

Sparse Projections for High-Dimensional Binary Codes

Yan Xia¹, Kaiming He², Pushmeet Kohli², Jian Sun²

¹University of Science and Technology of China. ²Microsoft Research.

Recent work on representation learning using deep neural networks has shown that features of thousands of dimensions or even more are useful for various recognition tasks. On the other hand, it has been noticed that for input signals of dimensionality d , the code-length b of the binary code required to achieve reasonable accuracy (compared with no encoding) is usually $O(d)$. Based on these facts, we consider the problem of learning long binary codes for high-dimensional arbitrary input signals.

We observe two key challenges arise while learning and using long binary codes: (1) lack of an effective regularizer for learned high-dimensional mapping and (2) high computational cost for computing long codes. In this paper, we overcome both these problems by introducing a sparsity encouraging regularizer that reduces the effective number of parameters involved in the learned projection operator. This regularizer can not only reduce overfitting but, due to the sparse nature of the projection matrix, also leads to a dramatic reduction in the computational cost.

To be specific, we use $X \in \mathbb{R}^{d \times n}$ to denote a matrix whose each column is a d -dimensional datum. Our target is to learn a projection matrix $R \in \mathbb{R}^{b \times d}$ for producing b -bit binary codes by $B = \text{sign}(RX) \in \{-1, 1\}^{b \times n}$. A widely considered objective function involves minimizing the distortion, as adopted in ITQ [2] and its variants [3, 6]. In ITQ, the objective function is only regularized by an orthogonality constraint, which is a weak regularizer. To further regularize the projection of high-dimensional data, we introduce a sparsity constraint. Our objective function is:

$$\begin{aligned} \min_{R, B} \|RX - B\|_F^2 \\ \text{s.t. } R^T R = I, \text{ and } |R|_0 \leq m. \end{aligned} \quad (1)$$

Here $\|\cdot\|_F$ denotes Frobenius norm, and $|\cdot|_0$ denotes the number of non-zero elements in matrix. Parameter m directly controls the sparsity of R .

However, it is challenging to optimize a matrix with both an orthogonality constraint and a sparsity constraint. To find a feasible solution, we adopt the variable-splitting and penalty techniques in optimization [1]. We introduce an auxiliary variable \bar{R} and relax the problem as:

$$\begin{aligned} \min_{\bar{R}, B} \|\bar{R}X - B\|_F^2 + \beta \|\bar{R}X - RX\|_F^2 \\ \text{s.t. } \bar{R}^T \bar{R} = I, \text{ and } |\bar{R}|_0 \leq m, \end{aligned} \quad (2)$$

where β is a penalty weight and we fix it to 1 in experiments. We solve (2) in an alternating manner: updating one variable with others fixed.

Step 1. Fix B and \bar{R} , update R . This sub-problem is:

$$\begin{aligned} \min_R \|RX - Z\|_F^2 \\ \text{s.t. } |R|_0 \leq m, \end{aligned} \quad (3)$$

where $Z = \bar{R}X$ is fixed. This sub-problem looks like a sparse coding problem under the m -sparsity constraint. But in commonly studied sparse coding problems, the variable R is a vector and RX is a vector-matrix multiplication. Here we develop a solution to the matrix variable R .

Instead of directly optimizing (3), we solve the following problem:

$$\min_{R, S} \phi(R, S) = \min_{R, S} \|RX - Z\|_F^2 + \|X\|_F^2 \|R - S\|_F^2 - \|RX - SX\|_F^2, \quad (4)$$

where ϕ is a surrogate objective function [5] and S is a matrix of the same size as R . The solution to (4) can be proven to give the solution to (3). We optimize (4) by an alternating algorithm as in [5]: fixing S and solving for R , and vice versa. In this way, the iterative solution to (3) is given as

$$R_{t+1} = \text{thr}_m \left(R_t + \frac{1}{\|X\|_F^2} (\bar{R} - R_t) X X^T \right), \quad (5)$$

This is an extended abstract. The full paper is available at the [Computer Vision Foundation webpage](#).

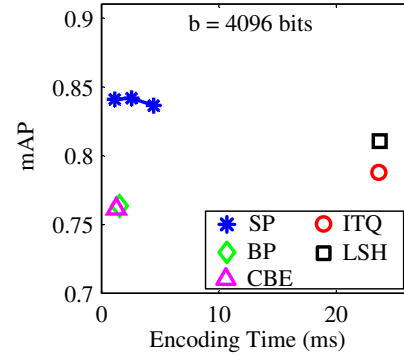


Figure 1: ANN search accuracy vs. encoding time on a dataset of one million 4096-d DNN features. Our sparse projection (SP) method, with 5%, 10%, and 15% non-zero elements, is more accurate and about 10× faster than ITQ [2] and LSH, and is more accurate than BP [3] and CBE [6].

where thr_m is an operator that only keeps the largest m entries. A natural initialization to (5) is to let $R_0 = \bar{R}$. Then the first iteration of (5) gives:

$$R = \text{thr}_m(\bar{R}). \quad (6)$$

In experiments, we find that the one-step solution gives comparable accuracy, but is much faster for training.

Step 2. Fix B and R , update \bar{R} . This sub-problem is an orthogonal procrustes problem. The solution is $\bar{R} = VU^T$, where U and V are from the SVD of XY^T . It is worth noting that according to modern studies [4], the above solution is valid for all $b \geq d$. Thus our method and ITQ are both naturally applicable for generating *higher*-dimensional binary codes.

Step 3. Fix R and \bar{R} , update B . This sub-problem is: $\min_B \|\bar{R}X - B\|_F^2 = \max_B \sum_{i,j} (\bar{R}X)_{ij} B_{ij}$. Since $B_{ij} \in \{-1, 1\}$, the solution is $B = \text{sign}(\bar{R}X)$.

We iteratively run in Step 1-3 (start from Step3) for 50 iterations, and finally a sparse matrix R is produced for projecting data.

We run comprehensive experiments of ANN search, image retrieval and image classification. Results show that our method not only leads to better accuracy than competitive dense projection methods (ITQ [2] and LSH) with the same code lengths, but is also over one order of magnitude faster. Furthermore, our method is more accurate than two recent methods (BP [3] and CBE [6]) that are designed for fast high-dimensional binary encoding. As an example, Figure 1 shows a comparison on one million 4096-d DNN features.

- [1] Richard Courant. Variational methods for the solution of problems of equilibrium and vibrations. *Bulletin of the American Mathematical Society*, 49(1):1–23, 1943.
- [2] Yunchao Gong and Svetlana Lazebnik. Iterative quantization: A procrustean approach to learning binary codes. In *CVPR*, 2011.
- [3] Yunchao Gong, Sanjiv Kumar, Henry A Rowley, and Svetlana Lazebnik. Learning binary codes for high-dimensional data using bilinear projections. In *CVPR*, 2013.
- [4] John C Gower and Garnt B Dijkstra. *Procrustes problems*, volume 3. Oxford University Press Oxford, 2004.
- [5] Kenneth Lange, David R Hunter, and Ilsoo Yang. Optimization transfer using surrogate objective functions. *Journal of computational and graphical statistics*, 2000.
- [6] Felix X Yu, Sanjiv Kumar, Yunchao Gong, and Shih-Fu Chang. Circulant binary embedding. In *ICML*, 2014.