# Efficient coordinate-wise leading eigenvector computation

Jialei Wang\* JIALEI@UCHICAGO.EDU

University of Chicago

Weiran Wang\* Weiranw@amazon.com

Amazon Alexa

Dan Garber Dangar@technion.ac.il

Technion - Israel Institute of Technology

Nathan Srebro NATI@TTIC.EDU

Toyota Technological Institute at Chicago

Editor:

#### Abstract

We develop and analyze efficient "coordinate-wise" methods for finding the leading eigenvector, where each step involves only a vector-vector product. We establish global convergence with overall runtime guarantees that are at least as good as Lanczos's method and dominate it for slowly decaying spectrum. Our methods are based on combining a shift-and-invert approach with coordinate-wise algorithms for linear regression.

#### 1. Introduction

Extracting the top eigenvalues/eigenvectors of a large symmetric matrix is a fundamental step in various problems including many machine learning applications. One prominent example of this problem is principal component analysis (PCA), in which we extract the top eigenvectors of the data covariance matrix, and there has been continuous effort in developing efficient stochastic/randomized algorithms for large-scale PCA (e.g., Garber et al., 2016; Shamir, 2015; Allen-Zhu and Li, 2016). The more general eigenvalue problems for large matrices without the covariance structure is relatively less studied. The method of choice for this problem has been the power method, or the faster but often relatively less known Lanczos algorithm (Golub and van Loan, 1996), which are based on iteratively computing matrix-vector multiplications with the input matrix until the component of the vector that lies in the tailing eigenspace vanishes. However, for very large-scale and dense matrices, even computing a single matrix-vector product is expensive.

An alternative is to consider much cheaper vector-vector products, i.e., instead of updating all the entries in the vector on each iteration by a full matrix-vector product, we consider the possibility of only updating one coordinate, by only computing the inner product of single row of the matrix with the vector. Such operations do not even require storing the entire matrix in memory. Intuitively, this may result in an overall significant speedup, in certain likely scenarios, since certain coordinates in the matrix-vector product are more valuable than others for making local progress towards converging to the leading eigenvector.

<sup>\*.</sup> Equal Contribution

<sup>©2017</sup> Jialei Wang, Weiran Wang, Dan Garber, and Nathan Srebro.

Indeed, this is the precise rational behind coordinate-descent methods that were extensively studied and are widely applied to convex optimization problems, see Wright (2015) for a comprehensive survey. Thus, given the structure of the eigenvalue problem which is extremely suitable for coordinate-wise updates, and the celebrated success of coordinate-descent methods for convex optimization, a natural question is whether such updates can be applied for eigenvector computation with provable global convergence guarantees, despite the inherent non-convexity of the problem.

Recently, Lei et al. (2016) have proposed two such methods, Coordinate-wise Power Method (CPM) and Symmetric Greedy Coordinate Descent (SGCD). Both methods update on each iteration only k entries in the vector, for some fixed k. CPM updates on each iteration the k coordinates that would change the most under one step of power method, while SGCD applies a greedy heuristic for choosing the coordinates to be updated. The authors show that CPM enjoys a global convergence rate, with rate similar to the classical power iterations algorithm provided that k, the number of coordinates to be updated on each iteration, is sufficiently large (or equivalently, the "noise" outside the k selected coordinate is sufficiently small). In principle, this might force k to be as large as the dimension d, and indeed in their experiments they set k to grow linearly with d, which is overall not significantly faster than standard power iterations, and does not truly capture the concept of coordinate updates proved useful to convex problems. The second algorithm proposed in Lei et al. (2016), SGCD, is shown to converge locally with a linear rate already for k=1(i.e., only a single coordinate is updated on each iteration), however this results assume that the method is initialized with a vector that is already sufficiently close (in a non-trivial way) to the leading eigenvector. The dependence of CPM and SGCD on the inverse relative eigengap<sup>1</sup> of the input matrix is similar to that of the power iterations algorithm, i.e. linear dependence, which in principal is suboptimal, since square-root dependence can be obtained by methods such as the Lanczos algorithm.

We present globally-convergent coordinate-wise algorithms for the leading eigenvector problem which resolves the abovementioned concerns and significantly improve over previous algorithms. Our algorithms update only a single entry at each step and enjoy linear convergence. Furthermore, for a particular variant, the convergence rate depends only on the square-root of the inverse relative eigengap, yielding a total runtime that dominates that of the standard power method and competes with that of the Lanczos's method. In Section 2, we discuss the basis of our algorithm, the shift-and-invert power method (Garber and Hazan, 2015; Garber et al., 2016), which transforms the eigenvalue problem into a series of convex least squares problems. In Section 3, we show the least squares problems can be solved using efficient coordinate descent methods that are well-studied for convex problems. This allow us to make use of principled coordinate selection rules for each update, all of which have established convergence guarantees.

We provide a summary of the time complexities of different globally-convergent methods in Table 1. In particular, it is observable that in cases where either the spectrum of the input matrix or the magnitude of its diagonal entries is slowly decaying, our methods can yield provable and significant improvements over previous methods. For example, for a spiked covariance model whose eigenvalues are  $\rho_1 > \rho_1 - \Delta = \rho_2 = \rho_3 = \dots$ , our algorithm

<sup>1.</sup> Defined as  $\frac{\rho_1(\mathbf{A})}{\rho_1(\mathbf{A})-\rho_2(\mathbf{A})}$  for a positive semidefinite  $\mathbf{A}$ , where  $\rho_i(\mathbf{A})$  is the *i*-th largest eigenvalue of  $\mathbf{A}$ .

Table 1: Time complexity (total number of coordinate updates) for finding an estimate  $\mathbf{w}$  satisfying  $(\mathbf{w}^{\top}\mathbf{p}_1)^2 \geq 1 - \epsilon$  with at least constant probability, where  $\mathbf{p}_1$  is the leading eigenvector of a positive semidefinite matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , with eigenvalues  $\rho_1, \rho_2, \ldots$  in descending order. Since all methods are randomized, we assume all are initialized with a random unit vector.

Method	Time complexity
Power method	$\mathcal{O}\left(\frac{\rho_1}{\rho_1-\rho_2}\cdot d^2\cdot\log\frac{1}{ \mathbf{w}_0^\top\mathbf{p}_1 \epsilon}\right)$
Lanczos	$\mathcal{O}\left(\sqrt{\frac{\rho_1}{\rho_1-\rho_2}}\cdot d^2\cdot\log\frac{1}{ \mathbf{w}_0^{\top}\mathbf{p}_1 \epsilon}\right)$
Ours (SI-GSL)	$\mathcal{O}\left(\frac{ ho_1 - \frac{1}{d}\sum_{i=1}^d  ho_i}{ ho_1 -  ho_2} \cdot d^2 \cdot \log \frac{1}{ \mathbf{w}_0^{T}\mathbf{p}_1 \epsilon}\right)$
Ours (SI-ACDM)	$\mathcal{O}\left(\frac{\frac{1}{d}\sum_{i=1}^{d}\sqrt{\rho_{1}-\mathbf{A}_{ii}}}{\sqrt{\rho_{1}-\rho_{2}}}\cdot d^{2}\cdot\log\frac{1}{ \mathbf{w}_{0}^{\top}\mathbf{p}_{1} \epsilon}\right)$

(SI-GSL) can have a runtime independent of the eigengap  $\Delta$ , while the runtime of Lanczos's method depends on  $\frac{1}{\sqrt{\Delta}}$ . We also verify this intuition empirically via numerical experiments.

**Notations** We use boldface uppercase letters (e.g., **A**) to denote matrices, and boldface lowercase letters (e.g., **x**) to denote vectors. For a positive definite matrix **M**, the vector norm  $\|\cdot\|_{\mathbf{M}}$  is defined as  $\|\mathbf{w}\|_{\mathbf{M}} = \sqrt{\mathbf{w}^{\top}\mathbf{M}\mathbf{w}} = \|\mathbf{M}^{\frac{1}{2}}\mathbf{w}\|$  for any **w**. We use  $\mathbf{A}[ij]$  to denote the element in row i and column j of the matrix **A**, and use  $\mathbf{x}[i]$  to denote the i-th element of the vector **x** unless stated otherwise. Additionally,  $\mathbf{A}[i:]$  and  $\mathbf{A}[:j]$  denote the i-th row and j-th column of **A** respectively.

**Problem formulation** We consider the task of extracting the top eigenvector of a symmetric positive definite matrix  $\mathbf{A} \in \mathbb{R}^d$ . Let the complete set of eigenvalues of  $\mathbf{A}$  be  $\rho_1 \geq \rho_2 \geq \cdots \geq \rho_d \geq 0$ , with corresponding eigenvectors  $\mathbf{p}_1, \ldots, \mathbf{p}_d$  which form an orthonormal basis of  $\mathbb{R}^d$ . Without loss of generality, we assume  $\rho_1 \leq 1$  (which can always be obtained by rescaling the matrix). Furthermore, we assume the existence of a positive eigenvalue gap  $\Delta := \rho_1 - \rho_2 > 0$  so that the top eigenvector is unique up to scaling.

#### 2. Shift-and-invert power method

In this section, we introduce the shift-and-invert approach to the eigenvalue problem and review its analysis, which will be the basis for our algorithms. The most popular iterative algorithm for the leading eigenvalue problem is the power method, which iteratively performs the following matrix-vector multiplications and normalization steps

$$\tilde{\mathbf{w}}_t \leftarrow \mathbf{A}\mathbf{w}_{t-1}, \quad \mathbf{w}_t \leftarrow \frac{\tilde{\mathbf{w}}_t}{\|\tilde{\mathbf{w}}_t\|}, \quad \text{for} \quad t = 1, \dots.$$

It can be shown that the iterates become increasingly aligned with the eigenvector corresponding to the largest eigenvalue in magnitude, and the number of iterations needed to achieve  $\epsilon$ -suboptimality in alignment is  $\mathcal{O}\left(\frac{\rho_1}{\Delta}\log\frac{1}{|\mathbf{w}_0^\top\mathbf{p}_1|\epsilon}\right)$  (Golub and van Loan, 1996)<sup>2</sup>.

<sup>2.</sup> We can always guarantee that  $|\mathbf{w}_0^{\mathsf{T}}\mathbf{p}_1| = \Omega(1/\sqrt{d})$  by taking  $\mathbf{w}_0$  to be a random vector on unit sphere.

We see that the computational complexity depends linearly on  $\frac{1}{\Delta}$ , and thus power method converges slowly if the gap is small.

Shift-and-invert (Golub and van Loan, 1996; Garber and Hazan, 2015; Garber et al., 2016) can be viewed as a pre-conditioning approach which improves the dependence of the time complexity on the eigenvalue gap. The main idea behind this approach is that, instead of running power method on  $\bf A$  directly, we can equivalently run power method on the matrix  $(\lambda \bf I - A)^{-1}$  where  $\lambda > \rho_1$  is a shifting parameter. Observe that  $(\lambda \bf I - A)^{-1}$  has exactly the same set of eigenvectors as  $\bf A$ , and its eigenvalues are

$$\beta_1 \ge \beta_2 \ge \dots \ge \beta_d > 0$$
, where  $\beta_i = \frac{1}{\lambda - \rho_i}$ .

If we have access to a  $\lambda$  that is slightly larger than  $\rho_1$ , and in particular if  $\lambda - \rho_1 = \mathcal{O}(1) \cdot \Delta$ , then the inverse relative eigengap of  $(\lambda \mathbf{I} - \mathbf{A})^{-1}$  is  $\frac{\beta_1}{\beta_1 - \beta_2} = O(1)$ , which means that the power method, applied to this shift and inverted matrix, will converge to the top eigenvector in only a poly-logarithmic number of iterations, in particular, without linear dependence on  $1/\Delta$ .

