

Semi-Supervised Hashing for Scalable Image Retrieval

¹Jun Wang, ²Sanjiv Kumar and ¹Shih-Fu Chang

¹Department of Electrical Engineering, Columbia University, New York, USA

²Google Research, New York, USA

Introduction and Motivation

- The emerging need of fast and accurate indexing for large scale content based image retrieval (CBIR)
 - Exhaustive search is infeasible due to computational and storage cost
 - Instead of exact nearest neighbor, approximate nearest neighbor (ANN) search is more practical for large databases
 - Compared with tree based ANN methods, hashing approach is more scalable and efficient
- Main issues:
 - Existing hashing methods mostly rely on random or principal projections, which are not very compact or accurate
 - Simple metrics are usually not enough to express semantic similarity
- Goal:** Given a large unlabeled set with a few pair-wise labeled points, learn data-dependent compact hash codes

Related Work

- Locality Sensitive Hashing (LSH) -- Indyk *et al.*, 98

$$h(x) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

\mathbf{w} - sampled from p distribution b - random shift

- Spectral Hashing (SH) -- Weiss *et al.* 08

$$h(\mathbf{x}) = \text{sgn}\left[\sin\left(\frac{\pi}{2} + \frac{k\pi}{b-a} \mathbf{w}^\top \mathbf{x}\right)\right]$$

k - spatial frequency $b-a$ - data range

\mathbf{w} - principal projections of data

Sinusoidal binarization over PCA projections

- Restricted Boltzmann Machines (RBMs) – Hinton *et al.* 06

➤ Learns deep belief networks to obtain compact binary codes

➤ Two training steps, unsupervised pre-training and supervised fine tuning

➤ Estimating a large number of weights requires lots of labeled data and long training time

Method	Projection Dependency	Learning Paradigm
LSH	data-independent	unsupervised
SH	data-dependent	unsupervised
RBM	–	unsupervised/supervised
SSH	data-dependent	semi-supervised

The conceptual comparison of the proposed SSH method with LSH, SH and RBMs.

Orthogonal Hash Codes

- Simplified objective function in matrix form

$$J(\mathbf{W}) = \frac{1}{2} \text{tr} [\mathbf{W}^\top \mathbf{X}_l \mathbf{S} \mathbf{X}_l^\top \mathbf{W}] + \frac{\eta}{2} \text{tr} [\mathbf{W}^\top \mathbf{X} \mathbf{X}^\top \mathbf{W}]$$

- By imposing orthogonality constraints, the objective function can be converted to a standard eigen problem

$$\arg \max_{\mathbf{W}} \frac{1}{2} \text{tr} \{ \mathbf{W}^\top \mathbf{M} \mathbf{W} \} \quad \text{"adjusted" covariance matrix} \\ \mathbf{M} = \mathbf{X}_l \mathbf{S} \mathbf{X}_l^\top + \eta \mathbf{X} \mathbf{X}^\top \\ \text{subject to: } \mathbf{W}^\top \mathbf{W} = \mathbf{I} \quad \text{supervised} \quad \text{unsupervised}$$

- Finding maximum variance projections over the data covariance matrix “adjusted” by incorporating the pair-wise labeled data

- Issues:** Most real datasets have the variance concentrated on the top several few projections. orthogonality constrains force to progressively pick those low-variance directions, substantially reducing the quality of hamming embedding

Non-Orthogonal Hash Codes

- Relax orthogonality constraints and combine them to the objective function as penalty term

$$J(\mathbf{W}) = \frac{1}{2} \text{tr} \{ \mathbf{W}^\top \mathbf{M} \mathbf{W} \} - \frac{\rho}{2} \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2$$

- Can be solved by a simple modification of the previous orthogonal solution

$$\mathbf{W} = [\mathbf{e}_1 \cdots \mathbf{e}_K] \quad \text{--- orthogonal solution}$$

principal projections of “adjusted” covariance matrix \mathbf{M}

- Choose proper coefficient ρ to guarantee \mathbf{Q} positive-definite

$$\mathbf{L} \mathbf{L}^\top = \mathbf{Q} = (\mathbf{I} + \frac{1}{\rho} \mathbf{M}) \quad \text{--- Cholesky factorization}$$

$$\mathbf{W}_{\text{nonorth}} = \mathbf{L} \cdot \mathbf{W} \quad \text{--- non-orthogonal solution}$$

Semi-Supervised Formulation

- Incorporate both metric similarity and semantic similarity for fast and accurate indexing

Given pair wise labeled data $(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{M}$: neighbor pair $(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{C}$ nonneighbor-pair

