# An Efficient Parallel Block Coordinate Descent Algorithm for Large-scale Precision Matrix Estimation using Graphics Processing Units

Seunghwan Lee

Department of Statistics
Inha University

May 28, 2021

*Joint work with Prof. Choi, Y.-G. and Prof. Yu, D.*

# Outline

- Precision Matrix Estimation

- Convex correlation selection method (CONCORD)

- Edge-Coloring in Graph Theory

- Parallel Coordinate Descent algorithm for CONCORD using Graphics Processing Units

- Numerical study

- Summary

## Estimation of Precision Matrix $\Omega$

▶ Observed data ($n \times p$ dimensional matrix):

$$\mathbf{X} = (X^1, X^2, \ldots, X^n)^T,$$

where $X^k = (X_1^k, \cdots, X_p^k)^T \sim \mathcal{N}(\mathbf{0}, \Sigma)$.

▶ *Conditional Independence*:

$$\omega_{ij} = 0 \iff X_i \perp\!\!\!\perp X_j \mid \{X_k : k \neq i, j\},$$

where $\omega_{ij}$ is the $(i, j)$-th element of the *precision matrix* $\Omega \equiv \Sigma^{-1}$.

# Estimation of Precision Matrix $\Omega$

▶ Loglikelihood of $\Omega$

$$\log L(\Omega) = \frac{n}{2} \log \det \Omega - \frac{1}{2} \sum_{k=1}^{n} (X^k)^{\mathrm{T}} \Omega \, X^k - \frac{np}{2} \log 2\pi$$

▶ Maximum likelihood estimator (MLE) of $\Omega$

$$\widehat{\Omega} = S^{-1}, \text{ where } S = \frac{1}{n} \sum_{k=1}^{n} X^k (X^k)^T$$

▶ In high-dimensional case $p > n$, $S$ is ill-conditioned and singular.

$\Rightarrow S^{-1}$ can not be obtained from $S$.

# $\ell_1$-regularization Methods

▶ Maximum likelihood approach

$$\widehat{\Omega} = \underset{\omega^{ij}, 1 \leq i,j \leq p}{\operatorname{argmin}} \Big\{ -\log\det(\Omega) + \operatorname{tr}(S\Omega) + \lambda \sum_{\substack{i \neq j \\ (\text{or } 1 \leq i,j \leq p)}} |\omega_{ij}| \Big\}$$

subject to $\Omega$ is positive definite and symmetric,
where $S = \frac{1}{n} \sum_{k=1}^{n} X^k (X^k)^T$ and $\lambda \geq 0$.

▶ Existing methods:
Yuan and Lin (2007);
Graphical lasso (Friedman et al., 2008; Witten et al., 2012;
Mazumerder and Hastie, 2012).

# $\ell_1$-regularization Methods

- $\ell_1$-minimization approach

$$\widehat{\Omega} = \underset{\Omega}{\operatorname{argmin}} \ \|\Omega\|_1 \text{ subject to } \|S\Omega - I_p\|_\infty \leq \lambda,$$

where $S = \frac{1}{n} \sum_{k=1}^{n} X^k (X^k)^T$, $\lambda \geq 0$, $\|A\|_1 = \sum_{i=1}^{n} \sum_{j=1}^{p} |a_{ij}|$ and $\|A\|_\infty = \max_{1 \leq i \leq n, 1 \leq j \leq p} |a_{ij}|$ for $A = (a_{ij}) \in \mathbb{R}^{n \times p}$.

- Existing methods:
  CLIME (Cai et al., 2011);
  ACLIME (Cai et al., 2016).

# $\ell_1$-regularization Methods

▶ Regression approach

$$X_i = \sum_{j \neq i} \beta_{ij} X_j + \epsilon_i = \sum_{j \neq i} \rho_{ij} \sqrt{\frac{\omega_{jj}}{\omega_{ii}}} X_j + \epsilon_i,$$

where $\beta_{ij} = -\omega_{ij}/\omega_{ii} = \rho_{ij} \sqrt{\omega_{jj}/\omega_{ii}}$, and $\epsilon_i$ is uncorrelated with $X_{-[i]} = \{X_j \mid j \neq i, \ 1 \leq j \leq p\}$ and has mean 0 and variance $1/\omega_{ii}$.

$(\omega_{ij} = 0 \iff \beta_{ij} = 0 \iff \rho_{ij} = 0)$

▶ Existing methods:

Neighborhood selection (Meinshausen and Bühlmann, 2006);

SPACE (Peng et al., 2009);

Matrix inversion with scaled lasso (Sun and Zhang, 2013);

CONCORD (Khare et al., 2015).

