

A stochastic Newton method for statistical inverse problems, with application to inverse scattering

Omar Ghattas

Joint work with:

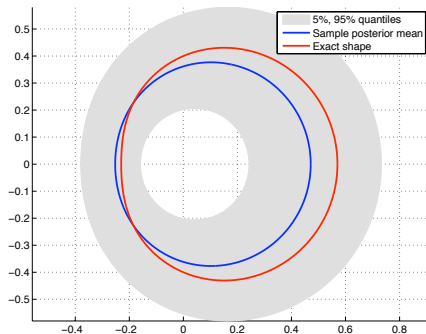
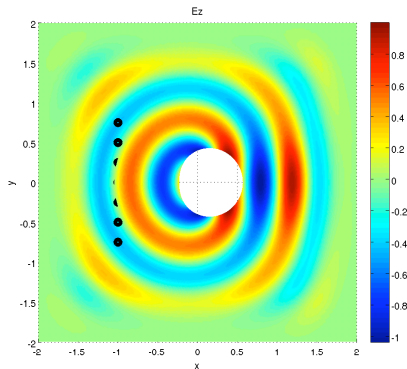
Tan Bui-Thanh Carsten Burstedde James Martin Lucas Wilcox

Institute for Computational Engineering and Sciences
Departments of Geological Sciences and Mechanical Engineering
The University of Texas at Austin

KAUST 22 March 2010

Bayesian inference for inverse scattering

Overall goals: Develop algorithms for solving statistical inverse problems that can scale to high dimensional probability spaces and expensive forward models, and tailor them to inverse shape and medium (acoustic, elastic, electromagnetic) scattering problems



Bayesian inference for inverse scattering, cont.

- Bayesian framework for statistical inverse problem: when data and/or model have uncertainties, solution of inverse problem expressed as a posterior probability density function
- Central challenge: for inverse problems characterized by high-dimensional parameter spaces, method of choice is to sample the posterior density using Markov chain Monte Carlo (MCMC)
- For inverse problems characterized by expensive forward simulations, contemporary MCMC methods become prohibitive
- Intractability of MCMC methods for large-scale statistical inverse problems can be traced to their black-box treatment of the parameter-to-observable map
- Goal: develop MCMC methods that exploit the structure of the parameter-to-observation map (including its derivatives), as has been done successfully in PDE-constrained optimization

- 1 Bayesian framework for statistical inverse problems
- 2 MCMC sampling
- 3 Langevin methods and stochastic Newton
- 4 Numerical results: 1D inverse elastic medium scattering
- 5 Numerical results: 2D inverse electromagnetic obstacle scattering

Bayesian formulation for statistical inverse problem

Given:

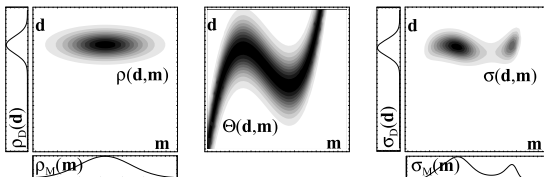
$\pi_{\text{pr}}(\mathbf{x})$:= prior p.d.f. of model parameters \mathbf{x}

$\pi_{\text{obs}}(\mathbf{y})$:= prior p.d.f. of the observables \mathbf{y}

$\pi_{\text{model}}(\mathbf{y}|\mathbf{x})$:= conditional p.d.f. relating \mathbf{y} and \mathbf{x}

Then *posterior p.d.f. of model parameters* is given by:

$$\begin{aligned}\pi_{\text{post}}(\mathbf{x}) &\stackrel{\text{def}}{=} \pi_{\text{post}}(\mathbf{x}|\mathbf{y}_{\text{obs}}) \\ &\propto \pi_{\text{pr}}(\mathbf{x}) \int_{\mathbf{y}} \frac{\pi_{\text{obs}}(\mathbf{y})\pi_{\text{model}}(\mathbf{y}|\mathbf{x})}{\mu(\mathbf{y})} d\mathbf{y} \\ &\propto \pi_{\text{pr}}(\mathbf{x}) \pi(\mathbf{y}_{\text{obs}}|\mathbf{x})\end{aligned}$$



From A. Tarantola, *Inverse Problem Theory*, SIAM, 2005

Gaussian additive noise

Given the parameter-to-observable map $\mathbf{y} = \mathbf{f}(\mathbf{x})$, a common noise model is Gaussian additive noise:

$$\mathbf{y}_{\text{obs}} = \mathbf{f}(\mathbf{x}) + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Gamma}_{\text{noise}})$$

If the prior is taken as Gaussian with mean \mathbf{x}_{pr} and covariance $\boldsymbol{\Gamma}_{\text{pr}}$, then the posterior can be written

$$\pi_{\text{post}}(\mathbf{x}) \propto \exp\left(-\frac{1}{2} \|\mathbf{f}(\mathbf{x}) - \mathbf{y}_{\text{obs}}\|_{\boldsymbol{\Gamma}_{\text{noise}}^{-1}}^2 - \frac{1}{2} \|\mathbf{x} - \mathbf{x}_{\text{pr}}\|_{\boldsymbol{\Gamma}_{\text{pr}}^{-1}}^2\right)$$

Note that “most likely” point is given by

$$\begin{aligned} \mathbf{x}_{\text{MAP}} &\stackrel{\text{def}}{=} \arg \max_{\mathbf{x}} \pi_{\text{post}}(\mathbf{x}) \\ &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{f}(\mathbf{x}) - \mathbf{y}_{\text{obs}}\|_{\boldsymbol{\Gamma}_{\text{noise}}^{-1}}^2 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_{\text{pr}}\|_{\boldsymbol{\Gamma}_{\text{pr}}^{-1}}^2 \end{aligned}$$

This is an (appropriately weighted) deterministic inverse problem!

