

---

# Supplementary Material: A Single-Pass Algorithm for Efficiently Recovering Sparse Cluster Centers of High-dimensional Data

---

**Jinfeng Yi**

JINFENGY@US.IBM.COM

IBM Thomas J. Watson Research Center, Yorktown Heights, NY 10598, USA

**Lijun Zhang**

ZHANGLJ@LAMDA.NJU.EDU.CN

National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China

**Jun Wang**

WANGJUN@US.IBM.COM

IBM Thomas J. Watson Research Center, Yorktown Heights, NY 10598, USA

**Rong Jin**

RONGJIN@CSE.MSU.EDU

**Anil K. Jain**

JAIN@CSE.MSU.EDU

Department of Computer Science and Engineering, Michigan State University, East Lansing, MI 48824 USA

**Theorem 1.** Let  $\epsilon = 1/(6m)$  be a parameter to control the success probability. Assume

$$\epsilon \leq \frac{c}{2^m} \quad ; \quad (1)$$

$$\frac{c}{2^m} \leq \frac{c}{2^m} \quad ; \quad (2)$$

$$T \leq \max \left( \frac{18}{\epsilon} \ln \frac{2K}{\epsilon}; \frac{3c}{\epsilon}; \left( \frac{6c}{\epsilon} \right) (\ln n + \ln d) \right) \quad (3)$$

where  $c$ ,  $c$  and  $c$  are some universal constants. Then, with a probability at least  $1 - 6m$ , we have

$$\epsilon^m = \max_{\leq i \leq K} k \hat{c}_i^m \quad c_i k \max \left( \epsilon; \frac{c}{2^m} \right) :$$

**Corollary 1.** The convergence rate for  $\epsilon$ , the maximum difference between the optimal cluster centers and the estimated ones, is  $O(\sqrt{(s \log d)/n})$  before reaching the optimal difference  $\epsilon$ .

## 1. Proof of Corollary 1

According to the assumption of  $\epsilon$  in (2), we know that  $\frac{c}{2^m} \leq \frac{\sqrt{s}}{T}$ . Since the value of  $T$  is dominated by the last term in the right side of (3), we have  $T \leq \frac{s}{\epsilon} \frac{d}{T}$ , which implies

$$n \leq 2^m T \leq 2^m \frac{s \log d}{\epsilon} :$$

Combining with the conclusion  $m \leq \frac{1}{\sqrt{m}}$ , we have

$$m \leq \sqrt{\frac{s \log d}{n}};$$

**Lemma 1.** Let  $\delta^t$  be the maximum difference between the optimal cluster centers and the ones estimated from iteration  $t$ , and  $\beta \in (0; 1)$  be the failure probability. Assume

$$\delta^t \leq \frac{1}{2} \sqrt{5 \ln(3K)}, \quad (4)$$

$$jS^t j \geq \frac{18}{\beta} \ln \frac{2K}{\beta}; \quad (5)$$

$$\delta^t \leq c \exp\left(\frac{(1 - 2^{-t})}{8(1 + 2^{-t})}\right) \left( \frac{c}{jS^t j} + \sqrt{\ln jS^t j} \right) + \frac{c}{jS^t j} + c \frac{\sqrt{\ln jS^t j} + \beta \ln d}{\sqrt{jS^t j}}; \quad (6)$$

for some constants  $c$ ,  $c$  and  $c$ . Then with a probability  $1 - \beta$ , we have

$$\delta^t \leq 2^{\frac{\rho}{s} t};$$

## 2. Proof of Lemma 1

For the simplicity of analysis, we will drop the superscript  $t$  through this analysis.

### 2.1. Preliminaries

We denote by  $C_k$  the support of  $\mathbf{c}_k$  and  $\bar{C}_k = [d] \setminus C_k$ . For any vector  $\mathbf{z}$ ,  $\mathbf{z}(C)$  is defined as  $[\mathbf{z}(C)]_i = z_i$  if  $i \in C$  and zero, otherwise.

