

A Proof of Theorem 3

The error analysis for the uniform+adaptive² algorithm relies on Lemma 7, which guarantees the error incurred by its uniform sampling step. The proof of Lemma 7 essentially follows Gittens (2011). We prove Lemma 7 using probability inequalities and some techniques of Boutsidis et al. (2011); Gittens (2011); Gittens and Mahoney (2013); Tropp (2012); the proof is in Appendix A.1.

Lemma 7 (Uniform Column Sampling). *Given an $m \times n$ matrix \mathbf{A} and a target rank k , let μ_k denote the matrix coherence of \mathbf{A} . By sampling*

$$c = \frac{\mu_k k \log(k/\delta)}{\theta \log \theta - \theta + 1},$$

columns uniformly without replacement to construct \mathbf{C} , the following inequality

$$\|\mathbf{A} - \mathcal{P}_{\mathbf{C},k}\mathbf{A}\|_F^2 \leq (1 + \delta^{-1}\theta^{-1})\|\mathbf{A} - \mathbf{A}_k\|_F^2.$$

holds with probability at least $1 - 2\delta$. Here $\delta \in (0, 0.5)$ and $\theta \in (0, 1)$ are arbitrary real numbers.

The error analysis for the two adaptive sampling steps of the uniform+adaptive² algorithm relies on Lemma 8, which follows immediately from (Wang and Zhang, 2013, Corollary 7 and Section 4.5).

Lemma 8. *Given an $m \times m$ symmetric matrix \mathbf{A} and a target rank k , we let \mathbf{C}_1 contain the c_1 columns of \mathbf{A} selected by a column sampling algorithm such that the following inequality holds:*

$$\|\mathbf{A} - \mathcal{P}_{\mathbf{C}_1}\mathbf{A}\|_F^2 \leq f\|\mathbf{A} - \mathbf{A}_k\|_F^2.$$

Then we select $c_2 = kf\epsilon^{-1}$ columns to construct \mathbf{C}_2 and $c_3 = (c_1 + c_2)\epsilon^{-1}$ columns to construct \mathbf{C}_3 , both using the adaptive sampling according to the residual $\mathbf{B}_1 = \mathbf{A} - \mathcal{P}_{\mathbf{C}_1}\mathbf{A}$ and $\mathbf{B}_2 = \mathbf{A} - \mathcal{P}_{[\mathbf{C}_1, \mathbf{C}_2]}\mathbf{A}$, respectively. Let $\mathbf{C} = [\mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3]$, we have that

$$\mathbb{P}\left\{\frac{\|\mathbf{A} - \mathbf{C}(\mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T) \mathbf{C}^T\|_F}{\|\mathbf{A} - \mathbf{A}_k\|_F} \geq 1 + s\epsilon\right\} \leq \frac{1 + \epsilon}{1 + s\epsilon},$$

where s is an arbitrary constant greater than 1.

Finally Theorem 3 is proved by combining Lemma 7 and Lemma 8. The proof is in Appendix A.2.

A.1 Proof of Lemma 7

Proof. We use uniform column sampling to select c column of \mathbf{A} to construct $\mathbf{C} = \mathbf{A}\mathbf{S}$. Here the $n \times c$ random matrix \mathbf{S} has one entry equal to one and the rest equal to zero in each column, and at most one nonzero entry in each row, and \mathbf{S} is uniformly distributed among $\binom{n}{c}$ such kind of

matrices. Applying Lemma 7 of Boutsidis et al. (2011), we get

$$\begin{aligned} & \|\mathbf{A} - \mathcal{P}_{\mathbf{C},k}\mathbf{A}\|_F^2 \\ & \leq \|\mathbf{A} - \mathbf{A}_k\|_F^2 + \|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_F^2 \|(\mathbf{V}_{\mathbf{A},k}^T \mathbf{S})^\dagger\|_2^2. \end{aligned} \quad (3)$$

Now we bound $\|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_2^2$ and $\|(\mathbf{V}_{\mathbf{A},k}^T \mathbf{S})^\dagger\|_2^2$ respectively using the techniques of Gittens (2011); Gittens and Mahoney (2013); Tropp (2012).

Let $\mathcal{I} \subset [n]$ be a random index set corresponding to \mathbf{S} . The support of \mathcal{I} is uniformly distributing among all the index sets in $2^{[n]}$ with cardinality c . According to Gittens and Mahoney (2013), the expectation of $\|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_F^2$ can be written as

$$\begin{aligned} \mathbb{E}\|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_F^2 &= \mathbb{E}\|(\mathbf{A} - \mathbf{A}_k)_{\mathcal{I}}\|_F^2 \\ &= c\mathbb{E}\|(\mathbf{A} - \mathbf{A}_k)_i\|_F^2 = \frac{c}{n}\|\mathbf{A} - \mathbf{A}_k\|_F^2. \end{aligned}$$

Applying Markov's inequality, we have that

$$\begin{aligned} \mathbb{P}\left\{\|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_F^2 \geq \frac{c}{n\delta}\|\mathbf{A} - \mathbf{A}_k\|_F^2\right\} \\ \leq \frac{\mathbb{E}\|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_F^2}{\frac{c}{n\delta}\|\mathbf{A} - \mathbf{A}_k\|_F^2} = \delta. \end{aligned} \quad (4)$$

Here $\delta \in (0, 0.5)$ is a real number defined later.