In shift-and-invert power method, the matrix-vector multiplications have the form  $\tilde{\mathbf{w}}_t \leftarrow (\lambda \mathbf{I} - \mathbf{A})^{-1} \mathbf{w}_{t-1}$ , which is equivalent to solving the convex least squares problem

$$\tilde{\mathbf{w}}_t \leftarrow \underset{\mathbf{w}}{\operatorname{arg\,min}} \ \frac{1}{2} \mathbf{w}^\top (\lambda \mathbf{I} - \mathbf{A}) \mathbf{w} - \mathbf{w}_{t-1}^\top \mathbf{w}. \tag{1}$$

Solving such least squares problems exactly could be costly if d is large. Fortunately, power method with approximate matrix-vector multiplications still converges, provided that the errors in each step is controlled; analysis of inexact power method together with applications to PCA and CCA can be found in Hardt and Price (2014); Garber et al. (2016); Ge et al. (2016); Wang et al. (2016).

We give the inexact shift-and-invert preconditioning power method in Algorithm 1, which consists of two phases. In Phase I, starting from an lower estimate of  $\tilde{\Delta}$  and an upper bound of  $\rho_1$ , namely  $\lambda_{(0)} = 1 + \tilde{\Delta}$  (recall that by assumption,  $\rho_1 \leq 1$ ), the **repeatuntil** loop locates an estimate of the eigenvalue  $\lambda_{(s)}$  such that  $0 < \lambda_{(s)} - \rho_1 = \Theta(1) \cdot \Delta$ , and it does so by estimating the top eigenvalue  $\beta_1$  of  $(\lambda \mathbf{I} - \mathbf{A})^{-1}$  (through a small number of power iterations in the **for** loop), and shrinking the gap between  $\lambda_{(s)}$  and  $\rho_1$ . In Phase II, the algorithm fixes the shift-and-invert parameter  $\lambda_{(s)}$ , and runs many power iterations in the last **for** loop to achieve an accurate estimate of  $\mathbf{p}_1$ .

We now provide convergence analysis of Algorithm 1 following that of Garber and Hazan (2015) closely, but improving their rate (removing an extra  $\log \frac{1}{\epsilon}$  factor) with the warm-start strategy for least squares problems by Garber et al. (2016), and making the initial versus final error ratio explicit in the exposition. We discuss the least-squares solver in Algorithm 1 in the next section.

Measure of progress Since  $\{\mathbf{p}_1, \dots, \mathbf{p}_d\}$  form an orthonormal basis of  $\mathbb{R}^d$ , we can write each normalized iterate as a linear combination of the eigenvectors as

$$\mathbf{w}_t = \sum_{i=1}^d \xi_{ti} \mathbf{p}_i, \text{ where } \xi_{ti} = \mathbf{w}_t^{\top} \mathbf{p}_i, \text{ for } i = 1, \dots, d,$$

#### **Algorithm 1** The shift-and-invert meta-algorithm for extracting top eigenvector of **A**.

```
Input: Data matrix A, an iterative least squares optimizer (LSO), a lower estimate \Delta for
    \Delta := \rho_1 - \rho_2 such that \tilde{\Delta} \in [c_1 \Delta, c_2 \Delta], and \underline{\sigma} = 1 + \frac{1 - c_2}{c_2} \tilde{\Delta}.
    Initialize \tilde{\mathbf{w}}_0 \in \mathbb{R}^d, \mathbf{w}_0 \leftarrow \frac{\tilde{\mathbf{w}}_0}{\|\tilde{\mathbf{w}}_0\|}
     // Phase I: locate a \lambda = \rho_1 + \mathcal{O}(1) \cdot \Delta
                              \lambda_{(0)} \leftarrow 1 + \Delta
    s \leftarrow 0,
    repeat
          s \leftarrow s + 1
          // Power method on (\lambda_{(s)}\mathbf{I} - \mathbf{A})^{-1} in crude regime
          for t = (s-1)m_1 + 1, \dots, sm_1 do
               Apply LSO to the problem \min_{\mathbf{w}} f_t(\mathbf{w}) := \frac{1}{2} \mathbf{w}^{\top} (\lambda_{(s-1)} \mathbf{I} - \mathbf{A}) \mathbf{w} - \mathbf{w}_{t-1}^{\top} \mathbf{w} from
               initialization \frac{\mathbf{w}_{t-1}}{\mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w}_{t-1}}, and output \tilde{\mathbf{w}}_t satisfying f_t(\tilde{\mathbf{w}}_t) \leq \min_{\mathbf{w}} f_t(\mathbf{w}) + \epsilon_t.
               Normalization: \mathbf{w}_t \leftarrow \frac{\tilde{\mathbf{w}}_t}{\|\tilde{\mathbf{w}}_t\|}
          end for
          // Estimate the top eigenvalue of (\lambda_{(s)}\mathbf{I} - \mathbf{A})^{-1}
         Apply LSO to the problem \min_{\mathbf{u}} l_s(\mathbf{u}) := \frac{1}{2} \mathbf{u}^{\top} (\lambda_{(s)} \mathbf{I} - \mathbf{A}) \mathbf{u} - \mathbf{w}_{sm_1}^{\top} \mathbf{u} from initialization \frac{\mathbf{w}_{sm_1}}{\mathbf{w}_{sm_1}^{\top} (\lambda \mathbf{I} - \mathbf{A}) \mathbf{w}_{sm_1}}, and output \mathbf{u}_s satisfying l_s(\mathbf{u}_s) \leq \min_{\mathbf{u}} l_s(\mathbf{u}) + \tilde{\epsilon}_s.
          Update: \Delta_s \leftarrow \frac{1}{2} \cdot \frac{1}{\mathbf{w}_{sm}^{\top} \mathbf{u}_s - \underline{\sigma}/8}, \ \lambda_{(s)} \leftarrow \lambda_{(s-1)} - \frac{\Delta_s}{2}
     until \Delta_{(s)} \leq \Delta
     \lambda_{(f)} \leftarrow \lambda_{(s)}
    // Power method on (\lambda_{(f)}\mathbf{I} - \mathbf{A})^{-1} in accurate regime
     for t = sm_1 + 1, \dots, sm_1 + m_2 do
          Apply LSO to the problem \min_{\mathbf{w}} f_t(\mathbf{w}) := \frac{1}{2} \mathbf{w}^{\top} (\lambda_{(f)} \mathbf{I} - \mathbf{A}) \mathbf{w} - \mathbf{w}_{t-1}^{\top} \mathbf{w} with initial-
         ization \frac{\mathbf{w}_{t-1}}{\mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w}_{t-1}}, and output an approximate solution \tilde{\mathbf{w}}_t satisfying f_t(\tilde{\mathbf{w}}_t) \leq
          \min_{\mathbf{w}} f_t(\mathbf{w}) + \epsilon_t.
          Normalization: \mathbf{w}_t \leftarrow \frac{\tilde{\mathbf{w}}_t}{\|\tilde{\mathbf{w}}_t\|}
     end for
```

and  $\sum_{i=1}^{d} \xi_{ti}^{2} = 1$ . Our goal is to have high alignment between the estimate  $\mathbf{w}_{t}$  and  $\mathbf{p}_{1}$ , i.e.,  $\xi_{t1} = \mathbf{w}_{t}^{\top} \mathbf{p}_{1} \geq 1 - \epsilon$  for  $\epsilon \in (0, 1)$ . Equivalently, we would like to have  $\rho_{1} - \mathbf{w}_{t}^{\top} \mathbf{A} \mathbf{w}_{t} \leq \rho_{1} \epsilon$  (see Lemma 5 in Appendix A).

### 2.1 Iteration complexity of inexact power method

**Output:**  $\mathbf{w}_{sm_1+m_2}$  is the approximate eigenvector.

Consider the inexact shift-and-invert power iterations:  $\tilde{\mathbf{w}}_t \approx \arg\min_{\mathbf{w}} f_t(\mathbf{w}) = \frac{1}{2}\mathbf{w}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w} - \mathbf{w}_{t-1}^{\top}\mathbf{w}$ ,  $\mathbf{w}_t \leftarrow \frac{\tilde{\mathbf{w}}_t}{\|\tilde{\mathbf{w}}_t\|}$ , where  $f_t(\tilde{\mathbf{w}}_t) \leq \min_{\mathbf{w}} f_t(\mathbf{w}) + \epsilon_t$ , for  $t = 1, \ldots$  We can divide the convergence behavior of this method into two regimes: in the crude regime, our goal is to estimate the top eigenvalue (as in Phase I of Algorithm 1) regardless of the existence of an eigengap; whereas in the accurate regime, our goal is to estimate the top eigenvector (as in Phase II of Algorithm 1).

Following Garber and Hazan (2015); Garber et al. (2016), we provide convergence guarantees for inexact shift-and-invert power iterations in Lemma 7 (Appendix B), quantifying the sufficient accuracy  $\epsilon_t$  in solving each least squares problem for the overall method to converge. The iteration complexity for the crude regime is independent of eigengap, while in the accurate regime the rate in which the error decreases does depend on the eigengap of  $(\lambda \mathbf{I} - \mathbf{A})^{-1}$ .

## 2.2 Bounding initial error for each least squares

For each least squares problem  $f_t(\mathbf{w})$  in inexact shift-and-invert power method, one can show that the optimal initialization based on the previous iterate is (Garber et al., 2016)

$$\mathbf{w}_t^{init} = \frac{\mathbf{w}_{t-1}}{\mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w}_{t-1}}.$$

We can bound the suboptimality of this initialization, defined as  $\epsilon_t^{init} := f_t(\mathbf{w}_t^{init}) - \min_{\mathbf{w}} f_t(\mathbf{w})$ . See full analysis in Appendix C.

**Lemma 1.** Initialize each least squares problem  $\min_{\mathbf{w}} f_t(\mathbf{w})$  in inexact shift-and-invert power method from  $\mathbf{w}_t^{init} = \frac{\mathbf{w}_{t-1}}{\mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w}_{t-1}}$ . Then the initial suboptimality  $\epsilon_t^{init}$  in  $f_t(\mathbf{w})$  can be bounded by the necessary final suboptimality  $\epsilon_t$  as follows.

- can be bounded by the necessary final suboptimality  $\epsilon_t$  as follows.

  1. In the crude regime,  $\epsilon_t^{init} \leq \frac{64\beta_1^2}{\epsilon^2 \beta_d^2} \left( \frac{(2\beta_1/\beta_d)^{T_1} 1}{(2\beta_1/\beta_d) 1} \right)^2 \cdot \epsilon_t$  where  $T_1$  is the total number of iterations used (see precise definition in Lemma 7).
- 2. In the accurate regime,  $\epsilon_t^{init} \leq \max(G_0, 1) \cdot \frac{16\beta_1^2}{(\beta_1 \beta_2)^2} \cdot \epsilon_t$  where  $G_0$  is the initial condition for shift-and-invert power iterations (see precise definition in Lemma 7).

#### 2.3 Iteration complexity of Algorithm 1

Based on the iteration complexity of inexact power method and the warm-start strategy, we derive the total number of iterations (least squares problems) needed by Algorithm 1.

**Lemma 2** (Iteration complexity of the **repeat-until** loop in Algorithm 1). Suppose that  $\tilde{\Delta} \in [c_1 \Delta, c_2 \Delta]$  where  $c_2 \leq 1$ . Set  $m_1 = \left \lceil 8 \log \left( \frac{16}{\xi_{01}^2} \right) \right \rceil$  where  $\xi_{01} = \mathbf{w}_0^{\top} \mathbf{p}_1$ , initialize each least squares problem as in Lemma 1, and maintain the ratio between initial and final error to be  $\frac{\epsilon_t^{init}}{\epsilon_t} = \frac{32 \cdot 10^{2m_1+1}}{\tilde{\Delta}^{2m_1}}$  and  $\frac{\tilde{\epsilon}_s^{init}}{\tilde{\epsilon}_s} = \frac{1024}{\tilde{\Delta}^2}$  throughout the **repeat-until** loop. Then for all  $s \geq 1$  it holds that  $\frac{1}{2}(\lambda_{(s-1)} - \rho_1) \leq \Delta_s \leq \lambda_{(s-1)} - \rho_1$ . Upon exiting this loop, the  $\lambda_{(f)}$  satisfies

$$\rho_1 + \frac{\tilde{\Delta}}{4} \le \lambda_{(f)} \le \rho_1 + \frac{3\tilde{\Delta}}{2},\tag{2}$$

and the number of iterations by **repeat-until** is  $\mathcal{O}\left(\log \frac{1}{\overline{\Delta}}\right)$ .

