# Hierarchical Prior Models and Krylov-Bayes

# Iterative Methods: Part 2

## Daniela Calvetti

## Case Western Reserve University

▶ D Calvetti, E Somersalo: Priorconditioners for linear systems. Inverse Problems 21 (2005) 13971418

▶ D Calvetti, F Pitolli, E Somersalo, B Vantaggi: Bayes meets Krylov: preconditioning CGLS for underdetermined systems http://arxiv.org/abs/1503.06844

# In the mood for Bayes

In the Bayesian framework for linear discrete inverse problems we need to solve a linear system of equations

- The linear solver apparatus (least-squares solvers, iterative methods, etc) can be used to update $z$ in the MAP calculation
- Conversely, the Bayesian framework can be very helpful to solve linear systems of equations
- More generally, the Bayesian approach can be used in the solution of nonlinear systems
- The Bayesian approach is particularly well-suited for under-determined linear systems

## General setting

Consider the problem of estimating $x \in \mathbb{R}^n$ from

$$b = \mathbf{F}(x) + \epsilon, \qquad \mathbf{F} : \mathbb{R}^n \longrightarrow \mathbb{R}^m,$$

Here we focus the attention on the special case where

$$\mathbf{F}(x) = \mathrm{A}x$$

with A is an $m \times n$ matrix of rank $m$, typically badly conditioned and of ill-determined rank.

We solve the linear system with a Krylov subspace iterative method.

# Bayesian solution of inverse problems

In the Bayesian framework for the solution of inverse problems,

- All unknown parameters are modeled as random variables and described in terms of their probability density functions;
- Here the unknowns are $\epsilon$ and $x$
- $\pi_{\mathrm{noise}}(\epsilon)$ describes what we know about the statistics of the noise and defines the *likelihood*;
- $\pi_{\mathrm{prior}}(x)$ expresses what we know about $x$ before taking into consideration the data and is called the *prior*
- The solution of the inverse problem is $\pi(x \mid b)$ and is called the *posterior*.

It follows from Bayes' formula that

$$\pi(x \mid b) \propto \pi_{\mathrm{prior}}(x)\pi_{\mathrm{noise}}(b - Ax).$$

# The noise in a Bayesian way

The linear discrete inverse problem that we consider is $b = Ax + \epsilon$.

- ► The noise term $\epsilon$ accounts for inaccuracies in the measurements as well as model uncertainties, i.e. discrepancy between reality and the model. (more on this tomorrow)
- ► We assume that $\epsilon \sim N(0, I_m)$.
- ► If $\epsilon = \epsilon_c \sim N(\mu_\epsilon, \Gamma)$, we can proceed as follows:
  - ► Compute a symmetric factorization of the precision matrix of $\epsilon$
    $\Gamma^{-1} = G^T G$
  - ► Make the change of variables

$$\epsilon = G(\epsilon_c - \mu_c) = G(b - Ax) - G\mu_c$$

- ► In the linear system

$$Gb = GAx + \epsilon$$

the noise is zero-mean white Gaussian.

# The unknown in a Bayesian way

Assume that $x \sim N(0, C)$, where $C$ is symmetric positive definite. It follows from Bayes formula that the posterior density is of the form

$$\pi(x \mid b) \propto \exp\left(-\frac{1}{2}\|Ax - b\|^2 - \frac{1}{2}x^T C^{-1}x\right).$$

Give a symmetric factorization of the precision matrix of $x$

$$C^{-1} = B^\mathsf{T} B$$

we can write the negative logarithm of the posterior, or *Gibbs energy* in the form

$$G(x) = \|Ax - b\|^2 + \|Bx\|^2 = \left\| \begin{bmatrix} A \\ B \end{bmatrix} x - \begin{bmatrix} b \\ 0 \end{bmatrix} \right\|^2.$$

# MAP estimate

The Maximum a Posteriori (MAP) estimate of $x$, $x_{\mathrm{MAP}}$ is the value of highest posterior probability, or equivalently, the minimizer of $G(x)$ .