Now we establish the bound for $\mathbb{E}\|\Omega_2^\dagger\|_2^2$ as follows. Let $\lambda_i(\mathbf{X})$ be the i -th largest eigenvalue of \mathbf{X} . Following the proof of Lemma 1 of Gittens (2011), we have

$$\begin{aligned} \|(\mathbf{V}_{\mathbf{A},k}^T \mathbf{S})^\dagger\|_2^2 &= \lambda_k^{-1}\left(\mathbf{V}_{\mathbf{A},k}^T \mathbf{S} \mathbf{S}^T \mathbf{V}_{\mathbf{A},k}\right) \\ &= \lambda_k^{-1}\left(\sum_{i=1}^c \mathbf{X}_i\right) \leq \lambda_{\min}^{-1}\left(\sum_{i=1}^c \mathbf{X}_i\right), \end{aligned} \quad (5)$$

where the random matrices $\mathbf{X}_1, \dots, \mathbf{X}_c$ are chosen uniformly at random from the set $\left\{(\mathbf{V}_{\mathbf{A},k}^T)_i (\mathbf{V}_{\mathbf{A},k}^T)_i^T\right\}_{i=1}^n$ without replacement. The random matrices are of size $k \times k$. We accordingly define

$$R = \max_i \lambda_{\max}(\mathbf{X}_i) = \max_i \|(\mathbf{V}_{\mathbf{A},k}^T)_i\|_2^2 = \frac{k}{n}\mu_k,$$

where μ_k is the matrix coherence of \mathbf{A} , and define

$$\begin{aligned} \beta_{\min} &= \lambda_{\min}\left(\mathbb{E}\sum_{i=1}^c \mathbf{X}_i\right) = \lambda_{\min}(c\mathbb{E}\mathbf{X}_1) \\ &= \lambda_{\min}\left(\frac{c}{n}\mathbf{V}_{\mathbf{A},k}^T \mathbf{V}_{\mathbf{A},k}\right) = \frac{c}{n}. \end{aligned}$$

Then we apply Lemma 9 and obtained the following inequality:

$$\mathbb{P}\left[\lambda_{\min}\left(\sum_{i=1}^c \mathbf{X}_i\right) \leq \frac{\theta c}{n}\right] \leq k \left[\frac{e^{\theta-1}}{\theta^\theta}\right]^{\frac{c}{k\mu_k}} \triangleq \delta, \quad (6)$$

where $\theta \in (0, 1]$ is a real number, and it follows that

$$c = \frac{\mu_k k \log(k/\delta)}{\theta \log \theta - \theta + 1}.$$

Applying (5) and (6), we have

$$\mathbb{P}\left\{\|(\mathbf{V}_{\mathbf{A},k}^T \mathbf{S})^\dagger\|_2^2 \geq \frac{n}{\theta c}\right\} \leq \delta. \quad (7)$$

Combining (4) and (7) and applying the union bound, we have the following inequality:

$$\begin{aligned} \mathbb{P}\left\{\|(\mathbf{A} - \mathbf{A}_k)\mathbf{S}\|_F^2 \geq \frac{c}{n\delta}\|\mathbf{A} - \mathbf{A}_k\|_F^2 \right. \\ \left. \text{or } \|(\mathbf{V}_{\mathbf{A},k}^T \mathbf{S})^\dagger\|_2^2 \geq \frac{n}{\theta c}\right\} \leq 2\delta. \quad (8) \end{aligned}$$

Finally, from (3) and (8) we have that the inequality

$$\|\mathbf{A} - \mathcal{P}_{\mathbf{C},k}\mathbf{A}\|_F^2 \leq (1 + \delta^{-1}\theta^{-1})\|\mathbf{A} - \mathbf{A}_k\|_F^2$$

holds with probability at least $1 - 2\delta$, by which the lemma follows. \square

Lemma 9 (Theorem 5.1.1 of Tropp (2012)). *We are given l independent random $d \times d$ SPDS matrices $\mathbf{X}_1, \dots, \mathbf{X}_l$ with the property*

$$\lambda_{\max}(\mathbf{X}_i) \leq R \quad \text{for } i = 1, \dots, l.$$

We define $\mathbf{Y} = \sum_{i=1}^l \mathbf{X}_i$ and $\beta_{\min} = \lambda_{\min}(\mathbb{E}\mathbf{Y})$. Then for any $\theta \in (0, 1]$, the following inequality holds:

$$\mathbb{P}\left\{\lambda_{\min}(\mathbf{Y}) \leq \theta \beta_{\min}\right\} \leq k \left[\frac{e^{\theta-1}}{\theta^\theta}\right]^{\frac{\beta_{\min}}{R}}.$$

A.2 Proof of the Theorem

Proof. The matrix \mathbf{C}_1 consists of c_1 columns selected by uniform sampling, and $\mathbf{C}_2 \in \mathbb{R}^{m \times c_2}$ and $\mathbf{C}_3 \in \mathbb{R}^{m \times c_3}$ are constructed by adaptive sampling. We set $\delta = 1/\sqrt{5}$ and $\theta = \sqrt{5}/4$ for Lemma 7, then we have

$$\begin{aligned} f &= 1 + \delta^{-1}\theta^{-1} = 5, \\ c_1 &= \frac{\mu_k k \log(k/\delta)}{\theta \log \theta - \theta + 1} = 8.7\mu_k k \log(\sqrt{5}k). \end{aligned}$$