## CONvex CORrelation selection methoD (CONCORD)

▶ Consider a log-pseudolikelihood function $L(\Omega; \lambda)$

$$L(\Omega; \lambda) = -\sum_{i=1}^{p} n \log \omega_{ii} + \frac{1}{2} \sum_{i=1}^{p} \sum_{k=1}^{n} \left( \omega_{ii} X_i^k + \sum_{j \neq i} \omega_{ij} X_j^k \right)^2 + \lambda \sum_{i<j} |\omega_{ij}|$$

▶ $L(\Omega; \lambda) = L(\omega_D, \omega_{-D}; \lambda)$ is convex for $(\omega_D, \omega_{-D})$,
where $\omega_D = (\omega_{11}, \omega_{22}, \ldots, \omega_{pp})^T$ and
$\omega_{-D} = (\omega_{12}, \omega_{13}, \ldots, \omega_{(p-1)p})^T$.

▶ It is shown that the cyclic coordinatewise minimization
guarantees the convergence to a global minimum of $L(\Omega; \lambda)$ if all
diagonal elements of the sample covariance matrix are positive.

## **Coordinate descent (CD) for CONCORD**

- ▶ Cyclic coordinatewise minimization

  **1.** Given the current solutions $\hat{\omega}_D^{(k)} = (\hat{\omega}_{ii}^{(k)}, 1 \leq i \leq p)$ at $k$-th iteration and the current updated solutions $(\hat{\omega}_{st})_{(s,t)\neq(i,j)}$,

  $$\hat{\omega}_{ij}^{(k+1)} = \underset{\omega_{ij}}{\mathrm{argmin}}\, L(\omega_{ij}; \hat{\omega}_D^{(k)}, (\hat{\omega}_{st})_{(s,t)\neq(i,j)}, \lambda)$$

  **2.** Given the current solution $\hat{\omega}_{-D}^{(k+1)} = (\hat{\omega}_{ij}^{(k+1)}, i < j)$,

  $$\hat{\omega}_D^{(k+1)} = \underset{\omega_D}{\mathrm{argmin}}\, L(\omega_D; \hat{\omega}_{-D}^{(k+1)}, \lambda)$$

  **3.** Repeat Steps 1 and 2 until convergence occurs.

# Coordinate descent (CD) for CONCORD

▶ Given the current solution $\hat{\omega}_{-D}^{(k)} = (\hat{\omega}_{ij}^{(k)}, i < j)$ at $k$-th iteration,

$$\hat{\omega}_{ii}^{(k+1)} = \frac{-\sum_{j \neq i} \hat{\omega}_{ij}^{(k)} T_{ij} + \sqrt{\left(\sum_{j \neq i} \hat{\omega}_{ij}^{(k)} T_{ij}\right)^2 + 4n T_{ii}}}{2 T_{ii}} \text{ for } i = 1, 2, \ldots, p, \quad (1)$$

where $T_{ij}$ denote the $(i, j)$th element of $\mathbf{X}^T \mathbf{X}$.

▶ Given the current solution $\hat{\omega}_D^{(k)} = (\hat{\omega}_{ii}^{(k)}, 1 \leq i \leq p)$ at $k$-th iteration,

$$\hat{\omega}_{ij}^{(k+1)} = \frac{\text{Soft}_\lambda(-\sum_{j' \neq j} \tilde{\omega}_{ij'} T_{jj'} - \sum_{i' \neq i} \tilde{\omega}_{i'j} T_{ii'})}{T_{ii} + T_{jj}} \text{ for } 1 \leq i < j \leq p, \quad (2)$$

where $\tilde{\omega}_{st}$ denotes the current iterative solution of $\omega_{st}$, $T_{ij}$ denotes the $(i, j)$th element of $\mathbf{X}^T \mathbf{X}$, $\text{Soft}_\tau(x) = \text{sign}(x)(|x| - \tau)_+$, and $(x)_+ = \max(0, x)$.

# Motivation: Parallelizable updating equations

▶ For the updating equations of $\hat{\omega}_D$, the current updated solution $\hat{\omega}_{ii}$ does not affect the other updating equations of $\hat{\omega}_{jj}$ for $j \neq i$.

▶ For the updating equations of $\hat{\omega}_{-D}$, $\hat{\omega}_{ij}$ depends on the current iterative solutions $(\hat{\omega}_{i1}, \ldots, \hat{\omega}_{i(j-1)}, \hat{\omega}_{i(j+1)}, \ldots, \hat{\omega}_{ip})$ and $(\hat{\omega}_{1j}, \ldots, \hat{\omega}_{(i-1)j}, \hat{\omega}_{(i+1)j}, \ldots, \hat{\omega}_{pj})$. Thus, the current updated solution $\hat{\omega}_{ij}$ does not affect other updating equations of $\hat{\omega}_{st}$ such that $s \neq i$ and $t \neq j$.
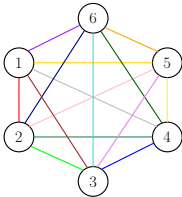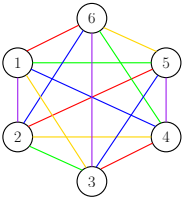
# Edge-Coloring

- Edge-Coloring:

  Edge-coloring is defined as an assignment of *colors to the edges* of a graph, in a way that any pair of adjacent edges (edges sharing at least one vertices) is colored in different colors.

