# 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* $\delta_1 = \delta/(6m)$ *be a parameter to control the success probability. Assume*

$$\Delta_* \geq \frac{\lambda_1}{\gamma} \geq \Delta_{max}; \tag{1}$$

$$\frac{1}{2}\rho^{-\frac{1}{2s}} \geq \lambda_1 \geq c\frac{1}{2}\rho^{-\frac{1}{2s}}; \tag{2}$$

$$T \geq \max\left(\frac{18}{\gamma_0}\ln\frac{2K}{\delta_1}; \frac{3c_2\gamma_0}{\gamma_1}; \left(\frac{6c_3}{\gamma_1}\right)^2(\ln n + \ln d)\right) \tag{3}$$

*where* $c$, $c_2$ *and* $c_3$ *are some universal constants. Then, with a probability at least* $1 - 6m\delta$, *we have*

$$\Delta^{m+1} = \max_{1\leq i\leq K}\|k\widehat{\mathbf{c}}_i^{m+1} - \mathbf{c}_i\|k \leq \max\left(\Delta_*; \frac{c}{\gamma}\rho^{-\frac{1}{2^m}}\right):$$

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

## 1. Proof of Corollary 1

According to the assumption of $\lambda_1$ in (2), we know that $\frac{1}{\lambda_1} \lesssim \frac{\sqrt{s}}{\gamma_1}$. Since the value of $T$ is dominated by the last term in the right side of (3), we have $T \lesssim \frac{s\log d}{\gamma_1 \cdot \gamma_1}$, which implies

$$n \lesssim 2^m T \lesssim 2^m \frac{s\log d}{\gamma_1\gamma_1}:$$

Combining with the conclusion $\epsilon_{m+1} \lesssim \frac{1}{\sqrt{2^m}}$, we have

$$\epsilon_{m+1} \lesssim \sqrt{\frac{s \log d}{n}}.$$

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

$$\epsilon^t \leq \frac{1}{2\sqrt{5 \ln(3K)}}, \quad \epsilon_{max}; \tag{4}$$

$$|S^t_j| \geq \frac{18}{\pi_0} \ln \frac{2K}{\delta}; \tag{5}$$

$$\epsilon^t \geq c_1 \exp\left( -\frac{(1 - 2\epsilon^t \sigma)^2}{8(1 + \epsilon^t)^2 \sigma^2} \right) (\sigma_0 + \sigma\sqrt{\ln|S^t_j|}) + \frac{c_2 \sigma_0}{|S^t_j|} + c_3 \sigma \frac{\sqrt{\ln|S^t_j|} + \sqrt[p]{\ln d}}{\sqrt{|S^t_j|}}; \tag{6}$$

*for some constants $c_1$, $c_2$ and $c_3$. Then with a probability $1 - 6\delta$, we have*

$$\epsilon^{t+1} \leq 2^p \bar{\sigma} \epsilon^t.$$

# 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 $\overline{C}_k = [d] \cap C_k$. 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 $\widehat{k}_i$ to denote index of the cluster assigned by the nearest neighbor search, i.e.,

$$\mathbf{x}_i = \mathbf{c}_{k_i} + \mathbf{g}_i \text{ and } \mathbf{g}_i \sim N(0, \sigma^2 I);$$
$$\widehat{k}_i = \arg\max_{j \in [K]} \widehat{\mathbf{c}}_j^\top \mathbf{x}_i.$$

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.,

$$S_k = \{\mathbf{x}_i \in S : \mathbf{x}_i = \mathbf{c}_k + \mathbf{g}_i \text{ and } \mathbf{g}_i \sim N(0, \sigma^2 I)\} \tag{7}$$

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

$$\widehat{S}_k = \{\mathbf{x}_i \in S : k = \arg\max_{j \in [K]} \widehat{\mathbf{c}}_j^\top \mathbf{x}_i\} \tag{8}$$