Then we set

$$\begin{aligned} c_2 &= kf\epsilon^{-1} = 5k\epsilon^{-1}, \\ c_3 &= (c_1 + c_2)\epsilon^{-1}, \end{aligned}$$

according to Lemma 8. Letting $s > 1$ be an arbitrary constant, we have that

$$\begin{aligned} &\mathbb{P}\left\{\frac{\|\mathbf{A} - \mathbf{C}\mathbf{U}\mathbf{C}^T\|_F}{\|\mathbf{A} - \mathbf{A}_k\|_F} \leq 1 + s\epsilon\right\} \\ &\geq \mathbb{P}\left\{\frac{\|\mathbf{A} - \mathbf{C}\mathbf{U}\mathbf{C}^T\|_F}{\|\mathbf{A} - \mathbf{A}_k\|_F} \leq 1 + s\epsilon \mid \frac{\|\mathbf{A} - \mathcal{P}_{\mathbf{C}_1}\mathbf{A}\|_F^2}{\|\mathbf{A} - \mathbf{A}_k\|_F^2} \leq f\right\} \\ &\quad \cdot \mathbb{P}\left\{\frac{\|\mathbf{A} - \mathcal{P}_{\mathbf{C}_1}\mathbf{A}\|_F^2}{\|\mathbf{A} - \mathbf{A}_k\|_F^2} \leq f\right\} \\ &\geq \left(1 - \frac{1 + \epsilon}{1 + s\epsilon}\right)(1 - 2\delta). \end{aligned}$$

where the last inequality follows from Lemma 7 and Lemma 8.

Repeating the sampling procedure for t times and letting $\mathbf{C}_{[i]}$ and $\mathbf{U}_{[i]}$ be the i -th sample, we obtain an upper error bound on the failure probability:

$$\begin{aligned} &\mathbb{P}\left\{\min_{i \in [t]} \left\{\frac{\|\mathbf{A} - \mathbf{C}_{[i]}\mathbf{U}_{[i]}\mathbf{C}_{[i]}^T\|_F}{\|\mathbf{A} - \mathbf{A}_k\|_F}\right\} \geq 1 + s\epsilon\right\} \\ &\leq \left(1 - \left(1 - \frac{1 + \epsilon}{1 + s\epsilon}\right)(1 - 2\delta)\right)^t \\ &= \left(1 + \frac{(s-1)(1-2\delta)}{\epsilon^{-1} + 1 + 2\delta(s-1)}\right)^{-t} \triangleq p. \end{aligned}$$

Taking logarithm of both sides of the equality and applying $\log(1+x) \approx x$ when x is small, we have

$$\begin{aligned} t &= \left[\log\left(1 + \frac{(1-2\delta)(s-1)}{\epsilon^{-1} + 1 + 2\delta(s-1)}\right)\right]^{-1} \log \frac{1}{p} \\ &\approx \frac{\epsilon^{-1} + 1 + 2\delta(s-1)}{(1-2\delta)(s-1)} \log \frac{1}{p}. \end{aligned}$$

Setting $s = 2$, we have that $t \approx (10\epsilon^{-1} + 18) \log(1/p)$.

Hence by sampling totally

$$c = (1 + \epsilon^{-1})(5k\epsilon^{-1} + 8.7\mu_k k \log(\sqrt{5}k))$$

columns and repeating the procedure for

$$t \geq (10\epsilon^{-1} + 18) \log(1/p)$$

times, the algorithm attains the upper error bound

$$\|\mathbf{A} - \mathbf{C}(\mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T) \mathbf{C}^T\|_F \leq (1 + 2\epsilon)\|\mathbf{A} - \mathbf{A}_k\|_F$$

with probability at least $1 - p$. Substituting 2ϵ by ϵ' yields the error bound in the theorem.

Time complexity and space complexity of Algorithm 1 is calculated as follows. The uniform sampling costs $\mathcal{O}(m)$ time; the first adaptive sampling round costs $\mathcal{O}(mc_1^2) + T_{\text{Multiply}}(m^2c_1)$ time; the second adaptive sampling round costs $\mathcal{O}(m(c_1 + c_2)^2) + T_{\text{Multiply}}(m^2(c_1 + c_2))$ time; computing the intersection matrix costs $\mathcal{O}(mc^2) +$

$T_{\text{Multiply}}(m^2c)$ time in general. So the total time complexity is $\mathcal{O}(mc^2) + T_{\text{Multiply}}(m^2c)$ without using Theorem 4, or $\mathcal{O}(m(c_1 + c_2)^2) + T_{\text{Multiply}}(m^2c)$ using Theorem 4. As for the space complexity, the Moore-Penrose inverse of an $m \times c$ matrix demands $\mathcal{O}(mc)$ space, and multiplying a $c \times m$ matrix \mathbf{C}^\dagger by an $m \times m$ matrix \mathbf{A} costs $\mathcal{O}(mc)$ space by partition \mathbf{A} into small blocks of size smaller than $m \times c$ and loading one block into RAM at a time to perform matrix multiplication. \square

B Proof of Theorem 4

Proof. Let $\mathbf{C} \in \mathbb{R}^{m \times c}$ consists of a subset of columns of \mathbf{A} . By row permutation \mathbf{C} can be expressed as

$$\mathbf{PC} = \begin{bmatrix} \mathbf{W} \\ \mathbf{A}_{21} \end{bmatrix}.$$

Then according to Lemma 10, the Moore-Penrose inverse of \mathbf{C} can be written as

$$\mathbf{C}^\dagger = \mathbf{W}^{-1}(\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \mathbf{P},$$

where $\mathbf{S} = \mathbf{A}_{21} \mathbf{W}^{-1}$. Then the intersection matrix of modified Nyström approximation to \mathbf{A} can be expressed as

$$\begin{aligned} \mathbf{U} &= \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \\ &= \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \mathbf{P} \mathbf{A} \mathbf{P}^T \\ &\quad \begin{bmatrix} \mathbf{I}_c \\ \mathbf{S} \end{bmatrix} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \mathbf{W}^{-1} \\ &= \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \\ &\quad \begin{bmatrix} \mathbf{W} & \mathbf{A}_{21}^T \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I}_c \\ \mathbf{S} \end{bmatrix} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \mathbf{W}^{-1} \\ &= \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \left(\mathbf{W} + \mathbf{A}_{21}^T \mathbf{S} + (\mathbf{A}_{21}^T \mathbf{S})^T \right. \\ &\quad \left. + \mathbf{S}^T \mathbf{A}_{22} \mathbf{S} \right) (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \mathbf{W}^{-1} \\ &\triangleq \mathbf{T}_1 (\mathbf{W} + \mathbf{T}_2 + \mathbf{T}_2^T + \mathbf{T}_3) \mathbf{T}_1^T. \end{aligned}$$

