



An Efficient Gauss-Newton Algorithm for Symmetric Low-Rank Product Matrix Approximations

Xin Liu

State Key Laboratory of Scientific and Engineering Computing
Institute of Computational Mathematics and Scientific/Engineering Computing
Academy of Mathematics and Systems Science
Chinese Academy of Sciences, China

(Joint work with **Zaiwen Wen**¹, **Yin Zhang**²)

2014 Workshop on Optimization for Modern Computation, BICMR, Peking University
September 3, 2014

¹Peking University, China

²Rice University, U.S.

Low-rank approximation: given $B \in \mathbb{R}^{n \times m}$ and $k < \min(n, m)$,

$$\min_{X, Z} \{ \|XZ^T - B\|_F : X \in \mathbb{R}^{n \times k}, Z \in \mathbb{R}^{m \times k} \}$$

Closely related to SVD

- k -Dominant (k -D) SVD of $n \times m \Rightarrow$ solution
- Solution + QR($n \times k$) and SVD ($k \times k$) \Rightarrow k -D SVD of $n \times m$

How about the symmetric case? for $A = A^T \in \mathbb{R}^{n \times n}$ (e.g., $A = BB^T$),

$$\min_X \{ \|XX^T - A\|_F : X \in \mathbb{R}^{n \times k} \}$$

Symmetric Low-Rank Product Optimization

A nonlinear, nonconvex least squares problem

$$\min_{X \in \mathbb{R}^{n \times k}} \|XX^T - A\|_F^2$$

Fundamental in low-rank matrix approximations

- Principal subspace of A :

$$\text{span}(X^*) = \text{span}\{q_1, q_2, \dots, q_k\}$$

where X^* is a global minimizer, and $\{q_j\}_{j=1}^k$ are k dominant eigenvectors of A .

- For $A = BB^T$, columns of X are “principal components” of B .

Most Established Eigensolvers

Why not just call `eigs` (or `svds`) in MATLAB? (ARPACK)

Why not use one of the existing eigensolvers?

Emerging applications demand new capacities.

- high efficiency at moderate accuracy
- high eigenspace dimensions
- high parallel scalability
- warm-start capacity

Established eigensolvers often lack in one or more aspects.

Advanced scientific computing and evolving computer architectures call for new algorithms (either of general or special purpose).

Block Methods

Block vs. Sequential (Lanczos-type Methods)

- Block SpMV: $AV = [Av_1 \ Av_2 \ \dots \ Av_k]$
- Sequential SpMv's: $Av \rightarrow A^2v \dots \rightarrow A^k v$
(+ inner products for orthogonalization)

As k increases, block methods are gaining advantages.

Block methods can be **warm-started** in an iterative setting.

Classic Block Method SSI: (power method)

$$X^{i+1} = \text{orth}(AX^i)$$

Other block algorithms:

- Block Jacobian-Davidson: Feast
- Trace minimization: LOBPCG, LMSVD

Research on block methods seems still largely unsettled.

Trace Minimization

$$\max_{X \in \mathbb{R}^{n \times k}} \text{tr}(X^T A X) \quad \text{s.t.} \quad X^T X = I.$$

- L.-Wen-Zhang, **L**imited **M**emory Block Krylov Subspace Optimization for Computing Dominant **S**ingular **V**alue **D**ecompositions, *SIAM Journal on Scientific Computing*, 35-3(2013);
- $A := BB^T$, main iteration of LMSVD: $X^{(i+1)} = \text{orth}(AY^{(i)})$, where $Y^{(i)} = \text{argmax}\{\phi(X) \mid X^T X = I, X \in \mathcal{S}_k\}$,
 $\mathcal{S}_k = \text{span}\{X_k, X_{k-1}, \dots, X_{k-p}\}$;
- **LMSVD** code is available at [MathWorks](#) (Google: LMSVD).

Two main types of operations: AX & $RR/orth$

As k increases, $AX \ll RR/orth \rightarrow$ bottleneck

Parallel Scalability

- $AX \rightarrow AX_1 \cup AX_2 \cup \dots \cup AX_k$. Higher.
- $RR/orth$ inherits sequentiality. Lower.

Avoid bottleneck?

- Do less $RR/orth$

No free lunch?

