# Coordinate-wise Power Method

Qi Lei, Kai Zhong[1], and Inderjit S. Dhillon[1][2]

[1]Institute for Computational Sciences and Engineering,  [2]Department of Computer Science, The University of Texas at Austin
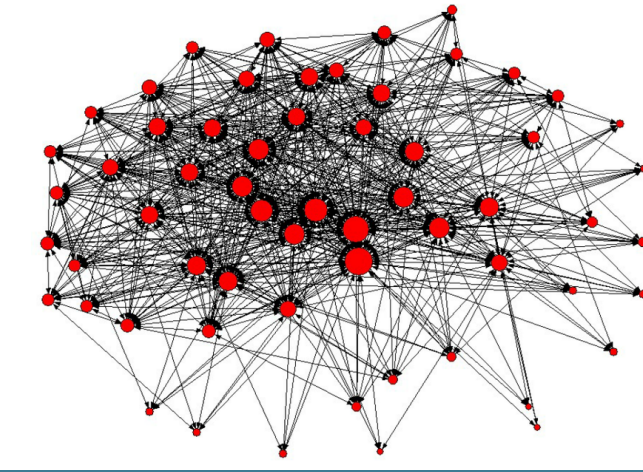
## Motivation

**Goal**: Given a matrix $A$, we seek to compute its dominant eigenvector $v_1$:

$$v_1 = \underset{\|x\|=1}{\operatorname{argmax}}\ x^T A^T A x \qquad (1)$$

Computing the dominant eigenvector of a given matrix/graph is meaningful for:

- Graph Centrality/PageRank
- Sparse PCA
- Spectral Clustering

The classic power method is still powerful in the sense of:

- Simplicity
- Small memory footprint
- Stable: being resistent to noise

We propose two coordinate-wise versions of the power method, from an optimization viewpoint.

## A brief review of the Power Method

- Given a matrix $A$, let its two dominant eigenvalues be $\lambda_1, \lambda_2$, and its dominant eigenvector is $v$. Power iteration conducts:

$$x^{(l+1)} \leftarrow \text{normalize}(A x^{(l)}) \qquad (2)$$

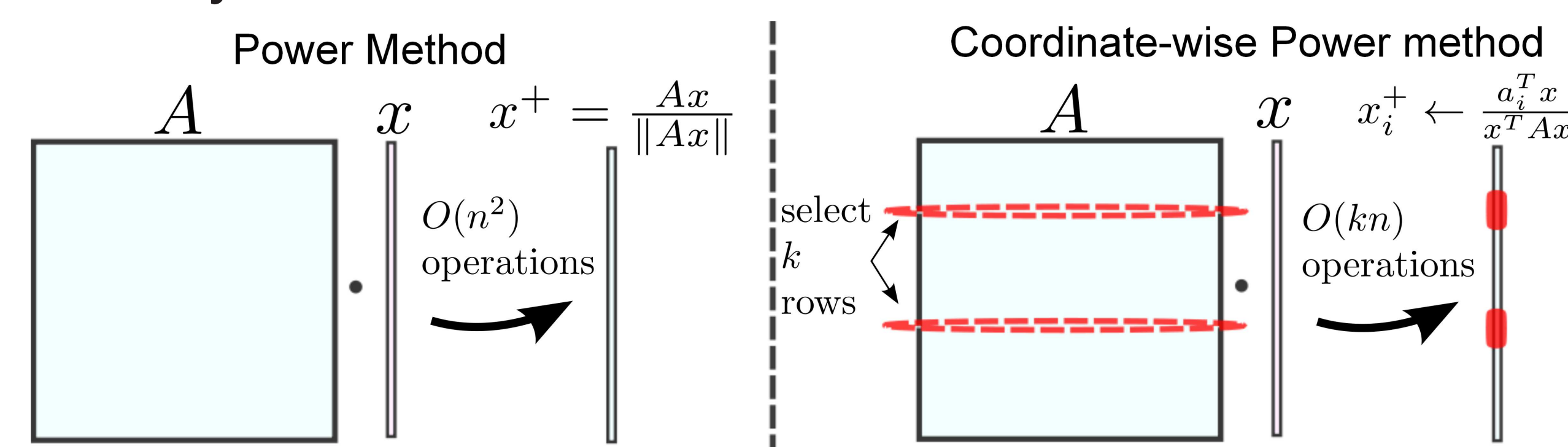- This is inefficient since some coordinates converge faster than others, e.g.,