Here the intermediate matrices are computed by

$$\begin{aligned} \mathbf{T}_0 &= \mathbf{A}_{21}^T \mathbf{A}_{21}, \\ \mathbf{T}_1 &= \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \\ &= \mathbf{W}^{-1} \left(\mathbf{I}_c + \mathbf{W}^{-1} \mathbf{T}_0 \mathbf{W}^{-1} \right)^{-1}, \\ \mathbf{T}_2 &= \mathbf{A}_{21}^T \mathbf{S} = \mathbf{A}_{21}^T \mathbf{A}_{21} \mathbf{W}^{-1} = \mathbf{T}_0 \mathbf{W}^{-1}, \\ \mathbf{T}_3 &= \mathbf{S}^T \mathbf{A}_{22} \mathbf{S} = \mathbf{W}^{-1} \left(\mathbf{A}_{21}^T \mathbf{A}_{22} \mathbf{A}_{21} \right) \mathbf{W}^{-1}. \end{aligned}$$

The matrix inverse operations are on $c \times c$ matrices which costs $\mathcal{O}(c^3)$ time. The matrix multiplication $\mathbf{A}_{21}^T \mathbf{A}_{22} \mathbf{A}_{21}$ requires time $T_{\text{Multiply}}((m-c)^2c)$. \square

Lemma 10 (The Moore Penrose Inverse of Partitioned Matrices (Ben-Israel and Greville, 2003, Page 179)). *Given a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ of rank of at least c which has a nonsingular $c \times c$ submatrix \mathbf{X}_{11} . By rearrangement of columns and rows by permutation matrices \mathbf{P} and \mathbf{Q} , the submatrix \mathbf{X}_{11} can be brought to the top left corner of \mathbf{X} , that is,*

$$\mathbf{P} \mathbf{X} \mathbf{Q} = \begin{bmatrix} \mathbf{X}_{11} & \mathbf{X}_{12} \\ \mathbf{X}_{21} & \mathbf{X}_{22} \end{bmatrix}.$$

Then the Moore-Penrose inverse of \mathbf{X} is

$$\begin{aligned} \mathbf{X}^\dagger &= \mathbf{Q} \begin{bmatrix} \mathbf{I}_c \\ \mathbf{T}^T \end{bmatrix} (\mathbf{I}_c + \mathbf{T} \mathbf{T}^T)^{-1} \mathbf{X}_{11}^{-1} \\ &\quad (\mathbf{I}_c + \mathbf{S} \mathbf{S}^T)^{-1} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \mathbf{P}, \end{aligned}$$

where $\mathbf{T} = \mathbf{X}_{11}^{-1} \mathbf{X}_{12}$ and $\mathbf{S} = \mathbf{X}_{21} \mathbf{X}_{11}^{-1}$.

C The Proof of Theorem 5

Proof. Suppose that $\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{A})$. We have that $\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{C}) = \text{rank}(\mathbf{A})$ because

$$\text{rank}(\mathbf{A}) \geq \text{rank}(\mathbf{C}) \geq \text{rank}(\mathbf{W}). \quad (9)$$

Thus there exists a matrix \mathbf{X} such that

$$\begin{bmatrix} \mathbf{A}_{21}^T \\ \mathbf{A}_{22} \end{bmatrix} = \mathbf{C} \mathbf{X}^T = \begin{bmatrix} \mathbf{W} \mathbf{X}^T \\ \mathbf{A}_{21} \mathbf{X}^T \end{bmatrix},$$

and it follows that $\mathbf{A}_{21} = \mathbf{X} \mathbf{W}$ and $\mathbf{A}_{22} = \mathbf{A}_{21} \mathbf{X}^T = \mathbf{X} \mathbf{W} \mathbf{X}^T$. Then we have that

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} \mathbf{W} & (\mathbf{X} \mathbf{W})^T \\ \mathbf{X} \mathbf{W} & \mathbf{X} \mathbf{W} \mathbf{X}^T \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} \\ \mathbf{X} \end{bmatrix} \mathbf{W} \begin{bmatrix} \mathbf{I} & \mathbf{X}^T \end{bmatrix}, \quad (10) \\ \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T &= \begin{bmatrix} \mathbf{W} \\ \mathbf{X} \mathbf{W} \end{bmatrix} \mathbf{W}^\dagger \begin{bmatrix} \mathbf{W} & (\mathbf{X} \mathbf{W})^T \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} \\ \mathbf{X} \end{bmatrix} \mathbf{W} \begin{bmatrix} \mathbf{I} & \mathbf{X}^T \end{bmatrix}. \quad (11) \end{aligned}$$

Here the second equality in (11) follows from $\mathbf{W} \mathbf{W}^\dagger \mathbf{W} = \mathbf{W}$. We obtain that $\mathbf{A} = \mathbf{C} \mathbf{W}^\dagger \mathbf{C}$. Then we show that $\mathbf{A} = \mathbf{C} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{C}^T$.

