}, \mathbf{Z}; \mu) = \sum_{n=1}^N \left( \|\mathbf{x}_n - \mathbf{f}(\mathbf{z}_n)\|^2 + \mu \|\mathbf{z}_n - \mathbf{h}(\mathbf{x}_n)\|^2 \right) \text{ s.t. } \mathbf{z}_n \in \{0, 1\}^L, n = 1, \dots, N. \quad (5)$$

Now we apply alternating optimization over  $\mathbf{Z}$  and  $(\mathbf{h}, \mathbf{f})$ . This results in the following two steps:

- Over  $\mathbf{Z}$  for fixed  $(\mathbf{h}, \mathbf{f})$ , the problem separates for each of the  $N$  codes. The optimal code vector for pattern  $\mathbf{x}_n$  tries to be close to the prediction  $\mathbf{h}(\mathbf{x}_n)$  while reconstructing  $\mathbf{x}_n$  well.
- Over  $(\mathbf{h}, \mathbf{f})$  for fixed  $\mathbf{Z}$ , we obtain  $L + 1$  independent problems for each of the  $L$  single-bit hash functions (which try to predict  $\mathbf{Z}$  optimally from  $\mathbf{X}$ ), and for  $\mathbf{f}$  (which tries to reconstruct  $\mathbf{X}$  optimally from  $\mathbf{Z}$ ).

We can now see the advantage of the auxiliary coordinates: the individual steps are (reasonably) easy to solve (although some work is still needed, particularly for the  $\mathbf{Z}$  step), and

```

input  $\mathbf{X}_{D \times N} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ ,  $L \in \mathbb{N}$ 
Initialize  $\mathbf{Z}_{L \times N} = (\mathbf{z}_1, \dots, \mathbf{z}_N) \in \{0, 1\}^{LN}$ 
for  $\mu = 0 < \mu_1 < \dots < \mu_\infty$ 
  for  $l = 1, \dots, L$  h step
     $h_l \leftarrow$  fit SVM to  $(\mathbf{X}, \mathbf{Z}_{:,l})$ 
  f  $\leftarrow$  least-squares fit to  $(\mathbf{Z}, \mathbf{X})$  f step
  for  $n = 1, \dots, N$  Z step
     $\mathbf{z}_n \leftarrow \arg \min_{\mathbf{z}_n \in \{0, 1\}^L} \|\mathbf{x}_n - \mathbf{f}(\mathbf{z}_n)\|^2 + \mu \|\mathbf{z}_n - \mathbf{h}(\mathbf{x}_n)\|^2$ 
  if no change in  $\mathbf{Z}$  and  $\mathbf{Z} = \mathbf{h}(\mathbf{X})$  then stop
return  $\mathbf{h}, \mathbf{Z} = \mathbf{h}(\mathbf{X})$ 

```

Figure 1. Binary autoencoder MAC algorithm.

besides they exhibit significant parallelism. We describe the steps in detail below. The resulting algorithm alternates steps over the encoder ( $L$  classifications) and decoder (one regression) and over the codes ( $N$  binary proximal operators; [25, 8]). During the iterations, we allow the encoder and decoder to be mismatched, since the encoder output does not equal the decoder input, but they are coordinated by  $\mathbf{Z}$  and as  $\mu$  increases the mismatch is reduced. The overall MAC algorithm to optimize a BA is in fig. 1.

Although a MAC algorithm can be shown to produce convergent algorithms as  $\mu \rightarrow \infty$  with a differentiable objective function, we cannot apply the theorem in [5] because of the binary nature of the problem. Instead, we show that our algorithm converges to a local minimum for a *finite*  $\mu$ , where “local minimum” is understood as in  $k$ -means: a point where  $\mathbf{Z}$  is globally minimum given  $(\mathbf{h}, \mathbf{f})$  and vice versa. The following theorem is valid for any choice of  $\mathbf{h}$  and  $\mathbf{f}$ , not just linear.

**Theorem 3.1.** *Assume the steps over  $\mathbf{h}, \mathbf{f}$  are solved exactly by finding their unique global minimum<sup>1</sup>. Then the MAC algorithm for the binary autoencoder stops at a finite  $\mu$ .*

*Proof.*  $\mu$  appears only in the  $\mathbf{Z}$  step. If  $\mathbf{Z}$  does not change there,  $\mathbf{f}$  and  $\mathbf{h}$  will not change either, since the  $\mathbf{h}$  and  $\mathbf{f}$  steps are exact. The  $\mathbf{Z}$  step over  $\mathbf{z}_n$  minimizes  $\|\mathbf{x}_n - \mathbf{f}(\mathbf{z}_n)\|^2 + \mu \|\mathbf{z}_n - \mathbf{h}(\mathbf{x}_n)\|^2$ , and from theorem 3.3 we have that  $\mathbf{h}(\mathbf{x}_n)$  is a global minimizer if  $\mu > \|\mathbf{x}_n - \mathbf{f}(\mathbf{h}(\mathbf{x}_n))\|^2$ . The statement follows from the fact that  $\|\mathbf{x}_n - \mathbf{f}(\mathbf{z})\|^2$  is bounded over all  $n, \mathbf{z}$  and  $\mathbf{f}$ . Let us prove this fact. Clearly this holds for fixed  $\mathbf{f}$  because  $(n, \mathbf{z}_n)$  take values on a finite set, namely  $\{1, \dots, N\} \times \{0, \dots, 2^L - 1\}$ . As for  $\mathbf{f}$ , even if the set of functions  $\mathbf{f}$  is infinite, the number of different functions  $\mathbf{f}$  that are possible is finite, because  $\mathbf{f}$  results from an exact fit to  $(\mathbf{Z}, \mathbf{X})$ , where  $\mathbf{X}$  is fixed and the set of possible  $\mathbf{Z}$  is finite (since each  $z_{nl}$  is binary).  $\square$