The value of $x_{\mathrm{MAP}}$ is the *least squares* solution of the linear system

$$\left[\begin{array}{c} A \\ B \end{array}\right] x = \left[\begin{array}{c} b \\ 0 \end{array}\right],$$

or, equivalently, the solution of the square linear system

$$(A^{\mathsf{T}}A + B^{\mathsf{T}}B)x = A^{\mathsf{T}}b.$$

## ... and Tikhonov regularization

The latter are the normal equations associated with the problem

$$x_{\mathrm{MAP}} = \mathrm{argmin}\{\|Ax - b\|^2 + \lambda\|Bx\|^2\} \qquad (1)$$

which is Tikhonov regularized solution with regularization parameter $\lambda = 1$ and linear regularization operator $B$.

- ▶ The computation of Tikhonov regularized solution with $B$ different from $I$ requires attention: moreover, a suitable value of the regularization parameter $\lambda$ must be chosen.
- ▶ This question has been studied extensively in the literature.

# Regularization operator beyond Tikhonov

- The operator $B$ brings into the solution additional information about $x$
- When the matrix A is underdetermined, the operator $B$ boosts the rank of the matrix of the linear system actually solved.
- In Tikhonov regularization when $n$ is large the introduction of $B$ may lead to a very large linear system

**Question:** How can we retain the benefits of $B$ while containing the computational costs?

# The alternative: Krylov subspace methods

As an alternative to Tikhonov regularization consider solving the linear system

$$b = Ax + \epsilon$$

with an iterative solver using the matrix $B$ as a right preconditioned. More specifically

- Consider the Conjugate Gradient for Least Squares method (CGLS)
- WLOG assume the initial approximate solution is $x_0 = 0$
- Define a termination rule based on the discrepancy

## Standard CGLS method

At the $k$th iteration step the approximate solution $x_k$ computed by the CGLS method satisfies

$$x_k = \operatorname{argmin}\{\|b - Ax\| \mid x \in \mathscr{K}_k\},$$

where the $k$th Krylov subspace is

$$\mathscr{K}_k = \operatorname{span}\{A^\mathsf{T} b, (A^\mathsf{T} A)A^\mathsf{T} b, \ldots, (A^\mathsf{T} A)^k A^\mathsf{T} b\}.$$

The noise is additive, zero-mean white Gaussian, thus

$$E\left\{\|\epsilon\|^2\right\} = m;$$

we stop iterating as soon as

$$\|Ax - b\|^2 < \tau m,$$

where $\tau = 1.2$. Typically, $k_{\text{last}} \ll m$.

# The question of the null space

- It follows from the canonical orthogonal decomposition in terms of fundamental subspaces that

$$\mathbb{R}^n = \mathscr{N}(\mathsf{A}) \oplus \mathscr{R}(\mathsf{A}^\mathsf{T})$$

- In standard CGLS any contribution to the solution from the null space must be added separately
- The right CGLS priorconditioner implicitly selects null space components based on the information contained in the data with the belief about $x$.

# CGLS with a whitened unknown

Assume that a prior we believe that $x \sim N(0, C)$. If

$$C^{-1} = B^T B$$

then

$$w = Bx, \qquad w \sim N(0, I_n).$$

Make the change of variable from $x$ to $w$ in the linear system

$$AB^{-1}w = b \qquad x = B^{-1}w, \tag{2}$$

let $\widetilde{A} = AB^{-1}$ and solve by CGLS for $w$. The $j$th iterate of the whitened problem satisfies

$$w_j = \mathrm{argmin}\{\|\widetilde{A}w - b\| \mid w \in \mathcal{K}_j(\widetilde{A}^T b, \widetilde{A}^T \widetilde{A})\}.$$

# Priorconditioning and the null space

The corresponding $j$th priorconditioned CGLS solution $\widetilde{x}_j = \mathsf{B}^{-1} w_j$ satisfies

$$\widetilde{x}_j \in \operatorname{span}\left\{ \mathsf{B}^{-1}(\widetilde{A}^{\mathsf{T}} \widetilde{A})^{\ell} \widetilde{A}^{\mathsf{T}} b \mid 0 \le \ell \le j-1 \right\}.$$

