

# On the Optimal Bit Complexity of Circulant Binary Embedding (supplementary material)

Saehoon Kim<sup>1,2</sup>, Jungtaek Kim<sup>1</sup>, and Seungjin Choi<sup>1</sup>

<sup>1</sup>CSE, POSTECH, Pohang, Korea    <sup>2</sup>Altrics, Seoul, Korea

{kshkawa,jtkim,seungjin}@postech.ac.kr

## 1 Outline

In this supplementary material, we present the proofs of Theorem 1, Corollary 1, and other lemmas in the submitted manuscript. Moreover, we introduce some experimental results for GIST1M and Flickr45K, which are not placed in the manuscript due to the space limit.

## 2 Proof of Theorem 1

**Theorem 1.** (*Theorem 1 in the manuscript*). Given  $\epsilon \in (0, 1)$  and any finite data set  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{S}^{d-1}$ , with probability at least  $1 - \exp(-c\epsilon^2 k)$ ,  $k = \mathcal{O}(\epsilon^{-2} \log n)$  implies that we have  $h : \mathcal{S}^{d-1} \rightarrow \{0, 1\}^k$  such that for all  $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}$

$$\left| d_H(h(\mathbf{x}_i), h(\mathbf{x}_j)) - \frac{\theta_{\mathbf{x}_i, \mathbf{x}_j}}{\pi} \right| \leq \epsilon,$$

where  $c > 0$  is a constant.

*Proof.* It can be easily proved by Hoeffding's inequality [1] and union bound. Hoeffding's inequality provides the upper-bound on the probability that the sum of independent random variables deviates from its expected value. When entries of  $\mathbf{G}$  are drawn independently from Gaussian with zero mean and unit variance, Bernoulli indicator variables  $\{\mathcal{I}_l(\mathbf{x}_i, \mathbf{x}_j)\}_{l=1}^k$  are marginally independent, given  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . Thus, the probability that the normalized Hamming distance (the sum of  $\mathcal{I}_l$ ) deviates from the angle distance (its mean), is given by

$$\mathbb{P} \left[ \left| d_H(h(\mathbf{x}_i), h(\mathbf{x}_j)) - \frac{\theta_{\mathbf{x}_i, \mathbf{x}_j}}{\pi} \right| \geq \epsilon \right] \leq 2 \exp \{-2k\epsilon^2\}. \quad (1)$$

Given a dataset  $\mathcal{D}$ , a  $\epsilon$ -distortion binary embedding involves the following probability:

$$\mathbb{P} \left[ \bigcap_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}} \left( \left| d_H(h(\mathbf{x}_i), h(\mathbf{x}_j)) - \frac{\theta_{\mathbf{x}_i, \mathbf{x}_j}}{\pi} \right| \leq \epsilon \right) \right]. \quad (2)$$

Strict positiveness of this probability guarantees  $\epsilon$ -distortion binary embedding. To this end, we use Boole's inequality (a.k.a. the union bound) which states that for any finite or countable set of events, the probability that at least one of the events happens is no greater than the sum of the probabilities of the individual events, i.e.,  $\mathbb{P}[\cup_i A_i] \leq \sum_i \mathbb{P}[A_i]$ , for a set of events  $\{A_i\}$ . Applying Boole's inequality

yields

$$\mathbb{P} \left[ \bigcap_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}} \left( \left| d_H(h(\mathbf{x}_i), h(\mathbf{x}_j)) - \frac{\theta_{\mathbf{x}_i, \mathbf{x}_j}}{\pi} \right| \leq \epsilon \right) \right] \quad (3)$$

$$= 1 - \mathbb{P} \left[ \bigcup_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}} \left( \left| d_H(h(\mathbf{x}_i), h(\mathbf{x}_j)) - \frac{\theta_{\mathbf{x}_i, \mathbf{x}_j}}{\pi} \right| \geq \epsilon \right) \right] \quad (4)$$

$$\geq 1 - \sum_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}} \mathbb{P} \left[ \left| d_H(h(\mathbf{x}_i), h(\mathbf{x}_j)) - \frac{\theta_{\mathbf{x}_i, \mathbf{x}_j}}{\pi} \right| \geq \epsilon \right] \quad (5)$$

$$\geq 1 - 2n^2 \exp \{-2k\epsilon^2\}. \quad (6)$$

The lower bound on the probability of event becomes interesting (i.e. greater than zero) when the condition  $k \geq \mathcal{O}(\epsilon^{-2} \log n)$  is given.  $\square$

