

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides

JEAN **TSHIMANGA** ILUNGA(FUNDP) joint work with

Serge Gratton (CERFACS) and Annick Sartenaer (FUNDP)

Sparse Days, Cerfacs, Toulouse, June 15-16, 2006

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 🕨 🐗 🖻 🕨 🗲



- General framework
- Preconditioning techniques considered
- Some properties
- Numerical experiments (data assimilation)
- Conclusions and perspectives



The goal of data assimilation

The goal of data assimilation is to find an initial state vector for which the trajector of the forecast model best fits an initial background and some observations. The vector found is the used to start a new simulation which gives an improved forecast.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 🕨 🐗 🖻 🕨 🗲



Data assimilation problem formulation

A variational formulation: nonlinear least-squares problem

$$\min_{\mathbf{x}\in\mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} ||\mathbf{x} - x_b||_{B^{-1}}^2 + \frac{1}{2} \sum_{j=0}^N ||\mathcal{H}_j(\mathcal{M}_j(\mathbf{x})) - y_j||_{R_j^{-1}}^2$$

- ► Size of real (operational) problems : $x, x_b \in \mathbb{R}^{10^6}$, $y_j \in \mathbb{R}^{10^5}$.
- The observations y_j and the background x_b are noisy.
- *M* is the model operator (nonlinear)
- H is the observation operator (nonlinear)



Solution strategy

- Typical approach : Incremental 4DVAR (i.e. inexact/truncated Gauss-Newton algorithm (GN)).
- GN leads to sequence of linear least-squares problems
- (Equivalently to) Sequence of linear symmetric positive definite systems to solve normal equations :

$$J_i^T J_i x = J_i^T r_i$$

• Whose matrix
$$A_i \equiv J_i^T J_i$$
 varies.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 🕨 4 🗟



The key Idea (in our study)

- ► We consider a symmetric and positive definite (spd) matrix A.
- Solve systems Ax = b₁, Ax = b₂, ..., Ax = br with RHS in sequence, by iterative methods: Conjugate Gradient (CG) or variants.
- Precondition the CG using informations obtained when solving the previous system.
- Extension of the idea to nonlinear process such as Gauss-Newton (GN) method. The matrix varies along the process.

 Outline
 Framework
 Techniques
 Properties
 Experiments
 Conclusions and Perspectives

 The CG algorithm (A is spd and large !)

 CG is an iterative method for solving
 I To the transmission of the transm

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x - b^T x \qquad \Leftrightarrow \qquad A x = b \text{ (stationary eq.)}$$

▶ Iterations: Given $x_0 \in \mathbb{R}^n$; $A \in \mathbb{R}^{n \times n}$; $b \in \mathbb{R}^n$

Set
$$r_0 \leftarrow Ax_0 - b_0$$
; $p_0 \leftarrow -r_0$; $i \leftarrow 0$

Loop on i

$$\begin{array}{rcl} \alpha_i & \leftarrow & (r_i^T r_i)/(p_i^T A p_i) \\ x_{i+1} & \leftarrow & x_i + \alpha_i p_i \\ r_{i+1} & \leftarrow & r_i + \alpha_i A p_i \\ \beta_{i+1} & \leftarrow & (r_{i+1}^T r_{i+1})/(r_i^T r_i) \\ p_{i+1} & \leftarrow & -r_{i+1} + \beta_{i+1} p_i \end{array}$$

▶ *r_i* are residuals; *p_i* are descent directions.

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 > < 🗇 > < 🗟 >

Outline Framework Techniques Properties Experiments Conclusions and Perspectives

The CG properties (in exact arithmetic !)

- Orthogonality of the residuals: $r_i^T r_j = 0$ if $i \neq j$.
- A-conjugacy of the descent directions: $p_i^T A p_j = 0$ if $i \neq j$.
- ► The distance of the iterate x_i to the solution x^{*} is related to the condition number of A, denoted by κ = λmax/λmin (≥ 1):

$$||x_i - x^*||_{\mathcal{A}} \le \eta_i ||x_0 - x^*||_{\mathcal{A}} \text{ with } \eta_i = 2\left(rac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}
ight)^i$$

 \Rightarrow The smaller $cond(A) \equiv \kappa$ is, the faster the convergence.

