

## A Appendix

### A.1 Proof of Lemma 4.1

We consider the difference between  $\mathbf{y}^{(l+1)}$  and  $\frac{A\mathbf{x}^{(l)}}{(\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)}}$  as noise, denoted by  $\mathbf{g}^{(l)}$ . To prove the results, we need to use Lemma A.1:

**Lemma A.1.** For any unit norm  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{c} \stackrel{\text{def}}{=} \mathbf{x} - \frac{A\mathbf{x}}{\mathbf{x}^T A\mathbf{x}}$  satisfies  $\mathbf{x}^T A\mathbf{x}\|\mathbf{c}\| \leq \sin \theta \lambda_1$ , where  $\theta$  is the angle between  $\mathbf{v}_1$  and  $\mathbf{x}$ .

*Proof.* Write  $\mathbf{x} = \cos \theta \mathbf{v}_1 + \sin \theta \mathbf{u}$ , where  $\mathbf{u} \perp \mathbf{v}_1$ . Then

$$\begin{aligned} & \|A\mathbf{x} - (\mathbf{x}^T A\mathbf{x})\mathbf{x}\|^2 \\ &= \|\cos \theta \lambda_1 \mathbf{v}_1 + \sin \theta A\mathbf{u} - (\cos^2 \theta \lambda_1 + \sin^2 \theta \mathbf{u}^T A\mathbf{u})(\cos \theta \mathbf{v}_1 + \sin \theta \mathbf{u})\|^2 \\ &= \|\cos \theta \sin^2 \theta (\lambda_1 - \mathbf{u}^T A\mathbf{u})\mathbf{v}_1 + \sin \theta (\cos^2 \theta ((\mathbf{u}^T A\mathbf{u})\mathbf{u} - \lambda_1 \mathbf{u}) + (A\mathbf{u} - (\mathbf{u}^T A\mathbf{u})\mathbf{u}))\|^2 \end{aligned}$$

Notice  $\mathbf{u}$ ,  $A\mathbf{u} - (\mathbf{u}^T A\mathbf{u})\mathbf{u}$  and  $\mathbf{v}_1$  are orthogonal to each other. Therefore,

$$\begin{aligned} & \|A\mathbf{x} - (\mathbf{x}^T A\mathbf{x})\mathbf{x}\|^2 \\ &= \cos^2 \theta \sin^2 \theta (\lambda_1 - \mathbf{u}^T A\mathbf{u})^2 + \sin^2 \theta \|A\mathbf{u} - (\mathbf{u}^T A\mathbf{u})\mathbf{u}\|^2 \\ &\leq \sin^2 \theta (\lambda_1 - \mathbf{u}^T A\mathbf{u})^2 + \sin^2 \theta \|A\mathbf{u}\|^2 \\ &\leq (\lambda_1 \sin \theta)^2 \end{aligned}$$

The last step makes use of the fact that  $\lambda_1 A - A^T A$  is positive semidefinite, so that  $\lambda_1 \mathbf{u}^T A\mathbf{u} \geq \mathbf{u}^T A^T A\mathbf{u} = \|A\mathbf{u}\|^2$  for any  $\mathbf{u}$ .  $\square$

Now we have the following corollary.

**Corollary A.1.1.** For  $\mathbf{g}^{(l)}$ ,  $\mathbf{x}^{(l)}$ ,  $\phi^{(l)}(k)$  defined for Algorithm 1,  $(\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)}\|\mathbf{g}^{(l)}\| \leq \sin \theta^{(l)} \lambda_1 \phi^{(l)}(k)$ .

This result is crucial to the following proof of Lemma 4.1.