## 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

$$
\begin{aligned}
&L_k(\mathbf{c}) \\
&= \|\mathbf{c}\|_1 + \|\mathbf{c} - \mathbf{c}_k\|^2 + \frac{1}{|\widehat{S}_k|} \sum_{\mathbf{x}_i \in \widehat{S}_k} \|\mathbf{x}_i - \mathbf{c}_k\|^2 - \frac{2}{|\widehat{S}_k|} \sum_{\mathbf{x}_i \in \widehat{S}_k} (\mathbf{c} - \mathbf{c}_k)^\top (\mathbf{x}_i - \mathbf{c}_k) \\
&= \|\mathbf{c}\|_1 + \|\mathbf{c} - \mathbf{c}_k\|^2 + \frac{1}{|\widehat{S}_k|} \sum_{\mathbf{x}_i \in \widehat{S}_k} \|\mathbf{x}_i - \mathbf{c}_k\|^2 \\
&\quad - \underbrace{2(\mathbf{c} - \mathbf{c}_k)^\top \frac{1}{|\widehat{S}_k|} \sum_{\mathbf{x}_i \in \widehat{S}_k \setminus S_k} (\mathbf{c}_{k_i} - \mathbf{c}_k)}_{A_k} - \underbrace{2(\mathbf{c} - \mathbf{c}_k)^\top \frac{1}{|\widehat{S}_k|} \sum_{\mathbf{x}_i \in \widehat{S}_k} \mathbf{g}_i}_{B_k}.
\end{aligned} \tag{9}
$$

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

$$
\begin{aligned}
& L_k(\mathbf{c}_k^*) - L_k(\mathbf{c}_k) \\
&= \|\mathbf{f}_k + \mathbf{c}_k\|_1 + \|\mathbf{f}_k\|^2 - 2\mathbf{f}_k^\top A_k - 2\mathbf{f}_k^\top B_k - \|\mathbf{c}_k\|_1 \\
& \quad - \|\mathbf{c}_k\|_1 - \|\mathbf{f}_k(C_k)\|_1 + \|\mathbf{f}_k(\overline{C}_k)\|_1 + \|\mathbf{f}_k\|^2 - 2\mathbf{f}_k^\top A_k - 2\mathbf{f}_k^\top B_k - \|\mathbf{c}_k\|_1 \\
& \quad - \|\mathbf{f}_k(C_k)\|_1 + \|\mathbf{f}_k(\overline{C}_k)\|_1 + \|\mathbf{f}_k\|^2 - 2\|\mathbf{f}_k\|_1 \|A_k\|_\infty - 2\|\mathbf{f}_k\|_1 \|B_k\|_\infty \\
&= (-1 + 2\|A_k\|_\infty + 2\|B_k\|_\infty)\|\mathbf{f}_k(C_k)\|_1 + (1 - 2\|A_k\|_\infty - 2\|B_k\|_\infty)\|\mathbf{f}_k(\overline{C}_k)\|_1 + \|\mathbf{f}_k\|^2 \\
& \quad -\sqrt{|C_k|}(-1 + 2\|A_k\|_\infty + 2\|B_k\|_\infty)\|\mathbf{f}_k(C_k)\| + (1 - 2\|A_k\|_\infty - 2\|B_k\|_\infty)\|\mathbf{f}_k(\overline{C}_k)\|_1 + \|\mathbf{f}_k\|^2;
\end{aligned}
$$

Thus, if

$$
2\|A_k\|_\infty + 2\|B_k\|_\infty;
$$

we have

$$
\|\mathbf{f}_k(C_k)\|^2 \geq \|\mathbf{f}_k\|^2 \geq (-1 + 2\|A_k\|_\infty + 2\|B_k\|_\infty)\sqrt{|C_k|}\|\mathbf{f}_k(C_k)\| \geq 2\sqrt{|C_k|}\|\mathbf{f}_k(C_k)\|) \geq \|\mathbf{f}_k(C_k)\| \geq 2\sqrt{|C_k|};
$$

and thus

$$
\|\mathbf{f}_k\|^2 \geq 2\sqrt{|C_k|}\|\mathbf{f}_k(C_k)\| \geq 4^2|C_k|) \geq \|\mathbf{f}_k\| \geq 2\sqrt{|C_k|};
$$

In summary, if

$$
2\|A_k\|_\infty + 2\|B_k\|_\infty; 8k \in 2[K]
$$

we have

$$
\max_{1 \leq k \leq K} \|\mathbf{c}_k^* - \mathbf{c}_k\| \leq 2^{\rho}\bar{s};
$$

In the following, we discuss how to bound $\|A_k\|_\infty$ and $\|B_k\|_\infty$.