The minimizers of  $E_Q(\mathbf{h}, \mathbf{f}, \mathbf{Z}; \mu)$  trace a path as a function of  $\mu \geq 0$  in the  $(\mathbf{h}, \mathbf{f}, \mathbf{Z})$  space. BA and BFA can be seen as the limiting cases of  $E_Q(\mathbf{h}, \mathbf{f}, \mathbf{Z}; \mu)$  when  $\mu \rightarrow \infty$

<sup>1</sup>This technical condition can always be guaranteed for linear  $\mathbf{h}$  or  $\mathbf{f}$  by adding a quadratic term with vanishing weight, for example.

and  $\mu \rightarrow 0^+$ , respectively (for BFA,  $\mathbf{f}$  and  $\mathbf{Z}$  can be optimized independently from  $\mathbf{h}$ , but  $\mathbf{h}$  must optimally fit the resulting  $\mathbf{Z}$ ). In practice, to learn the BFA model we set  $\mu$  to a small value and keep it constant while running the BA algorithm. As for BA itself, we increase  $\mu$  (times a constant factor, e.g. 2) and iterate the  $\mathbf{Z}$  and  $(\mathbf{h}, \mathbf{f})$  steps for each  $\mu$  value. Usually the algorithm stops in 10 to 20 iterations, when no further changes to the parameters occur.

### 3.1. $\mathbf{f}$ step

With a linear decoder this is a simple linear regression  $\min_{\mathbf{A}, \mathbf{b}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{A}\mathbf{z}_n - \mathbf{b}\|^2$  with data  $(\mathbf{Z}, \mathbf{X})$ , whose solution is (ignoring the bias for simplicity)  $\mathbf{A} = \mathbf{X}\mathbf{Z}^T(\mathbf{Z}\mathbf{Z}^T)^{-1}$  and can be computed in  $\mathcal{O}(NDL)$ . Note the constant factor in the  $\mathcal{O}$ -notation is small because  $\mathbf{Z}$  is binary, e.g.  $\mathbf{X}\mathbf{Z}^T$  involves only sums, not multiplications.

### 3.2. $\mathbf{h}$ step

This has the following form:

$$\begin{aligned} \min_{\mathbf{h}} \sum_{n=1}^N \|\mathbf{z}_n - \mathbf{h}(\mathbf{x}_n)\|^2 &= \min_{\mathbf{W}} \sum_{n=1}^N \|\mathbf{z}_n - \sigma(\mathbf{W}\mathbf{x}_n)\|^2 \\ &= \sum_{l=1}^L \min_{\mathbf{w}_l} \sum_{n=1}^N (z_{nl} - \sigma(\mathbf{w}_l^T \mathbf{x}_n))^2. \end{aligned} \quad (6)$$

Since  $\mathbf{Z}$  and  $\sigma(\cdot)$  are binary,  $\|\cdot\|^2$  is the Hamming distance and the objective function is the number of misclassified patterns, so it separates for each bit. So it is a classification problem for each bit, using as labels the auxiliary coordinates, where  $h_l$  is a linear classifier (a perceptron). However, rather than minimizing this, we will solve an easier, closely related problem: fit a linear SVM  $h_l$  to  $(\mathbf{X}, \mathbf{Z}_{\cdot l})$  where we use a high penalty for misclassified patterns but optimize the margin plus the slack. Besides being easier (by reusing well-developed codes for SVMs), this surrogate loss has the advantage of making the solution unique (no local optima) and generalizing better to test data (maximum margin). Also, although we used a quadratic penalty, the spirit of penalty methods is to penalize constraint violations  $(\mathbf{z}_n - \mathbf{h}(\mathbf{x}_n))$  increasingly. Since in the limit  $\mu \rightarrow \infty$  the constraints are satisfied exactly, the classification error using  $\mathbf{h}$  is zero, hence the linear SVM will find an optimum of the nested problem anyway. We use LIBLINEAR [9] with warm start (i.e., the SVM optimization is initialized from the previous iteration's SVM). Note the  $L$  SVMs and the decoder function  $\mathbf{f}$  can be trained in parallel.

### 3.3. $\mathbf{Z}$ step

From eq. (5), this is a binary optimization on  $NL$  variables, but it separates into  $N$  independent optimizations each on only  $L$  variables, with the form of a binary proximal operator [25, 8] (where we omit the index  $n$ ):

$$\min_{\mathbf{z}} e(\mathbf{z}) = \|\mathbf{x} - \mathbf{f}(\mathbf{z})\|^2 + \mu \|\mathbf{z} - \mathbf{h}(\mathbf{x})\|^2 \text{ s.t. } \mathbf{z} \in \{0, 1\}^L.$$

Thus, although the problem over each  $\mathbf{z}_n$  is binary and NP-complete, a good or even exact solution may be obtained, because practical values of  $L$  are small (typically 8 to 32 bits). Further, because of the intensive computation and large number of independent problems, this step can take much advantage of parallel processing.

We have spent significant effort into making this step efficient while yielding good, if not exact, solutions. Before proceeding, let us show how to reduce the problem, which as stated uses a matrix of  $D \times L$ , to an equivalent problem using a matrix of  $L \times L$ .