**Lemma 3** (Iteration complexity of the final **for** loop in Algorithm 1). Suppose that  $\tilde{\Delta} \in [c_1\Delta, c_2\Delta]$  where  $0 < c_1 < c_2 \le 1$ . Set  $m_2 = \left\lceil \frac{1}{2} \log_{\frac{9}{7}} \left( \frac{G_0^2}{\epsilon} \right) \right\rceil$  where  $G_0 := \frac{\sqrt{\sum_{i=2}^d (\lambda_{(f)} - \rho_i) \cdot (\mathbf{w}_0^\top \mathbf{p}_i)^2}}{\sqrt{(\lambda_{(f)} - \rho_1) \cdot (\mathbf{w}_0^\top \mathbf{p}_1)^2}}$ , initialize each least squares problem as in Lemma 1, and maintain the ratio between initial

and final error to be  $\frac{\epsilon_t^{init}}{\epsilon_t} = 100 \max(G_0, 1)$  throughout the final **for** loop. Then the output of Algorithm 1 satisfies  $\mathbf{w}_{sm_1+m_2}^{\top} \mathbf{p}_1 \geq 1 - \epsilon$ .

In the next section, we will be using linearly convergent solvers for the least squares problems, whose runtime depends on  $\log \frac{\epsilon_t^{init}}{\epsilon_t}$ . Lemma 2 implies that we need to solve  $sm_1 = \mathcal{O}\left(\log \frac{1}{\Delta}\right)$  least squares problems in Phase I, each with  $\log \frac{\epsilon_t^{init}}{\epsilon_t} = \mathcal{O}\left(\log \frac{1}{\Delta}\right)$ . And Lemma 3 implies that we need to solve  $m_2 = \mathcal{O}\left(\log \frac{1}{\epsilon}\right)$  least squares problems in Phase II, each with  $\log \frac{\epsilon_t^{init}}{\epsilon_t} = \mathcal{O}(1)$ .

## 3. Coordinate descent for least squares

Different from PCA, for general eigenvalue problems, the matrix **A** may not have the structure of data covariance, and fast stochastic gradient methods for finite-sums (such as SVRG Johnson and Zhang, 2013 used in Garber et al. 2016) does not apply. However, we can instead apply efficient coordinate descent (CD) methods in our setting. In this section, we review the CD methods and study their complexity for solving the least squares problems in Algorithm 1.

There is a rich literature on CD methods and they have attracted resurgent interests recently due to their simplicity and efficiency for big data problems; we refer the readers to Wright (2015) for a comprehensive survey. In each CD update, we pick one coordinate and take a step along the direction of negative coordinate-wise gradient. To state the convergence properties of CD methods, we need the definitions of a few key quantities. Recall that we would like to solve the least squares problems of the form  $\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{x} - \mathbf{y}^{\top}\mathbf{x}$ , with the optimal solution  $\mathbf{x}^* = (\lambda \mathbf{I} - \mathbf{A})^{-1}\mathbf{y}$ .

**Smoothness** The gradient  $\nabla f(\mathbf{x}) = (\lambda \mathbf{I} - \mathbf{A})\mathbf{x} - \mathbf{y}$  is coordinate-wise Lipschitz continuous: for all i = 1, ..., d,  $\mathbf{x} \in \mathbb{R}^d$ , and  $\alpha \in \mathbb{R}$ , we have

$$|\nabla_i f(\mathbf{x} + \alpha \mathbf{e}_i) - \nabla_i f(\mathbf{x})| \le L_i |\alpha|$$
 for  $L_i := \lambda - \mathbf{A}[ii]$ 

where  $\mathbf{e}_i$  is the *i*-th standard basis. Note that for least squares problems, the coordinate-wise gradient Lipschitz constant are the diagonal entries of its Hessian. Denote by  $L_{\max} := \max_{1 \leq i \leq d} L_i$  and  $\bar{L} = \frac{1}{d} \sum_{i=1}^d L_i$  the largest and average Lipschitz constant over all coordinates respectively. These two Lipschitz constants are to be distinguished from the "global" smoothness  $\tilde{L}$ , which satisfies

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \le \tilde{L} \cdot \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y}.$$

Since  $f(\mathbf{x})$  is quadratic,  $\tilde{L} = \sigma_{\text{max}} (\lambda \mathbf{I} - \mathbf{A})$ . Observe that (Wright, 2015)

$$1 \le \tilde{L}/L_{\text{max}} \le d,\tag{3}$$

with upper bound achieved by a Hessian matrix  $\mathbf{1}\mathbf{1}^{\top}$ .

**Strong-convexity** As mentioned before, the matrices  $(\lambda \mathbf{I} - \mathbf{A})$  in our least squares problems are always positive definite, with smallest eigenvalue  $\mu := \lambda - \rho_1 = \mathcal{O}(\Delta)$ . As a result,  $f(\mathbf{x})$  is  $\mu$ -strongly convex with respect to the  $\ell_2$  norm. Another strong convexity parameters

Table 2. Notation quick reference for CD methods.				
Notation	Description	Relevant properties for least squares		
$L_i$	smoothness parameter for <i>i</i> -th coordinate	$L_i := \lambda - \mathbf{A}[ii]$		
$ar{L}$	average smoothness of all coordinates	$\bar{L} = \lambda - \frac{1}{d} \sum_{i=1}^{d} \mathbf{A}[ii] = \lambda - \frac{1}{d} \sum_{i=1}^{d} \rho_d$		
$L_{ m max}$	largest smoothness of all coordinates	$L_{\max} := \max_{1 \le i \le d} L_i$		
$ ilde{L}$	global smoothness	$\tilde{L} = \lambda - \rho_d,  \bar{L} \le L_{\text{max}} \le \tilde{L} \le dL_{\text{max}}$		
$\mu$	strong-convexity parameter in $\ell_2$ -norm	$\mu = \lambda - \rho_1$		
$\mu_L$	strong-convexity parameter in $\ \cdot\ _L$ -norm	$\mu_L \ge \mu/(d\bar{L})$		

Table 2: Notation quick reference for CD methods.

we will need is  $\mu_L$  defined with respect to the norm  $\|\mathbf{z}\|_L = \sum_{i=1}^d \sqrt{L_i} |\mathbf{z}[i]|$ . It is shown by Nutini et al. (2015) that  $\mu_L \geq \mu/(d\bar{L})$ . We collect relevant parameters for CD methods in Table 2.

#### 3.1 Coordinate selection

The choice of coordinate to update is a central research topic for CD methods. We now discuss several coordinate selection rules, along with convergence rates, that are most relevant in our setting.

Gauss-Southwell-Lipschitz (GSL, Nutini et al., 2015): A greedy rule for selecting the coordinate where the coordinate with largest gradient magnitude (relative to the square root of Lipschitz constant) is chosen. The greedy rule tends to work well for high dimensional sparse problems (Dhillon et al., 2011).

Cyclic (Beck and Tetruashvili, 2013): Given an ordering of the coordinates, cyclic coordinate descent processes each coordinate exactly once according to the ordering in each pass. In the extension "essentially cyclic" rule, each coordinate is guaranteed to be chosen at least once in every  $\tau \geq d$  updates. This rule is particularly useful when the data can not fit into memory, as we can iteratively load a fraction of **A** and update the corresponding coordinates, and still enjoy convergence guarantee.

Accelerated randomized coordinate descent methods (ACDM): In each step, a coordinate is selected randomly, according to certain (possibly non-uniform) distribution based on the coordinate-wise Lipschitz constant. In accelerated versions of randomized CD (Nesterov, 2012; Lee and Sidford, 2013; Lu and Xiao, 2015; Allen-Zhu et al., 2016), one can maintain auxiliary sequences of iterates to better approximate the objective function and obtain faster convergence rate, at the cost of (slightly) more involved algorithm. We implement the variant NU-ACDM by Allen-Zhu et al. (2016), which samples each coordinate i with probability proportional to  $\sqrt{L_i}$  and has the fastest convergence rate to date.

We provide the pseudocode of coordinate descent of the above rules in Algorithm 2.<sup>3</sup> Note that the stepsize for the chosen coordinate j is  $\frac{1}{L_j}$ , the inverse Lipschitz constant; for least squares problems where  $f(\mathbf{x})$  is quadratic in each dimension, this stepsize exactly minimizes the function over  $\mathbf{x}[j]$  given the rest coordinates. In some sense, the GSL rule is the "optimal myopic coordinate update" (Nutini et al., 2015).

<sup>3.</sup> The pseudo code for accelerated randomized CD is more involved and we refer the readers to Allen-Zhu et al. (2016).

# Algorithm 2 Coordinate descent for minimizing $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^{\top}(\mathbf{I} - \mathbf{A})\mathbf{x} - \mathbf{y}^{\top}\mathbf{x}$ .

Input: Data  $\mathbf{A} \in \mathbb{R}^{d \times d}$ ,  $\mathbf{b} \in \mathbb{R}^d$ , initialization  $\mathbf{x}_0$ .

Compute gradient  $\mathbf{g} \leftarrow \lambda \mathbf{x} - \mathbf{A} \mathbf{x} - \mathbf{y}$ 

for  $t = 1, 2, ..., \tau$  do

Select a coordinate using one of the rules

GSL: 
$$j \leftarrow \underset{1 \leq i \leq d}{\operatorname{arg \, max}} \frac{|\mathbf{g}[i]|}{\sqrt{\lambda - \mathbf{A}[ii]}}$$

Cyclic:  $j \leftarrow \mod(t, d) + 1$ 

Random:  $j \leftarrow \text{random index } \in [1, d]$ 

Compute update:  $\delta \leftarrow -\frac{\mathbf{g}[j]}{\lambda - \mathbf{A}[jj]}$ Update coordinate:  $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + \delta \cdot \mathbf{e}_j$ 

Update gradient:

$$\mathbf{g} \leftarrow \mathbf{g} - \delta \cdot \mathbf{A}[:j], \quad \mathbf{g}[j] \leftarrow \mathbf{g}[j] + \lambda \delta$$

end for

**Output:**  $\mathbf{x}_{\tau}$  is the approximate solution.

Connection to the greedy selection rule of Lei et al. (2016) Given the current estimate  $\mathbf{x}$ , the coordinate-wise power method of Lei et al. (2016) selects the following coordinate to be updated:  $\arg\max_i |\mathbf{c}[i]|$  where  $\mathbf{c} = \frac{\mathbf{A}\mathbf{x}}{\mathbf{x}^{\top}\mathbf{A}\mathbf{x}} - \mathbf{x}$ , i.e., the coordinate that would be updated the most if a full power iteration were performed. Their selection rule is equivalent to choosing the largest element of  $((\mathbf{x}^{\top}\mathbf{A}\mathbf{x})\mathbf{I} - \mathbf{A})\mathbf{x}$ , where  $\mathbf{x}^{\top}\mathbf{A}\mathbf{x}$  is the current (lower) estimate of the eigenvalue  $\rho_1$ . On the other hand, the greedy Gauss-Southwell rule for our method chooses the largest element of the gradient  $(\lambda \mathbf{I} - \mathbf{A})\mathbf{x} - \mathbf{y}$ , where  $\mathbf{y}$  is an earlier estimate of the eigenvector and  $\lambda$  is an upper estimate of  $\rho_1$ .