## 2.3. Bound for $\|A_k\|_\infty$

From the definition of $A_k$ in (9), we have

$$
\|A_k\|_\infty \leq 2_0 \frac{|\widehat{S}_k \cap S_k|}{|\widehat{S}_k|};
$$

### 2.3.1. LOWER BOUND OF $|\widehat{S}_k|$

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 $_1, \ldots, _K$ are the weight of the Gaussian mixtures, and $_0 = \min\limits_{1 \leq i \leq K} _i$. According to the Chernoff bound (Angluin & Valiant, 1979) provided in Appendix A, we have, with a probability at least $1 -$

$$
|S_k| \geq _k|S|\left(1 - \sqrt{\frac{2}{_k|S|}\ln\frac{K}{}}\right) \overset{(5)}{\geq} \frac{2}{3}_k|S|; 8k \in 2[K]: \tag{10}
$$

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

$$
\begin{aligned}
\widehat{\mathbf{c}}_{k_i}^\top \mathbf{x}_i &= \widehat{\mathbf{c}}_{k_i}^\top(\mathbf{c}_{k_i} + \mathbf{g}_i) = \|\mathbf{c}_{k_i}\|^2 + [\widehat{\mathbf{c}}_{k_i} - \mathbf{c}_{k_i}]^\top \mathbf{c}_{k_i} + \widehat{\mathbf{c}}_{k_i}^\top \mathbf{g}_i \\
& \quad 1 - \|\widehat{\mathbf{c}}_{k_i} - \mathbf{c}_{k_i}\| - |\widehat{\mathbf{c}}_{k_i}^\top \mathbf{g}_i| \geq 1 - (1 + )\left|\mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_{k_i}}{\|\widehat{\mathbf{c}}_{k_i}\|}\right|; \\
\widehat{\mathbf{c}}_j^\top \mathbf{x}_i &= \widehat{\mathbf{c}}_j^\top(\mathbf{c}_{k_i} + \mathbf{g}_i) = \mathbf{c}_j^\top \mathbf{c}_{k_i} + [\widehat{\mathbf{c}}_j - \mathbf{c}_j]^\top \mathbf{c}_{k_i} + \widehat{\mathbf{c}}_j^\top \mathbf{g}_i \\
& \quad + \|\widehat{\mathbf{c}}_j - \mathbf{c}_j\| + |\widehat{\mathbf{c}}_j^\top \mathbf{g}_i| \leq + + (1 + )\left|\mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_j}{\|\widehat{\mathbf{c}}_j\|}\right|; j \neq k_i:
\end{aligned}
$$

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

$$1 \quad (1+ \ ) \left| \mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_{k_i}}{\|\widehat{\mathbf{c}}_{k_i}\|} \right| \quad + \quad +(1+ \ ) \left| \mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_j}{\|\widehat{\mathbf{c}}_j\|} \right| ; \ \forall j \neq k_i;$$

which leads to the following sufficient condition

$$\max_{1 \le j \le K} \left| \mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_j}{\|\widehat{\mathbf{c}}_j\|} \right| \quad \frac{1 \quad 2}{2(1+ \ )} \ , \quad g_0 \overset{(4)}{=} \frac{2 \ \sqrt{5 \ln(3K)}}{3} \quad \sqrt{2\ln(3K)}: \tag{11}$$

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

$$\Pr \left[ \max_{1 \le j \le K} \left| \mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_j}{\|\widehat{\mathbf{c}}_j\|} \right| \quad g_0 \right] \quad 1 \quad K \exp \left( \ \frac{g_0^2}{2 \ ^2} \right):$$

Define

$$= K \exp \left( \ \frac{g_0^2}{2 \ ^2} \right) \overset{(11)}{=} \frac{1}{3}: \tag{12}$$

In summary, we have proved the following lemma.