**Theorem 3.2.** *Let  $\mathbf{x} \in \mathbb{R}^D$  and  $\mathbf{A} \in \mathbb{R}^{D \times L}$ , with QR factorisation  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q}$  is of  $D \times L$  with  $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}$  and  $\mathbf{R}$  is upper triangular of  $L \times L$ , and  $\mathbf{y} = \mathbf{Q}^T\mathbf{x} \in \mathbb{R}^L$ . The following two problems have the same minima over  $\mathbf{z}$ :*

$$\min_{\mathbf{z} \in \{0, 1\}^L} \|\mathbf{x} - \mathbf{A}\mathbf{z}\|^2 \quad \min_{\mathbf{z} \in \{0, 1\}^L} \|\mathbf{y} - \mathbf{R}\mathbf{z}\|^2. \quad (7)$$

*Proof.* Let  $(\mathbf{Q} \ \mathbf{Q}_\perp)$  of  $D \times D$  be orthogonal, where the columns of  $\mathbf{Q}_\perp$  are an orthonormal basis of the nullspace of  $\mathbf{Q}^T$ . Then, since orthogonal matrices preserve Euclidean distances, we have:  $\|\mathbf{x} - \mathbf{A}\mathbf{z}\|^2 = \|(\mathbf{Q} \ \mathbf{Q}_\perp)^T(\mathbf{x} - \mathbf{A}\mathbf{z})\|^2 = \|\mathbf{Q}^T(\mathbf{x} - \mathbf{A}\mathbf{z})\|^2 + \|\mathbf{Q}_\perp^T(\mathbf{x} - \mathbf{A}\mathbf{z})\|^2 = \|\mathbf{y} - \mathbf{R}\mathbf{z}\|^2 + \|\mathbf{Q}_\perp^T\mathbf{x}\|^2$ , where the term  $\|\mathbf{Q}_\perp^T\mathbf{x}\|^2$  does not depend on  $\mathbf{z}$ .  $\square$

This achieves a speedup of  $2D/L$  (where the 2 factor comes from the fact that the new matrix is triangular), e.g. this is  $40\times$  if using 16 bits with  $D = 320$  GIST features in our experiments. Henceforth, we redefine the  $\mathbf{z}$  step as  $\min e(\mathbf{z}) = \|\mathbf{y} - \mathbf{R}\mathbf{z}\|^2 + \mu \|\mathbf{z} - \mathbf{h}(\mathbf{x})\|^2$  s.t.  $\mathbf{z} \in \{0, 1\}^L$ .

**Enumeration** For small  $L$ , this can be solved *exactly* by enumeration, at a worst-case runtime cost  $\mathcal{O}(L^{2^2L})$ , but with small constant factors in practice (see accelerations below).  $L = 16$  is perfectly practical in a workstation without parallel processing for the datasets in our experiments.

**Alternating optimization** For larger  $L$ , we use alternating optimization over groups of  $g$  bits (where the optimization over a  $g$ -bit group is done by enumeration and uses the same accelerations). This converges to a local minimum of the  $\mathbf{Z}$  step, although we find in our experiments that it finds near-global optima *if using a good initialization*. Intuitively, it makes sense to warm-start this, i.e., to initialize  $\mathbf{z}$  to the code found in the previous iteration's  $\mathbf{Z}$  step, since this should be close to the new optimum as we converge. However, empirically we find that the codes change a lot in the first few iterations, and that the following initialization works better (in leading to a lower objective value) in early iterations: we solve the relaxed problem on  $\mathbf{z}$  s.t.  $\mathbf{z} \in [0, 1]^L$  rather than  $\{0, 1\}^L$ . This is a strongly convex bound-constrained quadratic program (QP) in  $L$  variables for  $\mu > 0$  and its unique minimizer can be found efficiently.

We can further speed up the solution by noting that we have  $N$  QPs with some common, special structure. The objective is the sum of a term having the same matrix  $\mathbf{R}$  for all

QPs, and a term that is separable in  $\mathbf{z}$ . We have developed an ADMM algorithm [4] that is very simple, parallelizes or vectorizes very well, and reuses matrix factorizations over all  $N$  QPs. It is 10–100× faster than Matlab’s `quadprog`. We warm-start it from the continuous solution of the QP in the previous  $\mathbf{Z}$  step.

In order to binarize the continuous minimizer for  $\mathbf{z}_n$  we could simply round its  $L$  elements, but instead we apply a greedy procedure that is efficient and better (though still suboptimal). We optimally binarize from bit 1 to bit  $L$  by evaluating the objective function for bit  $l$  in  $\{0, 1\}$  with all remaining elements fixed (elements 1 to  $l - 1$  are already binary and  $l + 1$  to  $L$  are still continuous) and picking the best. Essentially, this is one pass of alternating optimization but having continuous values for some of the bits.

Finally, we pick the best of the binarized relaxed solution or the warm-start value and run alternating optimization. This ensures that the quadratic-penalty function (5) decreases monotonically at each iteration.

**Accelerations** Naively, the enumeration involves evaluating  $e(\mathbf{z})$  for  $2^L$  (or  $2^g$ ) vectors, where evaluating  $e(\mathbf{z})$  for one  $\mathbf{z}$  costs on average roughly  $L + 1$  multiplications and  $\frac{1}{4}L^2$  sums. This enumeration can be sped up or pruned while still finding a global minimum by using upper bounds on  $e$ , incremental computation of  $e$ , and necessary and sufficient conditions for the solution. Essentially, we need not evaluate every code vector, or every bit of every code vectors; we know the solution will be “near”  $\mathbf{h}(\mathbf{x})$ ; and we can recognize the solution when we find it.

Call  $\mathbf{z}^*$  a global minimizer of  $e(\mathbf{z})$ . An initial, good upper bound is  $e(\mathbf{h}(\mathbf{x})) = \|\mathbf{y} - \mathbf{R}\mathbf{h}(\mathbf{x})\|^2$ . In fact, we have the following sufficient condition for  $\mathbf{h}(\mathbf{x})$  to be a global minimizer. (We give it generally for any decoder  $\mathbf{f}$ .)