Exact solution found exactly in r iterations, where r ≤ n is the number of distinct eigenvalues of A ∈ ℝ^{n×n}.
 ⇒ The more clustered the eigenvalues are, the faster the

convergence.



Why to precondition ?

- Transform Ax = b in an equivalent system having a more favorable eigenvalues distribution for faster convergence.
- Use a preconditioning matrix H (which must be cheap to apply).
- ▶ Ideas to design preconditioner. *H* would :
 - approximates A^{-1} .
 - make cond(HA) < cond(A).</p>
 - ▶ make eigenvalues of *HA* more clustered than those of *A*.
- ▶ Note: when a preconditioning is used, residuals are:
 - Orthogonal if H is factored in LL^{T} .
 - Conjugate w.r.t. *H* if *H* is not factored.



Preconditioning techniques considered (I)

- We consider (second level preconditioning) techniques :
 - Solve $Ax = b_1$ and extract information *info*₁.
 - Use *info*₁ to solve $Ax = b_2$ and extract information *info*₂.
 - Use *info*₂ (and possibly *info*₁) to solve $Ax = b_3$ and ...
 - ▶ ...
- ► *Info_k* will be:
 - descent directions;
 - or other vectors such as eigenvectors of A ...

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 🕨 🗸 🗄



Preconditioning techniques considered (II)

We study and compare two approaches:

- Deflation [Frank, Vuik, 2001].
- Limited Memory Preconditioners (LMP): Preconditioners based on a set of A-conjugate directions. Generalization of known preconditioners:
 - spectral [Fisher, 1998],
 - L-BFGS [Nocedal, Morales, 2000], warm start [Gilbert, Lemaréchal, 1989].

We cover:

- Theoretical properties.
- Numerical experiments (data assimilation).



Deflation Techniques

- ► Given W ∈ ℝ^{n×k} (k ≪ n) formed with appropriate information obtained when solving the previous system.
- Consider the oblique projector $P = I AW(W^T AW)^{-1}W^T$.

• Split the solution vector as follows $x^* = \underbrace{(I - P^T)x^*}_{direct} + \underbrace{P^Tx^*}_{iterative}$.

- Compute $(I P^T)x^*$ with a direct method.
- Compute $P^T x^*$ with an iterative method.



Some Properties for Deflation

- Computation of $(I P^T)x^*$:
 - $(I P^T)x^* = W(W^T A W)^{-1} W^T A x^* = W(W^T A W)^{-1} W^T b.$
 - Note: $W^T A W \in \mathbb{R}^{k \times k}$ and $k \ll n$.
- Computation of $P^T x^*$:
 - Consider the compatible singular (but still symmetric) system PAy = Pb.
 - Use CG with $y_0 = 0$ to solve PAy = Pb.
 - Any solution y of PAy = Pb satisfies $P^T x^* = P^T y$.
 - Note: $cond(PA) \leq cond(A)$ and $PA = (PA)^T$.



Limited Memory Preconditioners (LMP)

General reformulation of LMP:

$$H_{k+1} = [I - \sum_{i=0}^{k} \frac{Aw_{i}w_{i}^{T}}{w_{i}^{T}Aw_{i}}]^{T}[I - \sum_{i=0}^{k} \frac{Aw_{i}w_{i}^{T}}{w_{i}^{T}Aw_{i}}] + \sum_{i=0}^{k} \frac{w_{i}w_{i}^{T}}{w_{i}^{T}Aw_{i}},$$

with $w_{i}^{T}Aw_{j}$ $\begin{cases} = 0 \quad \text{if } i \neq j \\ > 0 \quad \text{if } i = j \end{cases}$

Particular forms

- The w_i's are the descent directions obtained from CG: w_i = p_i ⇒ L-BFGS preconditioner.
- ► The w_i's are eigenvectors of A: w_i = v_i ⇒ spectral preconditioner.