- Do more $BLAS3$ (higher scalability than AX)

Orthogonal Free Models

Unconstrained Model: $\min_{X \in \mathbb{R}^{n \times k}} \frac{1}{4} \|X^T X\|_F^2 + \frac{1}{2} \text{tr}(X^T A X)$, Dai-Jiang-Cui, 2013

Trace-penalty Minimization

$$\min_{X \in \mathbb{R}^{n \times k}} f(X) := \frac{1}{2} \text{tr}(X^T A X) + \frac{\mu}{4} \|X^T X - I\|_F^2.$$

EIGPEN, Wen-Yang-L.-Zhang, 2013, available at "optimization online"

Good properties:

- Equivalence to Trace Minimization does NOT require $\mu \rightarrow \infty$
- No non-global local minimizer, less undesired saddle point
- RR/Orth mostly replaced by "big" BLAS3
- Efficient for moderate accuracy, numerically stable
- Parallel scalability appears promising

Algorithm

- Gradient method with Barzilai Borwein stepsize
- Rayleigh-Ritz restart

Why New Approach?

Questions to EIGPEN

- Gradient method + BB (How about high-order methods?)

- Condition number: $k = 1$: $\kappa(\nabla^2 f_\mu(\hat{X})) = \frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1}$; $k > 1$:

$\kappa(\nabla^2 f_\mu(\hat{X})) = 0$, (How about linear convergence rate?)

$$\begin{aligned} \kappa(\nabla^2 f_\mu(\hat{X})|_{Q_k^\perp}) &\triangleq \frac{\max_{S \in \mathbb{R}^{n \times k}} \{ \text{tr}(S^T \nabla^2 f_\mu(\hat{X})(S)) : \text{tr}(S^T S) = 1, S^T Q_k = 0 \}}{\min_{S \in \mathbb{R}^{n \times k}} \{ \text{tr}(S^T \nabla^2 f_\mu(\hat{X})(S)) : \text{tr}(S^T S) = 1, S^T Q_k = 0 \}} \\ &= \frac{\lambda_n - \lambda_1}{\lambda_{k+1} - \lambda_k}. \quad (\text{Explain why RR restart is useful.}) \end{aligned}$$

- μ should be tuned in properly. (How to avoid μ ?)

Solution

SLRP: $\min_{X \in \mathbb{R}^{n \times k}} \|XX^T - A\|_F^2$

- Let eigenvalues of A be in a descending order

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$$

- Eigenvalue decomposition:

$$A = Q_n \Lambda_n Q_n^T, \quad Q_n^T Q_n = I, \quad \Lambda_n \text{ diagonal}$$

- k -D principal eigenspace:

$$\text{span}(Q_k) \triangleq \text{span}\{q_1, q_2, \dots, q_k\}$$

- k -D principal eigenpair: (Q_k, Λ_k)

Equivalence

Assume that $A = A^T \in \mathbb{R}^{n \times n}$ such that $\lambda_k \geq 0$. Then $X \in \mathbb{R}^{n \times k}$ is a solution to $\min \|XX^T - A\|_F^2$ if and only if it has SVD:

$$X = Q_k \Lambda_k^{\frac{1}{2}} V^T,$$

where (Q_k, Λ_k) is a k -D principal eigenpair of A , $V \in \mathbb{R}^{k \times k}$ is orthogonal but otherwise arbitrary.

1st-order condition for SLRP:

$$AX = X(X^T X)$$

Stationary points span invariant subspaces.

Gauss-Newton Review

Nonlinear Least Squares Model:

$$\min_{x \in \mathbb{R}^n} f(x) \triangleq \frac{1}{2} r(x)^T r(x). \quad r(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m.$$

Linearize: $r(x + s) \approx r(x) + J(x)s$, where $J(x)$ is the Jacobian.

Normal equations + Line Search:
(minimize the linear approximation)

$$J(x)^T J(x)s = -J(x)^T r(x). \quad x = x + \alpha s.$$

Some properties:

- Fast for small residual. Slow for large residual.
- Local convergence may require $\alpha < 1$ all the time.

SLRP: Nonlinear Least Squares Model

$$\min_{X \in \mathbb{R}^{n \times k}} f(X) \triangleq \frac{1}{2} \|R(X)\|_F^2. \quad R(X) \triangleq XX^T - A.$$

Let $J(X) : \mathbb{R}^{n \times k} \rightarrow \mathbb{R}^{n \times n}$ be the Jacobian operator of $R(X)$ at X .

Normal equations: (size $nk \times nk$)

$$J(X)^T J(X)(S) = -J(X)^T (R(X)).$$

Infinitely many solutions since $J(X)$ is rank deficient.