**Theorem 3.3.** *Let  $e(\mathbf{z}) = \|\mathbf{x} - \mathbf{f}(\mathbf{z})\|^2 + \mu\|\mathbf{z} - \mathbf{h}(\mathbf{x})\|^2$ . Then: (1) A global minimizer  $\mathbf{z}^*$  of  $e(\mathbf{z})$  is at a Hamming distance from  $\mathbf{h}(\mathbf{x})$  of  $\frac{1}{\mu}\|\mathbf{x} - \mathbf{f}(\mathbf{h}(\mathbf{x}))\|^2$  or less. (2) If  $\mu > \|\mathbf{x} - \mathbf{f}(\mathbf{h}(\mathbf{x}))\|^2$  then  $\mathbf{h}(\mathbf{x})$  is a global minimizer.*

*Proof.*  $e(\mathbf{z}^*) = \|\mathbf{x} - \mathbf{f}(\mathbf{z}^*)\|^2 + \mu\|\mathbf{z}^* - \mathbf{h}(\mathbf{x})\|^2 \leq e(\mathbf{h}(\mathbf{x})) = \|\mathbf{x} - \mathbf{f}(\mathbf{h}(\mathbf{x}))\|^2 \Rightarrow \|\mathbf{z}^* - \mathbf{h}(\mathbf{x})\|^2 \leq \frac{1}{\mu}\|\mathbf{x} - \mathbf{f}(\mathbf{h}(\mathbf{x}))\|^2$ . (2) follows because the Hamming distance is integer.  $\square$

As  $\mu$  increases and  $\mathbf{h}$  improves, this bound becomes more effective and more of the  $N$  patterns are pruned. Upon convergence, the  $\mathbf{Z}$  step costs only  $\mathcal{O}(NL^2)$ . If we do have to search for a given  $\mathbf{z}_n$ , we keep a running bound (current best minimum)  $\bar{e} = \|\mathbf{y}_n - \mathbf{R}\bar{\mathbf{z}}\|^2 + \mu\|\bar{\mathbf{z}} - \mathbf{h}(\mathbf{x}_n)\|^2$ , and we scan codes in increasing Hamming distance to  $\mathbf{h}(\mathbf{x}_n)$  up to a distance of  $\bar{e}/\mu$ . Thus, we try first the codes that are more likely to be optimal, and keep refining the bound as we find better codes.

Second, since  $\|\mathbf{y} - \mathbf{R}\mathbf{z}\|^2$  separates over dimensions  $1, \dots, L$ , we evaluate it incrementally (dimension by dimension) and stop as soon as we exceed the running bound.

Finally, there exist global optimality necessary and sufficient conditions for binary quadratic problems that are easy to evaluate [2, 13]. This allows us to recognize the solution as soon as we reach it and stop the search (rather than do a linear search of all values, keeping track of the minimum). These conditions can also determine whether the continuous solution to the relaxed QP is a global minimizer of the binary QP.

### 3.4. Schedule for the penalty parameter $\mu$

The only user parameters in our method are the initialization for the binary codes  $\mathbf{Z}$  and the schedule for the penalty parameter  $\mu$  (sequence of values  $0 < \mu_1 < \dots < \infty$ ), since we use a penalty or augmented Lagrangian method. In general with these methods, setting the schedule requires some tuning in practice. Fortunately, this is simplified in our case for two reasons. 1) We need not drive  $\mu \rightarrow \infty$  because *termination occurs at a finite  $\mu$  and can be easily detected*: whenever  $\mathbf{Z}$  does not change compared to the previous  $\mathbf{Z}$  step, no further changes to the parameters can occur. This gives a practical stopping criterion. 2) In order to generalize well to unseen data, we stop iterating not when we (sufficiently) optimize  $E_Q(\mathbf{h}, \mathbf{f}; \mathbf{Z}; \mu)$ , but *when the precision in a validation set decreases*. This is a form of early stopping that guarantees that we improve (or leave unchanged) the initial  $\mathbf{Z}$ , and besides is faster. The initialization for  $\mathbf{Z}$  and further details about the schedule for  $\mu$  appear in section 4.

## 4. Experiments

We used three datasets in our experiments, commonly used as benchmarks for image retrieval. (1) CIFAR [15] contains 60 000  $32 \times 32$  color images in 10 classes. We ignore the labels in this paper and use  $N = 50\,000$  images as training set and 10 000 images as test set. We extract  $D = 320$  GIST features [24] from each image. (2) NUS-WIDE [7] contains  $N = 269\,648$  high-resolution color images and use  $N = 161\,789$  for training and 107 859 for test. We extract  $D = 128$  wavelet features [24] from each image. (3) SIFT-1M [12] contains  $N = 1\,000\,000$  training high-resolution color images and 10 000 test images, each represented by  $D = 128$  SIFT features.

We report precision and recall (%) in the test set using as true neighbors the  $K$  nearest images in Euclidean distance in the original space, and as retrieved neighbors in the binary space we either use the  $k$  nearest images in Hamming distance, or the images within a Hamming distance  $r$  (if no images satisfy the latter, we report zero precision).

Our experiments evaluate the effectiveness of our algorithm to minimize the BA objective and whether this translates into better hash functions (i.e., better image retrieval); its runtime and parallel speedup; and its precision and recall compared to representative state-of-the-art algorithms.

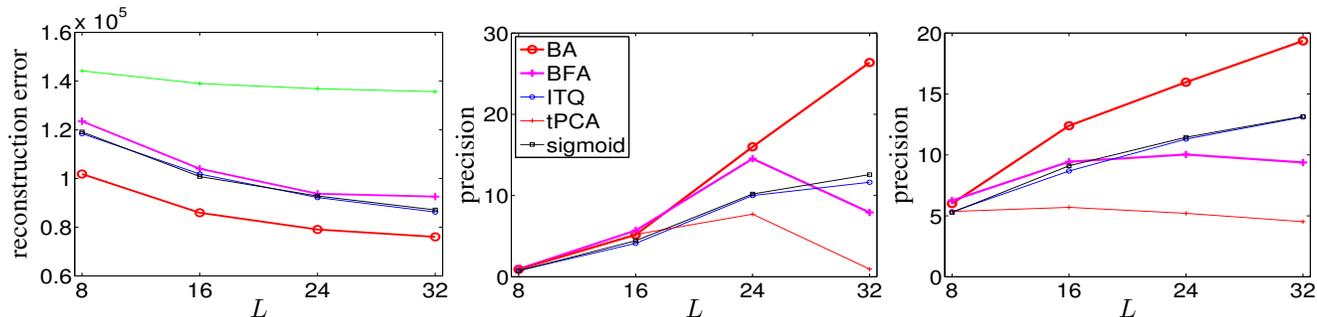


Figure 2. Wrapper vs filter optimization: BA objective function (*left*) and precision using neighbors within Hamming distance  $r = 2$  (*middle*) and using  $k = 50$  nearest neighbors (*right*).