$$A = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 3 & 0 \\ 1 & 0 & 2 \end{bmatrix}, \ x: \begin{bmatrix} 0.71 \\ 0.71 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0.53 \\ 0.80 \\ 0.27 \end{bmatrix} \rightarrow \begin{bmatrix} 0.45 \\ 0.81 \\ 0.36 \end{bmatrix} \rightarrow \begin{bmatrix} 0.42 \\ 0.82 \\ 0.39 \end{bmatrix} \rightarrow \begin{bmatrix} 0.41 \\ 0.82 \\ 0.40 \end{bmatrix}$$

Therefore we want to select and update important coordinates only.

- *One key question: how to select the coordinates?*
- *Another key problem: how to choose these coordinates without too much overhead?*

## Algorithm of Coordinate-wise Power Method (CPM)

**MAIN IDEA**: Choose $k$ coordinates with the most potential change and update them only.



1. Define auxiliary parameters:
   1.1 $z = Ax$ maintained for algorithm efficiency.
   1.2 Coordinate selection criterion: $c = \frac{z}{x^T z} - x$
2. Coordinate selection: let $\Omega$ be a set containing $k$ coordinates of $c$ with the largest magnitude.
3. Update the new iterate $x^+$:

$$y_i \leftarrow \begin{cases} \frac{z_i}{z^T z}, & i \in \Omega \\ x_i & i \notin \Omega \end{cases} \qquad x^+ \leftarrow \frac{y}{\|y\|}.$$

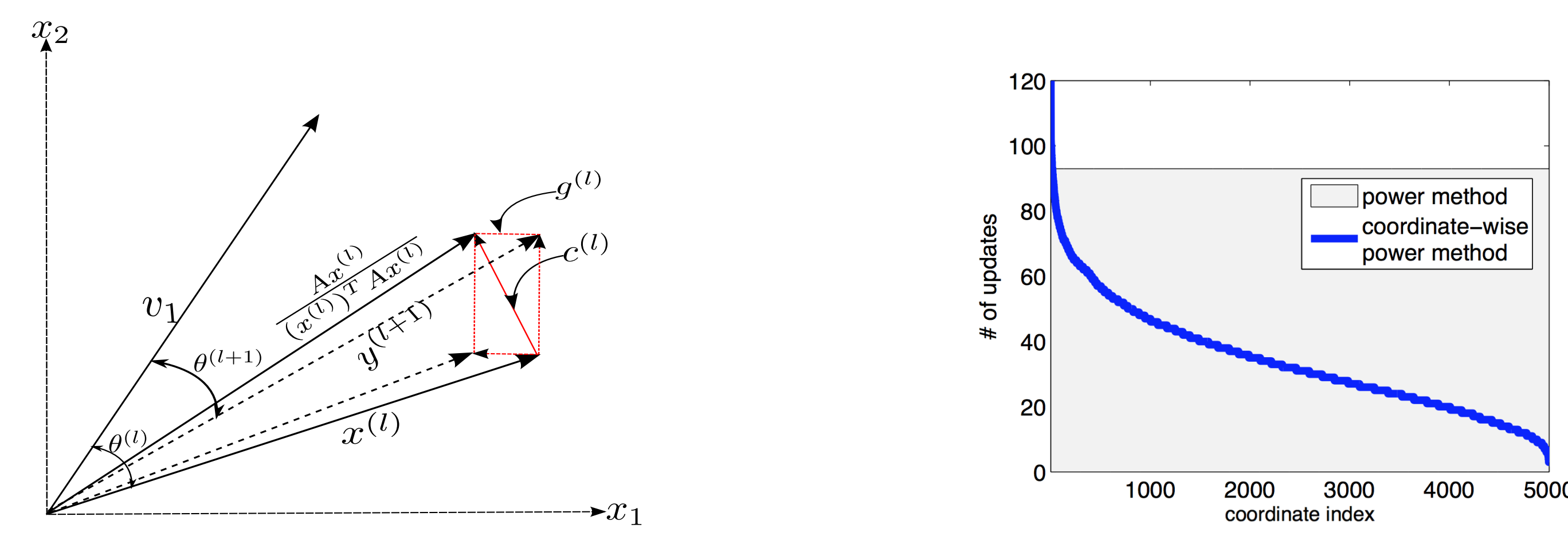4. Update the auxiliary parameters with the $k$ changes in $x$ with $O(kn)$ operations.