Gaussian additive noise, linear inverse problem

Suppose further the parameter-to-observable map is linear, i.e.

$$\mathbf{y} = \mathbf{F}\mathbf{x}$$

Then the posterior can be written

$$\pi_{\text{post}}(\mathbf{x}) \propto \exp\left(-\frac{1}{2} \|\mathbf{F}\mathbf{x} - \mathbf{y}_{\text{obs}}\|_{\mathbf{\Gamma}_{\text{noise}}^{-1}}^2 - \frac{1}{2} \|\mathbf{x} - \mathbf{x}_{\text{pr}}\|_{\mathbf{\Gamma}_{\text{pr}}^{-1}}^2\right)$$

The posterior is then Gaussian with

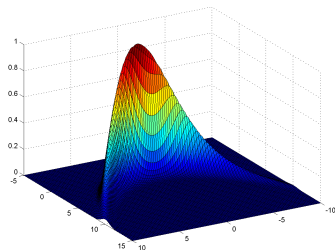
$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}_{\text{MAP}}, \mathbf{\Gamma}_{\text{post}})$$

The covariance is the inverse Hessian of the negative log posterior:

$$\begin{aligned}\mathbf{\Gamma}_{\text{post}}^{-1} &= \mathbf{F}^T \mathbf{\Gamma}_{\text{noise}}^{-1} \mathbf{F} + \mathbf{\Gamma}_{\text{pr}}^{-1} \\ &= \nabla_{\mathbf{x}}^2 (-\log \pi_{\text{post}})\end{aligned}$$

I.e., the covariance is given by the inverse Hessian of the regularized misfit function that is minimized by deterministic methods

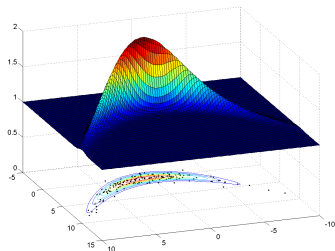
- 1 Bayesian framework for statistical inverse problems
- 2 MCMC sampling**
- 3 Langevin methods and stochastic Newton
- 4 Numerical results: 1D inverse elastic medium scattering
- 5 Numerical results: 2D inverse electromagnetic obstacle scattering



Example Probability Density

Given a probability density $\pi(\mathbf{x})$:

- How do we interrogate the distribution?
- Often high dimensional
- Computationally expensive



Sampled Probability Density

Given a probability density $\pi(\mathbf{x})$:

- How do we interrogate the distribution?
- Often high dimensional
- Computationally expensive

The MCMC Approach

- Replace $\pi(\mathbf{x})$ by a sample chain $\{\mathbf{x}_k\}$
- Compute using ergodic averages

$$\mathbb{E}[f(X)] = \int_{\mathbb{R}^n} f(x)\pi(dx) \approx \frac{1}{N} \sum_{j=1}^N f(x_k)$$

Metropolis-Hastings algorithm

- 1 $\mathbf{x}_k \leftarrow \mathbf{x}_0$
- 2 $k \leftarrow 0$
- 3 Choose a point \mathbf{y} from the proposal density $q(\mathbf{x}_k, \cdot)$
- 4 $\alpha \leftarrow \min \left(1, \frac{\pi(\mathbf{y})q(\mathbf{y}, \mathbf{x}_k)}{\pi(\mathbf{x}_k)q(\mathbf{x}_k, \mathbf{y})} \right)$
- 5 If $\alpha > \text{rand}([0, 1])$ Then
 Accept: $\mathbf{x}_{k+1} = \mathbf{y}$
 Otherwise
 Reject: $\mathbf{x}_{k+1} = \mathbf{x}_k$
 End If
- 6 $k \leftarrow k + 1$
- 7 Repeat from step 3

Some proposal functions

The best proposal function is just the pdf itself:

- $q(\mathbf{x}_k, \mathbf{y}) = \pi(\mathbf{y})$
- $\alpha(\mathbf{x}_k, \mathbf{y}) = \min\left(1, \frac{\pi(\mathbf{y})\pi(\mathbf{x}_k)}{\pi(\mathbf{x}_k)\pi(\mathbf{y})}\right) \equiv 1$
- Would like to use an approximation $\tilde{\pi}(\mathbf{y})$

Gaussian random walks:

- $q(\mathbf{x}_k, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Gamma})$
- Lots of freedom in choosing $\boldsymbol{\mu}$ and $\boldsymbol{\Gamma}$
- Both can depend on \mathbf{x}_k