### How much does respecting the binary constraints help?

We focus purely on the BA objective function (reconstruction error) and study the gain obtained by the MAC optimization, which respects the binary constraints, over the suboptimal, “filter” approach of relaxing the constraints (i.e., PCA) and then binarizing the result by thresholding at 0 (tPCA) or by optimal rotation (ITQ). We also try relaxing the step function to a sigmoid during training (with back-propagation using minibatches of 500 points) as in [26]. To compute the reconstruction error for tPCA and ITQ we find the optimal mapping  $\mathbf{f}$  given their binary codes. We use the NUS-WIDE-LITE subset of the NUS-WIDE dataset, containing  $N = 27\,807$  NUS-WIDE images for training and 27 807 images for test. We initialize BA from AGH [20] and BFA from tPCA and use alternating optimization in the  $\mathbf{Z}$  steps. We search for  $K = 50$  true neighbors and report results over a range of  $L = 8$  to 32 bits in fig. 2. We can see that BA dominates all other methods in reconstruction error, as expected, and also in precision, as one might expect. tPCA is consistently the worst method by a significant margin, while ITQ (very similar to the sigmoid) and BFA are intermediate. Hence, the more we respect the binary constraints during the optimization, the better the hash function. Further experiments below consistently show that the BA precision significantly increases over the (AGH) initialization and is leading or competitive over other methods.

### $\mathbf{Z}$ -step: alternating optimization and initialization

We study the MAC optimization if doing an inexact  $\mathbf{Z}$  step by using alternating optimization over groups of  $g$  bits. Specifically, we study the effect on the number of iterations and runtime of the group size  $g$  and of the initialization (warm-start vs relaxed QP). Fig. 3 shows the results in the CIFAR dataset using  $L = 16$  bits (so using  $g = 16$  gives an exact optimization), without using a validation-based stopping criterion (so we do optimize the training objective).

Surprisingly, the warm-start initialization leads to worse BA objective function values than the binarized relaxed one. Fig. 3(left) shows the dashed lines (warm-start for different  $g$ ) are all above the solid lines (relaxed for different  $g$ ). The reason is that, early during the optimization, the codes  $\mathbf{Z}$

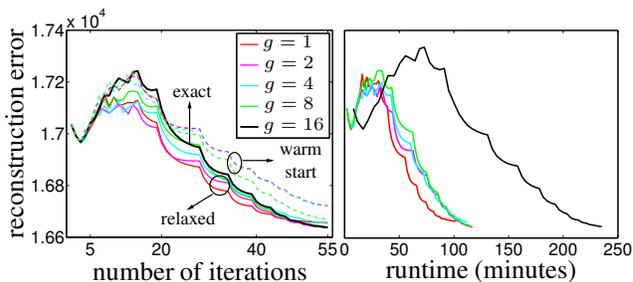


Figure 3. Iterations (*left*) and runtime (*right*) of the MAC optimization of the BA objective function (at each iteration, we run one  $\mathbf{Z}$  and  $(\mathbf{f}, \mathbf{h})$  step). In the  $\mathbf{Z}$  step, we use alternating optimization in  $g$ -bit groups ( $g = 16$  means exact optimization), and a warm-start vs relaxed initialization of  $\mathbf{Z}$ .

undergo drastic changes from one iteration to the next, so the warm-start initialization is farther from a good optimum than the relaxed one. Late in the optimization, when the codes change slowly, the warm-start does perform well. The relaxed initialization resulting optima are almost the same as using the exact binary optimization.

Also surprisingly, different group sizes  $g$  eventually converge to almost the same result as using the exact binary optimization if using the relaxed initialization. (If using warm-start, the larger  $g$  the better the result, as one would expect.) Likewise, in fig. 2, if using alternating optimization in the  $\mathbf{Z}$  step rather than enumeration, the curves for BA and BFA (not shown) barely vary. But, of course, the runtime per iteration grows exponentially on  $g$  (middle panel).

Hence, it appears that using faster, inexact  $\mathbf{Z}$  steps does not impair the model learnt, and we settle on  $g = 1$  with relaxed initialization as default for all our remaining experiments (unless we use  $L < 16$  bits, in which case we simply use enumeration).

### Parallel processing

Fig. 4 shows the BA training time speedup achieved with parallel processing, in CIFAR with  $L = 16$  bits. We use the Matlab Parallel Processing Toolbox with up to 12 processors and simply replace “for” with “parfor” loops so each iteration (over points in the  $\mathbf{Z}$  step, over bits in the  $\mathbf{h}$  step) is run in a different processor. We

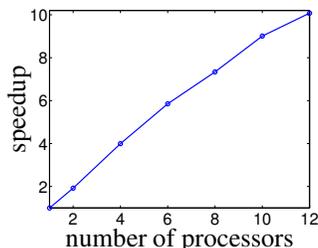


Figure 4. Parallel processing.

**Schedule of  $\mu$  and initial  $\mathbf{Z}$**  Since the objective is non-convex, our result does depend on the initial codes  $\mathbf{Z}$  (in the first iteration of the MAC algorithm), but we are guaranteed to improve or leave unchanged the precision (in the validation set) of the codes produced by any algorithm. We have observed that initializing from AGH [20] tends to produce best results overall, so we use this in all the experiments. Using ITQ [10] produces also good results (occasionally better but generally somewhat worse than AGH), and is a simpler and faster option if so desired. We initialize BFA from tPCA, since this seems to work best.