$$z^+ \leftarrow z + A^T_{\cdot,\Omega}(y_\Omega - x_\Omega). \qquad z^+ \leftarrow z^+/\|y\|.$$

$$c^+ \leftarrow \frac{z^+}{(x^+)^T z^+} - x^+.$$

5. Repeat $2 - 4$.

## Illustration on how CPM works



(a) Illustration on one update of CPM    (b) Number of updates of each coordinate

- (a) One iteration in CPM suffices similar result with the Power Method, but with less operations.
- (b) The unevenness of updates suggests that selecting important coordinates saves many useless updates in the Power method.

## Relation to Optimization & Coordinate selection rules

- Power method $\Longleftrightarrow$ Alternating minimization for Rank-1 matrix approximation:

$$\underset{x\in\mathbb{R},y\in\mathbb{R}}{\operatorname{argmin}} \left\{ f(x,y) = \|A - xy^T\|_F^2 \right\} \qquad (3)$$

- Updating rule for Alternation minimization:
$$x \leftarrow \operatorname{argmin}_\alpha f(\alpha, y) = \frac{Ay}{\|y\|^2}, \ y \leftarrow \operatorname{argmin}_\beta f(x,\beta) = \frac{A^T x}{\|x\|^2},$$

- The following coordinate selecting rules for (3) are equivalent:
  1. largest coordinate value change, denoted as $|\delta x_i|$;
  2. largest partial gradient (Gauss-Southwell rule), $|\nabla_i f(x)|$
  3. largest function value decrease, $|f(x + \delta x_i e_i) - f(x)|$

- A simple alternation of the objective function for Rank-1 matrix approximation for symmetric matrices:

| Algorithm | Compared to | Objective function |
|---|---|---|
| Power Method | Alternating Minimization | $f(x,y) = \|A - xy^T\|_F^2$ |
| CPM | Greedy Coordinate Descent | $f(x,y) = \|A - xy^T\|_F^2$ |
| SGCD | Greedy Coordinate Descent | $f(x) = \|A - xx^T\|_F^2$ |