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

$$\max_{1 \le j \le K} \left| \mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_j}{\|\widehat{\mathbf{c}}_j\|} \right| \quad g_0;$$

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

Define

$$S_k^1 = \left\{ \mathbf{x}_i \ 2 \ S_k : \max_{1 \le j \le K} \left| \mathbf{g}_i^\top \frac{\widehat{\mathbf{c}}_j}{\|\widehat{\mathbf{c}}_j\|} \right| \quad g_0 \right\} \quad \widehat{S}_k \setminus S_k: \tag{13}$$

Since each data point in $S_k$ has a probability at least $1 \quad$ to be assigned to set $S_k^1$, using the Chernoff bound again, we have, with a probability at least $1 \quad$,

$$|\widehat{S}_k| \quad |\widehat{S}_k \setminus S_k| \quad |S_k^1| \quad \mathbb{E}\left[ |S_k^1| \right] \left( 1 \quad \sqrt{\frac{2}{\mathbb{E}[|S_k^1|]} \ln \frac{K}{}} \right)$$

$$(1 \quad )|S_k| \left( 1 \quad \sqrt{\frac{2}{(1 \quad )|S_k|} \ln \frac{K}{}} \right)$$

$$\overset{(12)}{=} \frac{2}{3}|S_k| \left( 1 \quad \sqrt{\frac{3}{|S_k|} \ln \frac{K}{}} \right) \overset{(5), (10)}{=} \frac{1}{3}|S_k|; \forall k \ 2 \ [K]: \tag{14}$$

### 2.3.2. UPPER BOUND OF $|\widehat{S}_k \cap S_k|$

Define

$$O = \big[_{k=1}^{K} S_k^1 \quad S \text{ and } \overline{O} = \big[_{k=1}^{K} \left( \widehat{S}_k \cap S_k^1 \right) = S \cap O \quad S:$$

From Lemma 2, we know that with a probability at least $1 \quad$, each $\mathbf{x}_i \ 2 \ S_k$ belongs to the set $S_k^1 \quad O$. Thus, with probability at least $1 \quad$, each $\mathbf{x}_i \ 2 \ S$ belongs to $O$. In other words, with probability *at most* , each $\mathbf{x}_i \ 2 \ S$ belongs to $\overline{O}$. Based on the Chernoff bound, we have, with a probability at least $1 \quad$,

$$|\overline{O}| \quad 2\mathbb{E}\left[ |\overline{O}| \right] + 2 \ln \frac{1}{} \quad 2 \ |S| + 2 \ln \frac{1}{}: \tag{15}$$

Since $S_k^1 \quad S_k$, we have $\widehat{S}_k \cap S_k \quad \widehat{S}_k \cap S_k^1 \quad \overline{O}$. Therefore, with a probability at least $1 \quad$, we have

$$|\widehat{S}_k \cap S_k| \quad 2 \ |S| + 2 \ln \frac{1}{}; \forall k \ 2 \ [K]: \tag{16}$$

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

$$\|A_k\|_\infty \le 2\sigma_0 \frac{2|S| + 2\ln\frac{1}{\epsilon}}{\frac{2}{9}\alpha_k|S|} = \frac{18\sigma_0}{\alpha_k}\left(1 + \frac{1}{|S|}\ln\frac{1}{\epsilon}\right) = O(\sigma_0) + O\left(\frac{\sigma_0}{|S|}\right), \; \forall k \in [K]. \tag{17}$$

### 2.4. Bound for $\|B_k\|_\infty$

Notice that $\{\mathbf{g}_i : \mathbf{x}_i \in \widehat{S}_k\}$, determined by the estimated centers $\widehat{\mathbf{c}}_1, \ldots, \widehat{\mathbf{c}}_K$, is a specific subset of $\{\mathbf{g}_i : \mathbf{x}_i \in S\}$. Although $\mathbf{g}_i$ is drawn from the Gaussian distribution $N(0, \sigma^2 I)$, the distribution of elements in $\{\mathbf{g}_i : \mathbf{x}_i \in \widehat{S}_k\}$ is unknown. As a result, we cannot direct apply concentration inequality of Gaussian random vectors to bound $\|B_k\|_\infty$. Let $U_1 \in \mathbb{R}^{d \times K}$ be a matrix whose columns are basis vectors of the subspace spanned by $\widehat{\mathbf{c}}_1, \ldots, \widehat{\mathbf{c}}_K$, and $U_2 \in \mathbb{R}^{d \times (d-K)}$ be a matrix whose columns are basis vectors of the complementary subspace. We then divide each $\mathbf{g}_i$ as

$$\mathbf{g}_i = \mathbf{g}_i^{\|} + \mathbf{g}_i^{\perp},$$

where $\mathbf{g}_i^{\|} = U_1 U_1^\top \mathbf{g}_i$, and $\mathbf{g}_i^{\perp} = U_2 U_2^\top \mathbf{g}_i$.