- A special case of Baranyai's Theorem with a complete graph $\mathcal{K}_p$

### Theorem (Baranyai's Theorem)

*Suppose that $\mathcal{K}_p$ is an undirected complete graph with $p$ nodes. If $p$ is even, $\mathcal{K}_p$ is edge-colorable with at least $p - 1$ colors. If $p$ is odd, then $\mathcal{K}_p$ is edge-colorable with at least $p$ colors.*

# Example of Edge-Coloring

| Coloring scheme | with mutually distinct colors | with minimal number of colors |
|---|---|---|
| $\mathcal{K}_p$ with $p = 6$ |  |  |
| # of colors | $p(p-1)/2 \ = \ 15$ | $p-1 \ = \ 5$ |
| # of edge(s) for each color | 1 | $p/2 \ = \ 3$ |
| Collections of edges of the same colors $((i, j) \rightarrow ij)$ | {12}, {13}, {14}, {15}, {16}, {23}, {24}, {25}, {26}, {34}, {35}, {36}, {45}, {46}, {56} | {16, 25, 34}, {15, 23, 46}, {14, 26, 35}, {13, 24, 56}, {12, 36, 45} |

## Analogy between edge-coloring and update ordering

- ▶ Analogy between edge-coloring and update ordering:

    (A) Associate the edge-coloring of edge $(i, j)$ by $k$-th color
        with the update of $\hat{\omega}_{ij}$ as in (2) at the $k$-th step.

- ▶ Edge-coloring and CD (serial update):

    Coloring all edges by colors 1 through $p(p-1)/2$

- ▶ Edge-coloring and Parallel CD:

    Coloring multiple edges $(i, j)$ with the same $k$-th color, which means
    that the associated $\omega_{ij}$'s are simultaneously updated given the same
    current iterate.

# Circle method for edge-coloring in $\mathcal{K}_p$

- Partitioning an edge set of $\mathcal{K}_p$ as $p - 1$ ($p$) *equi-color subsets* in which edges have same color when $p$ is even (odd).

- To handle the differences between even $p$ and odd $p$, we consider a variable $p_{even}$, which is defined as $p_{even} = p$ if $p$ is even and $p_{even} = p + 1$ if $p$ is odd.

# Circle method for edge-coloring in $\mathcal{K}_p$

(i) Clockwise rotation for the round-robin table with the $(1, 1)$ element is fixed:

| 1 | 2 | 3 | $\cdots$ | $p_{even}/2$ |
|---|---|---|---|---|
| $p_{even}$ | $p_{even} - 1$ | $p_{even} - 2$ | $\cdots$ | $p_{even}/2 + 1$ |

$\rightarrow$

| 1 | $p_{even}$ | 2 | $\cdots$ | $p_{even}/2 - 1$ |
|---|---|---|---|---|
| $p_{even} - 1$ | $p_{even} - 2$ | $p_{even} - 3$ | $\cdots$ | $p_{even}/2$ |

$\rightarrow \cdots \rightarrow$

| 1 | 3 | 4 | $\cdots$ | $p_{even}/2 + 1$ |
|---|---|---|---|---|
| 2 | $p_{even}$ | $p_{even} - 1$ | $\cdots$ | $p_{even}/2 + 2$ |

(ii) Define matching pairs: for each table in (i), two indices in each columns of the given table defines a matching pair. For example, in the first table in (i), the matching pairs are
$\{(1, p_{even}), (2, p_{even} - 1), \ldots, (p_{even}/2, p_{even}/2 + 1)\}$.

(iii) Discard a pair containing $(p_{even})$ index in (ii) if $p$ is odd.

# Cyclic Reduction for $\|\Omega^{(k)} - \hat{\Omega}\|_\infty$

▶ $\|A\|_\infty = \max_{1 \le i \le n, 1 \le j \le p} |a_{ij}|$ for $A = (a_{ij}) \in \mathbb{R}^{n \times p}$.