It follows from

$$\mathsf{B}^{-1} \widetilde{A}^{\mathsf{T}} = \mathsf{B}^{-1} \mathsf{B}^{-\mathsf{T}} A^{\mathsf{T}} = \mathsf{C} A^{\mathsf{T}},$$

that

$$\mathsf{B}^{-1} (\widetilde{A}^{\mathsf{T}} \widetilde{A})^{\ell} \widetilde{A}^{\mathsf{T}} = (\mathsf{C} A^{\mathsf{T}} A)^{\ell} \mathsf{C} A^{\mathsf{T}}, \quad 0 \le \ell \le j-1.$$

Therefore

$$\widetilde{x}_j \in \mathsf{C}\big(\mathscr{N}(A)^{\perp}\big),$$

hence $\widetilde{x}_j$ is not necessarily orthogonal to the null space of A.

# Analizing the Krylov subspaces with the GSVD

**Theorem**
Given $(A, B)$ with $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times n}$, $m < n$, there is a factorization of the form

$$A = U \begin{bmatrix} 0_{m,n-m} & \Sigma_A \end{bmatrix} X^{-1}, \quad B = V \begin{bmatrix} I_{n-m} & \\ & \Sigma_B \end{bmatrix} X^{-1},$$

called the *generalized singular value decomposition*, where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices, $X \in \mathbb{R}^{n \times n}$ is invertible, and $\Sigma_A \in \mathbb{R}^{m \times m}$ and $\Sigma_B \in \mathbb{R}^{m \times m}$ are diagonal matrices.

The diagonal entries $s_1^{(A)}, \ldots, s_m^{(A)}$ and $s_1^{(B)}, \ldots, s_m^{(B)}$ of the matrices $\Sigma_A$ and $\Sigma_B$ are real, nonnegative and satisfy

$$
\begin{aligned}
s_1^{(A)} &\leq s_2^{(A)} \leq \ldots \leq s_m^{(A)} \\
s_1^{(B)} &\geq s_2^{(B)} \geq \ldots \geq s_m^{(B)} \\
(s_j^{(A)})^2 &+ (s_j^{(B)})^2 = 1, \qquad 1 \leq j \leq m. \quad (3)
\end{aligned}
$$

thus $0 < s_j^{(A)} \leq 1$ and $0 < s_j^{(B)} \leq 1$. The ratios $s_j^{(A)}/s_j^{(B)}$ for $1 \leq j \leq m$ are the generalized singular values of $(A, B)$. If A has full rank, the diagonal entries of $\Sigma_A$ are positive.

# C-orthogonality

**Theorem**

If we partition the matrix $X \in \mathbb{R}^{n \times n}$ in GSVD above as

$$X = \begin{bmatrix} X' & X'' \end{bmatrix}, \quad X' \in \mathbb{R}^{n \times (n-m)}, \, X'' \in \mathbb{R}^{n \times m},$$

it follows that

$$\mathrm{span}\{X'\} = \mathcal{N}(A),$$

and we can express $\mathbb{R}^n$ as a C-orthogonal direct sum,

$$\mathbb{R}^n = \mathrm{span}\{X'\} \oplus_C \mathrm{span}\{X''\} = \mathcal{N}(A) \oplus_C \mathrm{span}\{X''\}.$$

# Orthogonality and not

**Corollary 1**

$$\mathcal{N}(A)^{\perp} = \mathcal{R}(A^{\mathsf{T}}), \quad \mathcal{N}(A)^{\perp c} = \mathrm{span}\{X''\}.$$

**Corollary 2**
If $\mathcal{R}(A^{\mathsf{T}})$ is an invariant subspace of the covariance matrix C, then the iterates $\widetilde{x}_j$ are orthogonal to the null space of A.