In order to be able to use a fixed schedule, we make the data zero-mean and rescale it so the largest feature range is 1. This does not alter the ordering of the Euclidean distances and normalizes the scale; it is also recommended by LIBLINEAR to train the SVM faster [9]. We start with  $\mu_1 = 10^{-5}$  and double it after each iteration (one  $\mathbf{Z}$  and  $(\mathbf{f}, \mathbf{h})$  step). As noted in section 3, the algorithm will skip  $\mu$  values that do not improve the precision in the validation set, and will stop at a finite  $\mu$  value (past which no further changes occur).

It is of course possible to tweak all these settings ( $\mu$  schedule and initial  $\mathbf{Z}$ ) and obtain better results, but these defaults seem robust.

### Comparison with other algorithms in image retrieval

We compare BA and BFA with the following algorithms: thresholded PCA (tPCA), Iterative Quantization (ITQ) [10], Spectral Hashing (SH) [28], Kernelized Locality-Sensitive Hashing (KLSH) [17], AnchorGraph Hashing (AGH) [20], and Spherical Hashing (SPH) [11]. Note several of these learn nonlinear hash functions and use more sophisticated error functions (that better approximate the nearest neighbor ideal), while our BA uses a linear hash function and simply minimizes the reconstruction error. All experiments use the output of AGH and tPCA to initialize BA and BFA, respectively.

It is known that the retrieval performance of a given algorithm depends strongly on the size of the neighbor set used, so we report experiments with small and large number of points in the ground truth set. For NUS-WIDE, we

observe a nearly perfect scaling for this particular problem. As a rough indication of runtimes for BA, training the 50 000 CIFAR images and 161 789 NUS-WIDE images using  $L = 32$  bits with alternating

optimization in the  $\mathbf{Z}$  step takes 20' and 50', respectively (in a 4-core laptop).

considered as ground truth  $K = 100$  and  $K = 1\,500$  neighbors of the query point, and as set of retrieved neighbors, we retrieve either  $k$  nearest neighbors ( $k = 100, 1500$ ) or neighbors within Hamming distance  $r$  ( $r = 1, 2$ ). Figs. 5 and 6 show the results. For ANNSIFT-1M, we considered ground truth  $K = 10\,000$  neighbors and set of retrieved neighbors for  $k = 10\,000$  or  $r = 1$  to 3. Fig. 7 shows the results. All curves are the average over the test set.

Although the relative performance of the different methods varies depending on the reported set size, some trends are clear. Generally (though not always) BA beats all other methods, sometime by a significant margin. ITQ and SPH become close (sometimes comparable) to BA in CIFAR (not shown) and NUS-WIDE dataset, respectively. BFA is also quite competitive, but consistently worse than BA. The only situation when the precision of BA and BFA appears to decrease is when  $L$  is large and  $r$  is small. The reason is that many test images have no neighbors at a Hamming distance of  $r$  or less and we report zero precision for them. This suggests the hash function finds a way to avoid collisions as more bits are available. In practice, one would simply increase  $r$  to retrieve sufficient neighbors.

Finally, we repeated our experiments using centered and normalized data, i.e., cosine similarity instead of Euclidean distance to determine ground-truth neighbors. The only important change is that all methods see a small, consistent increase in precision compared to the Euclidean distance precision. Although the increase varies for each method, the method ranking remains nearly identical, with BA being competitive as in Euclidean distance. Fig. 7 shows this for ANNSIFT-1M. This is consistent with the fact that the BA objective does not directly measure neighbor preservation (in Euclidean or cosine distance), but instead preserves manifold structure in some sense (see section 5).

## 5. Discussion

One contribution of this paper is to reveal a connection between ITQ [10] (a popular, effective hashing algorithm) and binary autoencoders. ITQ can be seen as a fast, approximate optimization of the BA objective function, using a “filter” approach (relax the problem to obtain continuous codes, iteratively quantize the codes, then fit the hash function). Our BA algorithm is a corrected version of ITQ.

Admittedly, there are objective functions that are more suited for information retrieval than the autoencoder, by explicitly encouraging distances in the original and Hamming space to match in order to preserve nearest neighbors [28, 3, 20, 17, 19, 18]. However, autoencoders do result in good hash functions, as evidenced by the good performance of ITQ and our method (or of semantic hashing [26], using neural nets). The reason is that, with continuous codes, autoencoders can capture the data manifold in a smooth way and indirectly preserve distances, encouraging (dis)similar

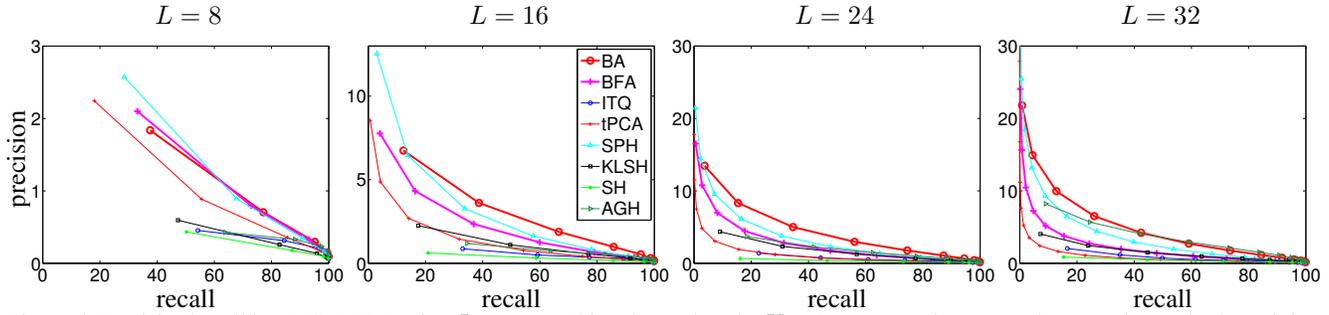


Figure 5. Precision/recall in NUS-WIDE using  $L = 8$  to 32 bits. Ground truth:  $K = 100$  nearest images to the query image in the training set. Retrieved neighbors: training images at Hamming distance  $\leq r$  of the query in binary space.