For any  $\mathbf{x}_i \in S$ , we use  $k_i$  to denote the index of the true cluster, and  $\hat{k}_i$  to denote index of the cluster assigned by the nearest neighbor search, i.e.,

$$\begin{aligned} \mathbf{x}_i &= \mathbf{c}_{k_i} + \mathbf{g}_i \text{ and } \mathbf{g}_i \sim N(0; I); \\ \hat{k}_i &= \arg \max_{j \in K} \hat{\mathbf{c}}_j^\top \mathbf{x}_i; \end{aligned}$$

Then, we can partition data points in  $S$  based on either the ground truth or the assigned cluster. Let  $S_k$  be the subset of data points in  $S$  that belong to the  $k$ -th cluster, i.e.,

$$S_k = \{\mathbf{x}_i \in S : \mathbf{x}_i = \mathbf{c}_k + \mathbf{g}_i \text{ and } \mathbf{g}_i \sim N(0; I)\} \quad (7)$$

Let  $\hat{S}_k$  be the subset of data points that are assigned to the  $k$ -th cluster based on the nearest neighbor search, i.e.,

$$\hat{S}_k = \{\mathbf{x}_i \in S : k = \arg \max_{j \in K} \hat{\mathbf{c}}_j^\top \mathbf{x}_i\} \quad (8)$$

### 2.2. The Main Analysis

Let  $L_k(\mathbf{c})$  be the objective function in Step 11 of Algorithm 1. We expand  $L_k(\mathbf{c})$  as

$$\begin{aligned} L_k(\mathbf{c}) &= kCk + kC \mathbf{c}_k k + \frac{1}{j\hat{S}_k j} \sum_{\mathbf{x}_i \in \hat{S}_k} k\mathbf{x}_i - \mathbf{c}_k k - \frac{2}{j\hat{S}_k j} \sum_{\mathbf{x}_i \in \hat{S}_k} (\mathbf{c} - \mathbf{c}_k)^\top (\mathbf{x}_i - \mathbf{c}_k) \\ &= kCk + kC \mathbf{c}_k k + \frac{1}{j\hat{S}_k j} \sum_{\mathbf{x}_i \in \hat{S}_k} k\mathbf{x}_i - \mathbf{c}_k k \\ &\quad - 2(\mathbf{c} - \mathbf{c}_k)^\top \underbrace{\frac{1}{j\hat{S}_k j} \sum_{\mathbf{x}_i \in \hat{S}_k \setminus S_k} (\mathbf{c}_{k_i} - \mathbf{c}_k)}_{A_k} - 2(\mathbf{c} - \mathbf{c}_k)^\top \underbrace{\frac{1}{j\hat{S}_k j} \sum_{\mathbf{x}_i \in \hat{S}_k} \mathbf{g}_i}_{B_k}; \end{aligned} \quad (9)$$

Let  $\mathbf{c}_k^*$  be the optimal solution that minimizes  $L_k(\mathbf{c})$ , and define  $\mathbf{f}_k = \mathbf{c}_k^* - \mathbf{c}_k$ . We have

$$\begin{aligned} & L_k(\mathbf{c}_k^*) - L_k(\mathbf{c}_k) \\ &= \|\mathbf{f}_k + \mathbf{c}_k\|_k + \|\mathbf{f}_k\|_k - 2\mathbf{f}_k^\top \mathbf{A}_k - 2\mathbf{f}_k^\top \mathbf{B}_k - \|\mathbf{c}_k\|_k \\ &= \|\mathbf{f}_k\|_k + \|\mathbf{c}_k\|_k + \|\mathbf{f}_k\|_k + \|\mathbf{f}_k\|_k + \|\mathbf{f}_k\|_k - 2\mathbf{f}_k^\top \mathbf{A}_k - 2\mathbf{f}_k^\top \mathbf{B}_k - \|\mathbf{c}_k\|_k \\ &= (\|\mathbf{f}_k\|_k + 2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty) \|\mathbf{f}_k\|_k + (\|\mathbf{f}_k\|_k + 2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty) \|\mathbf{f}_k\|_k + \|\mathbf{f}_k\|_k \\ & \quad \sqrt{jC_{kj}} (\|\mathbf{f}_k\|_k + 2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty) \|\mathbf{f}_k\|_k + (\|\mathbf{f}_k\|_k + 2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty) \|\mathbf{f}_k\|_k + \|\mathbf{f}_k\|_k : \end{aligned}$$