▶ Graphical representation of cyclic reduction:

$$
\begin{pmatrix}
\hat{\omega}_{11}^{(k+1)} - \hat{\omega}_{11}^{(k)} & \hat{\omega}_{12}^{(k+1)} - \hat{\omega}_{12}^{(k)} & \hat{\omega}_{13}^{(k+1)} - \hat{\omega}_{13}^{(k)} \\
& \hat{\omega}_{22}^{(k+1)} - \hat{\omega}_{22}^{(k)} & \hat{\omega}_{23}^{(k+1)} - \hat{\omega}_{23}^{(k)} \\
& & \hat{\omega}_{33}^{(k+1)} - \hat{\omega}_{33}^{(k)}
\end{pmatrix}
\implies (d_1, d_2, d_3, d_4, d_5, d_6)
$$

# Comparison of Computation times

▶ Dimension (# of nodes): 500, 1000, 2500, 5000

▶ Sample size (# of samples): 500, 1000, 2000

▶ Tuning parameter: $\lambda^* = \lambda/n = 0.1, 0.3$

▶ Network structure: AR(2) and Scale-free networks

▶ Generate samples

$$\mathbf{X}^1, \mathbf{X}^2, \ldots, \mathbf{X}^n \sim N(0, \Sigma),$$

where $\Sigma = \Omega^{-1}$ and $\Omega$ is from a network structure.

▶ We generate 10 data sets and report the averages of computation times ($\delta_{tol} = 10^{-5}$).

▶ System specs : Intel Xeon(R) W-2175 CPU (2.50GHz, Max Turbo (4.30GHz)) and 128 GB RAM with NVIDIA GeForce GTX 1080 Ti
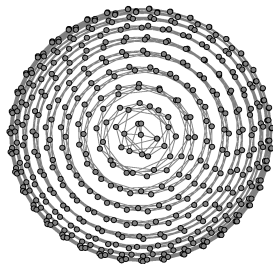
# Comparison of Computation times

- ▶ Implemented algorithms

  - ▶ CD-BLAS:

    Implemented CD algorithm using BLAS (`gconcord`)

  - ▶ CD-NAIVE:

    Implemented CD algorithm using C language only.

  - ▶ PCD-CPU:

    Implemented PCD algorithm without CUDA C (running on CPU)

  - ▶ PCD-GPU:

    Implemented PCD algorithm with CUDA C (running on GPU)

# AR(2) Network

AR(2) network has edges between $j$-th node and $(j-1)$-th node for
$j = 2, 3, \ldots, p$ and edges between $j$-th node and $(j-2)$-th node for
$j = 3, 4, \ldots, p$.

$$\Omega = (\omega_{ij})_{1 \le i,j \le p} = \begin{cases} 1 & \text{if } i = j \\ 0.45 & \text{if } |i - j| = 1 \\ 0.4 & \text{if } |i - j| = 2 \\ 0 & \text{otherwise} \end{cases}$$

## Scale-free Network

Degrees of nodes follow power law distribution having a form $P(k) \propto k^{-\alpha}$, where $P(k)$ is a fraction of nodes having $k$ connections and $\alpha$ is a preferential attachment parameter. We set $\alpha = 2.3$ and generate a scale-free network by using the Barabasi and Albert (BA) model.

$(i)\ \tilde{\Omega} = (\tilde{\omega}_{ij})_{1 \le i,j \le p} = \begin{cases} 1 & \text{if } i = j \\ U & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}$

$\qquad U \sim \mathcal{U}([-1, -0.5] \cup [0.5, 1]).$

$(ii)\ \Omega = (\omega_{ij})_{1 \le i,j \le p} = \frac{\tilde{\omega}_{ij}}{1.25 \sum_{k \ne i} |\tilde{\omega}_{ik}|}$

$(iii)\ \Omega = (\Omega + \Omega^T)/2$
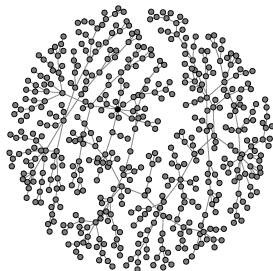
$(iv)\ \omega_{ii} = 1$ for $i = 1, 2, \ldots, p$