Since $\mathbf{C}^\dagger = (\mathbf{C}^T \mathbf{C})^\dagger \mathbf{C}^T$, we have that

$$\mathbf{C}^\dagger = (\mathbf{W}(\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W})^\dagger \mathbf{W} \begin{bmatrix} \mathbf{I} & \mathbf{X}^T \end{bmatrix},$$

and thus

$$\begin{aligned} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{W} &= (\mathbf{W}(\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W})^\dagger \mathbf{W} (\mathbf{I} + \mathbf{X}^T \mathbf{X}) \begin{bmatrix} \mathbf{W}(\mathbf{I} + \mathbf{X}^T \mathbf{X}) \\ \mathbf{W}(\mathbf{W}(\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W})^\dagger \mathbf{W} \end{bmatrix} \\ &= (\mathbf{W}(\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W})^\dagger \mathbf{W} (\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W}, \end{aligned}$$

where the second equality follows from Lemma 11 because $(\mathbf{I} + \mathbf{X}^T \mathbf{X})$ is positive definite. Similarly we have

$$\begin{aligned} & \mathbf{W} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{W} \\ &= \mathbf{W} (\mathbf{W} (\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W})^\dagger \mathbf{W} (\mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{W} = \mathbf{W}. \end{aligned}$$

Thus we have

$$\begin{aligned} \mathbf{C} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{C} &= \begin{bmatrix} \mathbf{I} \\ \mathbf{X} \end{bmatrix} \mathbf{W} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{W} \begin{bmatrix} \mathbf{I} & \mathbf{X}^T \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} \\ \mathbf{X} \end{bmatrix} \mathbf{W} \begin{bmatrix} \mathbf{I} & \mathbf{X}^T \end{bmatrix}. \end{aligned} \quad (12)$$

It follows from Equations (10) (11) (12) that $\mathbf{A} = \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T = \mathbf{C} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{C}^T$.

Conversely, when $\mathbf{A} = \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T$, we have that $\text{rank}(\mathbf{A}) \leq \text{rank}(\mathbf{W}^\dagger) = \text{rank}(\mathbf{W})$. By applying (9) we have that $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{W})$.

When $\mathbf{A} = \mathbf{C} \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^\dagger)^T \mathbf{C}^T$, we have $\text{rank}(\mathbf{A}) \leq \text{rank}(\mathbf{C})$. Thus there exists a matrix \mathbf{X} such that

$$\begin{bmatrix} \mathbf{A}_{21}^T \\ \mathbf{A}_{22} \end{bmatrix} = \mathbf{C} \mathbf{X}^T = \begin{bmatrix} \mathbf{W} \mathbf{X}^T \\ \mathbf{A}_{21} \mathbf{X}^T \end{bmatrix},$$

and therefore $\mathbf{A}_{21} = \mathbf{X} \mathbf{W}$. Then we have that

$$\mathbf{C} = \begin{bmatrix} \mathbf{W} \\ \mathbf{A}_{21} \end{bmatrix} = \begin{bmatrix} \mathbf{I} \\ \mathbf{X} \end{bmatrix} \mathbf{W},$$

so $\text{rank}(\mathbf{C}) \leq \text{rank}(\mathbf{W})$. Apply (9) again we have $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{W})$. \square

Lemma 11. $\mathbf{X}^T \mathbf{V} \mathbf{X} (\mathbf{X}^T \mathbf{V} \mathbf{X})^\dagger \mathbf{X}^T = \mathbf{X}^T$ for any positive definite matrix \mathbf{V} .

Proof. Since the positive definite matrix \mathbf{V} have a decomposition $\mathbf{V} = \mathbf{B}^T \mathbf{B}$ for some nonsingular matrix \mathbf{B} , so we have

$$\begin{aligned} & \mathbf{X}^T \mathbf{V} \mathbf{X} (\mathbf{X}^T \mathbf{V} \mathbf{X})^\dagger \mathbf{X}^T \\ &= (\mathbf{B} \mathbf{X})^T \left(\mathbf{B} \mathbf{X} ((\mathbf{B} \mathbf{X})^T (\mathbf{B} \mathbf{X}))^\dagger \right) (\mathbf{B} \mathbf{X})^T \mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} \\ &= (\mathbf{B} \mathbf{X})^T ((\mathbf{B} \mathbf{X})^T)^\dagger (\mathbf{B} \mathbf{X})^T (\mathbf{B}^T)^{-1} \\ &= (\mathbf{B} \mathbf{X})^T (\mathbf{B}^T)^{-1} \\ &= \mathbf{X}^T. \end{aligned} \quad \square$$

D Proof of Theorem 6

In Section D.1 we provide two key lemmas, and then in Section D.2 we prove Theorem 6 using the two lemmas.

D.1 Key Lemmas

Lemma 12. For an $m \times m$ matrix \mathbf{B} with diagonal entries equal to one and off-diagonal entries equal to α , the error incurred by the modified Nyström method is lower bounded by

$$\begin{aligned} & \|\mathbf{B} - \tilde{\mathbf{B}}_c^{\text{mod}}\|_F^2 \\ & \geq (1 - \alpha)^2 (m - c) \left(1 + \frac{2}{c} - (1 - \alpha) \frac{1 + o(1)}{\alpha c m / 2} \right). \end{aligned}$$

Proof. Without loss of generality, we assume the first c column of \mathbf{B} are selected to construct \mathbf{C} . We partition \mathbf{B} and \mathbf{C} as:

$$\mathbf{B} = \begin{bmatrix} \mathbf{W} & \mathbf{B}_{21}^T \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \quad \text{and} \quad \mathbf{C} = \begin{bmatrix} \mathbf{W} \\ \mathbf{B}_{21} \end{bmatrix}.$$

Here the matrix \mathbf{W} can be expressed by $\mathbf{W} = (1 - \alpha) \mathbf{I}_c + \alpha \mathbf{1}_c \mathbf{1}_c^T$. We apply the Sherman-Morrison-Woodbury formula