Thus, if

$$2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty;$$

we have

$$\|\mathbf{f}_k\|_k \leq \|\mathbf{f}_k\|_k \left( \|\mathbf{f}_k\|_k + 2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty \right) \sqrt{jC_{kj}} \|\mathbf{f}_k\|_k \leq 2\sqrt{jC_{kj}} \|\mathbf{f}_k\|_k \leq 2\sqrt{jC_{kj}};$$

and thus

$$\|\mathbf{f}_k\|_k \leq 2\sqrt{jC_{kj}} \|\mathbf{f}_k\|_k \leq 4\sqrt{jC_{kj}} \|\mathbf{f}_k\|_k \leq 2\sqrt{jC_{kj}};$$

In summary, if

$$2k\mathbf{A}_k k_\infty + 2k\mathbf{B}_k k_\infty \leq 8k \geq [K]$$

we have

$$\max_{1 \leq k \leq K} \|\mathbf{c}_k^* - \mathbf{c}_k\|_k \leq 2\sqrt{D_S};$$

In the following, we discuss how to bound  $k\mathbf{A}_k k_\infty$  and  $k\mathbf{B}_k k_\infty$ .

### 2.3. Bound for $k\mathbf{A}_k k_\infty$

From the definition of  $\mathbf{A}_k$  in (9), we have

$$k\mathbf{A}_k k_\infty \leq 2 \frac{j\widehat{S}_k n S_{kj}}{j\widehat{S}_{kj}}.$$

#### 2.3.1. LOWER BOUND OF $j\widehat{S}_{kj}$

First, we show that the size of  $S_k$  is lower-bounded, which means a significant amount of data points in  $S$  belong to the  $k$ -th cluster. Recall that  $w_1, \dots, w_K$  are the weight of the Gaussian mixtures, and  $w_k = \min_{1 \leq i \leq K} w_i$ . According to the Chernoff bound (Angluin & Valiant, 1979) provided in Appendix A, we have, with a probability at least 1

$$jS_{kj} \geq w_k jS_j \left( 1 - \sqrt{\frac{2}{k} \ln \frac{K}{k}} \right) \stackrel{(5)}{\geq} \frac{2}{3} w_k jS_j; \quad 8k \geq [K]; \quad (10)$$

Next, we prove that a larger amount of data points in  $S_k$  belong to  $\widehat{S}_k$ . We begin by analyzing the probability that the assigned cluster  $\widehat{K}_i$  of  $\mathbf{x}_i$  is the true cluster  $K_i$ . The similarity between  $\mathbf{x}_i$  and the estimated cluster centers can be bounded by