**Table 1:** Summary for AR(2) network with $\lambda^* = 0.1$

| $\lambda^*$ | $n$ | $p$ | Computation time (sec.) | | | | Iteration | | $|\hat{E}|$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | CD-BLAS | CD-NAIVE | PCD-CPU | PCD-GPU | CD | PCD | CD | PCD |
| 0.1 | 500 | 500 | 3.22 | 3.18 | 3.31 | 0.89 | 26.40 | 26.10 | 1976.70 | 1976.70 |
| | | | (0.02) | (0.02) | (0.02) | (0.01) | (0.16) | (0.18) | (9.52) | (9.52) |
| | | 1000 | 13.09 | 26.28 | 28.66 | 4.87 | 26.90 | 26.80 | 5078.90 | 5078.90 |
| | | | (0.09) | (0.17) | (0.18) | (0.04) | (0.18) | (0.20) | (14.19) | (14.19) |
| | | 2500 | 86.01 | 451.03 | 513.18 | 56.72 | 27.30 | 27.10 | 20327.30 | 20327.50 |
| | | | (0.83) | (4.30) | (3.28) | (0.38) | (0.26) | (0.18) | (45.87) | (45.90) |
| | | 5000 | 378.04 | 3646.43 | 4149.79 | 404.75 | 27.70 | 27.30 | 64307.20 | 64307.40 |
| | | | (2.94) | (19.49) | (51.93) | (2.29) | (0.15) | (0.15) | (69.26) | (69.13) |
| | 1000 | 500 | 2.07 | 3.16 | 3.29 | 0.88 | 25.70 | 25.50 | 1407.20 | 1407.20 |
| | | | (0.01) | (0.02) | (0.02) | (0.01) | (0.15) | (0.17) | (5.05) | (5.05) |
| | | 1000 | 25.90 | 25.84 | 28.27 | 4.76 | 26.20 | 26.10 | 2825.20 | 2825.20 |
| | | | (0.14) | (0.13) | (0.22) | (0.03) | (0.13) | (0.18) | (4.50) | (4.50) |
| | | 2500 | 167.90 | 428.93 | 490.24 | 54.91 | 26.20 | 26.20 | 7216.80 | 7216.80 |
| | | | (1.60) | (3.38) | (4.18) | (0.28) | (0.13) | (0.13) | (6.92) | (6.92) |
| | | 5000 | 694.62 | 3419.50 | 3862.95 | 385.66 | 26.20 | 26.00 | 14806.20 | 14806.20 |
| | | | (3.87) | (17.43) | (17.96) | (0.05) | (0.13) | (0.14) | (14.28) | (14.28) |
| | 2000 | 500 | 2.12 | 3.19 | 3.36 | 0.85 | 25.00 | 25.10 | 1393.10 | 1393.10 |
| | | | (0.00) | (0.00) | (0.01) | (0.00) | (0.00) | (0.10) | (3.74) | (3.74) |
| | | 1000 | 21.81 | 25.76 | 27.88 | 4.65 | 25.70 | 25.40 | 2803.10 | 2803.10 |
| | | | (0.39) | (0.16) | (0.26) | (0.03) | (0.15) | (0.16) | (4.70) | (4.70) |
| | | 2500 | 342.84 | 433.35 | 495.12 | 54.54 | 25.90 | 25.90 | 7006.90 | 7006.90 |
| | | | (1.26) | (1.65) | (1.93) | (0.21) | (0.10) | (0.10) | (9.79) | (9.79) |
| | | 5000 | 1389.95 | 3440.95 | 3929.09 | 386.23 | 26.00 | 26.00 | 14009.30 | 14009.40 |
| | | | (6.07) | (13.14) | (11.22) | (0.12) | (0.00) | (0.00) | (9.38) | (9.35) |

**Table 2:** Summary for AR(2) network with $\lambda^* = 0.3$

| $\lambda^*$ | $n$ | $p$ | Computation time (sec.) | | | | Iteration | | $|\hat{E}|$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | CD-BLAS | CD-NAIVE | PCD-CPU | PCD-GPU | CD | PCD | CD | PCD |
| 0.3 | 500 | 500 | 1.69 | 1.70 | 1.73 | 0.50 | 13.90 | 13.40 | 859.50 | 859.50 |
| | | | (0.06) | (0.05) | (0.06) | (0.02) | (0.46) | (0.52) | (5.00) | (5.00) |
| | | 1000 | 7.00 | 14.19 | 14.80 | 2.55 | 14.40 | 13.70 | 1722.20 | 1722.20 |
| | | | (0.13) | (0.26) | (0.45) | (0.07) | (0.27) | (0.40) | (5.87) | (5.87) |
| | | 2500 | 45.52 | 240.36 | 267.74 | 29.61 | 14.50 | 14.10 | 4297.80 | 4297.80 |
| | | | (0.72) | (3.57) | (3.37) | (0.38) | (0.22) | (0.18) | (7.12) | (7.12) |
| | | 5000 | 210.20 | 1937.84 | 2289.35 | 215.59 | 14.80 | 14.50 | 8624.10 | 8624.10 |
| | | | (4.05) | (38.32) | (23.56) | (2.47) | (0.29) | (0.17) | (17.08) | (17.08) |
| | 1000 | 500 | 1.06 | 1.59 | 1.62 | 0.46 | 12.40 | 12.10 | 853.60 | 853.60 |
| | | | (0.02) | (0.03) | (0.03) | (0.01) | (0.22) | (0.23) | (4.66) | (4.66) |
| | | 1000 | 12.06 | 12.31 | 13.05 | 2.22 | 12.20 | 11.80 | 1698.00 | 1698.00 |
| | | | (0.13) | (0.13) | (0.15) | (0.02) | (0.13) | (0.13) | (5.80) | (5.80) |
| | | 2500 | 78.69 | 203.18 | 225.83 | 25.51 | 12.50 | 12.10 | 4268.10 | 4268.10 |
| | | | (1.56) | (3.53) | (4.36) | (0.48) | (0.22) | (0.23) | (11.70) | (11.70) |
| | | 5000 | 350.79 | 1718.20 | 1901.72 | 182.99 | 13.00 | 12.30 | 8521.50 | 8521.50 |
| | | | (7.29) | (35.40) | (39.33) | (3.18) | (0.30) | (0.21) | (11.65) | (11.65) |
| | 2000 | 500 | 1.12 | 1.64 | 1.62 | 0.45 | 11.80 | 11.00 | 854.20 | 854.20 |
| | | | (0.01) | (0.02) | (0.00) | (0.01) | (0.13) | (0.00) | (3.14) | (3.14) |
| | | 1000 | 10.81 | 12.22 | 12.76 | 2.11 | 11.60 | 11.10 | 1717.00 | 1717.00 |
| | | | (0.21) | (0.16) | (0.11) | (0.02) | (0.16) | (0.10) | (4.96) | (4.96) |
| | | 2500 | 159.15 | 204.24 | 216.04 | 24.00 | 12.00 | 11.30 | 4291.10 | 4291.10 |
| | | | (0.09) | (0.27) | (1.86) | (0.32) | (0.00) | (0.15) | (5.94) | (5.94) |
| | | 5000 | 637.67 | 1597.47 | 1796.61 | 171.11 | 12.10 | 11.50 | 8578.20 | 8578.30 |
| | | | (5.89) | (15.28) | (32.73) | (2.46) | (0.10) | (0.17) | (14.48) | (14.51) |