### 3 Proof of Corollary 1

In this section, some details for the proof of Corollary 1 in the manuscript are introduced. Lemma 1 is the key for the proof, which needs Hanson-Wright inequality for sub-Gaussian random variables described as follows.

**Theorem 2.** (Hanson-Wright inequality for sub-Gaussian random variables [5]). Let  $\mathbf{X} = (X_1, \dots, X_d) \in \mathbb{R}^d$  be a random vector with independent sub-Gaussian random variables  $X_i$  which satisfy  $\mathbb{E}[X_i] = 0$  and  $\|X_i\|_{\psi_2} \leq K$ . Let  $\mathbf{A}$  be an  $d \times d$  matrix. Then, for every  $t \geq 0$ ,

$$\mathbb{P} \left( \left| \mathbf{X}^\top \mathbf{A} \mathbf{X} - \mathbb{E} [\mathbf{X}^\top \mathbf{A} \mathbf{X}] \right| > t \right) \leq 2 \exp \left[ -c \min \left( \frac{t^2}{K^4 \|\mathbf{A}\|_F^2}, \frac{t}{K^2 \|\mathbf{A}\|_2} \right) \right]. \quad (7)$$

**Lemma 1.** Letting two unit vectors be  $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{d-1}$ , suppose that  $\max\{\|\mathbf{x}\|_\infty, \|\mathbf{y}\|_\infty\} \leq \rho$  and  $\mathbf{x}^\top \mathbf{y} = 0$ . Given  $\mathbf{d}$  be i.i.d. Rademacher entries and  $\mathbf{D} = \text{diag}(\mathbf{d})$ , for all  $1 \leq i \neq j \leq d$ , the following is achieved.

$$\mathbb{P} (|\text{circ}(\mathbf{D}\mathbf{x})_{:,i}^\top \text{circ}(\mathbf{D}\mathbf{y})_{:,j}| > t) \leq 2 \exp(-ct^2 \rho^{-2}) \quad (8)$$

$$\mathbb{P} (|\text{circ}(\mathbf{D}\mathbf{x})_{:,i}^\top \text{circ}(\mathbf{D}\mathbf{x})_{:,j}| > t) \leq 2 \exp(-ct^2 \rho^{-2}), \quad (9)$$

where  $t \in (0, 1)$ ,  $c > 0$ , and  $\text{circ}(\cdot)$  is defined in Eq. 5 in the manuscript.

*Proof.* The proof is basically same with the proof of Lemma 4.5 in [4]. It introduces that  $\text{circ}(\mathbf{D}\mathbf{x})_{:,i}^\top \text{circ}(\mathbf{D}\mathbf{y})_{:,j}$  can be represented by  $\mathbf{d}^\top \mathbf{A} \mathbf{d}$  where  $(m+i, m+j)$ th entry of  $\mathbf{A}$  is  $x_{m+i} y_{m+j}$  for  $1 \leq m \leq d$  and remaining entries are 0. The matrix  $\mathbf{A}$  has the interesting properties:  $\|\mathbf{A}\|_F = \rho$ ,  $\|\mathbf{A}\|_2 = \rho^2$ , and  $\text{tr}(\mathbf{A}) = 0$ . Consequently, it is easy to show that

$$\mathbb{E} [\mathbf{d}^\top \mathbf{A} \mathbf{d}] = \mathbb{E} [\text{tr}(\mathbf{d} \mathbf{d}^\top \mathbf{A})] \quad (10)$$