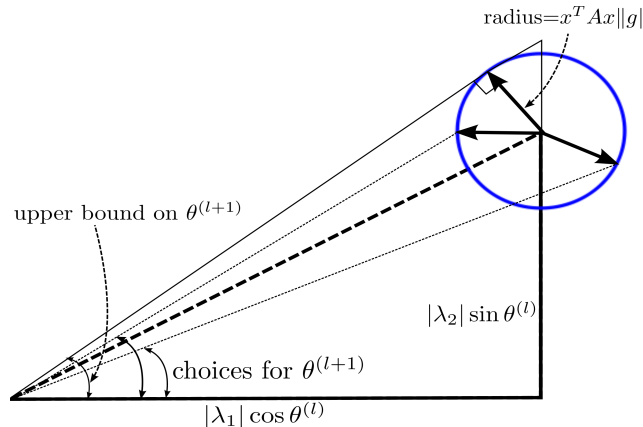


Figure 5: The central right triangle has a base-side of length  $\lambda_1 \cos \theta^{(l)}$  and height of at most  $\lambda_2 \sin \theta^{(l)}$ . The dashed line that ends in the center of the circle is  $A\mathbf{x}$  and the straight lines with an arrow are possible  $\mathbf{g}$ 's. Then Eq (11) can be represented by the tangent of the angle between the base-side and the dotted lines that ends on the circle of radius  $\mathbf{x}^T A\mathbf{x}\|\mathbf{g}\|$ . Therefore  $\tan \theta^{(l+1)} \leq \frac{\lambda_2 \sin \theta^{(l)} + \mathbf{x}^T A\mathbf{x}\|\mathbf{g}\| / \cos \theta^{(l+1)}}{\lambda_1 \cos \theta^{(l)}}$ .

Let  $U \in \mathbb{R}^{n \times (n-1)} = [\mathbf{v}_2 | \mathbf{v}_3 | \dots | \mathbf{v}_n]$  denote the orthonormal space of  $\mathbf{v}_1$ . The next iterate satisfies:

$$\begin{aligned}
& \tan \theta^{(l+1)} \\
&= \frac{\|U^T \mathbf{y}^{(l+1)}\|}{\mathbf{v}_1^T \cdot \mathbf{y}^{(l+1)}} \\
&= \frac{\|U^T \frac{A\mathbf{x}^{(l)}}{(\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)}} + U^T \mathbf{g}^{(l)}\|}{\mathbf{v}_1^T \frac{A\mathbf{x}^{(l)}}{(\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)}} + \mathbf{v}_1^T \mathbf{g}^{(l)}} \\
&\leq \frac{\sin \theta^{(l)} \lambda_2 + (\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)} \|U^T \mathbf{g}^{(l)}\|}{\cos \theta^{(l)} \lambda_1 + (\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)} \mathbf{v}_1^T \mathbf{g}^{(l)}} \tag{11}
\end{aligned}$$

$$\leq \frac{\sin \theta^{(l)} \lambda_2 + (\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)} \|\mathbf{g}^{(l)}\| / \cos \theta^{(l+1)}}{\cos \theta^{(l)} \lambda_1} \tag{12}$$

The logic from Eq (11) to Eq (12) is interpreted in Figure A.1.

Applying Lemma A.1 on Inequality (11), one gets

$$\tan \theta^{(l+1)} \leq \frac{\sin \theta^{(l)} \lambda_2 + \phi(k) \sin \theta^{(l)} \lambda_1}{\cos \theta^{(l)} \lambda_1 - \phi(k) \sin \theta^{(l)} \lambda_1} = \tan \theta^{(l)} \frac{\lambda_2 + \lambda_1 \phi(k)}{\lambda_1 (1 - \tan \theta^{(l)} \phi(k))}$$

Therefore with large enough  $k$  such that  $\phi(k) \leq \frac{\lambda_1 - \lambda_2}{\lambda_1 (1 + \tan \theta^{(l)})}$ , we could guarantee that  $\theta^{(l+1)} < \theta^{(l)}$ ,  $\frac{1}{\cos \theta^{(l+1)}} < \frac{1}{\cos \theta^{(l)}}$ . So continuing Eq. (12), we have

$$\begin{aligned}
\tan \theta^{(l+1)} &\leq \frac{\sin \theta^{(l)} \lambda_2 + (\mathbf{x}^{(l)})^T A\mathbf{x}^{(l)} \|\mathbf{g}^{(l)}\| / \cos \theta^{(l)}}{\cos \theta^{(l)} \lambda_1} \\
&\leq \frac{\sin \theta^{(l)} \lambda_2 + \phi(k) \sin \theta^{(l)} \lambda_1 / \cos \theta^{(l)}}{\cos \theta^{(l)} \lambda_1} \\
&= \tan \theta^{(l)} \left( \frac{\lambda_2}{\lambda_1} + \frac{\phi(k)}{\cos \theta^{(l)}} \right)
\end{aligned}$$