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{C}^{-1} + \mathbf{D} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{D} \mathbf{A}^{-1}$$

to compute \mathbf{W}^{-1} , yielding

$$\mathbf{W}^{-1} = \frac{1}{1 - \alpha} \mathbf{I}_c - \frac{\alpha}{(1 - \alpha)(1 - \alpha + c\alpha)} \mathbf{1}_c \mathbf{1}_c^T. \quad (13)$$

We expand the Moore-Penrose inverse of \mathbf{C} by Lemma 10 and obtain

$$\mathbf{C}^\dagger = \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix}$$

where

$$\mathbf{S} = \mathbf{B}_{21} \mathbf{W}^{-1} = \frac{\alpha}{1 - \alpha + c\alpha} \mathbf{1}_{m-c} \mathbf{1}_c^T.$$

It is easily verified that $\mathbf{S}^T \mathbf{S} = \left(\frac{\alpha}{1 - \alpha + c\alpha} \right)^2 (m - c) \mathbf{1}_c \mathbf{1}_c^T$.

Now we express the matrix constructed by the modified Nyström method in a partitioned form:

$$\begin{aligned} \tilde{\mathbf{B}}_c^{\text{mod}} &= \mathbf{C} \mathbf{C}^\dagger \mathbf{B} (\mathbf{C}^\dagger)^T \mathbf{C}^T \\ &= \begin{bmatrix} \mathbf{W} \\ \mathbf{B}_{21} \end{bmatrix} \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \mathbf{B} \\ &= \begin{bmatrix} \mathbf{I}_c \\ \mathbf{S} \end{bmatrix} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \mathbf{W}^{-1} \begin{bmatrix} \mathbf{W} \\ \mathbf{B}_{21} \end{bmatrix}^T \\ &= \begin{bmatrix} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \\ \mathbf{B}_{21} \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \mathbf{B} \\ &= \begin{bmatrix} \mathbf{I}_c \\ \mathbf{S} \end{bmatrix} \begin{bmatrix} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \\ \mathbf{B}_{21} \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} \end{bmatrix}^T. \end{aligned} \quad (14)$$

We then compute the submatrices $(\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1}$ and $\mathbf{B}_{21} \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1}$ respectively as follows. We apply

the Sherman-Morrison-Woodbury formula to compute $(\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1}$, yielding

$$\begin{aligned} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} &= \left(\mathbf{I}_c + \left(\frac{\alpha}{1 - \alpha + c\alpha} \right)^2 (m - c) \mathbf{1}_c \mathbf{1}_c^T \right)^{-1} \\ &= \mathbf{I}_c - \gamma_1 \mathbf{1}_c \mathbf{1}_c^T, \end{aligned} \quad (15)$$

where

$$\gamma_1 = \frac{m - c}{mc + \left(\frac{1 - \alpha}{\alpha} \right)^2 + \frac{2(1 - \alpha)c}{\alpha}}.$$

It follows from (13) and (15) that

$$\begin{aligned} \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} &= (\gamma_2 \mathbf{I}_c - \gamma_3 \mathbf{1}_c \mathbf{1}_c^T) (\mathbf{I}_c - \gamma_1 \mathbf{1}_c \mathbf{1}_c^T) \\ &= \gamma_2 \mathbf{I}_c + (\gamma_1 \gamma_3 c - \gamma_1 \gamma_2 - \gamma_3) \mathbf{1}_c \mathbf{1}_c^T \end{aligned} \quad (16)$$

where

$$\gamma_2 = \frac{1}{1 - \alpha} \quad \text{and} \quad \gamma_3 = \frac{\alpha}{(1 - \alpha)(1 - \alpha + c\alpha)}.$$

Then we have that

$$\begin{aligned} \mathbf{B}_{21} \mathbf{W}^{-1} (\mathbf{I}_c + \mathbf{S}^T \mathbf{S})^{-1} &= \alpha (\gamma_1 \gamma_3 c^2 - \gamma_3 c - \gamma_1 \gamma_2 c + \gamma_2) \mathbf{1}_{m-c} \mathbf{1}_c^T \\ &\triangleq \gamma \mathbf{1}_{m-c} \mathbf{1}_c^T, \end{aligned} \quad (17)$$

where

$$\begin{aligned} \gamma &= \alpha (\gamma_1 \gamma_3 c^2 - \gamma_3 c - \gamma_1 \gamma_2 c + \gamma_2) \\ &= \frac{\alpha (\alpha c - \alpha + 1)}{2\alpha c - 2\alpha - 2\alpha^2 c + \alpha^2 + \alpha^2 c m + 1}. \end{aligned} \quad (18)$$

Since $\mathbf{B}_{21} = \alpha \mathbf{1}_{m-c} \mathbf{1}_c^T$ and $\mathbf{B}_{22} = (1 - \alpha) \mathbf{I}_{m-c} + \alpha \mathbf{1}_{m-c} \mathbf{1}_{m-c}^T$, it is easily verified that

$$\begin{aligned} \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \mathbf{B} \begin{bmatrix} \mathbf{I}_c \\ \mathbf{S} \end{bmatrix} &= \begin{bmatrix} \mathbf{I}_c & \mathbf{S}^T \end{bmatrix} \begin{bmatrix} \mathbf{W} & \mathbf{B}_{21}^T \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I}_c \\ \mathbf{S} \end{bmatrix} \\ &= (1 - \alpha) \mathbf{I}_c + \lambda \mathbf{1}_c \mathbf{1}_c^T, \end{aligned} \quad (19)$$

where

$$\lambda = \frac{\alpha(3\alpha m - \alpha c - 2\alpha + \alpha^2 c - 3\alpha^2 m + \alpha^2 + \alpha^2 m^2 + 1)}{(\alpha c - \alpha + 1)^2}$$