#### 3.2 Convergence properties

**Lemma 4** (Iteration complexity of CD methods from the literature). The numbers of coordinate updates for Algorithm 2 to achieve  $f(\mathbf{x}_{\tau}) - f_* \leq \epsilon'$  where  $f_* = \min_{\mathbf{x}} f(\mathbf{x})$ , using different coordinate selection rules, are  $\mathcal{O}\left(\frac{1}{\mu_L} \cdot \log \frac{f(\mathbf{x}_0) - f_*}{\epsilon'}\right)$  for GSL,  $\mathcal{O}\left(\frac{dL_{\max}(1 + d\tilde{L}^2/L_{\max}^2)}{\mu} \cdot \log \frac{f(\mathbf{x}_0) - f_*}{\epsilon'}\right)$  for Cyclic, and  $\mathcal{O}\left(\frac{\sum_{i=1}^d \sqrt{L_i}}{\sqrt{\mu}} \cdot \log \frac{f(\mathbf{x}_0) - f_*}{\epsilon'}\right)$  for NU-ACDM.

Lemma 4 implies that the coordinate descent methods converges linearly for our strongly convex objective: the suboptimality decreases at different geometric rates for each method. Most remarkable is the time complexity of ACDM which has dependence on  $\frac{1}{\sqrt{\mu}}$ : as we mentioned earlier, the strong convexity parameter  $\mu$  for our least squares problems (in the accurate regime) are of the order  $\Delta$ , thus instantiating the shift-and-invert power method with ACDM leads to a total time complexity of which depends on  $\frac{1}{\sqrt{\Delta}}$ . In comparison, the total time complexity of standard power method depends on  $\frac{1}{\Delta}$ ). Thus we expect our algorithm to be much faster when the eigengap is small.

It is also illuminating to compare ACDM with full gradient descent methods. The iteration complexity to achieve the same accuracy is  $\mathcal{O}\left(\frac{\sqrt{\tilde{L}}}{\sqrt{\mu}}\right)$  for accelerated gradient descent (AGD, Nesterov, 2004). However, each full gradient calculation cost  $\mathcal{O}(d^2)$  (not taking into account the sparsity) and thus the total time complexity for AGD is  $\mathcal{O}\left(\frac{d^2\sqrt{\tilde{L}}}{\sqrt{\mu}}\right)$ . In contrast, each update of ACDM costs  $\mathcal{O}(d)$ , giving a total time complexity of  $\mathcal{O}\left(\frac{d\sum_{i=1}^{d}\sqrt{L_i}}{\sqrt{\mu}}\right)$ . In view of (3), we have  $\frac{d\sum_{i=1}^{d}\sqrt{L_i}}{\sqrt{\mu}} \leq \frac{d^2\sqrt{\tilde{L}_{\max}}}{\sqrt{\mu}} \leq \frac{d^2\sqrt{\tilde{L}}}{\sqrt{\mu}}$ , and thus ACDM is typically more efficient than AGD (and in the extreme case of all but one  $L_i$  being 0, the rate of ACDM is d times better). It is also observed empirically that ACDM outperforms AGD and the celebrated conjugate gradient method for solving least squares problems (Lee and Sidford, 2013; Allen-Zhu et al., 2016).

Although the convergence rates of GSL and cyclic rules are inferior than that of ACDM, we often observe competitive empirical performance from them. They are easier to implement and can be the method of choice in certain scenariors (e.g., cyclic coordinate descent may be employed when the data can not fit in memory). Finally, we remark that our general scheme and conclusion carries over to the case of block and parallel coordinate descent methods (Bradley et al., 2011; Richtarik and Takac, 2014), where more than one coordinate are updated in each step.

## 3.3 Putting it all together

Our proposed approach consists of instantiating Algorithm 1 with the different least squares optimizers (LSO) presented in Algorithm 2. Our analysis of the total time complexity thus combines the iteration complexity of shift-and-invert and the time complexity of the CD methods. First, by Lemma 2 the time complexity of Phase I is not dominant because it is independent of the final error  $\epsilon$ . Second, by Lemma 3 we need to solve  $\log \frac{d}{\epsilon}$  subproblem, each with constant ratio between initial and final suboptimality. Third, according to Lemma 4 and noting  $\lambda = \rho_1 + \mathcal{O}(1) \cdot \Delta$ , the number of coordinate updates for solving each subproblem is  $\mathcal{O}\left(\frac{1}{\mu_L}\right) = \mathcal{O}\left(\frac{d(\rho_1 - \frac{1}{d}\sum_{i=1}^d \rho_d)}{\Delta}\right)$  for GSL, and  $\mathcal{O}\left(\frac{\sum_{i=1}^d \sqrt{L_i}}{\sqrt{\mu}}\right) = \mathcal{O}\left(\frac{\frac{1}{d}\sum_{i=1}^d \sqrt{\rho_1 - A_{ii}}}{\sqrt{\Delta}}\right)$  for ACDM. Since each coordinate update costs  $\mathcal{O}(d)$  time, we obtain the time complexity of our algorithms given in Table 1.

#### 3.4 Extensions

Our algorithms can be easily extended to several other settings. Although our analysis have assumed that  $\mathbf{A}$  is positive semidefinite, given an estimate of the maximum eigenvalue, the shift-and-invert algorithm can always convert the problem into a series of convex least squares problems, regardless of whether  $\mathbf{A}$  is positive semidefinite, and so CD methods can be applied with convergence guarantee. To extract the tailing eigenvector of  $\mathbf{A}$ , one can equivalently extract the top eigenvector of  $-\mathbf{A}$ . Futhermore, extracting the top singular vectors of a non-symmetric matrix  $\mathbf{A}$  is equivalent to extracting the top eigenvectors of  $\begin{bmatrix} \mathbf{0} & \mathbf{A}^{\top} \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$ .

To extend our algorithms to extracting multiple top eigenvectors, we can use the "peeling" procedure: we iteratively extract one eigenvector at a time, removing the already ex-

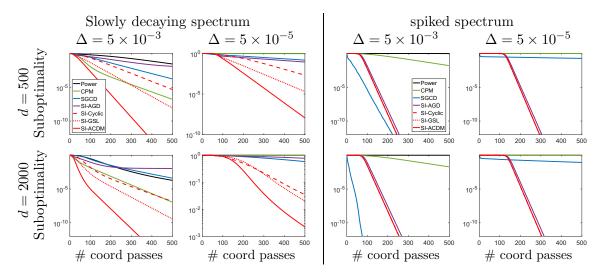


Figure 1: Comparison of the alignment suboptimality  $(1 - (\mathbf{w}_t^{\top} \mathbf{p}_1)^2)$  vs. number of coordinate passes) of various algorithms for computing the leading eigenvector on synthetic datasets with slowly decaying spectrum (left panel) and spiked spectrum (right panel).

tracted component from the data matrix before extracting the new direction. For example, to extract  $\mathbf{p}_2$ , we can equivalently extract the top eigenvector of  $\mathbf{A}' = (\mathbf{I} - \mathbf{p}_1 \mathbf{p}_1^\top) \mathbf{A} (\mathbf{I} - \mathbf{p}_1 \mathbf{p}_1^\top)$ . And note that we do not need to explicitly compute and store  $\mathbf{A}'$  which may be less sparse than  $\mathbf{A}$ ; all we need in CD methods are columns of  $\mathbf{A}'$  which can be evaluated by  $\mathcal{O}(d)$  vector operations, e.g.,  $\mathbf{A}'[:j] = (\mathbf{I} - \mathbf{p}_1 \mathbf{p}_1^\top) \cdot (\mathbf{A}[:,j] - \mathbf{p}_1[j] \cdot (\mathbf{A}\mathbf{p}_1))$  by storing  $\mathbf{A}\mathbf{p}_1 \in \mathbb{R}^d$ . We refer the readers to Allen-Zhu and Li (2016) for a careful analysis of the accuracy in solving each direction and the runtime.

#### 4. Experiments

We now demonstrate the efficiency of our algorithms, denoted as SI (shift-and-invert) + least squares solvers. We also include AGD as a solver since SI+AGD gives the same time complexity as Lanczos.

Synthetic datasets We first generate synthetic datasets to validate the fast convergence of our algorithms. Our test matrix has the form  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{U}^{\top}$ , where  $\mathbf{U}$  is a random rotation matrix and we tried two types of spectrum: i) the slowly decaying spectrum where  $\mathbf{S} = \text{diag}(1, 1 - \Delta, 1 - 2\Delta, \ldots)$ ; and ii) the spiked spectrum where  $\mathbf{S} = \text{diag}(1, 1 - \Delta, 1 - \Delta, 1 - \Delta, \ldots, 1 - \Delta)$ . The parameter  $\Delta$  (which is also the eigengap of  $\mathbf{A}$ ) controls the decay of the spectrum of  $\mathbf{A}$ . We use 4 passes of coordinate updates (each pass has d coordinate updates) for SI-Cyclic, SI-GSL and SI-ACDM, and 4 accelerated gradient updates for SI-AGD, to approximately solve the least squares problems (1), and set  $\tilde{\Delta} = 0.0001$  in Algorithm 1.

We test several settings of dimension d and  $\delta$ , and the results are summarized in Figure 1. Observe that in most cases, CPM/SGCD indeed converge faster than the standard power method, justifying the intuition that some coordinates are more important than others.

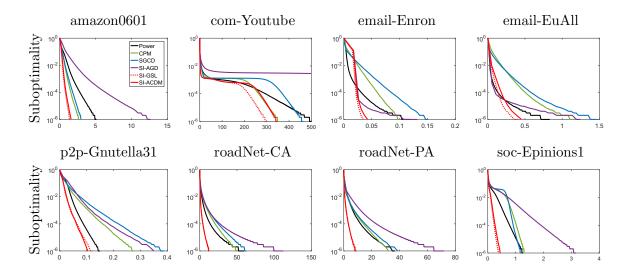


Figure 2: Comparison of the runtime efficiency  $(\rho_1 - \mathbf{w}_t^{\top} \mathbf{A} \mathbf{w}_t)$  vs. run time in seconds) of various algorithms for computing the leading eigenvector on real-world network datasets.

Furthermore, our algorithm significantly improve over CPM/SGCD for accurate estimate of the top eigenvector, especially when  $\Delta$  is small, validating the superior convergence rates of our methods. In the extreme case of the spiked covariance in Figure 1 (right panel) where other approaches converges extremely slowly when  $\Delta$  is small, while our methods still converges fast, and has milder dependence on  $\Delta$ .

Real-world datasets A major source of large-scale sparse eigenvalue problem in machine learning comes from graph-related applications (Shi and Malik, 2000; Ng et al., 2002; Fowlkes et al., 2004). We now examine the efficiency of the proposed algorithms on several large-scale networks datasets with up to a few million nodes<sup>4</sup>. Let **W** be the (binary) adjacency matrix of each network, our goal is to compute the top eigenvectors for the corresponding normalized graph Laplacian matrix (von Luxburg, 2007) defined as  $\mathbf{A} = \mathbf{D}^{-1/2}(\mathbf{D} - \mathbf{W})\mathbf{D}^{-1/2}$ , where **D** is a diagonal matrix containing the degrees of each node on the diagonal. This task can be extended to computing a low-rank approximation of the normalized graph Laplacian (Gittens and Mahoney, 2013).

The results are summarized in Figure 2.<sup>6</sup> We observe that the proposed methods usually improves over CPM and SGCD, with up to 5x speedup in runtime over CPM/SGD depending on the dataset.

 $<sup>4. \ \</sup>mathtt{http://snap.stanford.edu/data/index.html}$ 

<sup>5.</sup> We did experiment with the original task of Lei et al. (2016). For their task, our algorithms do not significantly improve over CPM/SGCD because the gap of the unnormalized adjacency matrix  $\mathbf{W}$  is quite large.

<sup>6.</sup> For real data, we do not have the ground truth leading eigenvector and it is hard to accurately estimate it when  $\Delta \to 0$ . So we measure the suboptimality in the objective which is more stable.