## A.2 Proof of Theorem 4.2

When  $\phi^{(l)}(k) \leq (\lambda_1 - \lambda_2) / (2\lambda_1 (1 + \tan \theta^{(l)}))$ , we obtain that,

$$\frac{\phi^{(l)}(k)}{\cos \theta^{(l)}} \leq \frac{\lambda_1 - \lambda_2}{2\lambda_1 (\cos \theta^{(l)} + \sin \theta^{(l)})} \leq (\lambda_1 - \lambda_2) / (2\lambda_1)$$

$$\tan \theta^{(l)} \leq \tan \theta^{(l-1)} (\lambda_2 / \lambda_1 + (\lambda_1 - \lambda_2) / (2\lambda_1)) \tag{13}$$

$$\leq \tan \theta^{(0)} \left( \frac{\lambda_1 + \lambda_2}{2\lambda_1} \right)^l \tag{14}$$

$$\leq \tan \theta^{(0)} e^{-l(\lambda_1 - \lambda_2) / (2\lambda_1)} \tag{15}$$

Therefore when  $l \geq 2 \frac{\lambda_1}{\lambda_1 - \lambda_2} \log \frac{\tan \theta^{(0)}}{\varepsilon}$ ,  $\tan \theta^{(l)} \leq \varepsilon$ .

## A.3 Proof of Corollary

To compare convergence rate between CPM and PM in comparable operations, one should notice one iteration of CPM costs around  $\frac{k}{n}$  percentage of operations as PM does. Therefore we should compare our convergence rate  $\frac{\lambda_1 + \lambda_2}{2\lambda_1}$  with  $(\frac{\lambda_2}{\lambda_1})^{\frac{k}{n}}$ . Therefore when

$$k < \frac{\log \left( \frac{\lambda_1 + \lambda_2}{2\lambda_1} \right)}{\log \frac{\lambda_2}{\lambda_1}} n,$$

our convergence rate is better than power method in terms of equivalent passes over data.

#### A.4 Proof of Theorem 4.3

**Lemma A.2.** In objective (5)  $f(\mathbf{x}) = \|A - \mathbf{x}\mathbf{x}^T\|_F^2$ , it can be shown that within area  $\mathbf{x} \in B_r(\sqrt{\lambda_1}\mathbf{v}_1) = \{\mathbf{y} \mid \|\mathbf{y} - \sqrt{\lambda_1}\mathbf{v}_1\| \leq r\}$ ,  $r = O(\sqrt{\lambda_1} - \frac{\lambda_2}{\sqrt{\lambda_1}})$ ,  $f(\mathbf{x})$  is strongly convex.

**Proof of A.2.** Notice for the objective function  $f$ ,  $\nabla f(\mathbf{x}) = -4(A\mathbf{x} - \|\mathbf{x}\|^2\mathbf{x})$ , Hessian matrix  $H(\mathbf{x}) = -4(A - \|\mathbf{x}\|^2I - 2\mathbf{x}\mathbf{x}^T)$ , and its stationary points are  $\mathbf{x}_i = \sqrt{\lambda_i}\mathbf{v}_i$ . Denote the eigenvalues  $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_r \geq 0 > \dots > \lambda_n$ , and with the assumption that the dominant eigenvalue is positive, we have  $\lambda_1 > |\lambda_n|$ .

At point  $\sqrt{\lambda_1}\mathbf{v}_1$ , the Hessian matrix of  $f$  is positive definite:

$$\begin{aligned} H(\sqrt{\lambda_1}\mathbf{v}_1) &= -4(A - \lambda_1 I - 2\lambda_1\mathbf{v}_1\mathbf{v}_1^T) \\ &= 4\lambda_1\mathbf{v}_1\mathbf{v}_1^T + 4\lambda_1 I - 4 \sum_{i=2}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T \end{aligned}$$

Therefore,  $H$  has the same eigenvectors as  $A$ :  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ , with respect to eigenvalues  $8\lambda_1, 4(\lambda_1 - \lambda_2), 4(\lambda_1 - \lambda_3), \dots, 4(\lambda_1 - \lambda_n)$ , which indicates that  $H$  is positive definite with its smallest eigenvector  $4(\lambda_1 - \lambda_2) > 0$ .