Many others:

- Hybrid Monte Carlo
- Gibbs sampling

Approaches to reducing the cost of MCMC

Basic difficulty: evaluating posterior density requires forward solve;
philosophies to circumvent:

- *Reduce-then-sample*:
 - Reduced model of the forward problem
 - POD (e.g. Wang and Zabaras, Willcox et al., Patera et al.)
 - Reduced model of the outputs
 - PC (e.g. Marzouk, Najm, and Rahn, Zabaras et al., Marzouk and Xiu)
 - Gaussian process model (e.g. O'Hagan and Kennedy, Higdon)
 - "Preconditioned" MCMC using reduced order models
 - Higdon, Lee, and Holloman
 - Christen and Fox
 - Efendiev, Hou, and Luo
 - Efendiev, Datta-Gupta, Ginting, Ma, and Mallick
- We will pursue: *Sample-then-reduce*
 - Exploit structure of $\pi_{\text{post}}(\boldsymbol{x})$, in particular local covariance related to Hessian of deterministic inverse problem
 - Build on fast algorithms from deterministic inverse solvers

- 1 Bayesian framework for statistical inverse problems
- 2 MCMC sampling
- 3 Langevin methods and stochastic Newton**
- 4 Numerical results: 1D inverse elastic medium scattering
- 5 Numerical results: 2D inverse electromagnetic obstacle scattering

Background: Langevin dynamics

Langevin dynamics

- Stochastic differential equation (continuous in time)
 - $\pi(\boldsymbol{x})$ is a stationary solution
 - \Rightarrow Trajectories sample $\pi(\boldsymbol{x})$
- Uses derivative information of $\pi(\boldsymbol{x})$
- Can be preconditioned for better performance

Discrete Langevin dynamics

- Discretization with timestep Δt introduces bias
- Use as proposal distribution for Metropolis-Hastings MCMC (U. Grenander and M. Miller, 1994)

Langevin-MCMC used for PDE-based inverse problems by A. Stuart and Y. Efendiev

Preconditioned Langevin MCMC

Given the target density $\pi(\mathbf{x})$, the associated Langevin SDE is given by:

$$d\mathbf{X}_t = -\mathbf{A}\nabla_{\mathbf{x}}(-\log \pi)dt + \sqrt{2}\mathbf{A}^{1/2}d\mathbf{W}_t$$

Discretize with a timestep Δt to derive Langevin proposal:

$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k - \mathbf{A}\nabla_{\mathbf{x}}(-\log \pi)\Delta t + \sqrt{2\Delta t}\mathbf{A}^{1/2}\mathcal{N}(\mathbf{0}, \mathbf{I})$$

Notes:

- Preconditioner \mathbf{A} must be symmetric positive definite
- Process is ergodic (convergence of time averages)
- \mathbf{W}_t is i.i.d. vector of standard Brownian motions
- \mathbf{W}_t has independent increments given by
 - $\mathbf{W}_{(t+\Delta t)} - \mathbf{W}_t \sim \mathcal{N}(\mathbf{0}, \Delta t \mathbf{I})$

Stochastic Newton's method

Langevin MCMC proposal given by:

$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k - \mathbf{A} \nabla_{\mathbf{x}}(-\log \pi) \Delta t + \sqrt{2\Delta t} \mathbf{A}^{1/2} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Take \mathbf{A} to be the inverse of the (local) Hessian (of the negative log posterior) and set $\Delta t = 1$:

$$\begin{aligned} \mathbf{A} &\equiv \mathbf{H}(\mathbf{x})^{-1} = \nabla_{\mathbf{x}}^2(-\log \pi(\mathbf{x}))^{-1} \\ &= (\mathbf{F}^T \mathbf{\Gamma}_{\text{noise}}^{-1} \mathbf{F} + \mathbf{\Gamma}_{\text{pr}}^{-1})^{-1} \quad (\text{e.g. for Gaussian noise and prior}) \end{aligned}$$

Then we have the stochastic equivalent of Newton's method:

$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k - \mathbf{H}^{-1} \nabla_{\mathbf{x}}(-\log \pi) + \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$$

Stochastic Newton: Optimal sampling of Gaussians

When the target density $\pi(\mathbf{x})$ is Gaussian, $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Gamma})$:

$$-\log \pi(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}\|_{\boldsymbol{\Gamma}^{-1}}$$

Apply Stochastic Newton:

$$\begin{aligned}\mathbf{x}_{k+1}^{\text{prop}} &= \mathbf{x}_k - \mathbf{H}^{-1} \nabla_{\mathbf{x}}(-\log \pi) + \mathbf{H}^{-1/2} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \mathbf{x}_k - \boldsymbol{\Gamma} \boldsymbol{\Gamma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) + \boldsymbol{\Gamma}^{1/2} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \boldsymbol{\mu} + \boldsymbol{\Gamma}^{1/2} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Gamma})\end{aligned}$$

Samples \mathbf{x}_k act like *independent* draws from the true pdf

Deterministic vs. Stochastic Newton

Deterministic Newton:

- Given a cost function $-\log \pi(\mathbf{x})$
- $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}^{-1} \nabla_{\mathbf{x}}(-\log \pi)$
- Minimizes local quadratic approximation at each step

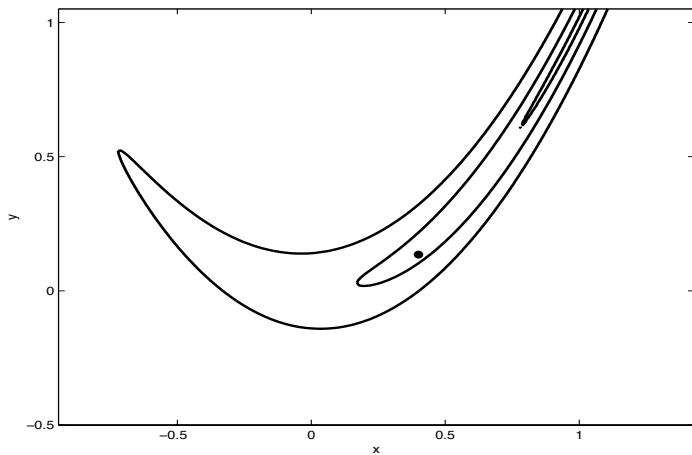
Stochastic Newton:

- Given a probability density $\pi(\mathbf{x})$
- $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}^{-1} \nabla_{\mathbf{x}}(-\log \pi) + \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$
- Samples local Gaussian approximation at each step

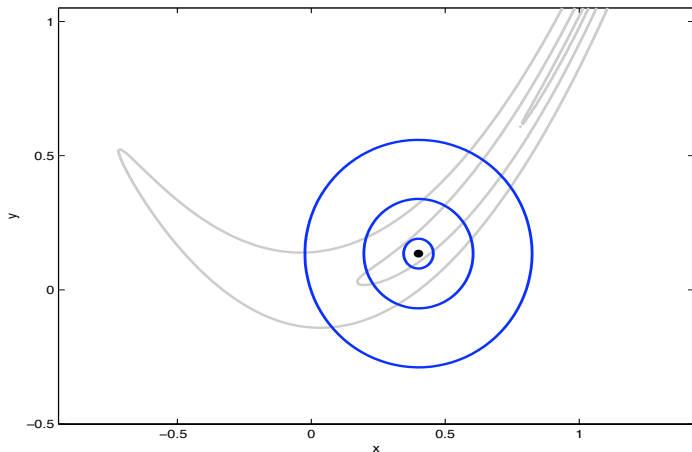
Unpreconditioned Langevin resembles steepest descent

- $\mathbf{x}_{k+1} = \mathbf{x}_k - \Delta t \nabla_{\mathbf{x}}(-\log \pi) + \sqrt{2\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I})$

Rosenbrock illustration

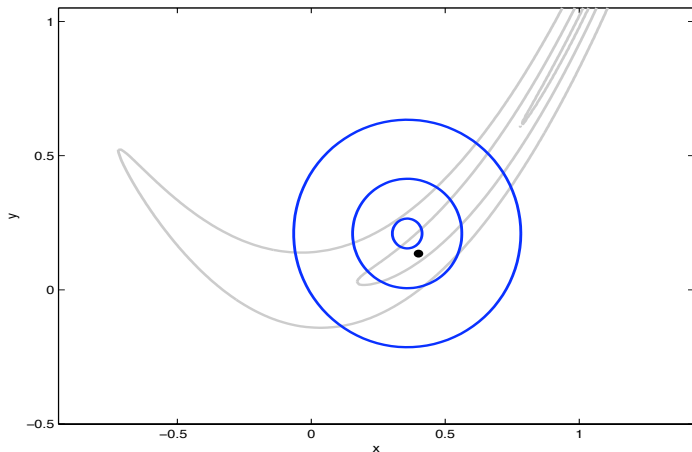


Rosenbrock illustration: Random walk



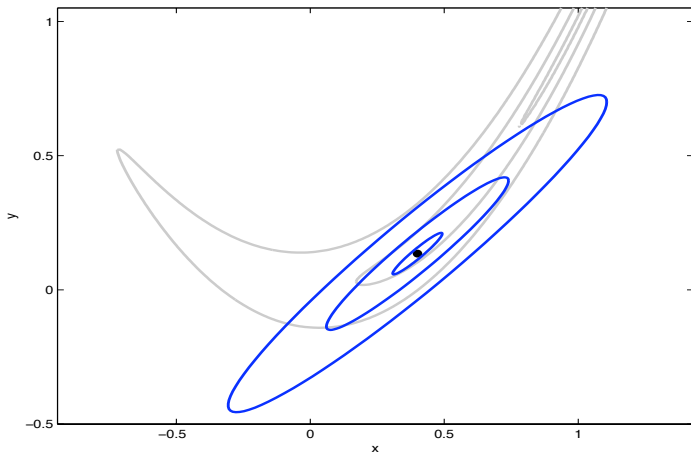
$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k + \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Rosenbrock illustration: Unpreconditioned Langevin