$$= \text{tr}(\mathbb{E} [\mathbf{d} \mathbf{d}^\top] \mathbf{A}) \quad (11)$$

$$= \text{tr}(\mathbf{I} \mathbf{A}) = \text{tr}(\mathbf{A}) = 0. \quad (12)$$

Then, Theorem 2 directly says that

$$\mathbb{P} (|\text{circ}(\mathbf{D}\mathbf{x})_{:,i}^\top \text{circ}(\mathbf{D}\mathbf{y})_{:,j}| > t) \leq 2 \exp \left[ -c \min \left( \frac{t^2}{\rho^2}, \frac{t}{\rho^2} \right) \right] \quad (13)$$

$$= 2 \exp[-ct^2 \rho^{-2}], \quad (14)$$

where the last equality is derived by  $t \in (0, 1)$  and  $K = 1$  for Bernoulli random variables. The same argument for Eq. 9 is developed to conclude the proof.  $\square$

**Corollary 1.** (Corollary 1 in the manuscript). Suppose that  $\mathbf{x}_i, \mathbf{x}_j$  from  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{S}^{d-1}$  satisfies  $\max\{\|\mathbf{x}_i\|_\infty, \|\mathbf{x}_j\|_\infty\} \leq \rho$ . Letting  $\epsilon \in (0, 1)$ , with probability at least  $1 - 4\exp(-\epsilon^2 k)$ , we achieve that

$$\max\{\|\hat{\mathbf{p}}_{i,m}\|_2, \|\hat{\mathbf{p}}_{j,m}\|_2\} \leq c_a \epsilon k \rho,$$

where  $c_a > 0$  is a constant.  $\{\hat{\mathbf{p}}_{i,m}, \hat{\mathbf{p}}_{j,m}\}_{m=1}^k$  are defined in Definition 4 in the manuscript.

*Proof.* The proof is concluded by switching Lemma 4.5 in the proof of Lemma 5.1 in [4] with Lemma 1, where some of the steps are twisted.  $\square$

## 4 Useful Properties

In this section, we present some details to understand the proof of Corollary 2 and Theorem 2 in the manuscript. Lemma 2 is introduced to describe the proof of Corollary 2 in the manuscript. Lemma 3 is described for the proof of Theorem 2 (Eq. (19)) in the manuscript.

**Lemma 2.** Given  $\|\mathbf{x}\|_2 \leq \delta$ , the following holds

$$\mathbb{P}[\|\mathbf{g}^\top \mathbf{x}\| \geq t\delta] \leq 2\exp(-0.5t^2), \quad (15)$$

where  $t \geq 0$  and  $\mathbf{g}$  follows  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .

*Proof.* It is easy to observe that  $\mathbf{g}^\top \mathbf{x}$  follows  $\mathcal{N}(0, \|\mathbf{x}\|_2^2)$ . Then,

$$\mathbb{P}\left[\frac{\mathbf{g}^\top \mathbf{x}}{\|\mathbf{x}\|_2} \geq t \frac{\delta}{\|\mathbf{x}\|_2}\right] \leq \mathbb{P}\left[\frac{\mathbf{g}^\top \mathbf{x}}{\|\mathbf{x}\|_2} \geq t\right] \quad (16)$$

$$\leq \exp(-0.5t^2), \quad (17)$$

where  $\frac{\mathbf{g}^\top \mathbf{x}}{\|\mathbf{x}\|_2}$  follows  $\mathcal{N}(0, 1)$  and the standard Chernoff inequality is applied. Since the same argument is applicable for  $\mathbb{P}[\mathbf{g}^\top \mathbf{x} \geq t\delta]$ , it concludes the proof.  $\square$