#### References

- Zeyuan Allen-Zhu and Yuanzhi Li. LazySVD: Even faster SVD decomposition yet without agonizing pain. In NIPS, 2016.
- Zeyuan Allen-Zhu, Zheng Qu, Peter Richtarik, and Yang Yuan. Even faster accelerated coordinate descent using non-uniform sampling. In *ICML*, 2016.
- Amir Beck and Luba Tetruashvili. On the convergence of block coordinate descent type methods. SIAM J. Imaging Sciences, 23(4):2037–2060, 2013.
- Joseph Bradley, Aapo Kyrola, Daniel Bickson, and Carlos Guestrin. Parallel coordinate descent for 11-regularized loss minimization. In *ICML*, 2011.
- Inderjit Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Nearest neighbor based greedy coordinate descent. In NIPS, 2011.
- Charless Fowlkes, Serge Belongie, Fan Chung, and Jitendra Malik. Spectral grouping using the Nyström method. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 26(2): 214–225, 2004.
- Dan Garber and Elad Hazan. Fast and simple PCA via convex optimization. arXiv:1509.05647 [math.OC], 2015.
- Dan Garber, Elad Hazan, Chi Jin, Sham M. Kakade, Cameron Musco, Praneeth Netrapalli, and Aaron Sidford. Faster eigenvector computation via shift-and-invert preconditioning. In *ICML*, 2016.
- Rong Ge, Chi Jin, Sham M. Kakade, Praneeth Netrapalli, and Aaron Sidford. Efficient algorithms for large-scale generalized eigenvector computation and canonical correlation analysis. In *ICML*, 2016.
- Alex Gittens and Michael Mahoney. Revisiting the Nystrom method for improved large-scale machine learning. In *ICML*, 2013.
- Gene H. Golub and Charles F. van Loan. *Matrix Computations*. third edition, 1996.
- Moritz Hardt and Eric Price. The noisy power method: A meta algorithm with applications. In NIPS, 2014.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In NIPS, 2013.
- Yin Tat Lee and Aaron Sidford. Efficient accelerated coordinate descent methods and faster algorithms for solving linear systems. In *Proc. of the 54th Annual Symposium on Foundations of Computer Science (FOCS 2013)*, 2013.
- Qi Lei, Kai Zhong, and Inderjit S. Dhillon. Coordinate-wise power method. In NIPS, 2016.
- Zhaosong Lu and Lin Xiao. On the complexity analysis of randomized block-coordinate descent methods. *Math. Prog.*, 152(1):615–642, 2015.

- Y. Nesterov. Introductory Lectures on Convex Optimization. A Basic Course. Number 87 in Applied Optimization. 2004.
- Yu. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. SIAM Journal on Optimization, 22(2):341–362, 2012.
- A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, 2002.
- Julie Nutini, Mark Schmidt, Issam H. Laradji, Michael Friedlander, and Hoyt Koepke. Coordinate descent converges faster with the gauss-southwell rule than random selection. In *ICML*, 2015.
- Peter Richtarik and Martin Takac. Parallel coordinate descent methods for big data optimization. *Math. Prog.*, pages 1–52, 2014.
- Ohad Shamir. A stochastic PCA and SVD algorithm with an exponential convergence rate. In *ICML*, 2015.
- Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Trans.* Pattern Analysis and Machine Intelligence, 22(8):888–905, 2000.
- Ulrike von Luxburg. A tutorial on spectral clustering. Statistics and Computing, 17(4): 395–416, December 2007.
- Weiran Wang, Jialei Wang, Dan Garber, and Nathan Srebro. Globally convergent stochastic optimization for canonical correlation analysis. In NIPS, 2016.
- Stephen J. Wright. Coordinate descent algorithms. Math. Prog., 2015.

## Appendix A. Auxiliary Lemmas

**Lemma 5.** Consider a positive semidefinite matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  with eigenvalues  $\rho_1 \geq \cdots \geq \rho_d \geq 0$  and corresponding eigenvectors  $\mathbf{p}_1, \ldots, \mathbf{p}_d$ . For a unit vector  $\mathbf{v}$  satisfying  $\mathbf{v}^\top \mathbf{p}_1 \geq 1 - \epsilon$  where  $\epsilon \in (0, 1)$ , we have

$$\mathbf{v}^{\top} \mathbf{A} \mathbf{v} \ge \rho_1 (1 - 2\epsilon).$$

*Proof.* By direct calculation, we obtain

$$\mathbf{v}^{\top} \mathbf{A} \mathbf{v} = \mathbf{v}^{\top} \left( \sum_{i=1}^{d} \rho_{i} \mathbf{p}_{i} \mathbf{p}_{i}^{\top} \right) \mathbf{v}$$

$$= \sum_{i=1}^{d} \rho_{i} \left( \mathbf{v}^{\top} \mathbf{p}_{i} \right)^{2}$$

$$\geq \rho_{1} (\mathbf{v}^{\top} \mathbf{p}_{1})^{2}$$

$$\geq \rho_{1} (1 - \epsilon)^{2}$$

$$\geq \rho_{1} (1 - 2\epsilon).$$

**Lemma 6.** For two nonzero vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ , we have

$$\left\| \frac{\mathbf{a}}{\|\mathbf{a}\|} - \frac{\mathbf{b}}{\|\mathbf{b}\|} \right\| \le 2 \frac{\|\mathbf{a} - \mathbf{b}\|}{\|\mathbf{a}\|}.$$

*Proof.* By direct calculation, we obtain

$$\begin{aligned} \left\| \frac{\mathbf{a}}{\|\mathbf{a}\|} - \frac{\mathbf{b}}{\|\mathbf{b}\|} \right\| &\le \left\| \frac{\mathbf{a}}{\|\mathbf{a}\|} - \frac{\mathbf{b}}{\|\mathbf{a}\|} \right\| + \left\| \frac{\mathbf{b}}{\|\mathbf{a}\|} - \frac{\mathbf{b}}{\|\mathbf{b}\|} \right\| \\ &= \frac{\|\mathbf{a} - \mathbf{b}\|}{\|\mathbf{a}\|} + \|\mathbf{b}\| \cdot \frac{\|\|\mathbf{a}\| - \|\mathbf{b}\|\|}{\|\mathbf{a}\| \|\mathbf{b}\|} \\ &\le \frac{\|\mathbf{a} - \mathbf{b}\|}{\|\mathbf{a}\|} + \frac{\|\mathbf{a} - \mathbf{b}\|}{\|\mathbf{a}\|} \\ &= 2\frac{\|\mathbf{a} - \mathbf{b}\|}{\|\mathbf{a}\|} \end{aligned}$$

where we have used the triangle inequality in the second inequality.

#### Appendix B. Analysis of inexact power method in different regimes

**Lemma 7** (Iteration complexity of inexact power method). Consider the following inexact shift-and-invert power iterations: for t = 1, ...,

$$\tilde{\mathbf{w}}_t \approx \underset{\mathbf{w}}{\operatorname{arg\,min}} \ f_t(\mathbf{w}) = \frac{1}{2} \mathbf{w}^{\top} (\lambda \mathbf{I} - \mathbf{A}) \mathbf{w} - \mathbf{w}_{t-1}^{\top} \mathbf{w},$$

$$\mathbf{w}_t \leftarrow \frac{\tilde{\mathbf{w}}_t}{\|\tilde{\mathbf{w}}_t\|},$$

where  $f_t(\tilde{\mathbf{w}}_t) \leq \min_{\mathbf{w}} f_t(\mathbf{w}) + \epsilon_t$ .

• (Crude regime) Define  $T_1 = \left\lceil \frac{2}{\epsilon} \log \left( \frac{4}{\epsilon \xi_{01}^2} \right) \right\rceil$ . Assume that for all  $t = 1, \ldots, T_1$ , the error in minimizing  $f_t(\mathbf{w})$  satisfies

$$\epsilon_t \le \frac{\epsilon^2 \beta_d^2}{128\beta_1} \left( \frac{(2\beta_1/\beta_d) - 1}{(2\beta_1/\beta_d)^{T_1} - 1} \right)^2.$$

Then we have  $\mathbf{w}_{T_1}^{\top}(\lambda - \mathbf{A})^{-1}\mathbf{w}_{T_1} \geq (1 - \epsilon)\beta_1$ .

• (Accurate regime) Define  $G_0 := \frac{\sqrt{\sum_{i=2}^d \xi_{0i}^2/\beta_i}}{\sqrt{\xi_{01}^2/\beta_1}}$ . Assume that for all  $t \ge 1$ , the error in minimizing  $f_t(\mathbf{w})$  satisfies

$$\epsilon_t \le \min\left(\sum_{i=2}^d \xi_{(t-1)i}^2/\beta_i, \ \xi_{(t-1)1}^2/\beta_1\right) \cdot \frac{(\beta_1 - \beta_2)^2}{32}.$$

Let  $\gamma = \frac{3\beta_1 + \beta_2}{\beta_1 + 3\beta_2} > 1$  and  $T_2 = \left\lceil \frac{1}{2} \log_{\gamma} \left( \frac{G_0^2}{\epsilon} \right) \right\rceil$ . Then we have  $\mathbf{w}_t^{\mathsf{T}} \mathbf{p}_1 \geq 1 - \epsilon$  for all  $t \geq T_2$ .

*Proof.* In the following, we use the shorthand  $\mathbf{M}_{\lambda} = (\lambda \mathbf{I} - \mathbf{A})^{-1}$ . Observe that

$$\mathbf{p}_i^{\top} \mathbf{M}_{\lambda} \mathbf{p}_i = \beta_i, \quad \text{for} \quad i = 1, \dots, d,$$
  
 $\mathbf{p}_i^{\top} \mathbf{M}_{\lambda} \mathbf{p}_j = 0, \quad \text{for} \quad i \neq j.$ 

Let the exact solution to the least squares problem  $\min_{\mathbf{w}} f_t(\mathbf{w})$  be  $\tilde{\mathbf{w}}_t^* = \mathbf{M}_{\lambda} \mathbf{w}_{t-1}$ . If we obtain an approximate solution  $\tilde{\mathbf{w}}_t$  such that  $f_t(\tilde{\mathbf{w}}_t) - f_t(\tilde{\mathbf{w}}_t^*) = \epsilon_t$ , it follows from the quadratic objective that

$$\epsilon_t = \frac{1}{2} \left( \tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_t^* \right)^{\top} \mathbf{M}_{\lambda}^{-1} \left( \mathbf{w}_t - \mathbf{w}_t^* \right) = \frac{1}{2} \left\| \tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_t^* \right\|_{\mathbf{M}_{\lambda}^{-1}}^2. \tag{4}$$

Furthermore, we have for the exact solution that

$$\tilde{\mathbf{w}}_t^* = \mathbf{M}_{\lambda} \mathbf{w}_{t-1} = \mathbf{M}_{\lambda} \sum_{i=1}^d \xi_{(t-1)i} \mathbf{p}_i = \sum_{i=1}^d \beta_i \xi_{(t-1)i} \mathbf{p}_i.$$
 (5)

**Crude regime** For the crude regime, we denote by n the number of eigenvalues of  $\mathbf{M}_{\lambda}$  that are greater than or equal to  $(1 - \frac{\epsilon}{4}) \beta_1$ .

We will relate the iterates of inexact power method to those of the exact power method

$$\tilde{\mathbf{v}}_t \leftarrow \mathbf{M}_{\lambda} \mathbf{v}_{t-1}, \qquad \mathbf{v}_t \leftarrow \frac{\tilde{\mathbf{v}}_t}{\|\tilde{\mathbf{v}}_t\|}, \qquad t = 1, \dots$$

from the same initialization  $\mathbf{v}_0 = \mathbf{w}_0$ .