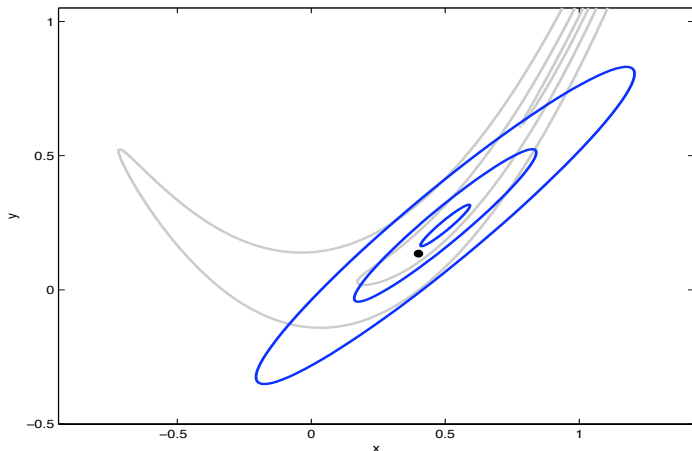
$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k - \Delta t \nabla_{\mathbf{x}} (-\log \pi) + \sqrt{2\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Rosenbrock illustration: Hessian-based Gaussian proposal



$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k + \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$$

Rosenbrock illustration: Hessian-preconditioned Langevin



$$\mathbf{x}_{k+1}^{\text{prop}} = \mathbf{x}_k - \mathbf{H}^{-1} \nabla_{\mathbf{x}} (-\log \pi) + \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$$

Stochastic Newton: Large-scale issues

At each MCMC step we need to

- solve systems of form $\mathbf{H}\mathbf{v} = \mathbf{b}$
- evaluate matvecs of form $\mathbf{H}^{-\frac{1}{2}}\mathbf{w}$

Key idea: **never** form \mathbf{H} ; instead:

- recognize that \mathbf{H} is sum of data misfit term, which is often equivalent to a compact operator, and (the inverse of) a smoothing prior, which is often equivalent to a differential operator:

$$\mathbf{F}^T \mathbf{\Gamma}_{\text{noise}}^{-1} \mathbf{F} + \mathbf{\Gamma}_{\text{pr}}^{-1}$$

- develop fast algorithms for low rank (in particular, truncated spectral decomposition) approximation of data misfit operator; often require constant number of forward/adjoint solves, independent of problem size
- combine with Sherman-Morrison-Woodbury to invert/factor (requires constant number of forward/adjoint solves)
- construct fast (multilevel) preconditioners for Hessian

- 1 Bayesian framework for statistical inverse problems
- 2 MCMC sampling
- 3 Langevin methods and stochastic Newton
- 4 Numerical results: 1D inverse elastic medium scattering
- 5 Numerical results: 2D inverse electromagnetic obstacle scattering

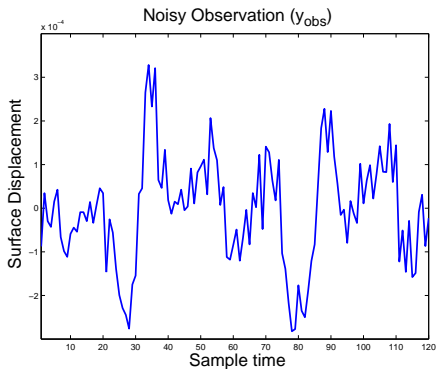
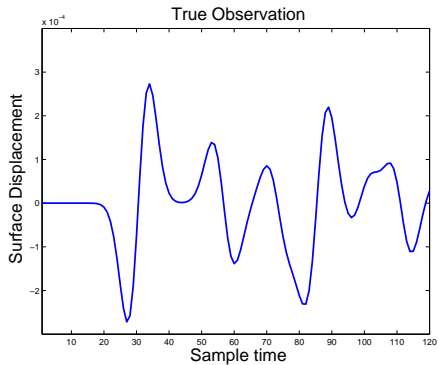
1D inverse elastic medium scattering

Forward problem: Given elastic modulus $a(x)$, solve 1D wave equation for Ricker wavelet source $g(t)$ to obtain observed waveform $y(0, t)$ at top surface:

$$\begin{aligned}\rho \frac{\partial^2 y}{\partial t^2} - \frac{\partial}{\partial x} \left(a(x) \frac{\partial y}{\partial x} \right) &= \delta(x - 0)g(t) \\ a \frac{\partial y}{\partial x} \Big|_{x=0} &= 0 \\ \sqrt{\rho a} \frac{\partial y}{\partial t} \Big|_{x=1} &= -a \frac{\partial y}{\partial x} \Big|_{x=1} \\ y|_{t=0} &= 0 \\ \dot{y}|_{t=0} &= 0\end{aligned}$$

Inverse problem: Given observed noisy waveform and layered parametrization of elastic medium modulus $a(x)$ with associated prior, recover layer moduli and associated uncertainty

Sample noisy observations



Prior and likelihood distributions

Gaussian smoothness prior between layers:

- Covariance matrix $\mathbf{\Gamma}$ between layers i and j :

$$\Gamma_{ij} = \theta_1 \exp\left(\frac{-(x_i - x_j)^2}{2\theta_2^2}\right)$$

- Prior mean is 5 for all layers

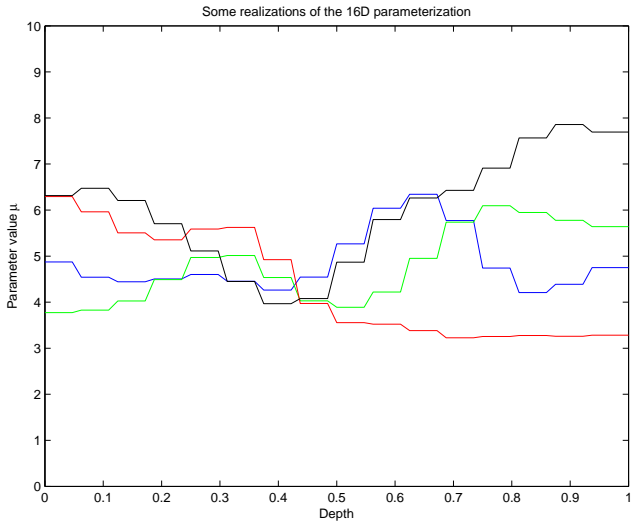
Gaussian likelihood function:

$$\pi_{\text{like}}(\mathbf{y}_{\text{obs}}|\mathbf{a}) = \exp\left(-\frac{1}{2}(\mathbf{f}(\mathbf{a}) - \mathbf{y}_{\text{obs}})^T \mathbf{\Gamma}_{\text{noise}}^{-1} (\mathbf{f}(\mathbf{a}) - \mathbf{y}_{\text{obs}})\right)$$

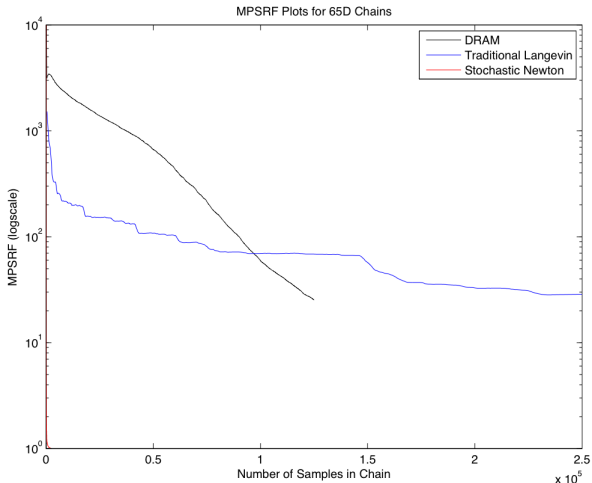
We wish to sample the posterior distribution:

$$\pi_{\text{post}}(\mathbf{a}|\mathbf{y}_{\text{obs}}) \propto \pi_{\text{pr}}(\mathbf{a})\pi_{\text{like}}(\mathbf{y}_{\text{obs}}|\mathbf{a})$$

Some realizations of 16-layer prior



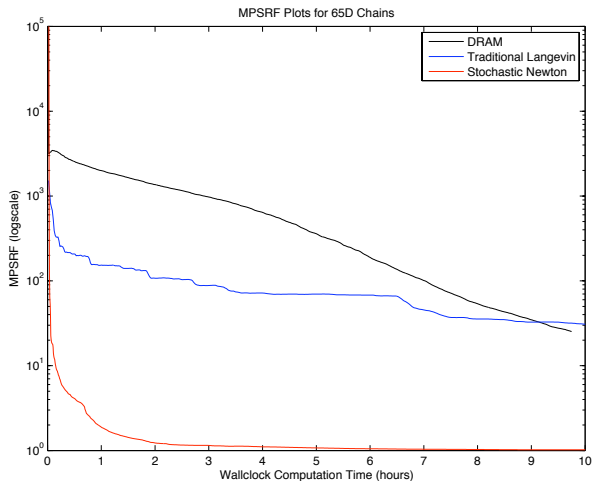
Convergence comparison for 65-layer problem



Multivariate potential scale reduction factor (MPSRF) convergence statistic for 65-layer problem

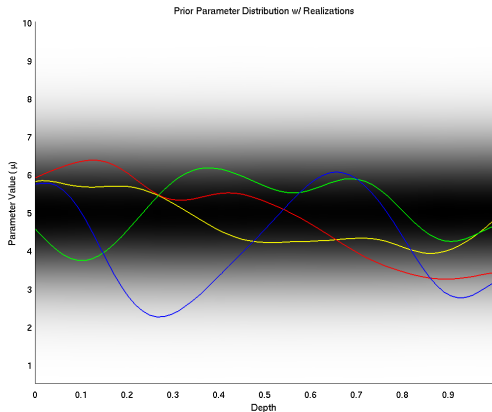
unpreconditioned Langevin vs. stochastic Newton vs. Adaptive Metropolis

Rescaled 65-layer MPSRF convergence



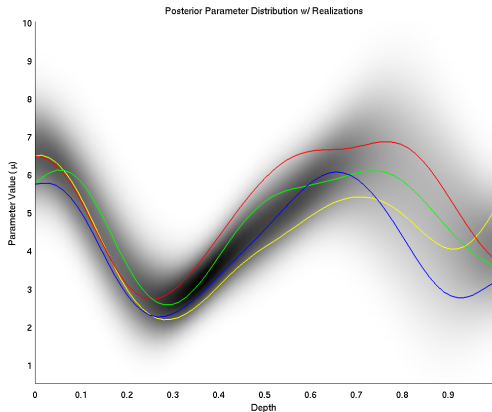
MPSRF statistic for 65-layer problem as function of wall clock time
(dense implementation – not recommended!)