**Table 3:** Summary for scale-free network with $\lambda^* = 0.1$

| $\lambda^*$ | $n$ | $p$ | Computation time (sec.) | | | | Iteration | | $|\hat{E}|$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | CD-BLAS | CD-NAIVE | PCD-CPU | PCD-GPU | CD | PCD | CD | PCD |
| 0.1 | 500 | 500 | 1.33 | 1.35 | 1.53 | 0.46 | 11.30 | 12.20 | 2348.30 | 2348.30 |
| | | | (0.02) | (0.02) | (0.04) | (0.01) | (0.15) | (0.29) | (13.88) | (13.88) |
| | | 1000 | 5.64 | 11.41 | 13.18 | 2.34 | 11.90 | 12.60 | 8090.90 | 8090.90 |
| | | | (0.15) | (0.30) | (0.28) | (0.05) | (0.31) | (0.27) | (22.41) | (22.41) |
| | | 2500 | 43.59 | 228.30 | 269.74 | 30.84 | 14.20 | 14.60 | 42601.50 | 42601.50 |
| | | | (1.59) | (8.18) | (9.14) | (1.05) | (0.51) | (0.50) | (34.34) | (34.34) |
| | | 5000 | 193.67 | 1872.72 | 2345.42 | 230.10 | 14.10 | 15.50 | 144508.00 | 144507.70 |
| | | | (3.84) | (37.21) | (39.78) | (3.98) | (0.28) | (0.27) | (58.82) | (58.88) |
| | 1000 | 500 | 0.96 | 1.44 | 1.61 | 0.48 | 11.60 | 12.50 | 598.20 | 598.20 |
| | | | (0.02) | (0.03) | (0.03) | (0.01) | (0.27) | (0.27) | (5.46) | (5.46) |
| | | 1000 | 10.76 | 11.00 | 12.49 | 2.19 | 11.20 | 11.70 | 1340.00 | 1340.00 |
| | | | (0.28) | (0.27) | (0.22) | (0.04) | (0.29) | (0.21) | (4.94) | (4.94) |
| | | 2500 | 92.85 | 235.54 | 269.49 | 29.01 | 13.90 | 13.70 | 4471.70 | 4471.70 |
| | | | (2.22) | (5.80) | (3.32) | (0.32) | (0.31) | (0.15) | (18.16) | (18.16) |
| | | 5000 | 369.05 | 1829.25 | 2121.88 | 209.48 | 13.80 | 14.10 | 12557.50 | 12557.60 |
| | | | (7.76) | (38.16) | (58.87) | (6.02) | (0.29) | (0.41) | (21.90) | (21.94) |
| | 2000 | 500 | 1.09 | 1.60 | 1.77 | 0.49 | 11.90 | 12.80 | 506.50 | 506.50 |
| | | | (0.01) | (0.02) | (0.02) | (0.01) | (0.18) | (0.20) | (0.50) | (0.50) |
| | | 1000 | 10.62 | 11.97 | 12.90 | 2.18 | 11.70 | 11.60 | 1011.00 | 1011.00 |
| | | | (0.25) | (0.20) | (0.32) | (0.05) | (0.21) | (0.31) | (1.02) | (1.02) |
| | | 2500 | 187.80 | 239.21 | 257.65 | 28.62 | 14.50 | 13.50 | 2517.50 | 2517.50 |
| | | | (5.47) | (6.88) | (4.20) | (0.47) | (0.43) | (0.22) | (1.56) | (1.56) |
| | | 5000 | 680.01 | 1705.53 | 2078.65 | 209.39 | 13.10 | 14.10 | 5045.90 | 5045.90 |
| | | | (9.11) | (23.24) | (26.84) | (2.66) | (0.18) | (0.18) | (2.74) | (2.74) |