On the one hand, we can show that exact power iterations converges quickly for eigenvalue estimation. Since  $\mathbf{v}_t = \frac{\mathbf{M}_{\lambda}^t \mathbf{v}_0}{\|\mathbf{M}_{\lambda}^t \mathbf{v}_0\|}$  and  $\mathbf{M}_{\lambda}^t \mathbf{v}_0 = \sum_{i=1}^d \beta_i^t \xi_{0i} \mathbf{p}_i$ , we have

$$\sum_{i=n+1}^{d} (\mathbf{v}_{t}^{\top} \mathbf{p}_{i})^{2} = \frac{\sum_{i=n+1}^{d} \beta_{i}^{2t} \xi_{0i}^{2}}{\sum_{i=1}^{d} \beta_{i}^{2t} \xi_{0i}^{2}} \leq \frac{\beta_{n+1}^{2t} \sum_{i=n+1}^{d} \xi_{0i}^{2}}{\beta_{1}^{2t} \xi_{01}^{2}}$$
$$\leq \frac{1}{\xi_{01}^{2}} \left(1 - \frac{\epsilon}{4}\right)^{2t} \leq \frac{e^{-\epsilon t/2}}{\xi_{01}^{2}}.$$

This indicates that after  $T_1 = \left\lceil \frac{2}{\epsilon} \log \left( \frac{4}{\epsilon \xi_{01}^2} \right) \right\rceil$  iterations, we have  $\sum_{i=n+1}^d (\mathbf{v}_t^\top \mathbf{p}_i)^2 \leq \frac{\epsilon}{4}$  for all  $t \geq T_1$ . And as a result, for all  $t \geq T_1$ , it holds that

$$\mathbf{v}_{t}^{\top} \mathbf{M}_{\lambda} \mathbf{v}_{t}$$

$$= \mathbf{v}_{t}^{\top} \left( \sum_{i=1}^{d} \beta_{i} \mathbf{p}_{i} \mathbf{p}_{i}^{\top} \right) \mathbf{v}_{t} = \sum_{i=1}^{d} \beta_{i} (\mathbf{v}_{t}^{\top} \mathbf{p}_{i})^{2}$$

$$\geq \beta_{n} \sum_{i=1}^{n} (\mathbf{v}_{t}^{\top} \mathbf{p}_{i})^{2} \geq \left( 1 - \frac{\epsilon}{4} \right) \beta_{1} \cdot \left( 1 - \sum_{i=n+1}^{d} (\mathbf{v}_{t}^{\top} \mathbf{p}_{i})^{2} \right)$$

$$\geq \left( 1 - \frac{\epsilon}{4} \right)^{2} \beta_{1} \geq \left( 1 - \frac{\epsilon}{2} \right) \beta_{1}.$$

On the other hand, we can lower bound  $\mathbf{w}_t^{\top} \mathbf{M}_{\lambda} \mathbf{w}_t$  using  $\mathbf{v}_t^{\top} \mathbf{M}_{\lambda} \mathbf{v}_t$ :

$$\mathbf{w}_{t}^{\top} \mathbf{M}_{\lambda} \mathbf{w}_{t} = (\mathbf{v}_{t} + \mathbf{w}_{t} - \mathbf{v}_{t})^{\top} \mathbf{M}_{\lambda} (\mathbf{v}_{t} + \mathbf{w}_{t} - \mathbf{v}_{t})$$

$$\geq \mathbf{v}_{t}^{\top} \mathbf{M}_{\lambda} \mathbf{v}_{t} + 2(\mathbf{w}_{t} - \mathbf{v}_{t})^{\top} \mathbf{M}_{\lambda} \mathbf{v}_{t}$$

$$\geq \mathbf{v}_{t}^{\top} \mathbf{M}_{\lambda} \mathbf{v}_{t} - 2\beta_{1} \|\mathbf{w}_{t} - \mathbf{v}_{t}\|$$

$$\geq \left(1 - \frac{\epsilon}{2}\right) \beta_{1} - 2\beta_{1} \|\mathbf{w}_{t} - \mathbf{v}_{t}\|.$$

Therefore, we will have  $\mathbf{w}_t^{\top} \mathbf{M}_{\lambda} \mathbf{w}_t \geq (1 - \epsilon) \beta_1$  as desired if we can guarantee  $\|\mathbf{w}_{T_1} - \mathbf{v}_{T_1}\| \leq \frac{\epsilon}{4}$ .

We now upper bound the difference  $S_t := \|\tilde{\mathbf{w}}_t - \tilde{\mathbf{v}}_t\|$  by induction. Due to the assumption of same initialization, we have  $S_0 = 0$ . For  $t \ge 2$ , we can decompose the error into two terms:

$$\begin{aligned} \|\tilde{\mathbf{w}}_t - \tilde{\mathbf{v}}_t\| &\leq \|\tilde{\mathbf{w}}_t - \mathbf{M}_{\lambda} \mathbf{w}_{t-1}\| + \|\mathbf{M}_{\lambda} \mathbf{w}_{t-1} - \tilde{\mathbf{v}}_t\| \\ &\leq \|\tilde{\mathbf{w}}_t - \mathbf{M}_{\lambda} \mathbf{w}_{t-1}\| + \|\mathbf{M}_{\lambda} \mathbf{w}_{t-1} - \mathbf{M}_{\lambda} \mathbf{v}_{t-1}\|. \end{aligned}$$

The first term concerns the error from inexact minimization of  $f_t(\mathbf{w})$  and can be bounded using (4):

$$\|\tilde{\mathbf{w}}_t - \mathbf{M}_{\lambda} \mathbf{w}_{t-1}\| \le \left\| \mathbf{M}_{\lambda}^{\frac{1}{2}} \right\| \cdot \|\tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_t^*\|_{\mathbf{M}_{\lambda}^{-1}} \le \sqrt{2\beta_1 \epsilon_t}.$$
 (6)

The second term concerns the error from inexact target:

$$\|\mathbf{M}_{\lambda}\mathbf{w}_{t-1} - \mathbf{M}_{\lambda}\mathbf{v}_{t-1}\| \le \|\mathbf{M}_{\lambda}\| \cdot \|\mathbf{w}_{t-1} - \mathbf{v}_{t-1}\|$$

$$= \beta_{1} \left\| \frac{\tilde{\mathbf{w}}_{t-1}}{\|\tilde{\mathbf{w}}_{t-1}\|} - \frac{\tilde{\mathbf{v}}_{t-1}}{\|\tilde{\mathbf{v}}_{t-1}\|} \right\| \le 2\beta_{1} \frac{\|\tilde{\mathbf{w}}_{t-1} - \tilde{\mathbf{v}}_{t-1}\|}{\|\tilde{\mathbf{v}}_{t-1}\|}$$
(7)

where the last inequality is due to Lemma 6.

Noting that  $\|\tilde{\mathbf{v}}_{t-1}\| = \|\mathbf{M}_{\lambda}\mathbf{v}_{t-2}\| \ge \beta_d \|\mathbf{v}_{t-2}\| = \beta_d$ , and combining (6) and (7), we obtain

$$S_t \le \sqrt{2\beta_1 \epsilon_t} + \frac{2\beta_1}{\beta_d} S_{t-1}.$$

Fixing  $\epsilon_t = \epsilon'$  for all t and unfolding the above inequality over t yield

$$S_{T_1} \le \sqrt{2\beta_1 \epsilon'} \cdot \sum_{t=0}^{T_1} \left(\frac{2\beta_1}{\beta_d}\right)^t = \sqrt{2\beta_1 \epsilon'} \cdot \frac{(2\beta_1/\beta_d)^{T_1} - 1}{(2\beta_1/\beta_d) - 1}.$$

By Lemma 6 again, we have

$$\|\mathbf{w}_{t} - \mathbf{v}_{t}\| = \left\| \frac{\tilde{\mathbf{w}}_{t}}{\|\tilde{\mathbf{w}}_{t}\|} - \frac{\tilde{\mathbf{v}}_{t}}{\|\tilde{\mathbf{v}}_{t}\|} \right\| \leq 2 \frac{\|\tilde{\mathbf{w}}_{t} - \tilde{\mathbf{v}}_{t}\|}{\|\tilde{\mathbf{v}}_{t}\|} \leq \frac{2S_{t}}{\beta_{d}}$$
$$\leq \frac{\sqrt{8\beta_{1}\epsilon'}}{\beta_{d}} \cdot \frac{(2\beta_{1}/\beta_{d})^{T_{1}} - 1}{(2\beta_{1}/\beta_{d}) - 1}.$$

Setting the RHS to be  $\frac{\epsilon}{4}$  gives

$$\epsilon' = \frac{\epsilon^2 \beta_d^2}{128\beta_1} \left( \frac{(2\beta_1/\beta_d) - 1}{(2\beta_1/\beta_d)^{T_1} - 1} \right)^2.$$

**Accurate regime** In the accurate regime, the potential function we use to evaluate the progress of each iteration is

$$G(\mathbf{w}_t) = \frac{\|\mathcal{P}_{\perp}\mathbf{w}_t\|_{\mathbf{M}_{\lambda}^{-1}}}{\|\mathcal{P}_{\parallel}\mathbf{w}_t\|_{\mathbf{M}_{\lambda}^{-1}}} = \frac{\sqrt{\sum_{i=2}^{d} \xi_{ti}^2/\beta_i}}{\sqrt{\xi_{t1}^2/\beta_1}},$$

where  $\mathcal{P}_{\perp}$  and  $\mathcal{P}_{\parallel}$  denote projections onto the subspaces that are orthogonal and parallel to  $\mathbf{p}_1$  respectively. Note that

$$|\sin \theta_t| = \sqrt{\sum_{i=2}^d \xi_{ti}^2} \le |\tan \theta_t| = \frac{\sqrt{\sum_{i=2}^d \xi_{ti}^2}}{\sqrt{\xi_{t1}^2}} \le G(\mathbf{w}_t)$$

where  $\theta_t$  is the angle between  $\mathbf{w}_t$  and  $\mathbf{p}_1$ .

Our goal in this regime is to have  $|\sin \theta_t| \leq \sqrt{\epsilon}$ , as this implies

$$\mathbf{w}_t^{\mathsf{T}} \mathbf{p}_1 = \cos \theta_t = \sqrt{1 - \sin^2 \theta_t} \ge 1 - \sin^2 \theta_t \ge 1 - \epsilon$$

as desired. We ensure  $|\sin \theta_t| \leq \sqrt{\epsilon}$  by requiring that  $G(\mathbf{w}_t) \leq \sqrt{\epsilon}$ .