65-layer prior



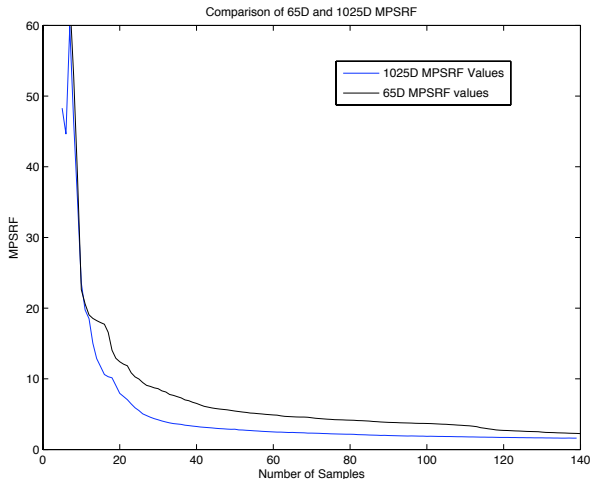
Density plot of marginal pdfs of prior of elastic moduli of 65 layers
Blue curve is “truth” modulus used to synthesize observations
Other colors are draws from prior

65-layer posterior



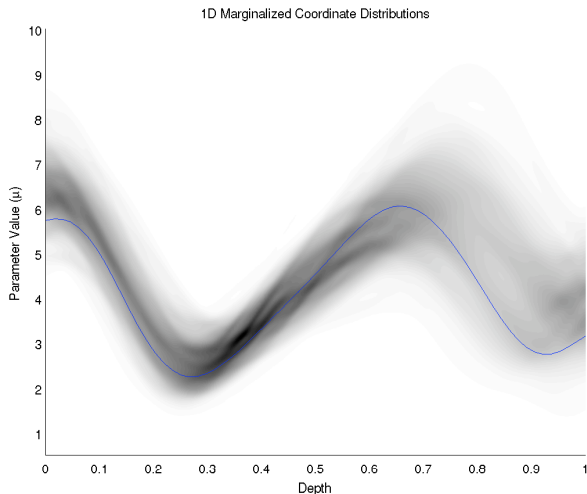
Density plot of marginal pdfs of posterior of elastic moduli of 65 layers
Blue curve is “truth” modulus used to synthesize observations
Other colors are draws from posterior

MPSRF convergence statistic for 1025-layer problem



MPSRF statistic for 1025-layer problem compared with 65-layer
(1025-layer results based on fast low-rank implementation)

1025 layer posterior



Density plot of marginal pdfs of posterior of elastic moduli of 1025 layers
Blue curve is “truth” modulus used to synthesize observations

Outline

- 1 Bayesian framework for statistical inverse problems
- 2 MCMC sampling
- 3 Langevin methods and stochastic Newton
- 4 Numerical results: 1D inverse elastic medium scattering
- 5 Numerical results: 2D inverse electromagnetic obstacle scattering

Statistical inverse electromagnetic obstacle scattering

Maxwell's equations:

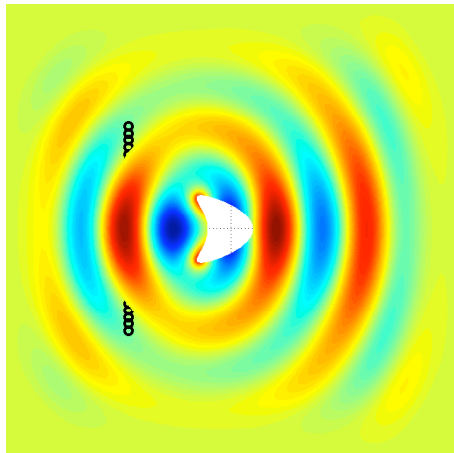
$$\nabla \times \mathbf{E} = -\mu \frac{\partial \mathbf{H}}{\partial t}$$

$$\nabla \times \mathbf{H} = \epsilon \frac{\partial \mathbf{E}}{\partial t}$$

$$\nabla \cdot \mathbf{E} = 0$$

$$\nabla \cdot \mathbf{H} = 0$$

- \mathbf{E} ... Electric field
- \mathbf{H} ... Magnetic field
- μ ... permeability
- ϵ ... permittivity



Scattered electric field E_z

Problem setup

- Forward code based on C version of Matlab code from J. Hesthaven and T. Warburton, *Nodal Discontinuous Galerkin Methods*
- Extended to include adjoint-based shape gradient and shape Hessian
- Discontinuous Galerkin 3rd-order spectral elements in space
- Fourth-order, five-stage explicit Runge Kutta scheme in time
- Prior favors shape with small area

$$\pi_{\text{pr}} = \exp\left(-\frac{\beta}{2} \int_0^{2\pi} r^2 d\theta\right), \quad \beta = 0.1$$