Assume data zero centered and use mean threshold: $h_k(\mathbf{x}) = \text{sgn}(\mathbf{w}_k^\top \mathbf{x})$

$$\text{Empirical fitness: } J(\mathbf{H}) = \sum_k \left(\sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{M}} h_k(\mathbf{x}_i) h_k(\mathbf{x}_j) - \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{C}} h_k(\mathbf{x}_i) h_k(\mathbf{x}_j) \right)$$

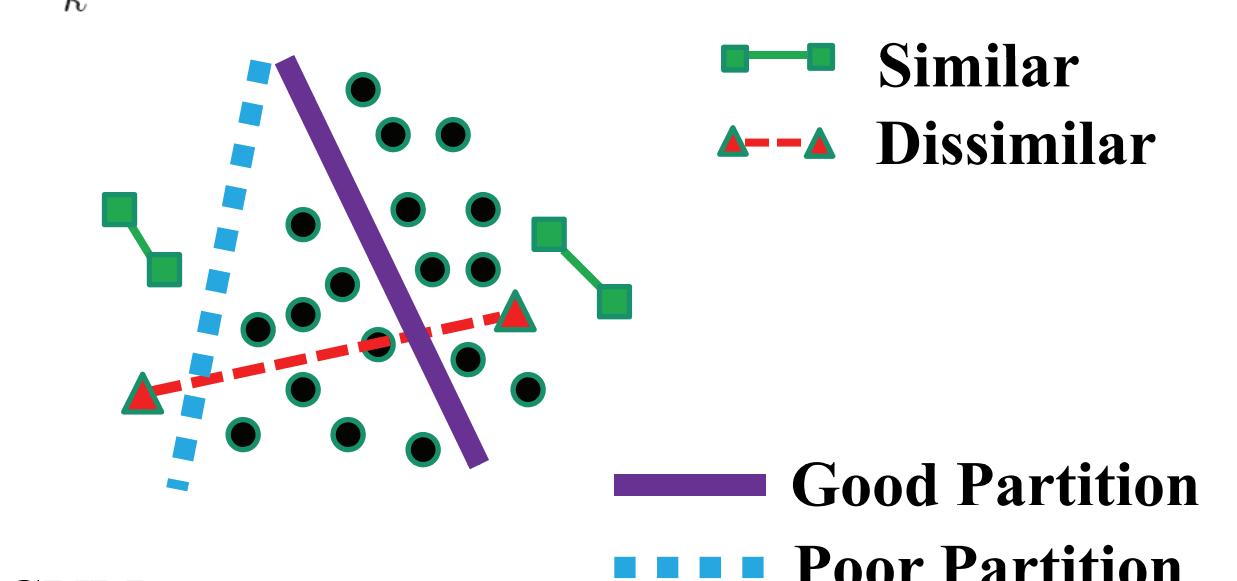
Regularization of maximum variance: $\sum_k E[\|h_k(\mathbf{x}) - \mu_k\|^2]$

Final objective function

$$J(\mathbf{W}) = \frac{1}{2} \text{tr} [\mathbf{W}^\top \mathbf{X}_l \mathbf{S} \mathbf{X}_l^\top \mathbf{W}] + \frac{\eta}{2} \sum_k E[\|w_k^\top \mathbf{x}\|^2]$$