In view of (4) and (5), we can bound the numerator and denominator of  $G(\mathbf{w}_t)$  respectively with regard to  $G(\mathbf{w}_{t-1})$ :

$$\begin{aligned} & \left\| \mathcal{P}_{\perp} \frac{\tilde{\mathbf{w}}_{t}}{\|\tilde{\mathbf{w}}_{t}\|} \right\|_{\mathbf{M}_{\lambda}^{-1}} \\ & \leq \frac{1}{\|\tilde{\mathbf{w}}_{t}\|} \left( \left\| \mathcal{P}_{\perp} \tilde{\mathbf{w}}_{t}^{*} \right\|_{\mathbf{M}_{\lambda}^{-1}} + \left\| \mathcal{P}_{\perp} \left( \tilde{\mathbf{w}}_{t} - \tilde{\mathbf{w}}_{t}^{*} \right) \right\|_{\mathbf{M}_{\lambda}^{-1}} \right) \\ & \leq \frac{1}{\|\tilde{\mathbf{w}}_{t}\|} \left( \left\| \sum_{i=2}^{d} \beta_{i} \xi_{(t-1)i} \mathbf{p}_{i} \right\|_{\mathbf{M}_{\lambda}^{-1}} + \left\| \tilde{\mathbf{w}}_{t} - \tilde{\mathbf{w}}_{t}^{*} \right\|_{\mathbf{M}_{\lambda}^{-1}} \right) \\ & = \frac{1}{\|\tilde{\mathbf{w}}_{t}\|} \left( \sqrt{\sum_{i=2}^{d} \beta_{i} \xi_{(t-1)i}^{2}} + \sqrt{2\epsilon_{t}} \right), \end{aligned}$$

and

$$\begin{aligned} & \left\| \mathcal{P}_{\parallel} \frac{\tilde{\mathbf{w}}_{t}}{\|\tilde{\mathbf{w}}_{t}\|} \right\|_{\mathbf{M}_{\lambda}^{-1}} \\ & \geq \frac{1}{\|\tilde{\mathbf{w}}_{t}\|} \left( \left\| \mathcal{P}_{\parallel} \tilde{\mathbf{w}}_{t}^{*} \right\|_{\mathbf{M}_{\lambda}^{-1}} - \left\| \mathcal{P}_{\parallel} \left( \tilde{\mathbf{w}}_{t} - \tilde{\mathbf{w}}_{t}^{*} \right) \right\|_{\mathbf{M}_{\lambda}^{-1}} \right) \\ & \geq \frac{1}{\|\tilde{\mathbf{w}}_{t}\|} \left( \left\| \beta_{1} \xi_{(t-1)1} \mathbf{p}_{1} \right\|_{\mathbf{M}_{\lambda}^{-1}} - \left\| \tilde{\mathbf{w}}_{t} - \tilde{\mathbf{w}}_{t}^{*} \right\|_{\mathbf{M}_{\lambda}^{-1}} \right) \\ & = \frac{1}{\|\tilde{\mathbf{w}}_{t}\|} \left( \sqrt{\beta_{1} \xi_{(t-1)1}^{2}} - \sqrt{2\epsilon_{t}} \right). \end{aligned}$$

Consequently, we have

$$G(\mathbf{w}_{t}) \leq \frac{\sqrt{\sum_{i=2}^{d} \beta_{i} \xi_{(t-1)i}^{2}} + \sqrt{2\epsilon_{t}}}{\sqrt{\beta_{1} \xi_{(t-1)1}^{2}} - \sqrt{2\epsilon_{t}}}$$

$$\leq \frac{\beta_{2} \sqrt{\sum_{i=2}^{d} \xi_{(t-1)i}^{2} / \beta_{i}} + \sqrt{2\epsilon_{t}}}{\beta_{1} \sqrt{\xi_{(t-1)1}^{2} / \beta_{1}} - \sqrt{2\epsilon_{t}}}$$

$$= G(\mathbf{w}_{t}) \cdot \frac{\beta_{2} + \frac{\sqrt{2\epsilon_{t}}}{\sqrt{\sum_{i=2}^{d} \xi_{(t-1)i}^{2} / \beta_{i}}}}{\beta_{1} - \frac{\sqrt{2\epsilon_{t}}}{\sqrt{\xi_{(t-1)1}^{2} / \beta_{1}}}}.$$

As long as

$$\sqrt{2\epsilon_t} \le \min\left(\sqrt{\sum_{i=2}^d \xi_{(t-1)i}^2/\beta_i}, \sqrt{\xi_{(t-1)1}^2/\beta_1}\right) \cdot \frac{\beta_1 - \beta_2}{4},$$

or equivalently

$$\epsilon_t \le \min\left(\sum_{i=2}^d \xi_{(t-1)i}^2/\beta_i, \ \xi_{(t-1)1}^2/\beta_1\right) \cdot \frac{(\beta_1 - \beta_2)^2}{32},$$

we are guaranteed that

$$G(\mathbf{w}_t) \le \frac{\beta_1 + 3\beta_2}{3\beta_1 + \beta_2} \cdot G(\mathbf{w}_{t-1}).$$

And when this holds for all  $t \geq 1$ , the sequence  $\{G(\mathbf{w}_t)\}_{t=0,...}$  decreases (at least) at a constant geometric rate of  $\frac{\beta_1+3\beta_2}{3\beta_1+\beta_2} < 1$ . Therefore, the number of iterations needed to achieve  $|\sin \theta_{T_2}| \leq \sqrt{\epsilon}$  is  $T_2 = \left\lceil \log_{\gamma} \left(\frac{G(\mathbf{w}_0)}{\sqrt{\epsilon}}\right) \right\rceil = \left\lceil \frac{1}{2} \log_{\gamma} \left(\frac{G_0^2}{\epsilon}\right) \right\rceil$  for  $\gamma = \frac{3\beta_1+\beta_2}{\beta_1+3\beta_2}$ .

## Appendix C. Bounding initial error for least squares

For each least squares problem  $f_t(\mathbf{w})$  in inexact shift-and-invert power method, we can use  $\tilde{\mathbf{w}}_{t-1}$  from the previous iteration as initialization, and the intuition behind is that as the power power method converges, the least squares problems become increasingly similar. However, an even better warm-start is to use as initialization  $\alpha_t \mathbf{w}_{t-1}$  for some  $\alpha_t$  and minimize the initial suboptimality over  $\alpha_t$  (Garber et al., 2016; Ge et al., 2016). Let the exact solution to the least squares problem  $\min_{\mathbf{w}} f_t(\mathbf{w})$  be  $\tilde{\mathbf{w}}_t^* = (\lambda \mathbf{I} - \mathbf{A})^{-1} \mathbf{w}_{t-1}$ , which gives the optimal objective function value  $f_t^* = -\frac{1}{2} \mathbf{w}_{t-1}^{\mathsf{T}} (\lambda \mathbf{I} - \mathbf{A})^{-1} \mathbf{w}_{t-1}$ . Observe that

$$\epsilon_t^{init} = f_t(\alpha_t \mathbf{w}_{t-1}) - f_t^*$$

$$= \frac{\alpha_t^2}{2} \cdot \mathbf{w}_{t-1}^{\top} (\lambda \mathbf{I} - \mathbf{A}) \mathbf{w}_{t-1} - \alpha_t \cdot \mathbf{w}_{t-1}^{\top} \mathbf{w}_{t-1} - f_t^*.$$

This is a quadratic function of  $\alpha_t$ , with minimum achieved at  $\alpha_t = \frac{1}{\mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w}_{t-1}}$ . With this choice of  $\alpha_t$ , we obtain

$$\epsilon_{t}^{init} \leq f_{t}(\beta_{1}\mathbf{w}_{t-1}) - f_{t}^{*} 
= \frac{\beta_{1}^{2}}{2} \cdot \mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})\mathbf{w}_{t-1} - \beta_{1} + \frac{1}{2}\mathbf{w}_{t-1}^{\top}(\lambda \mathbf{I} - \mathbf{A})^{-1}\mathbf{w}_{t-1} 
= \frac{\beta_{1}^{2}}{2} \sum_{i=1}^{d} \xi_{(t-1)i}^{2}/\beta_{i} - \beta_{1} \sum_{i=1}^{d} \xi_{(t-1)i}^{2} + \frac{1}{2} \sum_{i=1}^{d} \beta_{i} \xi_{(t-1)i}^{2} 
\leq \frac{1}{2} \sum_{i=1}^{d} (\beta_{1} - \beta_{i})^{2} \cdot \xi_{(t-1)i}^{2}/\beta_{i} 
\leq \frac{\beta_{1}^{2}}{2} \sum_{i=2}^{d} \xi_{(t-1)i}^{2}/\beta_{i}. \tag{8}$$

As we show in the next lemma, this initialization will allow the ratio between initial and final error to be a constant in the accurate regime, independent of the final suboptimality in alignment.

In the crude regime, we simply use the rough bound

$$\epsilon_t^{init} \le f_t(\mathbf{0}) - f_t^* = \frac{1}{2} \mathbf{w}_{t-1}^{\mathsf{T}} (\lambda \mathbf{I} - \mathbf{A})^{-1} \mathbf{w}_{t-1}$$
$$= \frac{1}{2} \sum_{i=1}^d \beta_i \xi_{(t-1)i}^2 \le \frac{\beta_1}{2}. \tag{9}$$

Substituting (8) and (9) into Lemma 7 yields Lemma 1.

#### Appendix D. Proof of Lemma 2

*Proof.* Observe that the eigenvalues of  $\mathbf{M}_{\lambda} = (\lambda \mathbf{I} - \mathbf{A})^{-1}$  are functions of  $\lambda$  as  $\lambda$  changes over iterations. Below we use  $\sigma_i(\mathbf{M}_{\lambda})$  to denote the i-th eigenvalue of  $\mathbf{M}_{\lambda}$  instead of  $\beta_1, \ldots, \beta_d$  to make this dependence explicit.

Define the condition number of  $\mathbf{M}_{\lambda}$  as

$$\kappa_{\lambda} := \frac{\sigma_1(\mathbf{M}_{\lambda})}{\sigma_d(\mathbf{M}_{\lambda})} = \frac{\frac{1}{\lambda - \rho_1}}{\frac{1}{\lambda - \rho_d}} = \frac{\lambda - \rho_d}{\lambda - \rho_1},$$

Suppose we have the upper and lower bound of all  $\sigma_1(\mathbf{M}_{\lambda_{(s)}})$  used in the **repeat-until** loop:

$$\overline{\sigma} \geq \sigma_1(\mathbf{M}_{\lambda_{(s)}}), \quad \underline{\sigma} \geq \sigma_1(\mathbf{M}_{\lambda_{(s)}}), \quad \text{for } s = 1, \dots,$$

whose specific values will be provided shortly.

Throughout the loop, we require for all iteration s that

$$\sqrt{2\overline{\sigma}\tilde{\epsilon}_s} \le \frac{\underline{\sigma}}{8}.\tag{10}$$

Let the power iterations for different s start from the same initialization  $\mathbf{w}_0$  with  $\xi_{01} = \mathbf{w}_0^{\mathsf{T}} \mathbf{p}_1$ , and apply Lemma 7 (crude regime) with  $\epsilon = \frac{1}{4}$ . By our choice of  $m_1$  and setting the ratio

$$\frac{\epsilon_t^{init}}{\epsilon_t} = 1024\kappa_{\lambda_{(s)}}^2 \left(\frac{(2\kappa_{\lambda_{(s)}})^{m_1} - 1}{2\kappa_{\lambda_{(s)}} - 1}\right)^2 \tag{11}$$

according to Lemma 1 (crude regime), we obtain

$$\mathbf{w}_{sm_1}^{\top} \mathbf{M}_{\lambda_{(s-1)}} \mathbf{w}_{sm_1} \ge \frac{3}{4} \sigma_1(\mathbf{M}_{\lambda_{(s-1)}}). \tag{12}$$

In view of the definition of the vector  $\mathbf{u}_s$ , and following the same argument in (6), we have

$$\left\|\mathbf{u}_{s}-\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}}\right\|\leq\sqrt{2\sigma_{1}(\mathbf{M}_{\lambda_{(s-1)}})\cdot\tilde{\epsilon}_{s}}.$$

Then for every iteration s of **repeat-until**, it holds that

$$\begin{split} &\mathbf{w}_{sm_{1}}^{\top}\mathbf{u}_{s} \\ &= \mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} + \mathbf{w}_{sm_{1}}^{\top}\left(\mathbf{u}_{s} - \mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}}\right) \\ &\in \left[\mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} - \sqrt{2\sigma_{1}(\mathbf{M}_{\lambda_{(s-1)}}) \cdot \tilde{\epsilon}_{s}}, \right. \\ & \left. \mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} + \sqrt{2\sigma_{1}(\mathbf{M}_{\lambda_{(s-1)}}) \cdot \tilde{\epsilon}_{s}}\right] \\ &\in \left[\mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} - \sqrt{2\overline{\sigma}\tilde{\epsilon}_{s}}, \right. \\ & \left. \mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} + \sqrt{2\overline{\sigma}\tilde{\epsilon}_{s}}\right] \\ &\in \left[\mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} - \frac{\underline{\sigma}}{8}, \ \mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} + \frac{\underline{\sigma}}{8}\right], \end{split}$$

where we have used the Cauchy-Schwarz inequality in the second step and (10) in the last step.