## Algorithm of Symmetric Greedy Coordinate Descent(SGCD)

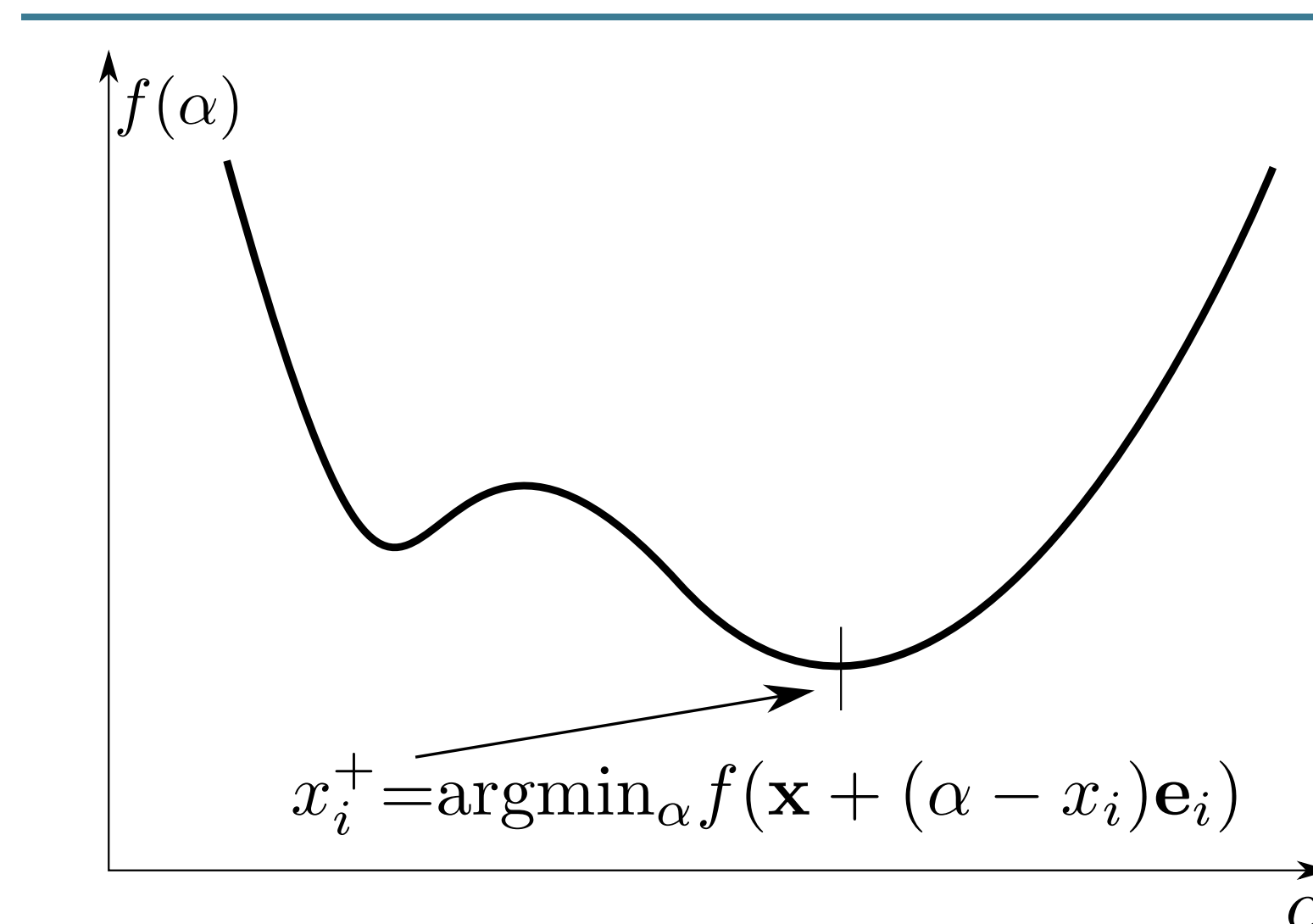- We also propose a new method we call Symetric Greedy Coordinate Descent (SGCD) for symmetric matrices.
- **MAIN IDEA**: use greedy and exact coordinate descent on $f(x) = \|A - xx^T\|_F^2$.
- Main differences:
  1. A different coordinate selection criterion: $c = \frac{Ax}{\|x\|^2} - x$ (parallel to the gradient of $f(x)$)
  2. A different update rule of $x^+$ in $\Omega$

$$x_i^+ = \begin{cases} \operatorname{argmin}_\alpha f(x + (\alpha - x_i)e_i), & \text{if } i \in \Omega, \\ x_i, & \text{if } i \notin \Omega. \end{cases}$$



- Exact update:
- Solve $x_i^+ = \alpha$ such that $\nabla f(x + (\alpha - x_i)e_i) = \alpha^3 + p\alpha + q = 0$, where
$p = \|x\|^2 - x_i^2 - a_{ii}$,
$q = -a_i^T x + a_{ii}x_i$.
- $O(n)$ operations

## Convergence guarantees for CPM and SGCD

- For Coordinate-wise Power Method (CPM), we prove global linear convergence for any positive semidefinite matrix $A$.