$$\begin{aligned} \widehat{\mathbf{c}}_{k_i}^\top \mathbf{x}_i &= \widehat{\mathbf{c}}_{k_i}^\top (\mathbf{c}_{k_i} + \mathbf{g}_i) = \|\mathbf{c}_{k_i}\|_k + [\widehat{\mathbf{c}}_{k_i} - \mathbf{c}_{k_i}]^\top \mathbf{c}_{k_i} + \widehat{\mathbf{c}}_{k_i}^\top \mathbf{g}_i \\ & \geq \|\mathbf{c}_{k_i}\|_k - \|\widehat{\mathbf{c}}_{k_i} - \mathbf{c}_{k_i}\|_k + \widehat{\mathbf{c}}_{k_i}^\top \mathbf{g}_i \geq \|\mathbf{c}_{k_i}\|_k - \|\widehat{\mathbf{c}}_{k_i} - \mathbf{c}_{k_i}\|_k + \|\mathbf{g}_i\|_k \left| \frac{\widehat{\mathbf{c}}_{k_i}^\top \mathbf{c}_{k_i}}{\|\mathbf{c}_{k_i}\|_k} \right|; \\ \widehat{\mathbf{c}}_j^\top \mathbf{x}_i &= \widehat{\mathbf{c}}_j^\top (\mathbf{c}_{k_i} + \mathbf{g}_i) = \widehat{\mathbf{c}}_j^\top \mathbf{c}_{k_i} + [\widehat{\mathbf{c}}_j - \mathbf{c}_j]^\top \mathbf{c}_{k_i} + \widehat{\mathbf{c}}_j^\top \mathbf{g}_i \\ & \leq \|\widehat{\mathbf{c}}_j - \mathbf{c}_j\|_k + \|\mathbf{c}_{k_i}\|_k + \|\widehat{\mathbf{c}}_j\|_k \|\mathbf{g}_i\|_k + \|\widehat{\mathbf{c}}_j\|_k \|\mathbf{g}_i\|_k \left| \frac{\widehat{\mathbf{c}}_j^\top \mathbf{c}_{k_i}}{\|\mathbf{c}_{k_i}\|_k} \right|; \quad j \neq k_i; \end{aligned}$$

Hence,  $\mathbf{x}_i$  will be assigned to cluster  $k_i$  if

$$1 - (1 + \epsilon) \left| \mathbf{g}_i^\top \frac{\hat{\mathbf{c}}_{k_i}}{k \hat{\mathbf{c}}_{k_i} k} \right| > (1 + \epsilon) \left| \mathbf{g}_i^\top \frac{\hat{\mathbf{c}}_j}{k \hat{\mathbf{c}}_j k} \right|; \forall j \neq k_i;$$

which leads to the following sufficient condition

$$\max_{\leq j \leq K} \left| \mathbf{g}_i^\top \frac{\hat{\mathbf{c}}_j}{k \hat{\mathbf{c}}_j k} \right| \leq \frac{1 - 2\epsilon}{2(1 + \epsilon)}, \quad g_i \stackrel{(4)}{\leq} \frac{2 \sqrt{5 \ln(3K)}}{3} \sqrt{2 \ln(3K)}. \quad (11)$$

It is easy to verify that for any fixed direction  $\hat{\mathbf{c}}$  with  $k \hat{\mathbf{c}} k = 1$ ,  $\mathbf{g}_i^\top \mathbf{c}$  is a Gaussian random variable with mean 0 and variance  $\frac{1}{2}$ . Based on the tail bound for the Gaussian distribution (Chang et al., 2011) provided in Appendix B, we have

$$\Pr \left[ \max_{\leq j \leq K} \left| \mathbf{g}_i^\top \frac{\hat{\mathbf{c}}_j}{k \hat{\mathbf{c}}_j k} \right| \geq g \right] \leq K \exp \left( -\frac{g^2}{2} \right);$$

Define

$$= K \exp \left( -\frac{g^2}{2} \right) \stackrel{(11)}{\leq} \frac{1}{3}. \quad (12)$$

In summary, we have proved the following lemma.

**Lemma 2.** Under the condition in (4), with a probability at least  $1 - \frac{1}{3}$ ,  $\mathbf{x}_i = \mathbf{c}_{k_i} + \mathbf{g}_i \geq S_{k_i}$  satisfies

$$\max_{\leq j \leq K} \left| \mathbf{g}_i^\top \frac{\hat{\mathbf{c}}_j}{k \hat{\mathbf{c}}_j k} \right| \leq g;$$

and is assigned to the correct cluster  $k_i$  based on the nearest neighbor search (i.e.,  $\hat{k}_i = k_i$ ).