**Table 4:** Summary for scale-free network with $\lambda^* = 0.3$

| $\lambda^*$ | $n$ | $p$ | Computation time (sec.) | | | | Iteration | | $|\hat{E}|$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | CD-BLAS | CD-NAIVE | PCD-CPU | PCD-GPU | CD | PCD | CD | PCD |
| 0.3 | 500 | 500 | 1.03 | 1.06 | 1.14 | 0.37 | 8.80 | 9.00 | 364.70 | 364.70 |
| | | | (0.02) | (0.02) | (0.02) | (0.00) | (0.13) | (0.15) | (1.69) | (1.69) |
| | | 1000 | 4.43 | 9.06 | 10.14 | 1.80 | 9.40 | 9.60 | 713.30 | 713.30 |
| | | | (0.08) | (0.16) | (0.22) | (0.04) | (0.16) | (0.22) | (2.13) | (2.13) |
| | | 2500 | 34.08 | 176.28 | 208.13 | 22.30 | 10.50 | 10.50 | 1755.40 | 1755.40 |
| | | | (0.78) | (3.97) | (3.55) | (0.46) | (0.27) | (0.22) | (5.31) | (5.31) |
| | | 5000 | 138.34 | 1343.89 | 1568.08 | 157.45 | 10.40 | 10.60 | 3569.40 | 3569.40 |
| | | | (2.87) | (27.90) | (43.39) | (4.52) | (0.22) | (0.31) | (6.12) | (6.12) |
| | 1000 | 500 | 0.77 | 1.15 | 1.23 | 0.37 | 9.00 | 9.30 | 367.80 | 367.80 |
| | | | (0.00) | (0.00) | (0.02) | (0.00) | (0.00) | (0.15) | (1.50) | (1.50) |
| | | 1000 | 8.65 | 8.94 | 9.75 | 1.71 | 9.00 | 9.00 | 715.00 | 715.00 |
| | | | (0.15) | (0.14) | (0.16) | (0.03) | (0.15) | (0.15) | (2.09) | (2.09) |
| | | 2500 | 68.93 | 176.52 | 198.19 | 21.84 | 10.50 | 10.50 | 1758.80 | 1758.80 |
| | | | (1.09) | (2.82) | (3.00) | (0.32) | (0.17) | (0.15) | (2.63) | (2.63) |
| | | 5000 | 267.18 | 1323.10 | 1530.87 | 150.17 | 9.90 | 10.10 | 3582.60 | 3582.60 |
| | | | (2.69) | (13.37) | (15.77) | (1.48) | (0.10) | (0.10) | (3.96) | (3.96) |
| | 2000 | 500 | 0.87 | 1.26 | 1.33 | 0.38 | 9.00 | 9.10 | 367.30 | 367.30 |
| | | | (0.00) | (0.00) | (0.01) | (0.00) | (0.00) | (0.10) | (1.24) | (1.24) |
| | | 1000 | 8.21 | 9.53 | 10.43 | 1.75 | 9.10 | 9.20 | 712.70 | 712.70 |
| | | | (0.07) | (0.09) | (0.14) | (0.02) | (0.10) | (0.13) | (1.73) | (1.73) |
| | | 2500 | 134.24 | 173.16 | 199.99 | 22.10 | 10.40 | 10.40 | 1760.00 | 1760.00 |
| | | | (2.07) | (2.64) | (3.10) | (0.34) | (0.16) | (0.16) | (3.50) | (3.50) |
| | | 5000 | 513.17 | 1294.45 | 1482.85 | 148.71 | 9.90 | 10.00 | 3585.10 | 3585.10 |
| | | | (5.14) | (13.08) | (1.64) | (0.00) | (0.10) | (0.00) | (4.13) | (4.13) |

# Summary

- We proposed the efficient parallel coordinate descent algorithm for CONCORD that simultaneously updates $p_{even}/2$ off-diagonal elements, which is $p/2$ for an even $p$ and $(p-1)/2$ for an odd $p$.

- We also showed that $p_{even}/2$ is the maximum of the number of simultaneously updatable elements in the CONCORD-CD algorithm by applying the theoretical results in the edge-coloring.

- Comprehensive numerical studies show that the proposed CONCORD-PCD algorithm is adequate for GPU-parallel computation and more efficient than the original CONCORD-CD algorithm for large datasets.

# References.

– Cai, T. T., Liu, W., and Luo, X. (2011). A constrained l1 minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association*, **106**(494), 594-607.

– Cai, T. T., Liu, W., and Zhou, H. H. (2016). Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation. *The Annals of Statistics*, **44**(2), 455-488.

– Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, **9** (3), 432-441.