Special structure of normal equations allows low-cost solution:

$$SX^T X + XS^T X = AX - X(X^T X)$$

GN Direction

Let $X \in \mathbb{R}^{n \times k}$ be full rank, and $\mathcal{P}_X = X(X^T X)^{-1} X^T$. Then

$$S_C = \left(I - \frac{1}{2} \mathcal{P}_X \right) \left(AX(X^T X)^{-1} - X \right) + XC,$$

where $C^T = -C$, satisfies the normal equations. In particular, for $C = 0$,

$$S_0 = \left(I - \frac{1}{2} \mathcal{P}_X \right) \left(AX(X^T X)^{-1} - X \right)$$

is a minimum weighted-norm Gauss-Newton direction at X .

SLRPGN (Theoretical Version)

Gauss-Newton (GN):

- While not “converged”, do
 - ① If $\sigma_{\min}(X) < \delta$, set $X = X + P$; — **Correction Step**
 - ② Select $\alpha = \min\left(1, \sigma_{\min}^3(X) / \|\nabla f(X)\|_F\right)$, set $X = X + \alpha S_0$.

Calculation GN step:

- $Y = X(X^T X)^{-1}$, $G = AY - X$
- $S_0 = G - X(Y^T G)/2$

Computational cost:

- 1 block SpMV: AY
- 3 dense matrix multiplications
- 1 $k \times k$ linear system with n rhs

SLRPGN (Practical Implementation)

So far, in practice

- $\alpha = 1$ appears always to work well;
- Correction step can hardly be invoked.

$[X, Y] = GN(A, X)$

- While not “converged”, do
 - 1 $Y = X(X^T X)^{-1}$
 - 2 $X = AY - X(Y^T AY - I)/2$

Simple Algorithm

- Two-liner with no parameters
- No orthogonalization
- No Rayleigh-Ritz (unless eigenpairs are required)

Step Size and Correction Step

Full Rankness: $\sigma_{\min}(X^{i+1}) \geq 0.75 \sigma_{\min}(X^i)$

Correction Step:

- $\delta \leq \left(\frac{\lambda_n/\lambda_1}{4+\sqrt{20}} \right) \sqrt{\frac{\lambda_n}{k}}$
- $X_c = X + P$ ($:= \sqrt{\frac{\lambda_n}{\rho}} UV_p^T$, where $U^T X = 0$ and $U^T U = I$)

Key properties:

- $\sigma_{\min}(X_c) \geq \delta$
- $f(X_c) < f(X) - \frac{1}{4} \lambda_n^2$

Convergence Results

Theorem (Global Convergence)

Suppose that $A > 0$. Let $\{X^i\}$ be generated by **SLRPGN(TH)** starting from a full-rank initial point. Then after finite number of iterations, step-size $\alpha = 1$ will always be taken, **no more correction step**, and $\nabla f(X_j) \rightarrow 0$.

$f(X)$ does not have any local (non-global) minimum. It is unlikely that the iterates get trapped at a saddle point. **Better local convergence result holds if we further assume $\lambda_k > \lambda_{k+1}$.**

Theorem (Q-Linear Rate)

Suppose $A > 0$ and $\lambda_k > \lambda_{k+1}$. Then $\{X^i\}$, a sequence generated by **SLRPGN(PR)** starting from a full-rank initial point $X^0 \in \mathcal{L}(\underline{\gamma})$, **globally converges to $\mathcal{L}(f^*)$** , where $\mathcal{L}(\underline{\gamma}) := \{X \mid f(X) \leq \underline{\gamma}\}$ denotes the level set, f^* denotes the global minimum of SLRP and $\underline{\gamma} > f^*$ is a constant. Moreover, the gradient sequence $\{\nabla f(X^i)\}$ converges to zero at a **Q-linear rate $\frac{\lambda_{k+1}}{\lambda_k}$** .

Low-rank Approximation without SVD

Recall for $B \in \mathbb{R}^{n \times m}$,

$$\min_{X, Z} \{ \|XZ^T - B\|_F : X \in \mathbb{R}^{n \times k}, Z \in \mathbb{R}^{m \times k} \}$$

SVD:

$$XZ^T = (U_k \Sigma_k) V_k^T,$$

which is the principal part of $B = U \Sigma V^T$.