Now to show  $f$  is strongly convex within the neighborhood  $B_r(\sqrt{\lambda_1}\mathbf{v}_1)$ , we denote  $\mathbf{x} = \sqrt{\lambda_1}\mathbf{v}_1 + \mathbf{h}$ ,  $\|\mathbf{h}\| \leq r$ , and introduce

$$G(\mathbf{h}, \mathbf{g}) \stackrel{\text{def}}{=} \frac{\mathbf{g}^T H(\sqrt{\lambda_1}\mathbf{v}_1 + \mathbf{h}) \mathbf{g}}{\mathbf{g}^T \mathbf{g}}$$

which could represent the range of eigenvalues to  $H(\sqrt{\lambda_1}\mathbf{v}_1 + \mathbf{h})$ . Notice

$$\begin{aligned} \nabla_{\mathbf{h}} G(\mathbf{h}, \mathbf{g}) &= 8\sqrt{\lambda_1}\mathbf{v}_1^T + 8\mathbf{h} + 16(\sqrt{\lambda_1}\mathbf{v}_1^T \frac{\mathbf{g}}{\|\mathbf{g}\|} + \mathbf{h}^T \frac{\mathbf{g}}{\|\mathbf{g}\|}) \frac{\mathbf{g}}{\|\mathbf{g}\|} \\ , \text{ and } \|\nabla_{\mathbf{h}} G(\mathbf{h}, \mathbf{g})\| &\leq 8\sqrt{\lambda_1} + 8\|\mathbf{h}\| + 16(\sqrt{\lambda_1} + \|\mathbf{h}\|) \\ &= 24(\sqrt{\lambda_1} + \|\mathbf{h}\|) \\ &\leq 24(\sqrt{\lambda_1} + r), \forall \mathbf{h} \in B_r(0) \end{aligned}$$

By mean-value theorem,

$$\begin{aligned} |G(\mathbf{h}, \mathbf{g}) - G(0, \mathbf{g})| &\leq \left( \sup_{\mathbf{h} \in B_r(0)} \|\nabla_{\mathbf{h}} G(\mathbf{h}, \mathbf{g})\| \right) \|\mathbf{h}\| \\ &\leq 24(\sqrt{\lambda_1} + r)r, \forall \mathbf{h} \in B_r(0), \forall \mathbf{g} \in \mathbb{R}^n \end{aligned}$$

With some proper relaxation, when  $r = \frac{1}{30} \frac{\lambda_1 - \lambda_2}{\sqrt{\lambda_1}}$ , we have  $|G(\mathbf{h}, \mathbf{g}) - G(0, \mathbf{g})| \leq \lambda_1 - \lambda_2$ .

Recall  $G(0, \mathbf{g}) = \frac{\mathbf{g}^T H(\sqrt{\lambda_1}\mathbf{v}_1) \mathbf{g}}{\|\mathbf{g}\|^2} \geq 4(\lambda_1 - \lambda_2)$ ,  $\forall \mathbf{g} \in \mathbb{R}^n$ .

$$\begin{aligned} G(\mathbf{h}, \mathbf{g}) &\geq G(0, \mathbf{g}) - |G(\mathbf{h}, \mathbf{g}) - G(0, \mathbf{g})| \\ &\geq 3(\lambda_1 - \lambda_2), \forall \mathbf{g} \in \mathbb{R}^n, \|\mathbf{h}\| < r, \\ &\text{i.e.} \\ H(\sqrt{\lambda_1}\mathbf{v}_1 + \mathbf{h}) &\succeq 3(\lambda_1 - \lambda_2), \forall \mathbf{h}, \|\mathbf{h}\| \leq \frac{1}{30} \frac{\lambda_1 - \lambda_2}{\sqrt{\lambda_1}} \end{aligned}$$

Therefore the cost function is  $3(\lambda_1 - \lambda_2)$ -strongly convex within the area  $\mathbf{x} \in B_r(\sqrt{\lambda_1}\mathbf{v}_1)$ .  $\square$

**Lemma A.3.** In area  $B_r(\sqrt{\lambda_1}\mathbf{v}_1)$ , where  $r = \frac{\lambda_1 - \lambda_2}{30\sqrt{\lambda_1}}$ ,  $\nabla_i f$  satisfies coordinate-wise Lipschitz continuous with parameter  $L \leq 14\lambda_1 - 2\lambda_2 + 4\max_i |a_{ii}|$ .