J. TSHIMANGA I.

14/36



Spectral Properties for LMP (I)

- . The matrix A is fixed.
 - The matrix A to precondition is the same as the previous one (only the RHS changes).
 - ► Theorem : Assume that λ₁,..., λ_n is the spectrum of A, then the spectrum μ₁,..., μ_n of the preconditioned matrix H_{k+1}A satisfies:

$$\begin{cases} \mu_j = 1, & \text{for } j = 1, \dots, k \\ \lambda_{j-k}(A) \le \mu_j \le \lambda_j(A), & \text{for } j = k+1, \dots, n, \end{cases}$$

where $\lambda_j(A)$ is the *j*-th eigenvalue of A (increasing order assumed).

J. TSHIMANGA I.

15/36

Outline Framework Techniques **Properties** Experiments Conclusions and Perspectives

Spectral properties for LMP (II)



- Eigenvalues translated to 1.
- The rest of the spectrum is not expanded compared to the spectrum of A.



Spectral Properties for LMP (III)

The matrix A varies.

- ▶ $\tilde{A} = A + \tau E$ with $||E||_F = 1$ and $\tau \in R^+$ (and small enough : first order perturbation)
- Theorem : Assume that τ is small, then the k + 1 first eigenvalues of the preconditioned matrix H_{k+1}Ã satisfies:

$$\tilde{\mu}_j \leq 1 + \lambda_{max}(W^T E W) * \tau + o(\tau), \quad \text{for } j = 1, \dots, k+1$$

λ_{max}(W^TEW) * τ is a measure of the default of conjugacy of columns vectors of W w.r.t. the new matrice Ã

J. TSHIMANGA I.

17/36

 Outline
 Framework
 Techniques
 Properties
 Experiments
 Conclusions and Perspectives

 Existence of a factored form for the LMP (not the

Cholesky factor and obtained by construction!)

L-BFGS:

• A possible factored form is $H_{k+1} = L_{k+1}L_{k+1}^T$ where:

$$L_{k+1} = \prod_{i=0}^{k} \left(I - \frac{\mathbf{s}_i \mathbf{y}_i^{T}}{\mathbf{y}_i^{T} \mathbf{s}_i} + \frac{\mathbf{s}_i}{\sqrt{\mathbf{y}_i^{T} \mathbf{s}_i}} \frac{\mathbf{r}_i^{T}}{\|\mathbf{r}_i\|} \right),$$

with $s_i = x_{i+1} - x_i$ and $y_i = r_{i+1} - r_i$.

- Same cost in memory and CPU as the unfactored form.
- Spectral:
 - A possible factored form is $H_{k+1} = L_{k+1}^2$ where:

$$L_{k+1} = I + \sum_{i=1}^{k+1} \left(\frac{1}{\sqrt{\lambda_i}} - 1 \right) \frac{v_i v_i^{T}}{v_i^{T} v_i}$$

Same cost in memory as the unfactored form.

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 🕨 < 🗇 🕨 < 🗄 🕨



Why looking for a factored form $H = LL^T$?

- ▶ With a non factored form, we use CG preconditioned by *H*.
- ► With a factored form, we solve L^TALu = L^Tb; x = Lu. Advantages:
 - More appropriate if reorthogonalization of the residuals is used.
 - ► Least-squares min_x ||A_x b|| or AA^T_x = A^T_b: LSQR (or CGLS) is more accurate than CG in presence of rounding errors but works with (A, A^T, L, L^T, b) instead of (A^TA, A^Tb, H).
 - When accumulating preconditioners, symmetry and positiveness are still maintained:

$$L_1^T A L_1 y_1 = L_1^T b_1, \quad L_2^T (L_1^T A L_1) L_2 y_2 = L_2^T L_1^T b_2, \quad \dots$$



Why to reorthogonalize the residuals ?

- In finite precision, residuals often loose orthogonality (or conjugacy) and theoretical convergence is then slowed down.
- Reorthogonalization of residuals in CG is terribly successful when matrix-vector product is very expensive compared to other computations in CG (see example in the next slide).
- Note: to restore orthogonality or conjugacy, working with L^TAL and the canonical inner-product is better (memory, CPU, error propagation) than working on A preconditioned by the no factored matrix H.

Example of reorthogonalisation effect : CERFACS data assimilation system (1 000 000 unknowns)





J. TSHIMANGA I.

21/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 🕨 🔬 ী 🕨 💐

 Outline
 Framework
 Techniques
 Properties
 Experiments
 Conclusions and Perspectives

Experiments with unpreconditioned LSQR



LSQR is better than CG !