GN:

- $A = @ (X) B * (B' * X)$;
- $[X, Y] = SLRPGN(A, X)$;
- $Z^T = Y^T B$

Platform

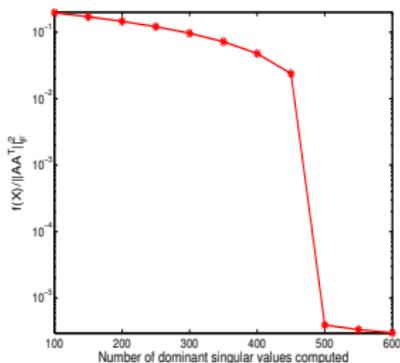
All the experiments were preformed on a linux workstation with 2 Intel Xeon E5-2697 CPUs (2.70GHz, 12 cores) and 128GB of memory running Ubuntu 12.04 and MATLAB 2013b.

Tested Methods

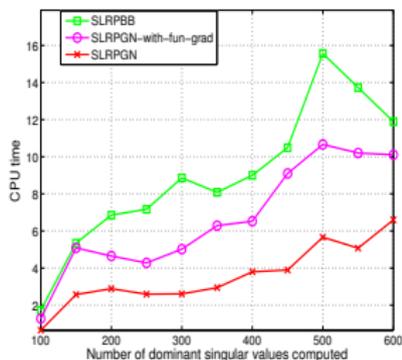
- MATLAB **EIGS** — Lanczos-based (ARPACK, Sorensen et.al.)
- **LANSVD** — Lanczos-based (PROPACK, R. M. Larsen)
- **LMSVD** — block subspace method (**L.-Wen-Zhang, SISC, 2013**)
- **SLRPBB** – BB + gradient (EIGPEN) (**Wen-Yang-L.-Zhang**)
- **SLRPGN** – proposed GN algorithm

Required Accuracy: moderate

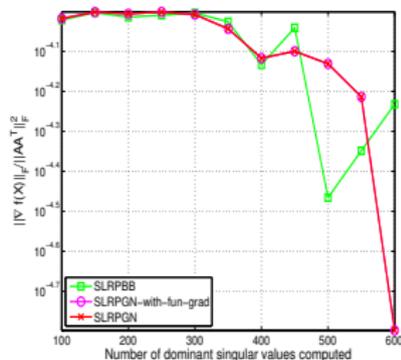
Comparison Between BB and GN



(a) Optimal residual



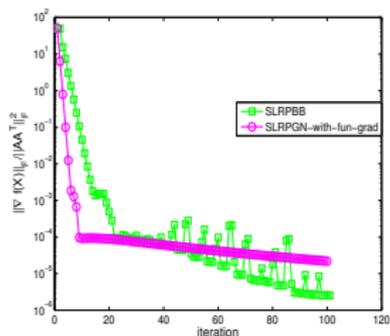
(b) Runtime



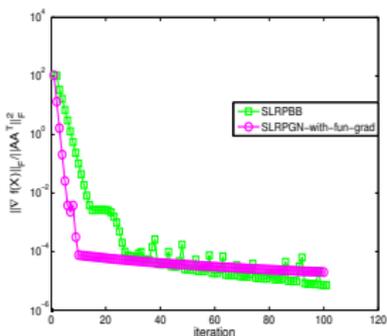
(c) gradient norm

(SLRPBB and SLRPGN with varying number of computed singular values on the random example “randcolu”, size: 10000×10000 , $\text{tol} = 10^{-4}$)

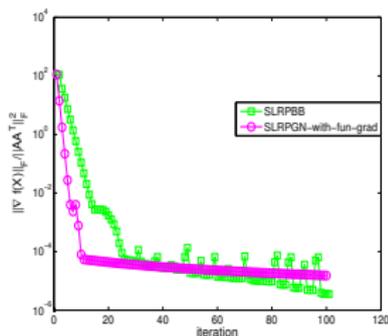
Comparison Between BB and GN (Cont'd)