### Theorem 1

Convergence rate: require $T = O(\frac{\lambda_1}{\lambda_1-\lambda_2}\log(\frac{1}{\varepsilon}))$ to achieve $\tan\theta_{x^{(t)},v_1} \le \epsilon$ provided the "noise rate" $\frac{\|c_{[n]-\Omega}\|}{\|c\|} \lesssim \frac{\lambda_1-\lambda_2}{\lambda_1}$.

- For the method of Symmetric Greedy Coordinate Descent (SGCD), we prove local linear convergence:

### Theorem 2

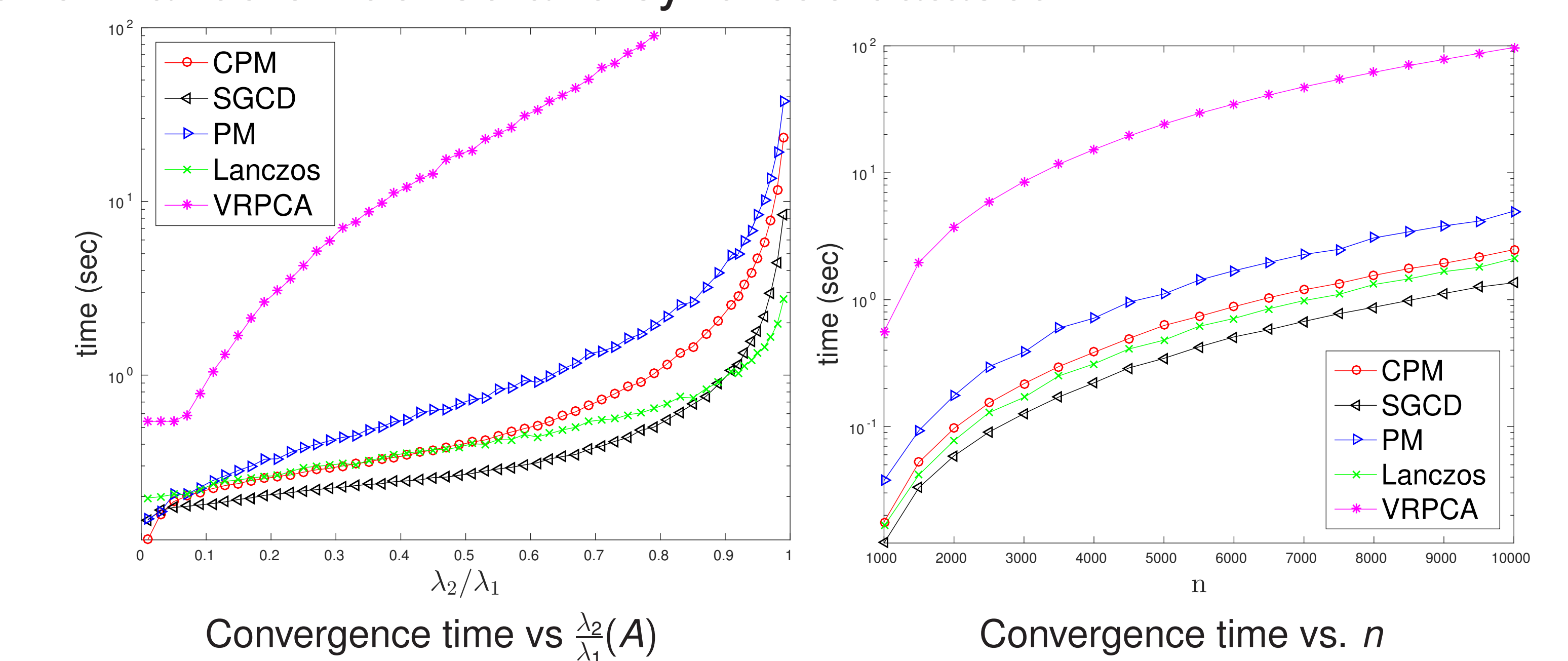Convergence rate: require $T = O(\frac{\lambda_1}{\lambda_1-\lambda_2}\log(\frac{1}{\varepsilon}))$ to achieve $f(x^{(t)}) - f(v) \le \epsilon$ provided $x^{(0)}$ sufficiently close to $v_1$: $\|x^{(0)} - v_1\| \lesssim \frac{\lambda_1-\lambda_2}{\sqrt{\lambda_1}}$.
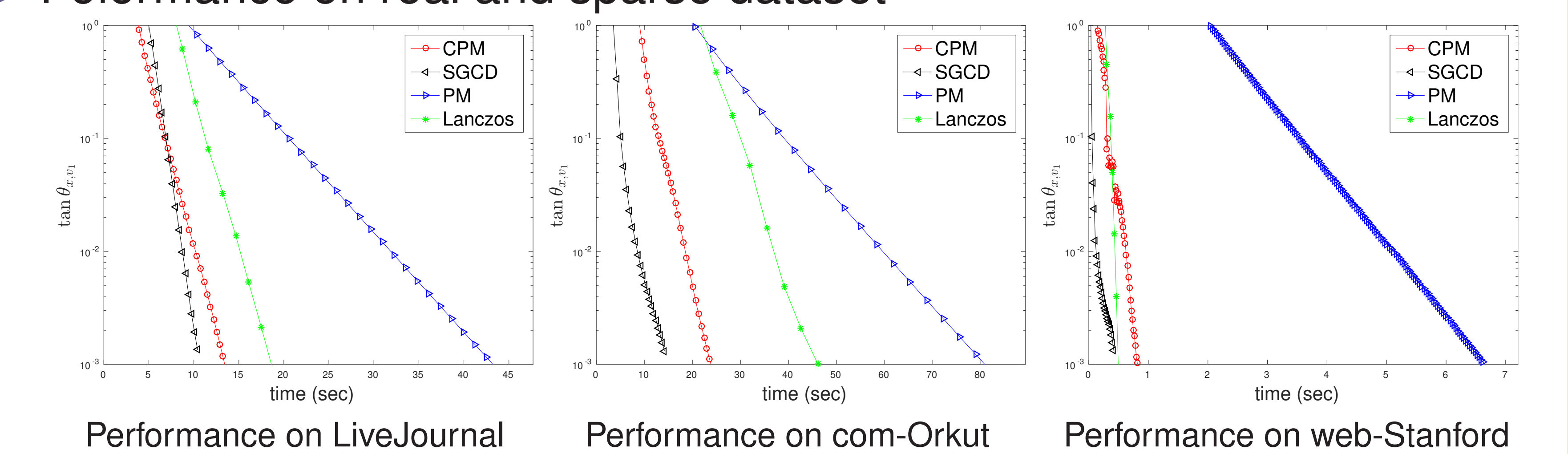
## Experimental Results