**Proof of Lemma A.3:** Our goal is to find  $L$  that satisfies  $|\nabla_i f(\mathbf{x} + \alpha \mathbf{e}_i) - \nabla_i f(\mathbf{x})| \leq L|\alpha|$ ,  $\forall \mathbf{x}, \alpha$  s.t.  $\mathbf{x}, \mathbf{x} + \alpha \mathbf{e}_i \in B_r(\sqrt{\lambda_1}\mathbf{v}_1)$ .

Notice that  $r = \frac{\lambda_1 - \lambda_2}{30\sqrt{\lambda_1}}$ , and  $\|\mathbf{x}\| \leq \sqrt{\lambda_1} + r$ ,  $|\alpha| \leq 2r$ .

Now

$$\begin{aligned}
& |\nabla_i f(\mathbf{x} + \alpha \mathbf{e}_i) - \nabla_i f(\mathbf{x})| \\
&= 4\|\mathbf{x} + \alpha \mathbf{e}_i\|^2(x_i + \alpha) - a_{ii}\alpha - \|\mathbf{x}\|^2 x_i \\
&\leq 4\|\mathbf{x} + \alpha \mathbf{e}_i\|^2 \alpha + \alpha^2 x_i + 2\alpha x_i^2 + 4|a_{ii}\alpha| \\
&\leq 4|\alpha|((\sqrt{\lambda_1} + r)^2 + 2r(\sqrt{\lambda_1} + r) + 2(\sqrt{\lambda_1} + r)^2) + 4|a_{ii}\alpha| \\
&= 4|\alpha|(3\lambda_1 + 10\sqrt{\lambda_1}r + 5r^2) + 4|a_{ii}\alpha| \\
&\leq [12\lambda_1 + 2(\lambda_1 - \lambda_2) + 4|a_{ii}|]|\alpha|
\end{aligned}$$

□

Remark:  $L = 14\lambda_1 - 2\lambda_2 + 4\max_i |a_{ii}|$ , for real application like social network,  $a_{ii} = 0$  and  $L = 14\lambda_1 - 2\lambda_2$ .

With the Lipschitz continuous and strongly convex properties, we show convergence by quoting the result of [13]:

**Lemma A.4.** When  $f$  is strongly convex as  $\nabla^2 f \succeq \mu I$ , and  $\nabla f$  satisfies coordinate-wise  $L$ -Lipschitz continuous, meaning

$$|\nabla_i f(\mathbf{x} + \alpha \mathbf{e}_i) - \nabla_i f(\mathbf{x})| \leq L|\alpha|,$$

$\forall i = 1, 2, \dots, n, \forall \mathbf{x} \in \text{convex set } \mathbb{S}$ , and  $\forall \alpha$  such that  $\mathbf{x} + \alpha \mathbf{e}_i \in \mathbb{S}$ . Then with Gauss-Southwell rule the optimization on  $f$  satisfies linear convergence:

$$f(\mathbf{x}^{(l+1)}) - f(\mathbf{x}^*) \leq (1 - \frac{\mu_1}{L})[f(\mathbf{x}^{(l)}) - f(\mathbf{x}^*)]. \quad (16)$$

Here  $\mu_1 = \inf_{\mathbf{x}, \mathbf{y} \in \mathbb{S}} \frac{\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_\infty}{\|\mathbf{x} - \mathbf{y}\|_1} \in [\frac{\mu}{n}, \mu]$

Therefore, the convergence rate for updating one coordinate at a time with Gauss-Southwell rule becomes  $(1 - \frac{\mu}{L})^n$ ,  $\mu = \inf_{\mathbf{x}, \mathbf{y}} \frac{\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_\infty}{\|\mathbf{x} - \mathbf{y}\|_1} \in [\frac{3(\lambda_1 - \lambda_2)}{n}, 3(\lambda_1 - \lambda_2)]$ ,  $L = 14\lambda_1 - 2\lambda_2 + 4\max_i |a_{ii}|$ .