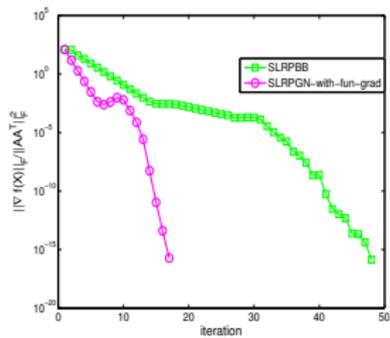
(d) $k = 100$



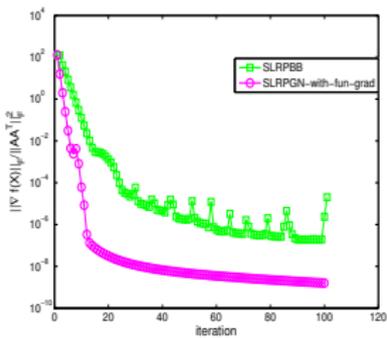
(e) $k = 400$



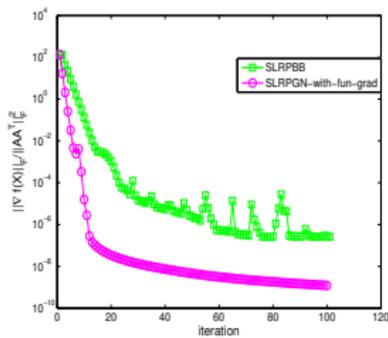
(f) $k = 450$



(g) $k = 500$



(h) $k = 550$



(i) $k = 600$

Robust Principal Component Analysis

Data matrix is

$$M = L_0 + S_0 + \omega \in \mathbb{R}^{m \times n},$$

where L_0 is low-rank, S_0 is sparse and ω is small noise.

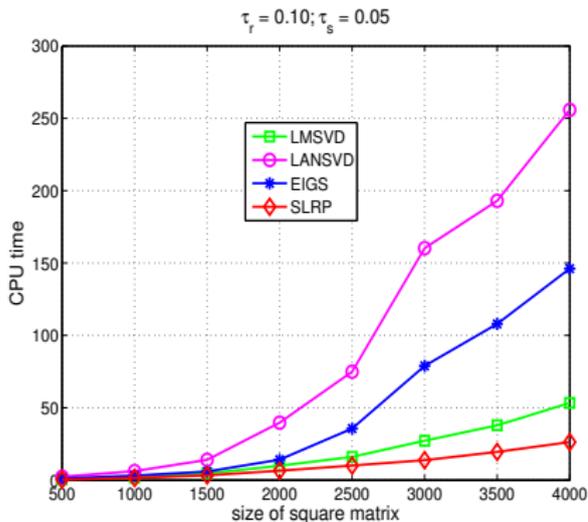
Given M , find L_0 and S_0 approximately by solving:

$$\min_{L,S} \|L\|_* + \mu \|S\|_1, \quad \text{s.t. } L + S = M.$$

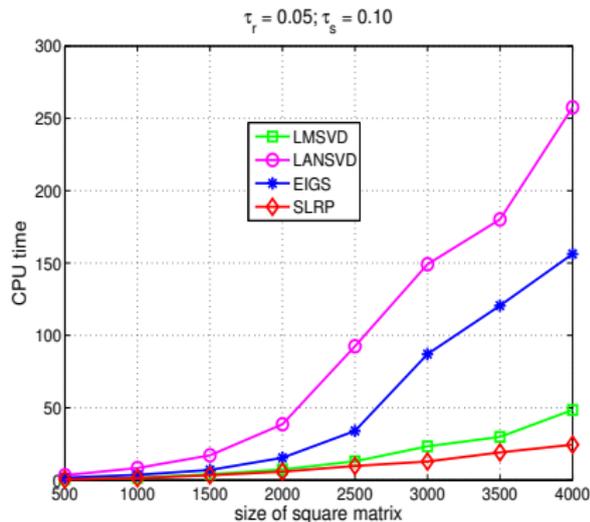
MATLAB Code: IALM (Lin *et al*)

- Alternating Direction Multiplier Method (ADMM)
- Calls SVD at every iteration (warm-start desired)
- Test cases: random instances

CPU Time in Seconds



(a) rank = 0.10n Sample 5%



(b) rank = 0.05n Sample 10%

(All achieved similar accuracy)

Find a low-rank matrix from a sampled set of its entries

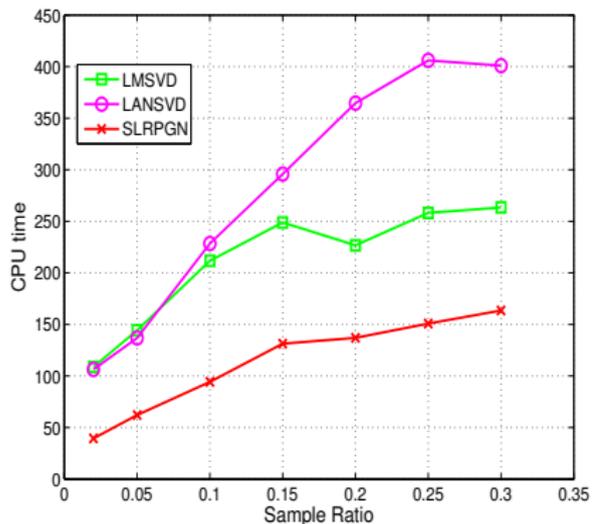
Given the entries of $M + \omega$ in Ω , find $X \approx M$ by solving:

$$\min_X \|X\|_*, \quad \text{s.t. } X_{ij} = M_{ij}, \quad \forall (i, j) \in \Omega.$$