J. TSHIMANGA I.

22/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 + 4 🗇 + 4 🗟

 Outline
 Framework
 Techniques
 Properties
 Experiments
 Conclusions and Perspectives

Experiments with LSQR preconditioned with factored L-BFGS



LSQR is again better than CG !

J. TSHIMANGA I.

23/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 + 4 🗟 + 4 🗟

Experiments description

Algorithmic variants tested:

- Use CG to solve the normal equations.
- Compare 3 preconditioning techniques:
 - Deflation technique (using spectral information).
 - Spectral preconditioner (using spectral info. but differently).
 - L-BFGS preconditioner (using descent directions).
- Where spectral information is needed, use Ritz (vectors) as approximations of the eigenvectors.
- Ritz vectors are obtained by mean of a variant of CG: the Lanczos algorithm which combines linear and eigen solvers.



Ranking of the preconditioners using the basic strategies.

- Diagnostics are made on the non linear optimisation problem
- The number of CG iterations is fixed.

 Outline
 Framework
 Techniques
 Properties
 Experiments
 Conclusions and Perspectives

Results L-BFGS - Deflation



Deflation is better than L-BFGS !

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 + 4 🗇 + 4 🗟

Outline Framework Techniques Properties Experiments Conclusions and Perspectives

Results Noprecond - L-BFGS



L-BFGS is better than Noprecond !

J. TSHIMANGA I.

27/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 + 4 🗇 + 4 🗟

 Outline
 Framework
 Techniques
 Properties
 Experiments
 Conclusions and Perspectives

Results Spectral - Noprecond



Noprecond is better than Spectral !

J. TSHIMANGA I.

28/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 + 4 🗟 + 4 🗟

Is it better to wait until the system does not vary much for the preconditioner to be efficient ?

- Our main theory assumes that the matrix A is fixed. In practice (data assimilation) the matrix varies.
- It is Known that the GN method may not be locally convergent on some problems. But in our experiments the GN process converges and the steps thus become smaller.
- This means that the matrix does not change so much, when approaching the solution.

Outline Framework Techniques Properties Experiments Conclusions and Perspectives

Results Noprecond - L-BFGS(2nd)



L-BFGS(2nd) is better than Noprecond !

J. TSHIMANGA I.

30/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 + 🗸 🗇 + 💐

Outline	Framework	Techniques	Properties	Experiments	Conclusions and Perspectives

Results L-BFGS(1rst) - L-BFGS(2nd)



L-BFGS(2nd) is better than L-BFGS(1rst) !

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides 4 🗆 🕨 4 🗇 🕨 4

Outline Framework Techniques Properties Experiments Conclusions and Perspectives

Results Noprecond - deflation(4th)



Deflation(4th) is better than Noprecond !

J. TSHIMANGA I.

32/36

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 + 🗸 🗇 + 💐



Results deflation(4th) - deflation(1rst)



Deflation(1rst) is a little better than Deflation(4th) !

J. TSHIMANGA I.

On Some Preconditioning Techniques for Linear Least-Squares Problems with Multiple Right-Hand Sides < 🗆 🕨 🔶 🗟



Remarks on our system !

- Spectral preconditioner:
 - Does not work in our case.
- L-BFGS preconditioner:
 - Requires no large changes in the matrix.
 - Based on by-products of CG.
 - More efficient than the spectral preconditioner or than no preconditioner.
- Deflation:
 - Is stable even when the matrix changes.
 - May be expensive $(W^T A W)$ in CPU time.
 - More efficient than the other techniques.



- Properties of LMP preconditioners understood (when the matrix A is fixed or changes a "bit").
- Existence of factored forms for particular instances of LMP.
- Preliminary tests show weakness of spectral compared to deflation and L-BFGS in a data assimilation experiment.
- One technique (deflation) not yet used before in data assimilation has been tested in this field.



Perpectives and work in progress

Experiment sampling techniques to select information:

- Which information to capture ?
- Purge the preconditioner or not ?
- When and how to decide to apply the preconditioner ?

► Make tests in a more realistic data assimilation environment:

- CERFACS 2D Shallow Water or any.
- CERFACS operational system or any.