**Lemma 3.** As in the manuscript, suppose that  $E_{(i,j),r}$  is defined as

$$E_{(i,j),r} \triangleq \left| \frac{1}{k} \sum_{l=1}^k \mathcal{I}_{(i,j)}^{r,l} - \text{ang}(\mathbf{x}_i, \mathbf{x}_j) \right| \leq \epsilon, \quad (18)$$

where  $\mathcal{I}_{(i,j)}^{r,l}$  is the union of the following events:

$$E_{margin}^1 \triangleq \mathcal{I}[\mathbf{g}^\top \hat{\mathbf{x}}_{i,m} > \epsilon \text{ and } \mathbf{g}^\top \hat{\mathbf{x}}_{j,m} < -\epsilon] \quad (19)$$

$$E_{margin}^2 \triangleq \mathcal{I}[\mathbf{g}^\top \hat{\mathbf{x}}_{i,m} < -\epsilon \text{ and } \mathbf{g}^\top \hat{\mathbf{x}}_{j,m} > \epsilon]. \quad (20)$$

Given Condition 1 in the manuscript, the following holds

$$\mathbb{P}[E_{(i,j),r}^c] \leq 2\exp(-2\epsilon^2 k), \quad (21)$$

where  $\epsilon \in (0, 1)$  and  $k > 0$  is the binary code length.

*Proof.* According Corollary 2 in the manuscript and the third condition (i.e.  $c_3 \rho k < \epsilon$ ) in Condition 1, with at least probability  $1 - 4\exp(-\epsilon^2 k)$ , we have that  $\max\{\|\hat{\mathbf{p}}_{i,m}\|_2, \|\hat{\mathbf{p}}_{j,m}\|_2\} < \epsilon$ . Then, it is easy to see that  $\mathbf{g}^\top \hat{\mathbf{x}}_{i,m}$  and  $\mathbf{g}^\top \hat{\mathbf{x}}_{j,m}$  follow a uni-variate Gaussian distribution with zero mean and the variance within  $(1 - \epsilon^2, 1]$ . It is trivial to show the following:

$$\mathbb{P}[\mathcal{I}_{(i,j)}^{r,l}] \leq \text{ang}(\hat{\mathbf{x}}_{i,m}, \hat{\mathbf{x}}_{j,m}) \leq \text{ang}(\mathbf{x}_i, \mathbf{x}_j) + c_a \epsilon k \rho \leq \text{ang}(\mathbf{x}_i, \mathbf{x}_j) + \epsilon^2, \quad (22)$$

where the second inequality is derived by Lemma 1 in the manuscript. As discussed in [6],  $\{\mathcal{I}_{(i,j)}^{r,l}\}_{l=1}^k$  are marginally independent. Hoeffding's inequality leads that

$$\mathbb{P}\left(\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}-\mathbb{E}\left[\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}\right]\geq\epsilon\right)\leq\exp(-2\epsilon^2k) \quad (23)$$

The lower bound on the left-hand side is derived as follows.

$$\mathbb{P}\left(\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}-\mathbb{E}\left[\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}\right]\geq\epsilon\right)=\mathbb{P}\left(\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}\geq\epsilon+\frac{1}{k}\sum_{i=1}^k\mathbb{P}\left[\mathcal{I}_{(i,j)}^{r,l}\right]\right) \quad (24)$$

$$\geq\mathbb{P}\left(\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}\geq\epsilon(1+\epsilon)+\text{ang}(\mathbf{x}_i,\mathbf{x}_j)\right), \quad (25)$$

where the last inequality is introduced by Eq. 22 and  $\epsilon(1+\epsilon)\approx\epsilon$  because  $\epsilon\ll 1$ . It concludes that

$$\mathbb{P}\left(\frac{1}{k}\sum_{i=1}^k\mathcal{I}_{(i,j)}^{r,l}-\text{ang}(\mathbf{x}_i,\mathbf{x}_j)\geq\epsilon\right)\leq\exp(-2\epsilon^2k). \quad (26)$$