Define

$$S_k = \left\{ \mathbf{x}_i \geq S_k : \max_{\leq j \leq K} \left| \mathbf{g}_i^\top \frac{\hat{\mathbf{c}}_j}{k \hat{\mathbf{c}}_j k} \right| \leq g \right\} \quad \hat{S}_k \setminus S_k; \quad (13)$$

Since each data point in  $S_k$  has a probability at least  $\frac{1}{3}$  to be assigned to set  $S_k$ , using the Chernoff bound again, we have, with a probability at least  $1 - \frac{1}{3}$ ,

$$\begin{aligned} j \hat{S}_k \setminus S_k &\leq j \hat{S}_k - j S_k \leq \mathbb{E}[j S_k] \left( 1 + \sqrt{\frac{2}{\mathbb{E}[j S_k]} \ln \frac{K}{j S_k}} \right) \\ &\leq (1 + \epsilon) j S_k \left( 1 + \sqrt{\frac{2}{(1 - \epsilon) j S_k} \ln \frac{K}{j S_k}} \right) \\ &\stackrel{(12)}{\leq} \frac{2}{3} j S_k \left( 1 + \sqrt{\frac{3}{j S_k} \ln \frac{K}{j S_k}} \right) \stackrel{(5), (10)}{\leq} \frac{1}{3} j S_k; \forall k \geq [K]. \end{aligned} \quad (14)$$

### 2.3.2. UPPER BOUND OF $j \hat{S}_k \setminus S_k$

Define

$$O = \bigcup_k S_k \quad S \text{ and } \bar{O} = \bigcup_k (\hat{S}_k \setminus S_k) = S \cap \bar{O} \quad S;$$

From Lemma 2, we know that with a probability at least  $1 - \frac{1}{3}$ , each  $\mathbf{x}_i \geq S_k$  belongs to the set  $S_k \subseteq O$ . Thus, with probability at least  $1 - \frac{1}{3}$ , each  $\mathbf{x}_i \geq S$  belongs to  $O$ . In other words, with probability at most  $\frac{1}{3}$ , each  $\mathbf{x}_i \geq S$  belongs to  $\bar{O}$ . Based on the Chernoff bound, we have, with a probability at least  $1 - \frac{1}{3}$ ,

$$j \bar{O} \leq 2 \mathbb{E}[j \bar{O}] + 2 \ln \frac{1}{1 - \frac{1}{3}} \leq 2 j S + 2 \ln \frac{1}{1 - \frac{1}{3}}. \quad (15)$$

Since  $S_k \subseteq S$ , we have  $\hat{S}_k \setminus S_k \subseteq \bar{O}$ . Therefore, with a probability at least  $1 - \frac{1}{3}$ , we have

$$j \hat{S}_k \setminus S_k \leq 2 j S + 2 \ln \frac{1}{1 - \frac{1}{3}}; \forall k \geq [K]. \quad (16)$$

Combining (10), (14) and (16), we have, with probability at least  $1 - 3\epsilon$

$$kA_k k_\infty \leq \frac{2}{\epsilon} \frac{jS_j + 2 \ln \frac{1}{\epsilon}}{k_j S_j} = \frac{18}{k} \left( \frac{1}{jS_j} + \ln \frac{1}{\epsilon} \right) = O\left(\frac{1}{jS_j}\right) + O\left(\ln \frac{1}{\epsilon}\right); \forall k \geq [K]; \quad (17)$$