– Khare, K., Oh, S.-Y. and Rajaratnam, B. (2015). A convex pseudolikelihood framework for high dimensional partial correlation estimation with convergence guarantees. *Journal of the Royal Statistical Society:Series B*, **77**, 803-825.

– Mazumder, R. and Hastie, T. (2012). The graphical lasso: New insights and alternatives. *Electronic Journal of Statistics*, **6**(August), 2125-2149.

– Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, **34**(3), 1436-1462.

– Peng, J., Wang, P., Zhou, N., and Zhu, J. (2009). Partial correlation estimation by joint sparse regression models. *Journal of the American Statistical Association*, **104**(486), 735-746.

– Sun, T. and Zhang, C. H. (2013). Sparse matrix inversion with scaled lasso. *Journal of Machine Learning Research*, **14**, 3385-3418.

– Witten, D. M., Friedman, J. H., and Simon, N. (2011). New insights and faster computations for the graphical lasso. *Journal of Computational and Graphical Statistics*, **20**(4), 892-900.

– Yuan, M., and Lin, Y. (2007). Model Selection and Estimation in the Gaussian Graphical Model. *Biometrika*, **94**(1), 19-35.

# Cyclic Reduction for $\|\Omega^{(k)} - \hat{\Omega}\|_\infty$

▶ Cyclic reduction

Let $\theta = (\theta_1, \ldots, \theta_m) = \text{vech}(\Omega)$, which is a half-vectorization for the parameter $\Omega$, and $\mathbf{d} = (d_j)_{1 \leq j \leq m} = \theta^{(k)} - \hat{\theta}$.

  ▶ Initialization:

  for $q = z - 1$,
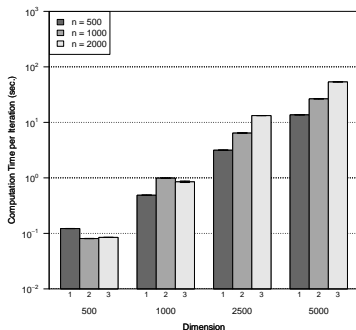  $$d_j \leftarrow \begin{cases} \max(|d_j|, |d_{j+2^q}|) & \text{if } j + 2^q \leq m \\ d_j & \text{if } j + 2^q > m \end{cases} \text{ for } j = 1, \ldots, 2^q,$$

  where $z = \lceil \log_2(m) \rceil$, where $\lceil x \rceil$ is the smallest integer greater than or equal to $x$.

  ▶ Parallel Update:
  for $q = z - 2, \ldots, 0$,

  $$d_j \leftarrow \max(|d_j|, |d_{j+2^q}|) \text{ for } j = 1, \ldots, 2^q,$$

# Computation times per iteration for AR(2)



(a) CD-BLAS, $\lambda^* = 0.1$        (b) PCD-GPU, $\lambda^* = 0.1$

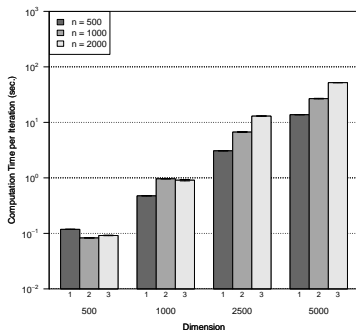# Computation times per iteration for AR(2)
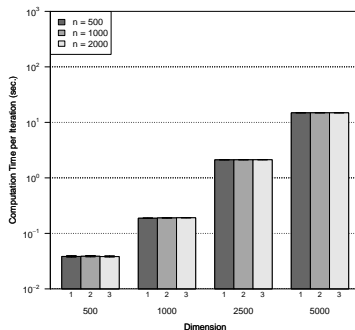


(c) CD-BLAS, $\lambda^* = 0.3$      (d) PCD-GPU, $\lambda^* = 0.3$

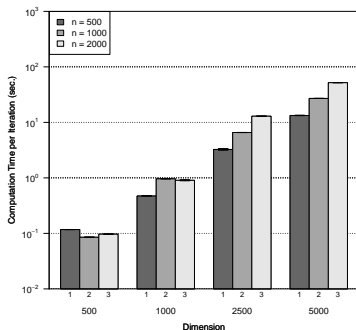# Computation times per iteration for Scale-free



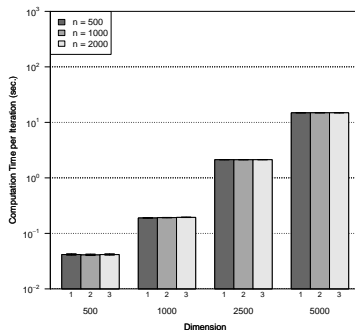(a) CD-BLAS, $\lambda^* = 0.1$

(b) PCD-GPU, $\lambda^* = 0.1$

# Computation times per iteration for Scale-free



(c) CD-BLAS, $\lambda^* = 0.3$

(d) PCD-GPU, $\lambda^* = 0.3$