**Corollary 3**
When $C(\mathcal{R}(A^{\mathsf{T}}))$ is not C-orthogonal to $\mathcal{N}(A)$, $\widetilde{x}_j$ may have a component in the null space of A. This component is invisible to the data.

# Priorconditionting and the Lanczos process

The first $k$ residual vectors computed by CGLS normalized to have unit length $v_0, v_1, \ldots, v_{k-1}$ form an orthonormal basis for the Krylov subspace $\mathcal{K}_k(A^\mathsf{T} b, A^\mathsf{T} A)$.

It can be shown that

$$A^\mathsf{T} A V_k = V_k T_k - \frac{\sqrt{\beta_{k-1}}}{\alpha_{k-1}} v_k e_k^\mathsf{T}., \qquad V_k = [v_0, v_1, \ldots, v_{k-1}].$$

It follows from the orthogonality of the $v_j$ that the tridiagonal matrix $T_k$ is the projection of $A^\mathsf{T} A$ onto the Krylov subspace $\mathcal{K}_k(A^\mathsf{T} b, A^\mathsf{T} A)$.

$$V_k^\mathsf{T}(A^\mathsf{T} A) V_k = T_k.$$

# The Lanczos tridiagonal matrix

The $k$th CGLS iterate can be expressed as

$$x_k = V_k y_k,$$

where $y_k$ solves the $k \times k$ linear system

$$T_k y = \|r_0\| e_1.$$

Thus the $k$th CGLS iterate $x_k$ is the lifting of $y_k$ via $V_k$.

The eigenvalues of $T_k$ are the Ritz values of $A^T A$ and approximate of the eigenvalues of $A^T A$.

# Ritz values and convergence rate

**Theorem**
For all $k$, $1 \leq k \leq r$, where $r$ is the rank of A, there exists
$\xi_k, \lambda_1 \leq \xi_k \leq \lambda_r$ such that the norm of the residual vector satisfies

$$\|r_k\|^2 = \frac{1}{\xi_k^{2k+1}} \sum_{i=1}^{n} \left[ \prod_{j=1}^{k} \left( \lambda_i - \theta_j^{(k)} \right)^2 \right] \left( r_0^\mathsf{T} q_i \right)^2,$$

where $q_i$ is the eigenvector of $A^\mathsf{T}A$ corresponding to the eigenvalue
$\lambda_i$, and $\theta_j^{(k)}$ is the $j$th eigenvalue of the tridiagonal matrix $T_k$.

The quality of the eigenvalues approximations in the projected
problem affects the number of iterations needed to meet the
stopping rule.

# A simple deconvolution problem

*Forward model:* Deconvolution problem with few data,

$$g(t) = \int_0^1 a(t-s)f(s)ds, \quad a(t) = \left( \frac{J_1(\kappa t)}{\kappa t} \right)^2,$$

Discretize:

$$g(t) \approx \frac{1}{n} \sum_{k=1}^n a(t-s_k)f(s_k), \quad 1 \le j \le n,$$

Discrete noisy observations at $t_1, \ldots, t_m$, $m \ll n$.

$$b_\ell = g(t_\ell) + \varepsilon_\ell, \quad 1 \le \ell \le m,$$

or, in matrix notation, $A \in \mathbb{R}^{m \times n}$,

$$b = Ax + \varepsilon, \quad x_k = f(s_k).$$

# Computed examples: Deconvolution

*Prior:* Define the precision matrix $C^{-1}$ as

$$C^{-1} = B^T B, \quad B = \beta \begin{bmatrix} \alpha & & & & \\ -1 & 2 & -1 & & \\ & & \ddots & & \\ & & -1 & 2 & -1 \\ & & & & \alpha \end{bmatrix},$$

where $\alpha > 0$ is chosen so that prior variance is as uniform as possible over the interval.

*Parameters:* Set $n = 150$, $m = 6$.

# Basis vectors



The six basis vectors that span $\mathscr{R}(\mathsf{A}^\mathsf{T})$ (dashed line), and the vectors that span $\mathsf{C}\big(\mathscr{R}(\mathsf{A}^\mathsf{T})\big)$ (solid line).