#### 2.4. Bound for $kB_k k_\infty$

Notice that  $f\mathbf{g}_i : \mathbf{x}_i \in \hat{S}_k g$ , determined by the estimated centers  $\hat{\mathbf{c}}_1; \dots; \hat{\mathbf{c}}_K$ , is a specific subset of  $f\mathbf{g}_i : \mathbf{x}_i \in S g$ . Although  $\mathbf{g}_i$  is drawn from the Gaussian distribution  $N(0; I)$ , the distribution of elements in  $f\mathbf{g}_i : \mathbf{x}_i \in \hat{S}_k g$  is unknown. As a result, we cannot directly apply concentration inequality of Gaussian random vectors to bound  $kB_k k_\infty$ . Let  $U \in \mathbb{R}^{d \times K}$  be a matrix whose columns are basis vectors of the subspace spanned by  $\hat{\mathbf{c}}_1; \dots; \hat{\mathbf{c}}_K$ , and  $U \in \mathbb{R}^{d \times (d-K)}$  be a matrix whose columns are basis vectors of the complementary subspace. We then divide each  $\mathbf{g}_i$  as

$$\mathbf{g}_i = \mathbf{g}_i^\parallel + \mathbf{g}_i^\perp;$$

where  $\mathbf{g}_i^\parallel = U U^\top \mathbf{g}_i$ , and  $\mathbf{g}_i^\perp = U^\perp U^{\perp\top} \mathbf{g}_i$ .

First, we upper bound  $kB_k k_\infty$  as

$$kB_k k_\infty \leq \underbrace{\left\| \frac{1}{j\hat{S}_k j} \sum_{\mathbf{x}_i \in \hat{S}_k} \mathbf{g}_i^\perp \right\|_\infty}_{\hat{B}_k^1} + \underbrace{\frac{j\hat{S}_k n S_{k,j}}{j\hat{S}_k j} \left\| \frac{1}{j\hat{S}_k n S_{k,j}} \sum_{\mathbf{x}_i \in \hat{S}_k \setminus S_k^1} \mathbf{g}_i^\parallel \right\|_\infty}_{\hat{B}_k^2} + \underbrace{\frac{jS_{k,j}}{j\hat{S}_k j} \left\| \frac{1}{jS_{k,j}} \sum_{\mathbf{x}_i \in S_k^1} \mathbf{g}_i^\parallel \right\|_\infty}_{\hat{B}_k^3}; \quad (18)$$

In the following, we discuss how to bound each term in the right hand side of (18).

##### 2.4.1. UPPER BOUND OF $\hat{B}_k$

Following the property of Gaussian random vector,  $\sum_{\mathbf{x}_i \in \hat{S}_k} U^\top \mathbf{g}_i = \left( \sqrt{j\hat{S}_k j} \right)$  can be treated as a  $(d - K)$ -dimensional Gaussian random vector. As a result, each element of  $U \sum_{\mathbf{x}_i \in \hat{S}_k} U^\top \mathbf{g}_i = \left( \sqrt{j\hat{S}_k j} \right)$  is a Gaussian random variable with variance smaller than 1. Based on the tail bound for the Gaussian distribution (Chang et al., 2011) provided in Appendix B and the union bound, with a probability at least  $1 - \epsilon$ , we have

$$\left\| \sum_{\mathbf{x}_i \in \hat{S}_k} \mathbf{g}_i^\perp = \left( \sqrt{j\hat{S}_k j} \right) \right\|_\infty = \left\| U \sum_{\mathbf{x}_i \in \hat{S}_k} U^\top \mathbf{g}_i = \left( \sqrt{j\hat{S}_k j} \right) \right\|_\infty \leq \sqrt{2 \ln \frac{Kd}{\epsilon}}; \forall k \geq [K];$$

which implies

$$\hat{B}_k \leq \sqrt{\frac{2 \ln \frac{Kd}{\epsilon}}{j\hat{S}_k j}} \stackrel{(10), (14)}{\leq} \sqrt{\frac{2 \ln \frac{Kd}{\epsilon}}{2 \frac{k_j S_j}{9}}} = O\left(\sqrt{\frac{\ln d}{jS_j}}\right); \forall k \geq [K]; \quad (19)$$