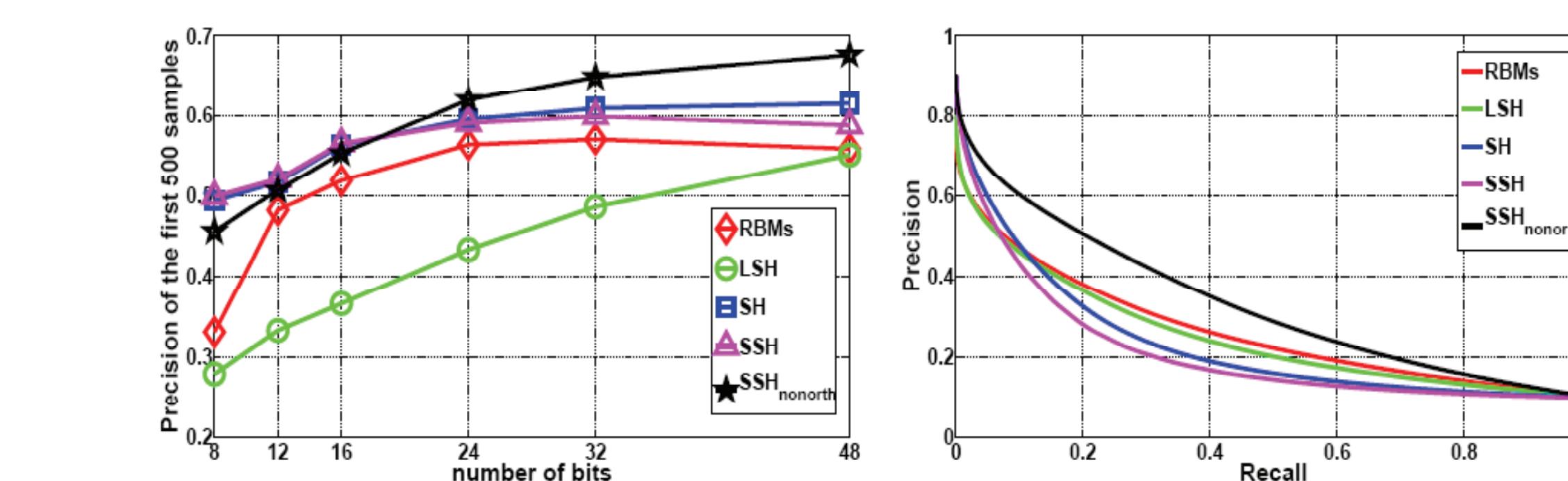
$$S_{ij} = \begin{cases} 1 & : (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{M} \\ -1 & : (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{C} \\ 0 & : \text{otherwise.} \end{cases}$$

AD The other semi-supervised hashing work in this year's CVPR:
Y. Mu, et al., Weakly-Supervised Hashing in Kernel Space

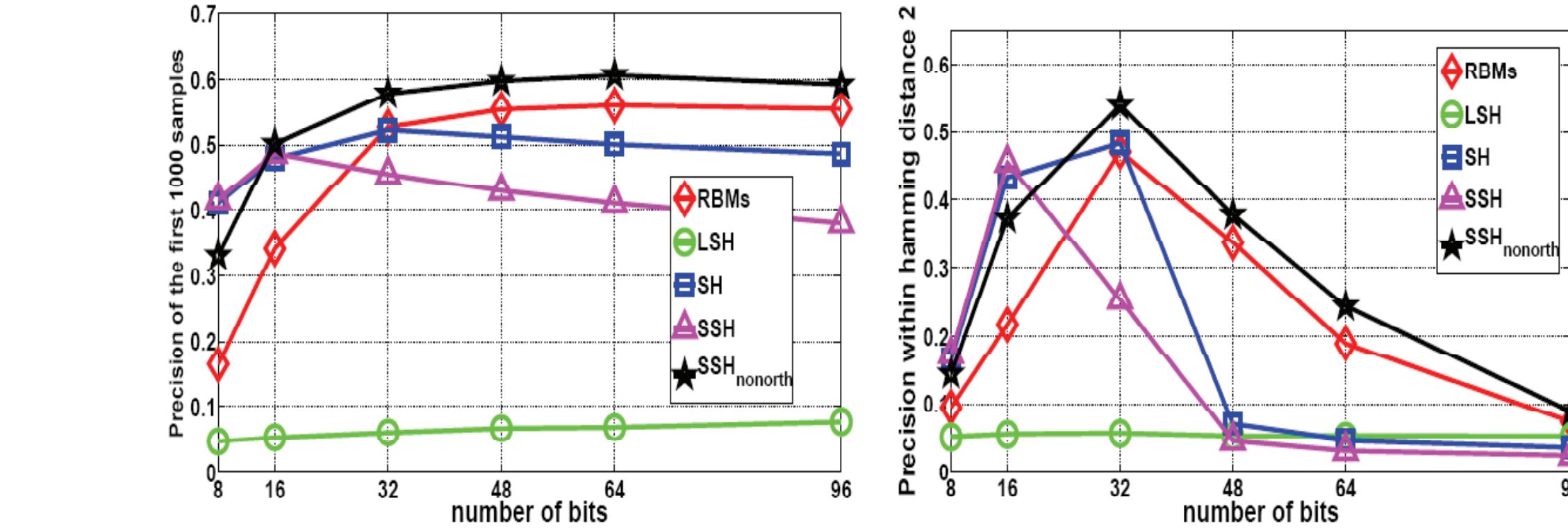


Experiments and Conclusions

- MNIST (70K, 2K training labels, 1K test)



- One million GIST (10K training labels, 2K test)



- Conclusions

- A semi-supervised formulation for hash function learning
- Orthogonal and non-orthogonal solutions by simple SVD
- Easy implementation, fast training, and superior performance