## A.5 Greedy Coordinate Descent and Coordinate Selection Rules

For an arbitrary matrix  $A \in \mathbb{R}^{n \times d}$ , we can formulate rank-1 matrix approximation:

$$\arg \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^d} f(\mathbf{x}, \mathbf{y}) = \|A - \mathbf{x}\mathbf{y}^T\|_F^2 \quad (17)$$

Notice that  $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) = 2(\|\mathbf{y}\|^2 \mathbf{x} - A\mathbf{y})$ . When fixing  $\mathbf{y}$ , we obtain the optimal solution of  $\mathbf{x}$  to be  $\mathbf{x} = \frac{A\mathbf{y}}{\|\mathbf{y}\|^2}$  and vice versa,  $\mathbf{y} = \frac{A^T \mathbf{x}}{\|\mathbf{x}\|^2}$ . And for symmetric matrices, this alternating minimization algorithm is exactly power method apart from the normalization constant.

Recall our coordinate-wise power method. At each iteration we only update the coordinates with the largest changes. Nevertheless here we can formally interpret this rule as the well-studied Gauss-Southwell rule [12], where the coordinates that maximize the gradient norm is selected. As  $\nabla_{x_i} f(\mathbf{x}, \mathbf{y}) = 2(\|\mathbf{y}\|^2 x_i - \mathbf{a}_i^T \mathbf{y}) = 2\|\mathbf{y}\|^2(x_i - \frac{\mathbf{a}_i^T \mathbf{y}}{\|\mathbf{y}\|^2})$ , Gauss-Southwell gives the same choice of coordinates as our coordinate-wise power method.

Meanwhile, specifically for quadratic objectives, Gauss-Southwell rule actually select the coordinates based on the decrease in the objective function, leading to optimal updates, i.e.,

$$\Delta f_i := f(\mathbf{x}', \mathbf{y}) - f(\mathbf{x}, \mathbf{y}) = -\|\mathbf{y}\|_2^2(x_i - \frac{\mathbf{a}_i^T \mathbf{y}}{\|\mathbf{y}\|_2^2})^2 = -\frac{(\nabla_{x_i} f)^2}{4\|\mathbf{y}\|^2}$$

where  $\mathbf{x}' = \mathbf{x} + (x'_i - x_i)\mathbf{e}_i$ , and  $x'_i = \frac{\mathbf{a}_i^T \mathbf{y}}{\|\mathbf{y}\|^2}$  is the updated coordinate.

Here we summarize the three coordinate selection rules: **(a)** largest coordinate value change,  $|x'_i - x_i|$ , where  $x'_i$  is the next iterate; **(b)** largest partial gradient (Gauss-Southwell),  $|\nabla_i f(\mathbf{x})|$ ; **(c)** largest

function value decrease,  $|f(\mathbf{x}') - f(\mathbf{x})|$ , where  $\mathbf{x}'_i = \mathbf{x} + (x'_i - x_i)\mathbf{e}_i$ . With the good property of quadratic function Eq. (4), for each alternating minimization step, the three selection rules are equivalent. Therefore now with the aid of the objective function, our coordinate selection strategy in CPM, similar as in (a), is now consistent with the rule (c) with its nature in choosing the most "important" coordinates.

Given the optimization interpretation, the extension of CPM to computing the top- $r$  eigenvectors of a symmetric matrix is straightforward. For the objective function  $f(X, Y) = \|A - XY^T\|_F^2$ , where  $X, Y \in \mathbb{R}^{n \times r}$ , the partial gradient of  $f(X, Y)$  with respect to matrices  $X, Y$  becomes  $2(XY^T Y - AY)$  and  $2(YX^T X - AX)$ . By evaluating the norm in each rows of the gradient, we could select and update row by row for  $X$  and  $Y$  by  $\mathbf{a}_i^T Y(Y^T Y)^{-1}$  and  $\mathbf{a}_i^T X(X^T X)^{-1}$ . Although the algorithm is well-defined and can speedup power method for computing top- $r$  eigenvectors, power method (a.k.a. subspace iteration) is typically not used for computing the dominant  $r$  (especially for large  $r$ ) eigenvectors [16]. Therefore we don't expand the discussion of this direction here.

### A.6 Choice of $k$

The choice of  $k$  could be viewed as choosing the block size for greedy block coordinate descent, which is usually tuned in practice or determined by objective's separable property.