##### 2.4.2. UPPER BOUND OF $\hat{B}_k$

First, we have

$$\left\| \frac{1}{j\hat{S}_k n S_{k,j}} \sum_{\mathbf{x}_i \in \hat{S}_k \setminus S_k^1} \mathbf{g}_i^\parallel \right\|_\infty = \left\| \frac{1}{j\hat{S}_k n S_{k,j}} \sum_{\mathbf{x}_i \in \hat{S}_k \setminus S_k^1} U U^\top \mathbf{g}_i \right\|_\infty = \left\| \frac{1}{j\hat{S}_k n S_{k,j}} \sum_{\mathbf{x}_i \in \hat{S}_k \setminus S_k^1} U^\top \mathbf{g}_i \right\|_\infty \quad (20)$$

Since  $U^\top \mathbf{g}_i =$  can be treated as a  $K$ -dimensional Gaussian random vector, based on the tail bound for the distribution (Laurent & Massart, 2000), we have with a probability at least  $1 - \epsilon$ ,

$$\|U^\top \mathbf{g}_i\| \leq \left( \sqrt{\frac{K}{\epsilon}} + \sqrt{2 \log \frac{1}{\epsilon}} \right)$$

Applying the union bound again, with a probability at least  $1 - \delta$ , we have

$$\max_{\leq i \leq |S|} \|U^\top \mathbf{g}_i\| \leq \left( \rho_{\overline{K}} + \sqrt{2 \log \frac{jS_j}{\delta}} \right) \quad (21)$$

Combining (20) and (21), we have

$$\widehat{B}_k \leq \frac{9}{k} \left( \rho_{\overline{K}} + \frac{1}{jS_j} \ln \frac{1}{\delta} \right) \left( \rho_{\overline{K}} + \sqrt{2 \log \frac{jS_j}{\delta}} \right) = O\left( \sqrt{\ln jS_j} \right) + O\left( \frac{\sqrt{\ln jS_j}}{jS_j} \right); \delta k \geq [K] \quad (22)$$

### 2.4.3. UPPER BOUND OF $\widehat{B}_k$

First, we have

$$\left\| \frac{1}{jS_{k^j}} \sum_{\mathbf{x}_i \in S_k^1} \mathbf{g}_i \right\|_\infty = \left\| U \frac{1}{jS_{k^j}} \sum_{\mathbf{x}_i \in S_k^1} U^\top \mathbf{g}_i \right\|_\infty = \left\| \frac{1}{jS_{k^j}} \sum_{\mathbf{x}_i \in S_k^1} U^\top \mathbf{g}_i \right\| := u_k \quad (23)$$

Recall the definition of  $S_k$  in (13). Due to the fact that the domain is symmetric, we have  $E[U^\top \mathbf{g}_i] = 0$ . Under the condition in (21), we can invoke the following lemma to bound  $u_k$ .

**Lemma 3.** (Lemma 2 from (Smale & Zhou, 2007)) Let  $H$  be a Hilbert space and  $\mathbf{g}_i$  be a random variable on  $(Z; \mathcal{H})$  with values in  $H$ . Assume  $k \leq M < 1$  almost surely. Denote  $\mu = E(k)$ . Let  $\{Z_i, \mathbf{g}_i^m\}$  be independent random drawers of  $\mathbf{g}_i$ . For any  $0 < \delta < 1$ , with confidence  $1 - \delta$ ,

$$\left\| \frac{1}{m} \sum_i^m (\mathbf{g}_i - E[\mathbf{g}_i]) \right\| \leq \frac{2M \ln(2/\delta)}{m} + \sqrt{\frac{2 \mu \ln(2/\delta)}{m}}$$

From Lemma 3 and the union bound, with a probability at least  $1 - \delta$ , we have

$$u_k \leq \left( \rho_{\overline{K}} + \sqrt{2 \log \frac{jS_j}{\delta}} \right) \left( \frac{2 \ln(2K/\delta)}{jS_{k^j}} + \sqrt{\frac{2 \ln(2K/\delta)}{jS_{k^j}}} \right); \delta k \geq [K] \quad (24)$$