In view of (10) and (12), it follows that

$$\mathbf{w}_{sm_{1}}^{\top}\mathbf{u}_{s} - \underline{\sigma}/8$$

$$\in \left[\mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}} - \frac{\underline{\sigma}}{4}, \mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}}\right]$$

$$\in \left[\frac{3}{4}\sigma_{1}(\mathbf{M}_{\lambda_{(s-1)}}) - \frac{\underline{\sigma}}{4}, \mathbf{w}_{sm_{1}}^{\top}\mathbf{M}_{\lambda_{(s-1)}}\mathbf{w}_{sm_{1}}\right]$$

$$\in \left[\frac{1}{2}\sigma_{1}(\mathbf{M}_{\lambda_{(s-1)}}), \sigma_{1}(\mathbf{M}_{\lambda_{(s-1)}})\right].$$

By the definition of  $\Delta_s$  in Algorithm 1 and the fact that  $\sigma_1(\mathbf{M}_{\lambda_{(s-1)}}) = \frac{1}{\lambda_{(s-1)}-\rho_1}$ , we have

$$\Delta_{s} = \frac{1}{2} \cdot \frac{1}{\mathbf{w}_{sm_{1}}^{\top} \mathbf{u}_{s} - \underline{\sigma}/8}$$

$$\in \left[ \frac{1}{2} \left( \lambda_{(s-1)} - \rho_{1} \right), \ \lambda_{(s-1)} - \rho_{1} \right]. \tag{13}$$

And as a result,

$$\lambda_{(s)} = \lambda_{(s-1)} - \frac{\Delta_s}{2} \ge \lambda_{(s-1)} - \frac{1}{2} \left( \lambda_{(s-1)} - \rho_1 \right)$$
$$= \frac{\lambda_{(s-1)} + \rho_1}{2},$$

and thus by induction (note  $\lambda_{(0)} \geq \rho_1$ ) we have  $\lambda_{(s)} \geq \rho_1$  throughout the **repeat-until** loop.

From (13) we also obtain

$$\lambda_{(s)} - \rho_1 = \lambda_{(s-1)} - \rho_1 - \frac{\Delta_s}{2}$$

$$\leq \lambda_{(s-1)} - \rho_1 - \frac{1}{4} \left( \lambda_{(s-1)} - \rho_1 \right)$$

$$= \frac{3}{4} \left( \lambda_{(s-1)} - \rho_1 \right).$$

To sum up,  $\lambda_{(s)}$  approaches  $\rho_1$  from above and the gap between  $\lambda_{(s)}$  and  $\rho_1$  reduces at the geometric rate of  $\frac{3}{4}$ . Thus after at most  $T_3 = \left\lceil \log_{3/4} \left( \frac{\tilde{\Delta}}{\lambda_{(0)} - \rho_1} \right) \right\rceil = \mathcal{O}\left( \log\left(\frac{1}{\tilde{\Delta}}\right) \right)$  iterations, we reach a  $\lambda_{(T_3)}$  such that  $\lambda_{(T_3)} - \rho_1 \leq \tilde{\Delta}$ . And in view of (13), the **repeatuntil** loop exits in the next iteration. Hence, the overall number of iterations is at most  $T_3 + 1 = \mathcal{O}\left(\frac{1}{\tilde{\Delta}}\right)$ .

We now analyze  $\lambda_{(f)}$  and derive the interval it lies in. Note that  $\Delta_f \leq \tilde{\Delta}$  and  $\Delta_{f-1} > \tilde{\Delta}$  by the exiting condition. In view of (13), we have

$$\lambda_{(f)} - \rho_1 = \lambda_{(f-1)} - \rho_1 - \frac{\Delta_f}{2} \le 2\Delta_f - \frac{\Delta_f}{2}$$
$$= \frac{3\Delta_f}{2} \le \frac{3\tilde{\Delta}}{2}.$$

On the other hand, we have

$$\lambda_{(f)} - \rho_1 = \lambda_{(f-1)} - \rho_1 - \frac{\Delta_f}{2}$$

$$\geq \lambda_{(f-1)} - \rho_1 - \frac{1}{2} \left( \lambda_{(f-1)} - \rho_1 \right)$$

$$= \frac{1}{2} \left( \lambda_{(f-1)} - \rho_1 \right). \tag{14}$$

If f = 1, then by our choice of  $\lambda_{(0)}$  we have that  $\lambda_{(f)} - \rho_1 \geq \tilde{\Delta}$ . Otherwise, by unfolding (14) one more time, we have that

$$\lambda_{(f)} - \rho_1 \ge \frac{1}{4} \left( \lambda_{(f-2)} - \rho_1 \right) \ge \frac{\Delta_{f-1}}{4} \ge \frac{\tilde{\Delta}}{4}.$$

Thus in both case, we have that  $\lambda_{(f)} - \rho_1 \geq \frac{\tilde{\Delta}}{4}$  holds.

Since the  $\lambda_{(s)}$  values are monotonically non-increasing and lower-bounded by  $\rho_1 + \frac{\Delta}{4}$ , we have

$$\max_{s} \ \sigma_{1}(\mathbf{M}_{\lambda_{(s)}}) = \sigma_{1}(\mathbf{M}_{\lambda_{(f)}}) = \frac{1}{\lambda_{(f)} - \rho_{1}} \le \frac{4}{\tilde{\Delta}} =: \overline{\sigma},$$

and

$$\min_{s} \sigma_{1}(\mathbf{M}_{\lambda_{(s)}}) = \sigma_{1}(\mathbf{M}_{\lambda_{(0)}}) = \frac{1}{\lambda_{(0)} - \rho_{1}} = \frac{1}{1 + \tilde{\Delta} - \rho_{1}}$$

$$\geq \frac{1}{1 + c_{2}\Delta - \Delta} \geq 1 + (1 - c_{2})\Delta$$

$$\geq 1 + \frac{1 - c_{2}}{c_{2}}\tilde{\Delta} =: \underline{\sigma},$$

where the first inequality holds since by definition of  $\Delta$  it follows that  $\rho_1 = \rho_2 + \Delta \geq \Delta$ . Then according to (10), we can set

$$\tilde{\epsilon}_s = \frac{\underline{\sigma}^2}{128\overline{\sigma}} = \frac{\left(1 + \frac{1 - c_2}{c_2}\tilde{\Delta}\right)^2}{128 \cdot \frac{4}{\tilde{\Delta}}} \ge \frac{1}{128 \cdot \frac{4}{\tilde{\Delta}}} = \frac{\tilde{\Delta}}{512}.$$

Because the initialization for minimizing  $l_s(\mathbf{u})$  gives an initial suboptimality of  $\tilde{\epsilon}_s^{init} \leq$  $\frac{\sigma_1\left(\mathbf{M}_{\lambda_{(s-1)}}\right)}{2} \leq \frac{\overline{\sigma}}{2} = \frac{2}{\tilde{\Delta}}, \text{ we achieve } \tilde{\epsilon}_s\text{-suboptimality by requiring } \frac{\tilde{\epsilon}_s^{init}}{\tilde{\epsilon}_s} = \frac{1024}{\tilde{\Delta}^2}.$  It remains to fulfill the requirement (11) for **repeat-until**. Note that the condition

numbers are bounded throughout:

$$\kappa_{\lambda_{(s)}} = \frac{\lambda_{(s)} - \rho_d}{\lambda_{(s)} - \rho_1} \le \frac{\lambda_{(s)}}{\lambda_{(s)} - \rho_1} = 1 + \frac{\rho_1}{\lambda_{(s)} - \rho_1} \\
\le 1 + \rho_1 \cdot \sigma_1(\mathbf{M}_{\lambda_{(s)}}) \le 1 + \frac{4}{\tilde{\Delta}} \le \frac{5}{\tilde{\Delta}}.$$
(15)

We can then bound the ratio in (11) for all least squares subproblems:

$$\frac{\epsilon_t^{init}}{\epsilon_t} \le \frac{1024 \cdot 25}{\tilde{\Delta}^2} \left( \frac{(10/\tilde{\Delta})^{m_1}}{9/\tilde{\Delta}} \right)^2 \le \frac{32 \cdot 10^{2m_1 + 1}}{\tilde{\Delta}^{2m_1}}.$$

## Appendix E. Proof of Lemma 3

*Proof.* Observe that for  $\lambda = \rho_1 + c(\rho_1 - \rho_2)$ , we have

$$\frac{\sigma_1(\mathbf{M}_{\lambda})}{\sigma_1(\mathbf{M}_{\lambda}) - \sigma_2(\mathbf{M}_{\lambda})} = \frac{\frac{1}{\lambda - \rho_1}}{\frac{1}{\lambda - \rho_1} - \frac{1}{\lambda - \rho_2}} = \frac{\lambda - \rho_2}{\rho_1 - \rho_2}$$
$$= \frac{\rho_1 + c(\rho_1 - \rho_2) - \rho_2}{\rho_1 - \rho_2} = c + 1$$

and

$$\frac{3\sigma_1(\mathbf{M}_{\lambda}) + \sigma_2(\mathbf{M}_{\lambda})}{\sigma_1(\mathbf{M}_{\lambda}) + 3\sigma_1(\mathbf{M}_{\lambda})} = \frac{\frac{3}{\lambda - \rho_1} + \frac{1}{\lambda - \rho_2}}{\frac{1}{\lambda - \rho_1} + \frac{3}{\lambda - \rho_2}} = \frac{4c + 3}{4c + 1}.$$

In view of (2), we have  $\lambda_{(f)} - \rho_1 \leq \frac{3}{2}\tilde{\Delta} \leq \frac{3c_2}{2}\Delta \leq \frac{3}{2}\Delta$ , i.e.,  $c \leq \frac{3}{2}$  for  $\lambda_{(f)}$ . As a result, we can apply Lemma 7 (accurate regime) with  $\gamma = \frac{9}{7}$ , and we are guaranteed to achieve the desired alignment with the specified  $m_2$ . The choice of  $\frac{\epsilon_t^{init}}{\epsilon_t}$  follows from an application of Lemma 1 (accurate regime) with  $\frac{\beta_1^2}{(\beta_1 - \beta_2)^2} \leq \frac{25}{4}$ .

## Appendix F. More details of experiments on networks

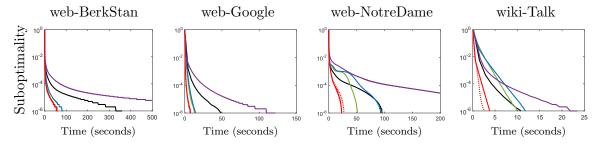


Figure 3: More results on comparison of the runtime efficiency  $(\rho_1 - \mathbf{w}_t^{\top} \mathbf{A} \mathbf{w}_t)$  vs. run time in seconds) of various algorithms for computing the leading eigenvector on real-world network datasets.

The statistics of tested real world datasets are summarized in Table 3.

For power method, CPM and SGCD, we adopt the C++ implementation of Lei et al. (2016) and implement our algorithms alongside theirs, and run all experiments on a Linux machine with Intel i7 CPU of frequency 3.20GHz. We use 4 passes of coordinate updates for solving least squares problems and set  $\tilde{\Delta}=0.05$  in Algorithm 1.

Table 3: List of network datasets used in the experiments.

Name	#Nodes	#Edges
amazon0601	403,394	3,387,388
com-Youtube	1,134,890	2,987,624
email-Enron	36,692	183,831
email-EuAll	$265,\!214$	420,045
p2p-Gnutella31	$62,\!586$	147,892
roadNet-CA	1,965,206	2,766,607
roadNet-PA	1,379,917	1,921,660
soc-Epinions1	$75,\!879$	508,837
web-BerkStan	$685,\!230$	7,600,595
web-Google	875,713	5,105,039
web-NotreDame	$325{,}729$	1,497,134
wiki-Talk	2,394,385	5,021,410