However, it would be better if  $k$  could be prescribed and only depend on  $n$ , as we don't know other properties like  $\frac{\lambda_2}{\lambda_1}$  beforehand. In Corollary 4.2.1 it shows the upper bound of  $k$  ranges from  $6\%n$  to  $50\%n$  when  $\frac{\lambda_2}{\lambda_1}$  ranges from  $10^{-5}$  to  $1 - 10^{-5}$ . Meanwhile, experiments also show that the performance of our algorithms isn't too sensitive to the choice of  $k$ . See Figure 6 a large range of  $k$  guarantees good performances. Thus we chose  $k = \frac{n}{20}$  through out experiments in this paper, which is a theoretically and experimentally favorable choice.

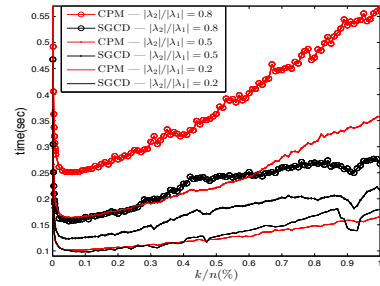


Figure 6: Convergence time with different  $k$  for different  $\lambda_2/\lambda_1$

### A.7 Out-of-core Algorithm

Here we formally present the algorithm for the out-of-core case.

---

#### Algorithm 3 PM, CPM, SGCD for out-of-core matrix $A$

---

- 1: **Initialization:** Separate and save matrix  $A \in \mathbb{R}^{n \times n}$  into  $m$  files, each containing  $n/m$  rows of  $A$  and being able to fit into memory. Initialize random unit vector  $\mathbf{x}^{(0)}$ .
  - 2: **for**  $l = 1$  **to**  $L$  **do**
  - 3:   **for**  $i = 1$  **to**  $m$  **do**
  - 4:     Set  $\Omega = (\frac{(i-1)n}{m} + 1) : \frac{in}{m}$ .
  - 5:     For PM, calculate  $A_{\Omega, :} \mathbf{x}^{(l-1)}$
  - 6:     For CPM, do Step 4 in Algorithm 1 for  $t$  times.
  - 7:     For SGCD, do Step 4 in Algorithm 2 for  $t$  times.
  - 8:   Update  $\mathbf{x}^{(l)}$ .
  - 9: **Output:** Approximate dominant eigenvector  $\mathbf{x}^{(L)}$
- 

### A.8 Extension of Coordinate-wise Mechanism on the Jacobi method

For coordinate-wise Jacobi method for solving  $A\mathbf{x} = \mathbf{b}$ , the algorithm is included here:

And for each iteration, it takes  $O(nnz(R) + n)$  operations for naive Jacobi, and  $O(\frac{k}{n}nnz(R) + n)$  for coordinate-wise Jacobi. This coordinate-wise mechanism also reminds us of Gauss-Seidel method. Recall that Gauss-Seidel:

**Initialize:**  $A = L + U$ , where  $L$  is lower triangular matrix and  $U$  is upper triangular matrix  
**Iterations:**  $\mathbf{x}^+ \leftarrow L^{-1}(\mathbf{b} - U\mathbf{x})$ .

---

**Algorithm 4** Coordinate-wise Jacobi Method

---

- 1: **Input:** Symmetric diagonal dominant matrix  $A \in \mathbb{R}^{n \times n}$ , vector  $\mathbf{b} \in \mathbb{R}^n$ , number of selected coordinates,  $k$ , and number of iterations,  $L$ .
- 2: Initialize  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  and set  $A = D + R$ , where  $D$  is diagonal component of  $A$  and  $R$  is the remainder part.  $\mathbf{z}^{(0)} = R\mathbf{x}^{(0)}$ . Set coordinate selecting criterion  $\mathbf{c}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)} = \mathbf{b} - D\mathbf{x}^{(0)} - \mathbf{z}^{(0)}$ .
- 3: **for**  $l = 1$  **to**  $L$  **do**
- 4: Let  $\Omega^{(l)}$  be a set containing  $k$  coordinates of  $\mathbf{c}^{(l-1)}$  with the largest magnitude. Execute the following updates:

$$\begin{aligned} x_j^{(l)} &= \begin{cases} (b_j - z_j^{(l-1)})/D_{jj}, & j \in \Omega^{(l)} \\ x_j^{(l-1)}, & j \notin \Omega^{(l)} \end{cases} \\ \mathbf{z}^{(l)} &= \mathbf{z}^{(l-1)} + R(\mathbf{x}_{\Omega^{(l)}}^{(l)} - \mathbf{x}_{\Omega^{(l)}}^{(l-1)}) \\ \mathbf{c}^{(l)} &= \mathbf{b} - D\mathbf{x}^{(l)} - \mathbf{z}^{(l)} \end{aligned}$$