A symmetric argument is easily developed to conclude the proof.  $\square$

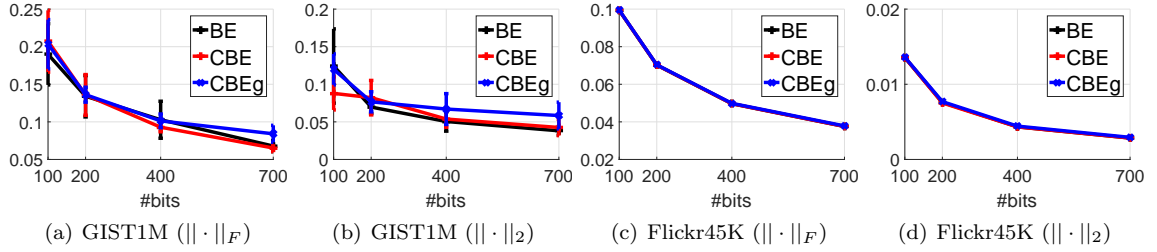


Figure 1: Plots for the relative error on approximating the angle between vectors measured by Frobenious and spectral norms for GIST1M and Flickr45K datasets, showing that CBE is almost identical to the standard binary embedding.

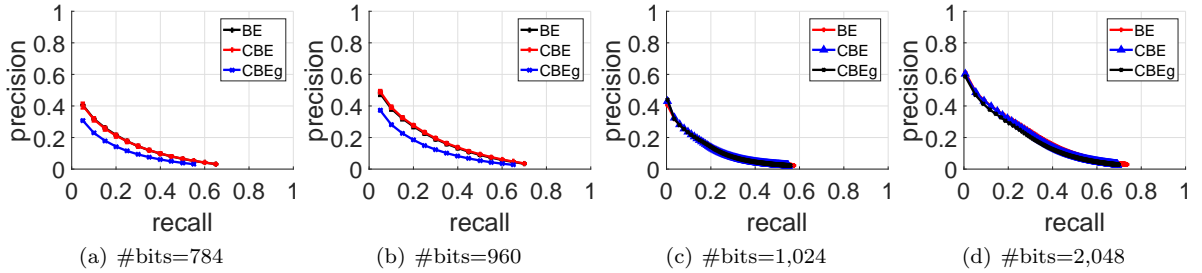


Figure 2: Precision-recall curves from Hamming ranking evaluation to compare CBE with BE and CBEg on GIST1M(the first two subfigures) and Flickr45K(the last two subfigures).

## 5 Experiments

In this section, we introduce experimental results for angle preservation and hamming ranking evaluation on GIST1M and Flickr45K datasets. Flickr45K consists of 45,000 images randomly chosen from 1 million

Flickr images [2]. VLAD [3] makes use of 500 cluster centers to represent images which are  $200 \times 128 = 25,600$ -dimensional vectors. To evaluate the angle preservation, we measured the relative error on approximating the angle with the normalized Hamming distance as in the manuscript. To perform Hamming ranking evaluation, we randomly selected 1,000 queries from test datasets and computed 100 nearest neighbors for ground-truths, where angular distance is used.

Clearly, Figure 1 and 2 suggest that CBE preserves the angular distance as well as BE, which has been observed in the manuscript. Unlike MNIST, CIFAR-10, and GIST1M, the performance of CBEg is indistinguishable from the ones of CBE and BE for Flickr45K, which means that the distortion induced by Gaussian random sequences is not serious in case of a large data dimension.

## References

- [1] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [2] M. J. Huiskes and M. S. Lew. The mir flickr retrieval evaluation. In *MIR '08: Proceedings of the 2008 ACM International Conference on Multimedia Information Retrieval*, New York, NY, USA, 2008. ACM.
- [3] H. Jégou, M. Douze, C. Schmid, and P. Pérez. Aggregating local descriptors into a compact image representation. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, San Francisco, CA, 2010.
- [4] S. Oymak. Near-optimal sample complexity bounds for circulant binary embedding, 2016. arXiv preprint arXiv:1603.03178v2.
- [5] M. Rudelson and R. Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electron. Commun. Probab*, 18, 2013.
- [6] F. X. Yu, A. Bhaskara, S. Kumar, Y. Gong, and S.-F. Chang. On binary embedding using circulant matrices, 2015. arXiv preprint arXiv:1511.06480v2.