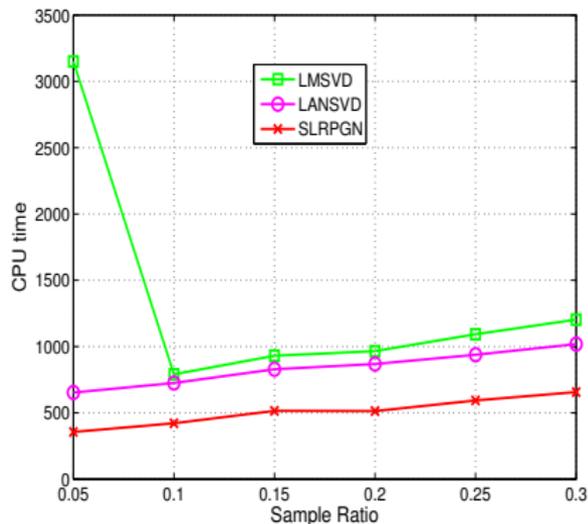
MATLAB Code: SVT and NNLS

- Singular Value Thresholding
- Calls SVD at every iteration (warm-start desired)
- Test cases: random instances

CPU Time in Seconds



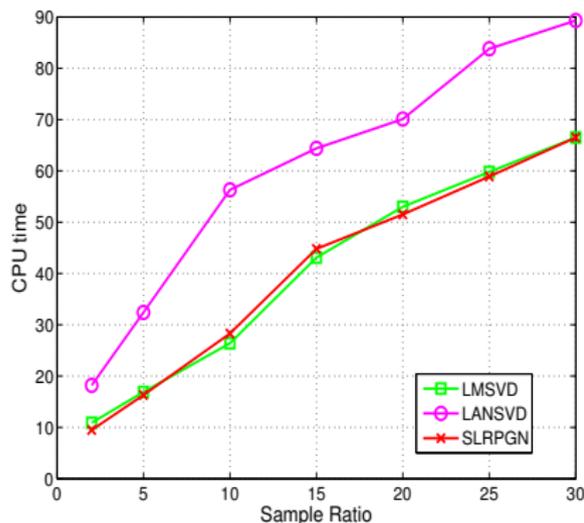
(a) rank = 10



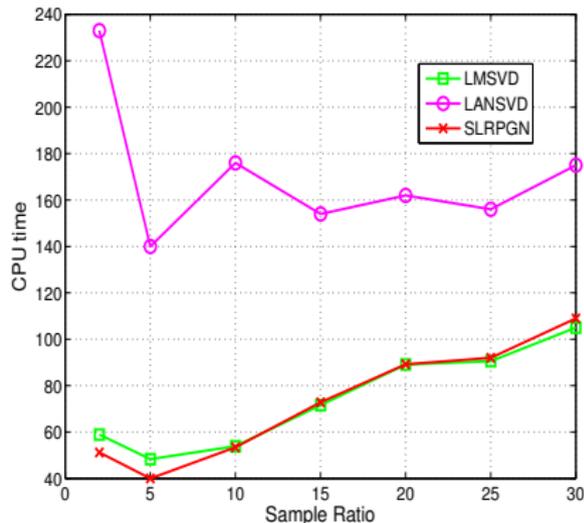
(b) rank = 50

(All achieved similar accuracy)

CPU Time in Seconds³



(a) rank = 10



(b) rank = 50

(All achieved similar accuracy)

³The sparse-dense matrix multiplication uses MKL

Summery and Remarks

SLRP: $\min \|XX^T - A\|_F^2$. Output (X, Y)

GN: $Y = X(X^T X)^{-1}$; $X = AY - X(Y^T AY - I)/2$

- SLRPGN: simple and parameter-free
- Principal subspace without SVD, nor Rayleigh-Ritz
- Benefit of concurrency already seen in plain MATLAB
- Global convergence and local Q-linear convergence rate
- Effective for small residuals and low-moderate accuracy (so far)

Further Works

- Strategically placed Rayleigh-Ritz will improve accuracy
- Other eigen-techniques (poly-filtering, deflation, ...) help too
- SLRPGN + (a few RR): Potential as eigensolver worth investigating
- Parallel scalability to be exploited



Thank you for your attention!