First, we upper bound $\|B_k\|_\infty$ as

$$\|B_k\|_\infty \le \underbrace{\left\|\frac{1}{|\widehat{S}_k|}\sum_{\mathbf{x}_i \in \widehat{S}_k}\mathbf{g}_i^{\perp}\right\|_\infty}_{\widehat{B}_k^1} + \underbrace{\frac{|\widehat{S}_k \cap S_k^1|}{|\widehat{S}_k|}\left\|\frac{1}{|\widehat{S}_k \cap S_k^1|}\sum_{\mathbf{x}_i \in \widehat{S}_k \setminus S_k^1}\mathbf{g}_i^{\|}\right\|_\infty}_{\widehat{B}_k^2} + \underbrace{\frac{|S_k^1|}{|\widehat{S}_k|}\left\|\frac{1}{|S_k^1|}\sum_{\mathbf{x}_i \in S_k^1}\mathbf{g}_i^{\|}\right\|_\infty}_{\widehat{B}_k^3}. \tag{18}$$

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

#### 2.4.1. UPPER BOUND OF $\widehat{B}_k^1$

Following the property of Gaussian random vector, $\sum_{\mathbf{x}_i \in \widehat{S}_k} U_2^\top \mathbf{g}_i \sim \left(\sigma\sqrt{|\widehat{S}_k|}\right)$ can be treated as a $(d - K)$-dimensional Gaussian random vector. As a result, each element of $U_2 \sum_{\mathbf{x}_i \in \widehat{S}_k} U_2^\top \mathbf{g}_i \sim \left(\sigma\sqrt{|\widehat{S}_k|}\right)$ 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 $1 - \delta$, we have

$$\left\|\sum_{\mathbf{x}_i \in \widehat{S}_k}\mathbf{g}_i^{\perp} \sim \left(\sigma\sqrt{|\widehat{S}_k|}\right)\right\|_\infty = \left\|U_2\sum_{\mathbf{x}_i \in \widehat{S}_k} U_2^\top \mathbf{g}_i \sim \left(\sigma\sqrt{|\widehat{S}_k|}\right)\right\|_\infty \le \sigma\sqrt{2\ln\frac{Kd}{\delta}}, \; \forall k \in [K],$$

which implies

$$\widehat{B}_k^1 \le \sigma\sqrt{\frac{2\ln\frac{Kd}{\epsilon}}{|\widehat{S}_k|}} \overset{(10),(14)}{\le} \sigma\sqrt{\frac{2\ln\frac{Kd}{\epsilon}}{2\alpha_k|S|/9}} = O\left(\sigma\sqrt{\frac{\ln d}{|S|}}\right), \; \forall k \in [K]. \tag{19}$$

#### 2.4.2. UPPER BOUND OF $\widehat{B}_k^2$

First, we have

$$\left\|\frac{1}{|\widehat{S}_k \cap S_k^1|}\sum_{\mathbf{x}_i \in \widehat{S}_k \setminus S_k^1}\mathbf{g}_i^{\|}\right\|_\infty = \left\|\frac{1}{|\widehat{S}_k \cap S_k^1|}\sum_{\mathbf{x}_i \in \widehat{S}_k \setminus S_k^1} U_1 U_1^\top \mathbf{g}_i\right\|_\infty \le \left\|\frac{1}{|\widehat{S}_k \cap S_k^1|}\sum_{\mathbf{x}_i \in \widehat{S}_k \setminus S_k^1} U_1^\top \mathbf{g}_i\right\| \tag{20}$$

Since $U_1^\top \mathbf{g}_i \sim \sigma$ can be treated as a $K$-dimensional Gaussian random vector, based on the tail bound for the $\chi^2$ distribution (Laurent & Massart, 2000), we have with a probability at least $1 - \delta$,

$$\|U_1^\top \mathbf{g}_i\| \le \sigma\left(\sqrt{K} + \sqrt{2\log\frac{1}{\delta}}\right)$$

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

$$\max_{1 \le i \le |\mathcal{S}|} \|U_1^\top \mathbf{g}_i\| \le \left( \sqrt{\frac{p}{K}} + \sqrt{2 \log \frac{|S|}{\delta}} \right) \tag{21}$$