Combining (23) and (24), we have

$$\begin{aligned} \widehat{B}_k &\leq \left( \rho_{\overline{K}} + \sqrt{2 \log \frac{jS_j}{\delta}} \right) \left( \frac{2}{jS_{k^j}} \ln \frac{2K}{\delta} + \sqrt{\frac{2}{jS_{k^j}} \ln \frac{2K}{\delta}} \right) \\ &\stackrel{(10), (14), (5)}{=} \left( \rho_{\overline{K}} + \sqrt{2 \log \frac{jS_j}{\delta}} \right) 2 \sqrt{\frac{9}{k jS_j} \ln \frac{2K}{\delta}} = O\left( \sqrt{\frac{\ln jS_j}{jS_j}} \right); \delta k \geq [K] \end{aligned} \quad (25)$$

In summary, under the condition that (10), (14) and (15) are true, with a probability at least  $1 - 3\delta$ ,

$$k B_k k_\infty \leq O\left( \sqrt{\ln jS_j} \right) + O\left( \frac{\sqrt{\ln jS_j} + \rho_{\overline{\ln d}}}{\sqrt{jS_j}} \right); \delta k \geq [K] \quad (26)$$

## A. Chernoff Bound

**Theorem 2** (Multiplicative Chernoff Bound (Angluin & Valiant, 1979)). Let  $X_1, X_2, \dots, X_n$  be independent binary random variables with  $\Pr[X_i = 1] = p_i$ . Denote  $S = \sum_{i=1}^n X_i$  and  $\mu = E[S] = \sum_{i=1}^n p_i$ . We have

$$\Pr[S \leq (1 - \delta)\mu] \leq \exp\left(-\frac{\delta^2 \mu}{2}\right); \text{ for } 0 < \delta < 1;$$

$$\Pr[S \geq (1 + \delta)\mu] \leq \exp\left(-\frac{\delta^2 \mu}{2 + \delta}\right); \text{ for } \delta > 0;$$

Therefore,

$$\Pr \left[ S \leq \left( 1 - \sqrt{\frac{2}{\delta} \ln \frac{1}{\delta}} \right) \right] \leq \exp \left( -\frac{\delta}{2} \right) < \delta < 1;$$

$$\Pr \left[ S \geq \left( 1 + \frac{\ln \frac{1}{\delta} + \sqrt{2 \ln \frac{1}{\delta}}}{\delta} \right) \right] \leq \delta < 1;$$

## B. Tail bounds for the Gaussian distribution

**Theorem 3** (Chernoff-type upper bound for the  $Q$ -function (Chang et al., 2011)). *The  $Q$ -function defined as*

$$Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^{\infty} \exp \left( -\frac{t^2}{2} \right) dt$$

is the tail probability of the standard Gaussian distribution. When  $x > 0$ , we have

$$Q(x) \leq \frac{1}{x} \exp \left( -\frac{x^2}{2} \right).$$

Let  $X \sim N(0; 1)$  be a Gaussian random variable. According to Theorem 3, we have

$$\Pr [ |X_j| \geq \sqrt{\frac{2}{\delta}} ] \leq \exp \left( -\frac{\delta}{2} \right); \text{ or}$$

$$\Pr \left[ |X_j| \geq \sqrt{2 \ln \frac{1}{\delta}} \right] \leq \delta.$$

## References

- Angluin, D. and Valiant, L.G. Fast probabilistic algorithms for hamiltonian circuits and matchings. *Journal of Computer and System Sciences*, 18(2):155–193, 1979.
- Chang, Seok-Ho, Cosman, Pamela C., and Milstein, Laurence B. Chernoff-type bounds for the gaussian error function. *IEEE Transactions on Communications*, 59(11):2939–2944, 2011.
- Laurent, B. and Massart, P. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302–1338, 2000.
- Smale, Steve and Zhou, Ding-Xuan. Learning theory estimates via integral operators and their approximations. *Constructive Approximation*, 26(2):153–172, 2007.