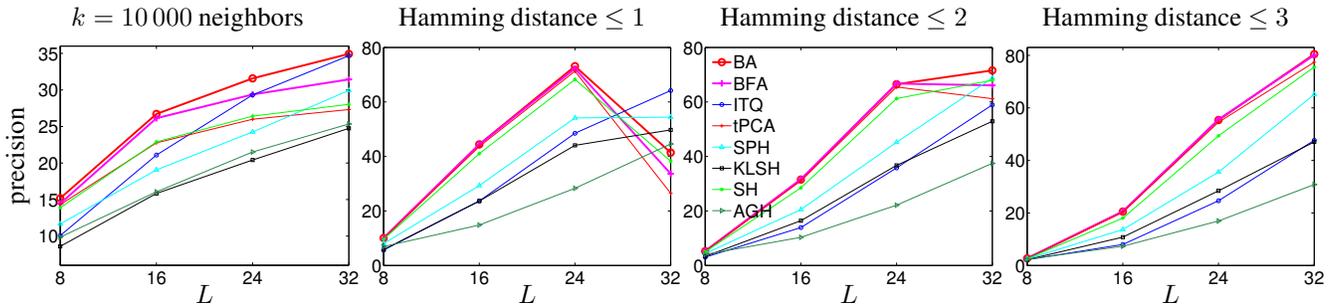


Figure 7. Precision in ANNSIFT-1M using  $L = 8$  to 32 bits. Like fig. 6 but with ground truth  $K = 10\,000$  and retrieved neighbors  $k = 10\,000$  or Hamming distance  $\leq r = 1$  to 3. In this figure we used the cosine similarity instead of Euclidean distance as ground truth.

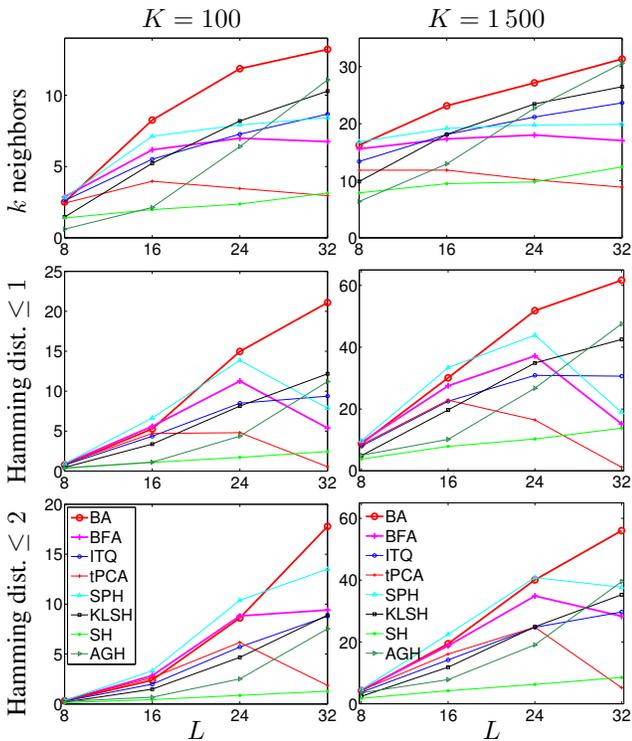


Figure 6. Precision in NUS-WIDE using  $L = 8$  to 32 bits. Ground truth:  $K = 100$  (left) and 1 500 (right column) nearest images to the query image in the training set. Retrieved neighbors:  $k = 100$  and 1 500 nearest images to, or images at Hamming distance  $\leq r = 1$  and 2 of the query, searching the training set binary codes.

images to have (dis)similar codes—even if this is worsened to some extent because of the quantization introduced with discrete codes. Autoencoders are also faster and easier to optimize and scale up better to large datasets.

## 6. Conclusion and future work

Up to now, many hashing approaches have essentially ignored the binary nature of the problem and have approximated it through relaxation and truncation, possibly disregarding the hash function when learning the binary codes. The inspiration for this work was to capitalize on the decoupling introduced by the method of auxiliary coordinates to be able to break the combinatorial complexity of optimizing with binary constraints, and to introduce parallelism into the problem. Armed with this algorithm, we have shown that respecting the binary nature of the problem during the optimization is possible in an efficient way and that it leads to better hash functions, competitive with the state-of-the-art. This was particularly encouraging given that the autoencoder objective is not the best for retrieval, and that we focused on linear hash functions.

The algorithm has an intuitive form (alternating classification, regression and binarization steps) that can reuse existing, well-developed code. The extension to nonlinear hash and reconstruction mappings is straightforward and it will be interesting to see how much these can improve over the linear case. This paper is a step towards constructing better hash functions using the MAC framework. We believe it may apply more widely to other objective functions.

## Acknowledgments

Work supported by NSF award IIS-1423515. We thank Ming-Hsuan Yang, Yi-Hsuan Tsai and Mehdi Alizadeh (UC Merced) for useful discussions.