- Scalability experiments between our methods compared to power method, Lanczos method and VRPCA (Ohad Shamir, 2015) conducted with C++ with Eigen library on one machine with 16G memory:
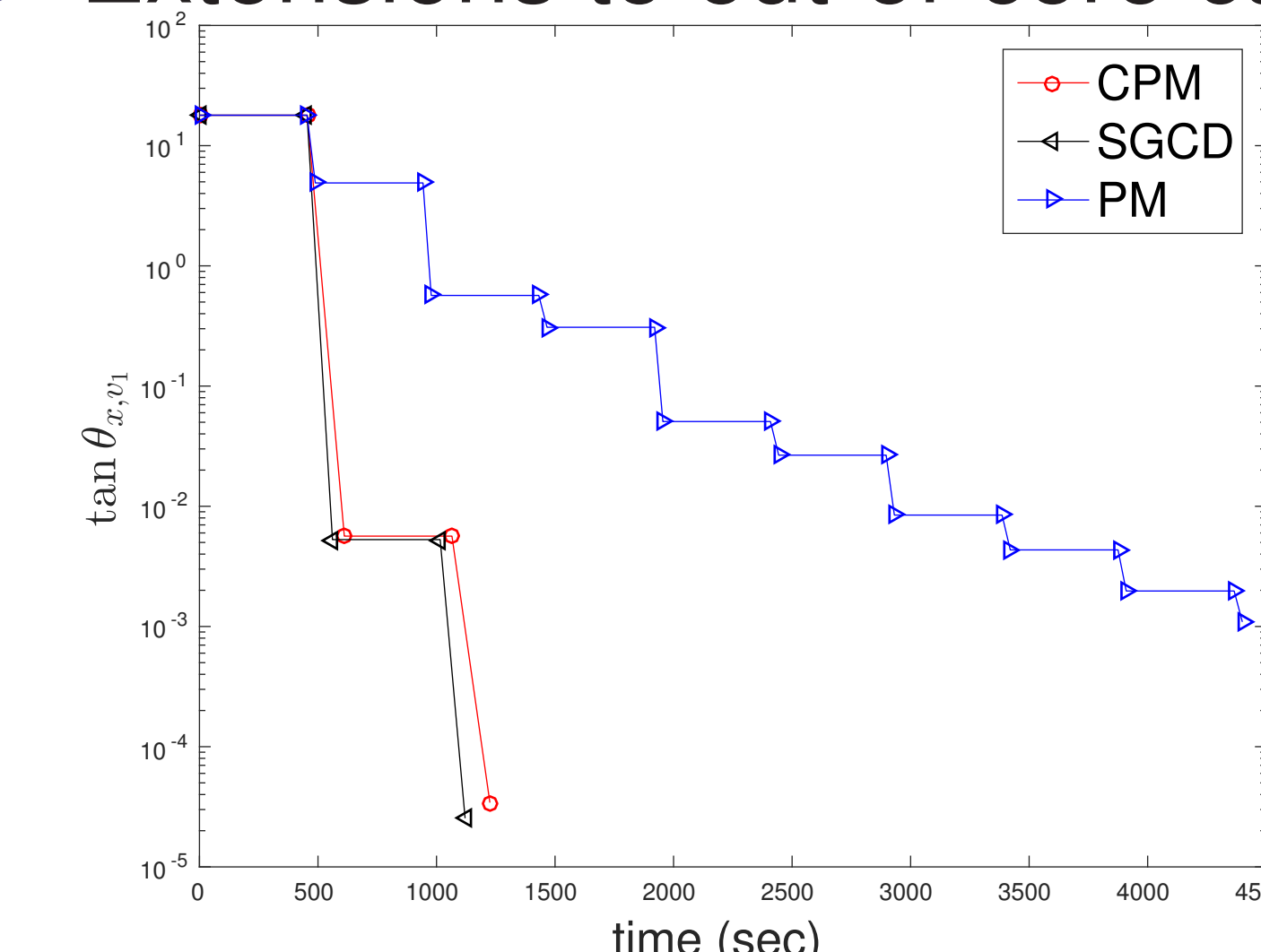
- Performance on dense and synthetic dataset



Convergence time vs $\frac{\lambda_2}{\lambda_1}(A)$    Convergence time vs. $n$

- Peformance on real and sparse dataset



Performance on LiveJournal    Performance on com-Orkut    Performance on web-Stanford

| Dataset | LiveJournal | com-Orkut | web-Stanford |
|---|---|---|---|
| # nodes | 4,847,571 | 3,072,626 | 281,903 |
| # nonzero | 86,220,856 | 234,370,166 | 3,985,272 |

- Extensions to out-of-core case:



- Existing methods can't be easily applied to out-of-core dataset.
- Our methods indicate that updating only $k$ coordinates of iterate $x$ still enhance the target direction
- we can choose a $k$ such that $k$ rows of data fit in memory and then fully update the corresponding coordinates