It follows from (14), (15), (17), and (19) that

$$\begin{aligned} \tilde{\mathbf{B}}_c^{\text{mod}} &= \begin{bmatrix} \mathbf{I}_c - \gamma_1 \mathbf{1}_c \mathbf{1}_c^T \\ \gamma \mathbf{1}_{m-c} \mathbf{1}_c^T \end{bmatrix} \left((1 - \alpha) \mathbf{I}_c + \lambda \mathbf{1}_c \mathbf{1}_c^T \right) \begin{bmatrix} \mathbf{I}_c - \gamma_1 \mathbf{1}_c \mathbf{1}_c^T \\ \gamma \mathbf{1}_{m-c} \mathbf{1}_c^T \end{bmatrix}^T \\ &\triangleq \begin{bmatrix} \tilde{\mathbf{B}}_{11} & \tilde{\mathbf{B}}_{21}^T \\ \tilde{\mathbf{B}}_{21} & \tilde{\mathbf{B}}_{22} \end{bmatrix}, \end{aligned}$$

where

$$\begin{aligned} \tilde{\mathbf{B}}_{11} &= (1 - \alpha) \mathbf{I}_c + [(1 - \gamma_1 c) \\ &\quad (\lambda - \lambda \gamma_1 c - (1 - \alpha) \gamma_1) - (1 - \alpha) \gamma_1] \mathbf{1}_c \mathbf{1}_c^T \\ &= (1 - \alpha) \mathbf{I}_c + \eta_1 \mathbf{1}_c \mathbf{1}_c^T, \\ \tilde{\mathbf{B}}_{21} &= \tilde{\mathbf{A}}_{12}^T = \gamma (1 - \gamma_1 c) (1 - \alpha + \lambda c) \mathbf{1}_{m-c} \mathbf{1}_c^T \\ &= \eta_2 \mathbf{1}_{m-c} \mathbf{1}_c^T, \\ \tilde{\mathbf{B}}_{22} &= \gamma^2 c (1 - \alpha + \lambda c) \mathbf{1}_{m-c} \mathbf{1}_{m-c}^T \\ &= \eta_3 \mathbf{1}_{m-c} \mathbf{1}_{m-c}^T, \end{aligned}$$

where

$$\begin{aligned} \eta_1 &= (1 - \gamma_1 c) (\lambda - \lambda \gamma_1 c - (1 - \alpha) \gamma_1) - (1 - \alpha) \gamma_1, \\ \eta_2 &= \gamma (1 - \gamma_1 c) (1 - \alpha + \lambda c), \\ \eta_3 &= \gamma^2 c (1 - \alpha + \lambda c), \end{aligned}$$

By dealing with the four blocks of $\tilde{\mathbf{B}}_c^{\text{mod}}$ respectively, we finally obtain that

$$\begin{aligned} &\|\mathbf{B} - \tilde{\mathbf{B}}_c^{\text{mod}}\|_F^2 \\ &= \|\mathbf{W} - \tilde{\mathbf{B}}_{11}\|_F^2 + 2\|\mathbf{B}_{21} - \tilde{\mathbf{B}}_{21}\|_F^2 + \|\mathbf{B}_{22} - \tilde{\mathbf{B}}_{22}\|_F^2 \\ &= c^2 (\alpha - \eta_1)^2 + 2c(m - c) (\alpha - \eta_2)^2 \\ &\quad + (m - c)(m - c - 1) (\alpha - \eta_3)^2 + (m - c)(1 - \eta_3)^2 \\ &= (m - c) (\alpha - 1)^2 (\alpha^4 c^2 m^2 - 4\alpha^4 c^2 m + 4\alpha^4 c^2 \\ &\quad + 2\alpha^4 c m^2 - 4\alpha^4 c m + \alpha^4 c + \alpha^4 m - \alpha^4 + 4\alpha^3 c^2 m \\ &\quad - 8\alpha^3 c^2 + 2\alpha^3 c m + 2\alpha^3 c - 2\alpha^3 m + 2\alpha^3 + 4\alpha^2 c^2 \\ &\quad + 2\alpha^2 c m - 7\alpha^2 c + \alpha^2 m + 4\alpha c - 2\alpha + 1) / (2\alpha c \\ &\quad - 2\alpha - 2\alpha^2 c + \alpha^2 + \alpha^2 c m + 1)^2 \\ &= (m - c) (\alpha - 1)^2 \left(1 + \frac{2}{c} - \frac{(1 - \alpha)}{c} (6\alpha c - 6\alpha \right. \\ &\quad - 12\alpha^2 c + 6\alpha^3 c + 6\alpha^2 - 2\alpha^3 + 3\alpha^2 c^2 - 3\alpha^3 c^2 \\ &\quad + 2\alpha^3 c^2 m + 3\alpha^2 c m - 3\alpha^3 c m + 2) / (2\alpha c - 2\alpha \\ &\quad \left. - 2\alpha^2 c + \alpha^2 + \alpha^2 c m + 1)^2 \right) \\ &= (m - c) (\alpha - 1)^2 \left(1 + \frac{2}{c} - (1 + o(1)) \frac{1 - \alpha}{\alpha c m / 2} \right). \end{aligned}$$

□

Lemma 13 (Lemma 19 of Wang and Zhang (2013)). *Given m and k , we let \mathbf{B} be an $\frac{m}{k} \times \frac{m}{k}$ matrix whose diagonal entries equal to one and off-diagonal entries equal to $\alpha \in [0, 1]$. We let \mathbf{A} be an $m \times m$ block-diagonal matrix*

$$\mathbf{A} = \text{diag}(\underbrace{\mathbf{B}, \dots, \mathbf{B}}_{k \text{ blocks}}). \quad (20)$$

Let \mathbf{A}_k be the best rank- k approximation to the matrix \mathbf{A} , then we have that

$$\|\mathbf{A} - \mathbf{A}_k\|_F = (1 - \alpha) \sqrt{m - k}.$$

D.2 Proof of the Theorem

Now we prove Theorem 6 using Lemma 12 and Lemma 13.