# Approximate solutions



Iterations with low additive noise ($\sigma = 5 \times 10^{-5}$) without prior conditioner (left) and with preconditioner.

# Observations

- Every vector whose support consists of points where all six basis functions of $\mathscr{R}(A^\mathsf{T})$ vanish is in the null space $\mathscr{N}(A)$
- Consequently, plain CGLS produces approximate solutions that are zero at those points
- The basis functions of $C\big(\mathscr{R}(A^\mathsf{T})\big)$ are non-zero everywhere
- Consequently, priorconditioned CGLS has no blind spots
- The price to pay is that priorconditioned CGLS requires more iterations

# Spectral approximation



Spectral approximation: Plain CGLS (left) and preconditioned CGLS (right). The grey band on the right is the spectral interval of the non-preconditioned matrix $A^T A$.

# Convergence history and null space contributions



Left: Convergence rates of the two algorithms. The dashed line marks the stopping criterion. Right: Component of the computed solution in the null space measured as
$\nu_k = \frac{\|P\widetilde{x}_k\|}{\|\widetilde{x}_k\|}, \quad P : \mathbb{R}^n \longrightarrow^{\perp} \mathcal{N}(A).$

# Computed examples: X-ray tomography



Image size: $N = 160 \times 160$ pixels. 20 illumination angles, 60 parallel beams per illumination angle.

## Correlation priors

Matèrn-Whittle correlation priors: Define the precision matrix as

$$C^{-1} = -I_n \otimes D - D \otimes I_n + \frac{1}{\lambda^2} I_N,$$

where $D \in \mathbb{R}^{n \times n}$ is the three-point finite difference approximation of the one-dimensional Laplacian with Dirichlet boundary conditions,

$$D = \frac{1}{n^2} \begin{bmatrix} -2 & 1 & & \\ 1 & -2 & \ddots & \\ & \ddots & & 1 \\ & & 1 & -2 \end{bmatrix},$$

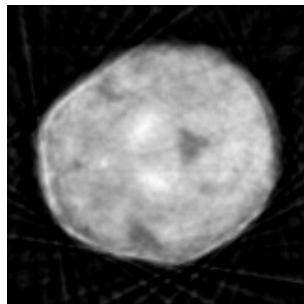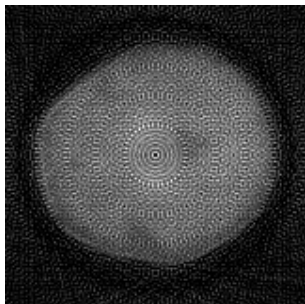and $\lambda > 0$ is the correlation length.

# Basis functions



Basis vector with no priorconditioning (upper left) and with
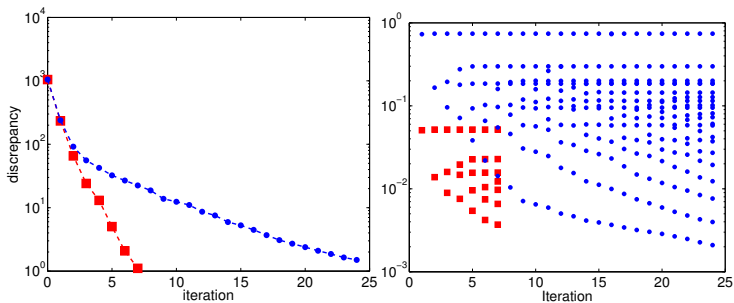priorconditioning. Correlation length 2, 4, 8, 16 and 32 pixels.

# Observations

- Every image whose support is on pixels not touched by a ray is in the null space of A
- $\Rightarrow$ plain CGLS iterates will be zero at those pixels
- Preconditioning makes the rays fuzzy, illuminating the dark pixels
- Reconstruction will be slightly blurred, but has fewer geometric artifacts
- Number of iterations needed will increase.

# Computed solutions



Reconstructions with plain CGLS (left) and priorconditioned CGLS (right).

# Converge history and spectral approximation



Convergence and spectral approximation.