Combining (20) and (21), we have

$$\widehat{B}_k^2 \le \frac{9}{k} \left( 1 + \frac{1}{|S|} \ln \frac{1}{\delta} \right) \left( \sqrt{\frac{p}{K}} + \sqrt{2 \log \frac{|S|}{\delta}} \right) = O(\sqrt{\ln |S|}) + O\left( \frac{\sqrt{\ln |S|}}{|S|} \right), \; \forall k \in [K]. \tag{22}$$

### 2.4.3. UPPER BOUND OF $\widehat{B}_k^3$

First, we have

$$\left\| \frac{1}{|S_k^1|} \sum_{\mathbf{x}_i \in \mathcal{S}_k^1} \mathbf{g}_i^\| \right\|_\infty = \left\| U_1 \frac{1}{|S_k^1|} \sum_{\mathbf{x}_i \in \mathcal{S}_k^1} U_1^\top \mathbf{g}_i \right\|_\infty \le \left\| \frac{1}{|S_k^1|} \sum_{\mathbf{x}_i \in \mathcal{S}_k^1} U_1^\top \mathbf{g}_i \right\| := u_k \tag{23}$$

Recall the definition of $S_k^1$ in (13). Due to the fact that the domain is symmetric, we have $\mathrm{E}\left[ U_1^\top \mathbf{g}_i \right] = 0$. 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 $\xi$ be a random variable on $(Z, \rho)$ with values in $H$. Assume $\|\xi\| \le M < \infty$ almost surely. Denote $\sigma^2(\xi) = \mathrm{E}(\|\xi\|^2)$. Let $\{z_i\}_{i=1}^m$ be independent random drawers of $\rho$. For any $0 < \delta < 1$, with confidence $1 - \delta$,*

$$\left\| \frac{1}{m} \sum_{i=1}^m (\xi_i - \mathrm{E}[\xi_i]) \right\| \le \frac{2M \ln(2/\delta)}{m} + \sqrt{\frac{2 \sigma^2(\xi) \ln(2/\delta)}{m}}$$

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

$$u_k \le \left( \sqrt{\frac{p}{K}} + \sqrt{2 \log \frac{|S|}{\delta}} \right) \left( \frac{2 \ln(2K/\delta)}{|S_k^1|} + \sqrt{\frac{2 \ln(2K/\delta)}{|S_k^1|}} \right); \; \forall k \in [K]. \tag{24}$$

Combining (23) and (24), we have

$$\widehat{B}_k^3 \le \left( \sqrt{\frac{p}{K}} + \sqrt{2 \log \frac{|S|}{\delta}} \right) \left( \frac{2}{|S_k^1|} \ln \frac{2K}{\delta} + \sqrt{\frac{2}{|S_k^1|} \ln \frac{2K^2}{\delta}} \right)$$

$$\overset{(10),(14),(5)}{\le} \left( \sqrt{\frac{p}{K}} + \sqrt{2 \log \frac{|S|}{\delta}} \right) 2 \sqrt{\frac{9}{k|S|} \ln \frac{2K}{\delta}} = O\left( \sqrt{\frac{\ln |S|}{|S|}} \right); \; \forall k \in [K]. \tag{25}$$

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

$$\|B_k\|_\infty \le O(\sqrt{\ln |S|}) + O\left( \frac{\sqrt{\ln |S|} + \sqrt{\frac{p}{\ln d}}}{\sqrt{|S|}} \right); \; \forall k \in [K]. \tag{26}$$

## A. Chernoff Bound

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

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

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

Therefore,

$$\Pr\left[S \geq \left(1 - \sqrt{\frac{2}{}\ln\frac{1}{}}\right)\right] \geq \quad ; \text{ for } \exp\left(-\frac{}{2}\right) < \; < 1;$$

$$\Pr\left[S \leq 2 + 2\ln\frac{1}{}\left(1 + \frac{\ln\frac{1}{\delta} + \sqrt{2 \ln\frac{1}{\delta}}}{}\right)\right] \geq \quad ; \text{ for } 0 < \; < 1:$$

## 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*

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

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

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

$$\Pr\left[|X| \geq \sqrt{2 \ln\frac{1}{}}\right] \leq \quad :$$

## 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.