- Shape is parametrized by 6 cosine modes

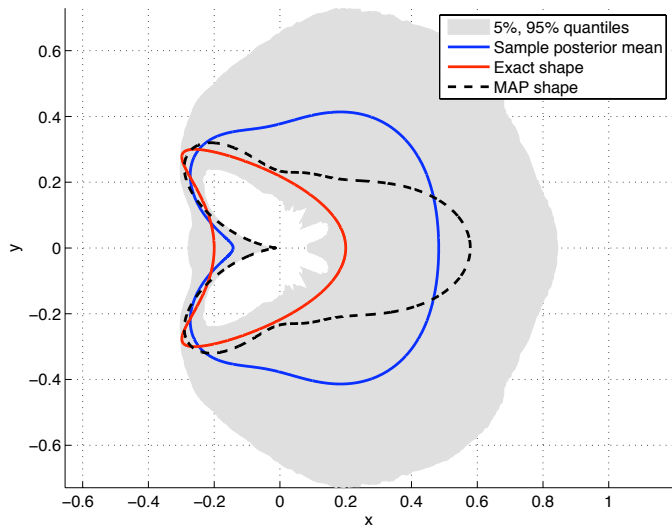
$$r(\theta) = \sum_{i=0}^5 a_i \cos(i\theta)$$

- Computational domain $\Omega = \{(x, y) : -1 \leq x, y \leq 1\}$
- PML domain $\Omega_{PML} = \{(x, y) : 1 \leq |x|, |y| \leq 2\}$
- Kite shape to generate synthetic observations:

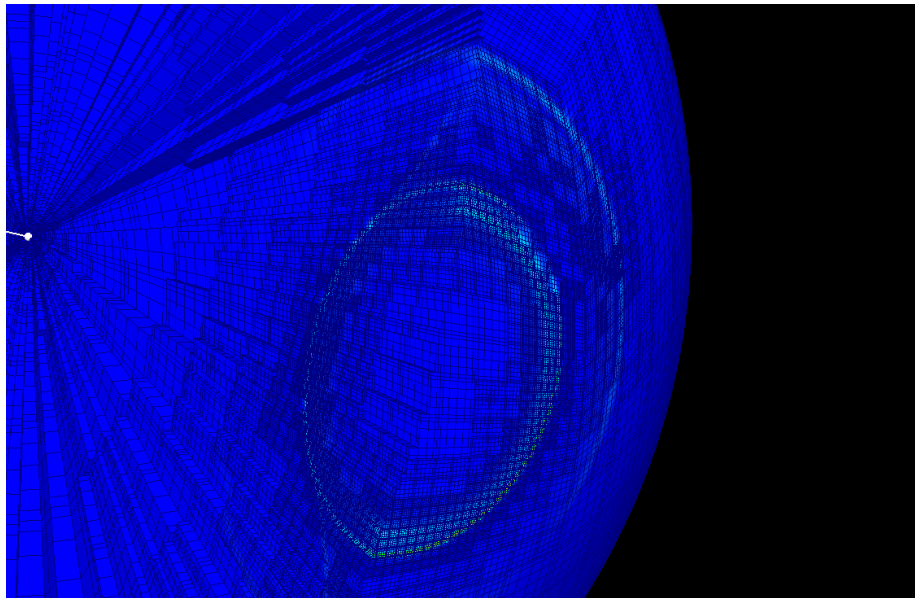
$$x = 0.2 [\cos(\theta) + 0.65(\cos(2\theta) - 1)], \quad y = 0.3 \sin(\theta)$$

- Incident wave $E_z^I = \cos(8(t - x))$, $H_x = 0$, $H_y = 0$ from left
- 31 observation points: $x = -0.9$, $y = \text{linspace}(-0.9, 0.9, 31)$
- E_z, H_x, H_y are observed in $0 \leq t \leq \pi$ with 5% Gaussian noise
- $\Delta t = 10^{-3}$ with 3324 time steps
- Mesh size $h_{\min} = 0.05$ ($\sim 135,000$ DOF)

Results using stochastic Newton for 6-parameter problem



3D Inverse elastic-acoustic wave propagation



Conclusions

- Target: PDE-based Bayesian inverse problems
- Stochastic Newton (inverse-Hessian-preconditioned Langevin MCMC)
 - motivated by connection to deterministic Newton method
 - exactly samples a Gaussian posterior
 - dense implementation shows vast improvement over black-box MCMC (adaptive Metropolis)
 - low rank implementation able to solve 1025-dimension inverse scattering problem
- Current work aimed at capitalizing on advances in deterministic PDE-based optimization and inverse methods to improve stochastic Newton
 - inexact Newton (Steihaug-Eisenstat-Walker ideas)
 - trust region methods
 - exploit “compact + differential” structure of Hessians (e.g. low rank approximations, Fredholm-multigrid type preconditioners)
- Exploiting deterministic PDE inverse problem structure should play an important role in scaling MCMC to high dimensions and expensive forward problems