- 5: **Output:**  $\mathbf{x}^{(L)}$
- 

And taking advantage of triangular form, the procedure could be simplified as the following version,

$$x_i^{(l+1)} \leftarrow \frac{1}{a_{ii}} (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(l+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(l)})$$

which is very similar to Jacobi method, but uses a forward substitution on newly computed  $x_i$ . Therefore our method is also like a greedy block version of Gauss-Seidel method. While Gauss-Seidel is like a cyclic coordinate version of our method.

We use Matlab to do some simple experiments on some synthetic data to measure the convergence time until the error is less than  $1e-5$ . Here error is measured by  $A$ -quadratic norm between current iteration  $\mathbf{x}^{(l)}$  from ground truth  $\mathbf{x}^*$ , namely,  $\|\mathbf{x}^{(l)} - \mathbf{x}^*\|_A = \sqrt{(\mathbf{x}^{(l)} - \mathbf{x}^*)^T A (\mathbf{x}^{(l)} - \mathbf{x}^*)}$ .

Table 2: Comparison between Jacobi method and Coordinate-wise Jacobi method. N/A denotes the algorithm doesn't converge.

Dataset	n	$\frac{\lambda_2}{\lambda_1}(A)$	$\sigma(D^{-1}R)$	Flops(/ $n^2$ )			Speedup	
				Jacobi	C-Jacobi	Gauss-Seidel	on Jacobi	on G-S
1	1000	0.7803	0.6870	35.035	<b>4.794</b>	7.007	<b>7.308</b>	<b>1.462</b>
2	1000	0.5565	0.9524	254.254	<b>4.284</b>	9.009	<b>59.350</b>	<b>2.103</b>
3	1000	0.5224	0.9942	2115.113	<b>4.488</b>	9.009	<b>471.282</b>	<b>2.007</b>
4	1000	0.5206	0.9986	8505.50	<b>4.08</b>	9.009	<b>2084.68</b>	<b>2.2081</b>
5	1000	0.495	<b>1.11</b>	N/A	<b>4.386</b>	9.009	N/A	<b>2.054</b>
6	5000	0.7792	0.6948	40.01	<b>5.321</b>	8.002	<b>7.519</b>	<b>1.504</b>
7	5000	0.5443	0.9529	290.058	<b>4.317</b>	9.002	<b>67.187</b>	<b>2.085</b>
8	5000	0.5146	0.9949	2703.54	<b>5.622</b>	10.002	<b>480.852</b>	<b>1.779</b>
9	5000	0.5111	0.9992	19760.0	<b>6.256</b>	10.002	<b>3158.76</b>	<b>1.599</b>
10	5000	0.5063	<b>1.02</b>	N/A	<b>6.256</b>	10.002	N/A	<b>1.599</b>

And from Table A.8, we can see that coordinate-wise Jacobi method shows significant speedup over the naive Jacobi method. And even when the matrix is no longer diagonal dominant, (see table when  $\sigma(D^{-1}R) > 1$ ), but still positive definite, coordinate-wise Jacobi method still converges. And this trait meets the convergence requirement for Gauss-Seidel method.

Although in this comparison coordinate-wise Jacobi doesn't beat up Gauss-Seidel that much, Gauss-Seidel has the disadvantage that it can not be done in parallel, while our method could be more flexible on that. For example, we could greedily update coordinates in each worker, rather than choosing globally the most greedy coordinates.

However, since for symmetric diagonal dominant matrices, Jacobi or Gauss-Seidel is not the state-of-the-art method for solving linear system, we will need to compare with other more powerful methods. And this algorithm lacks theoretical support at this point, so we consider this as an expansion of our current work on coordinate-wise power method. But still, it's worth mentioning that the coordinate-wise mechanism could be powerful applying to Jacobi method and maybe to other iterative methods in linear algebra too. Therefore in the future, we may continue exploiting the theory behind, and analyze why and how greediness impacts on Jacobi method or other iterative methods in linear algebra.