## References

- [1] A. Andoni and P. Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Comm. ACM*, 51(1):117–122, Jan. 2008. [2](#)
- [2] A. Beck and M. Teboulle. Global optimality conditions for quadratic optimization problems with binary constraints. *SIAM Journal on Optimization*, 11(1):179–188, 2000. [5](#)
- [3] M. Á. Carreira-Perpiñán. The elastic embedding algorithm for dimensionality reduction. In J. Fürnkranz and T. Joachims, editors, *Proc. of the 27th Int. Conf. Machine Learning (ICML 2010)*, pages 167–174, Haifa, Israel, June 21–25 2010. [7](#)
- [4] M. Á. Carreira-Perpiñán. An ADMM algorithm for solving a proximal bound-constrained quadratic program. arXiv:1412.8493 [math.OA], Dec. 29 2014. [5](#)
- [5] M. Á. Carreira-Perpiñán and W. Wang. Distributed optimization of deeply nested systems. arXiv:1212.5921 [cs.LG], Dec. 24 2012. [1](#), [3](#)
- [6] M. Á. Carreira-Perpiñán and W. Wang. Distributed optimization of deeply nested systems. In S. Kaski and J. Corander, editors, *Proc. of the 17th Int. Conf. Artificial Intelligence and Statistics (AISTATS 2014)*, pages 10–19, Reykjavik, Iceland, Apr. 22–25 2014. [1](#), [3](#)
- [7] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y.-T. Zheng. NUS-WIDE: A real-world web image database from National University of Singapore. In *Proc. ACM Conf. Image and Video Retrieval (CIVR'09)*, Santorini, Greece, July 8–10 2009. [5](#)
- [8] P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. In H. H. Bauschke, R. S. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz, editors, *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, Springer Series in Optimization and Its Applications, pages 185–212. Springer-Verlag, 2011. [3](#), [4](#)
- [9] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. *J. Machine Learning Research*, 9:1871–1874, Aug. 2008. [4](#), [7](#)
- [10] Y. Gong, S. Lazebnik, A. Gordo, and F. Perronnin. Iterative quantization: A Procrustean approach to learning binary codes for large-scale image retrieval. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 35(12):2916–2929, Dec. 2013. [2](#), [7](#)
- [11] J.-P. Heo, Y. Lee, J. He, S.-F. Chang, and S.-E. Yoon. Spherical hashing. In *Proc. of the 2012 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'12)*, pages 2957–2964, Providence, RI, June 16–21 2012. [7](#)
- [12] H. Jégou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 33(1):117–128, Jan. 2011. [5](#)
- [13] V. Jeyakumar, A. M. Rubinov, and Z. Y. Wu. Non-convex quadratic minimization problems with quadratic constraints: Global optimality conditions. *Math. Prog.*, 110(3):521–541, Sept. 2007. [5](#)
- [14] R. Kohavi and G. H. John. The wrapper approach. In H. Liu and H. Motoda, editors, *Feature Extraction, Construction and Selection. A Data Mining Perspective*. Springer-Verlag, 1998. [1](#)
- [15] A. Krizhevsky. Learning multiple layers of features from tiny images. Master's thesis, Dept. of Computer Science, University of Toronto, Apr. 8 2009. [5](#)
- [16] B. Kulis and T. Darrell. Learning to hash with binary reconstructive embeddings. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 22, pages 1042–1050. MIT Press, Cambridge, MA, 2009. [2](#)
- [17] B. Kulis and K. Grauman. Kernelized locality-sensitive hashing. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 34(6):1092–1104, June 2012. [2](#), [7](#)
- [18] G. Lin, C. Shen, Q. Shi, A. van den Hengel, and D. Suter. Fast supervised hashing with decision trees for high-dimensional data. In *Proc. of the 2014 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'14)*, pages 1971–1978, Columbus, OH, June 23–28 2014. [2](#), [7](#)
- [19] G. Lin, C. Shen, D. Suter, and A. van den Hengel. A general two-step approach to learning-based hashing. In *Proc. 14th Int. Conf. Computer Vision (ICCV'13)*, pages 2552–2559, Sydney, Australia, Dec. 1–8 2013. [2](#), [7](#)
- [20] W. Liu, J. Wang, S. Kumar, and S.-F. Chang. Hashing with graphs. In L. Getoor and T. Scheffer, editors, *Proc. of the 28th Int. Conf. Machine Learning (ICML 2011)*, pages 1–8, Bellevue, WA, June 28 – July 2 2011. [2](#), [6](#), [7](#)
- [21] B. Neyshabur, N. Srebro, R. Salakhutdinov, Y. Makarychev, and P. Yadollahpour. The power of asymmetry in binary hashing. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 26, pages 2823–2831. MIT Press, Cambridge, MA, 2013. [2](#)
- [22] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer-Verlag, New York, second edition, 2006. [3](#)
- [23] M. Norouzi and D. Fleet. Minimal loss hashing for compact binary codes. In L. Getoor and T. Scheffer, editors, *Proc. of the 28th Int. Conf. Machine Learning (ICML 2011)*, Bellevue, WA, June 28 – July 2 2011. [2](#)
- [24] A. Oliva and A. Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *Int. J. Computer Vision*, 42(3):145–175, May 2001. [5](#)
- [25] R. T. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM J. Control and Optim.*, 14(5):877–898, 1976. [3](#), [4](#)
- [26] R. Salakhutdinov and G. Hinton. Semantic hashing. *Int. J. Approximate Reasoning*, 50(7):969–978, July 2009. [2](#), [6](#), [7](#)
- [27] A. Torralba, R. Fergus, and Y. Weiss. Small codes and large image databases for recognition. In *Proc. of the 2008 IEEE*

*Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'08)*, Anchorage, AK, June 23–28 2008. [2](#)

- [28] Y. Weiss, A. Torralba, and R. Fergus. Spectral hashing. In D. Koller, Y. Bengio, D. Schuurmans, L. Bottou, and A. Culotta, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 21, pages 1753–1760. MIT Press, Cambridge, MA, 2009. [2](#), [7](#)
- [29] P. Whittle. On principal components and least square methods of factor analysis. *Skand. Aktur. Tidskr.*, 36:223–239, 1952. [2](#)
- [30] S. X. Yu and J. Shi. Multiclass spectral clustering. In *Proc. 9th Int. Conf. Computer Vision (ICCV'03)*, pages 313–319, Nice, France, Oct. 14–17 2003. [2](#)
- [31] D. Zhang, J. Wang, D. Cai, and J. Lu. Self-taught hashing for fast similarity search. In *Proc. of the 33rd ACM Conf. Research and Development in Information Retrieval (SIGIR 2010)*, pages 18–25, Geneva, Switzerland, July 19–23 2010. [2](#)