Proof. Let \mathbf{C} consist of c column sampled from \mathbf{A} and $\hat{\mathbf{C}}_i$ consist of c_i columns sampled from the i -th block diagonal matrix in \mathbf{A} . Without loss of generality, we assume $\hat{\mathbf{C}}_i$ consists of the first c_i columns of \mathbf{B} . Then the intersection matrix \mathbf{U} is computed by

$$\begin{aligned} \mathbf{U} &= \mathbf{C}^\dagger \mathbf{A} (\mathbf{C}^T)^\dagger \\ &= [\text{diag}(\hat{\mathbf{C}}_1, \dots, \hat{\mathbf{C}}_k)]^\dagger \mathbf{A} [\text{diag}(\hat{\mathbf{C}}_1^T, \dots, \hat{\mathbf{C}}_k^T)]^\dagger \\ &= \text{diag}(\hat{\mathbf{C}}_1^\dagger \mathbf{B} (\hat{\mathbf{C}}_1^\dagger)^T, \dots, \hat{\mathbf{C}}_k^\dagger \mathbf{B} (\hat{\mathbf{C}}_k^\dagger)^T). \end{aligned}$$

The modified Nyström approximation to \mathbf{A} is

$$\begin{aligned} \tilde{\mathbf{A}}_c^{\text{mod}} &= \mathbf{C} \mathbf{U} \mathbf{C}^T \\ &= \text{diag}(\hat{\mathbf{C}}_1 \hat{\mathbf{C}}_1^\dagger \mathbf{B} (\hat{\mathbf{C}}_1^\dagger)^T \hat{\mathbf{C}}_1^T, \dots, \hat{\mathbf{C}}_k \hat{\mathbf{C}}_k^\dagger \mathbf{B} (\hat{\mathbf{C}}_k^\dagger)^T \hat{\mathbf{C}}_k^T), \end{aligned}$$

and thus the approximation error is

$$\begin{aligned} \|\mathbf{A} - \tilde{\mathbf{A}}_c^{\text{mod}}\|_F^2 &= \sum_{i=1}^k \|\mathbf{B} - \hat{\mathbf{C}}_i \hat{\mathbf{C}}_i^\dagger \mathbf{B} (\hat{\mathbf{C}}_i^\dagger)^T \hat{\mathbf{C}}_i^T\|_F^2 \\ &\geq (1 - \alpha)^2 \sum_{i=1}^k (p - c_i) \left(1 + \frac{2}{c_i} - (1 - \alpha) \left(\frac{1 + o(1)}{\alpha c_i p / 2} \right) \right) \\ &= (1 - \alpha)^2 \left(\sum_{i=1}^k (p - c_i) \right. \\ &\quad \left. + \sum_{i=1}^k \frac{2(p - c_i)}{c_i} \left(1 - \frac{(1 - \alpha)(1 + o(1))}{\alpha p} \right) \right) \\ &\geq (1 - \alpha)^2 (m - c) \left(1 + \frac{2k}{c} \left(1 - \frac{k(1 - \alpha)(1 + o(1))}{\alpha m} \right) \right), \end{aligned}$$

where the former inequality follows from Lemma 12, and the latter inequality follows by minimizing over c_1, \dots, c_k . Finally we apply Lemma 13, and the theorem follows by setting $\alpha \rightarrow 1$. \square

E Supplementary Experiments

We have mentioned in Remark 1 that the resulting approximation accuracy is insensitive to the parameter μ in Algorithm 1, and setting μ to be exactly the matrix coherence does not in general give rise to the highest accuracy. To demonstrate this point of view, we conduct experiments on an RBF kernel matrix of the Letters Dataset with $\sigma = 0.2$, and we set $k = 10$.

We compare the uniform+adaptive² algorithm with different settings of μ ; we also employ the adaptive-full algorithm of Kumar et al. (2012), the near-optimal+adaptive algorithm of Wang and Zhang (2013), and the uniform sampling algorithm for comparison. The experiment

settings are the same to Section 6. Here the adaptive-full algorithm also has three steps: one uniform sampling and two adaptive sampling steps, and we set $c_1 = c_2 = c_3 = c/3$ according to Kumar et al. (2012). We plot the approximation errors in Figure 5.

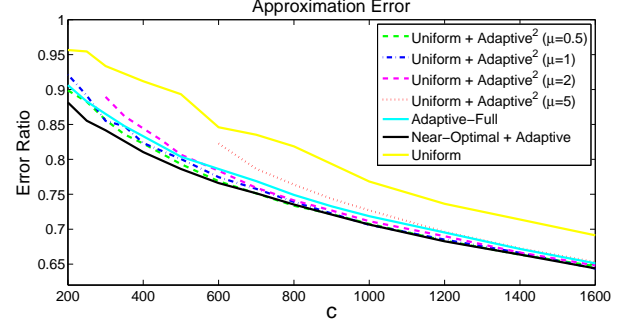


Figure 5: Effect of the parameter μ in Algorithm 1.

We can see from Figure 5 that different settings of μ does not have big influence on the approximation accuracy. We can also see that it is unnecessary to set μ to be exactly the matrix coherence; in this set of experiments, the uniform+adaptive² algorithm achieves the higher accuracy when $\mu = 0.5$ (the actual matrix coherence is $\mu_{10} = 62.05$).