The Joint Power of Deterministic and Bayesian Approaches for (Seismic) Inverse Problems

Yunan Yang, Daniel Appelö (CU Boulder), Matt Dunlop (NYU), Björn Engquist (UT-Austin), Kui Ren (Columbia), Alex Townsend (Cornell) April 27, 2020

This work is partially supported by NSF DMS-1913129.

Cornell SCAN Seminar over Zoom

General Large-Scale Inverse Problem



Given (many) $\{X, Y\}$, find *m*.

The Model *F*(*m*)

Model (m) (PDE) (NNets)

F is given; we just find m (e.g., PDEs). OR F is not known; m depends on F.

F is given; we just find m (e.g., PDEs).

- Pro: We know the best (exact) forward problem!
- Con: The forward and inverse problems are so nonlinear!

OR

F is not known; we are free to choose (e.g., XXX-net).

- Pro: The freedom to modify it to a "better" map (Over-Parametrization, ReLu)
- Con: Trial and error to build the model

Seismic Inversion: Earthquake Source, Hydrocarbons, etc.

Seismic inversion is one of the inherently more difficult families of large-scale nonlinear inverse problems.



Seismic Inversion



Waveform measurements from receivers at the surface



Subsurface properties (i.e. wave velocity or material density)

Forward Problem

 $\mathcal{F}: m
ightarrow u|_{\Gamma}$, $\Gamma \subseteq \partial \Omega$ or Ω

Inverse Problem

 $\mathcal{G}: u|_{\Gamma} \to m$

 ${\mathcal F}$ and ${\mathcal G}$ are often nonlinear.

 $m^* = \underset{m}{\operatorname{argmin}} J(f(m), g)$ $f(m) = u|_{\Gamma}$

J is an objective function measuring the difference between f and g.

Forward Wave Propagation

$$\begin{cases} m(\mathbf{x}) \frac{\partial^2 u(\mathbf{x}, t)}{\partial t^2} - \triangle u(\mathbf{x}, t) = s(\mathbf{x}, t) \\ \text{Zero i.c. in half-space } \Omega \\ \text{Neumann b.c. on } \partial \Omega \end{cases}$$

 $m(\mathbf{x}) = \frac{1}{c(\mathbf{x})^2}$, $c(\mathbf{x})$ is the wave velocity



m

Important Components in the deterministic approach



The objective function



Traditional Least-Squares (L² norm) Objective Function

$$J(m) = \frac{1}{2} \sum_{r} \int |f(x_r, t; m) - g(x_r, t)|^2 dt,$$
 (1)

- observed data g,
- simulated data $f(m) = u|_{\Gamma}$,
- receiver x_r,
- the model parameter m,
- Regularization is often added in (1).

Main Challenges

- 1. Local minima trapping
- 2. Sensitive to noise

The shift and dilation are typical effects from variations in velocity parameter m(x) = m (constant). For example:

$$\begin{cases} m \frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2}, & x > 0, t > 0, \\ u = 0, \quad \frac{\partial u}{\partial t} = 0, & x > 0, t = 0, \\ u = f(t), & x = 0, t > 0. \end{cases}$$

The solution to the equation is $u(x, t; m) = f(t - \sqrt{mx})$.

For fixed *x*, variation in *m* relates **shifts** in the signal.

For fixed *t*, variation in *m* generates the **dilation** in data.

[Engquist, Froese & Y, 2016]

Motivation of Using the Wasserstein Distance (EMD)



The Quadratic Wasserstein Distance

For $f, g \in \mathcal{P}(\Omega)$ ($f, g \ge 0$ and $\int f = \int g = 1$), the quadratic Wasserstein distance is formulated as

$$W_2(f,g) = \left(\inf_{T \in \mathcal{M}} \int |x - T(x)|^p f(x) dx\right)^{\frac{1}{2}}$$
(2)

 \mathcal{M} : the set of all maps that rearrange the distribution f into g.

[Monge, 1781]



Synthetic data f (left) and observed data g (right)

[Monge, 1781]



Synthetic data f (left) and observed data g (right)

Optimal Transport



Synthetic data f (left) and observed data g (right)

Optimal Transport



Synthetic data f (left) and observed data g (right)

Let $\{e_k\}_{k=1}^d$ be standard basis of the Euclidean space \mathbb{R}^d . Assume $s_k \in \mathbb{R}, \lambda_k \in \mathbb{R}^+, k = 1, \dots, d$ and $A = \text{diag}(1/\lambda_1, \dots, 1/\lambda_d)$. We define f_{Θ} as jointly the translation and dilation of g:

$$f_{\Theta}(\mathbf{x}) = \det(\mathbf{A})g(\mathbf{A}(\mathbf{x} - \sum_{k=1}^{d} s_k e_k)), \Theta = \{s_1, \dots, s_d, \lambda_1, \dots, \lambda_d\}.$$

Theorem (Convexity of W₂ in translation and dilation)

The optimal map between $f_{\Theta}(x)$ and g(y) is $y = T_{\Theta}(x)$ where $\langle T_{\Theta}(x), e_k \rangle = \frac{1}{\lambda_k} (\langle x, e_k \rangle - s_k), k = 1, \dots, d.$

Moreover, $I(\Theta) = W_2^2(f_{\Theta}(x), g)$ is a convex function of Θ .

[Y, 2019]

Tackling Nonconvexity







[Y-Engquist-Sun-Hamfeldt, 2016]

2. More Robust w.r.t. Noise (Perturbation)

Given strictly positive probability density $f = d\nu$, we can define a Laplace-type linear operator

$$L = -\Delta +
abla (-\log f) \cdot
abla$$

which satisfies the fundamental integration by parts formula:

$$\begin{split} \int_{\mathbb{R}^d} (Lh_1)h_2 d\nu &= \int_{\mathbb{R}^d} h_1 (Lh_2) d\nu = \int_{\mathbb{R}^d} \nabla h_1 \cdot \nabla h_2 d\nu. \\ \|h\|_{L^2(f)}^2 &= \int_{\mathbb{R}^d} h^2 d\nu, \quad \|h\|_{\dot{H}^1(f)}^2 = \int_{\mathbb{R}^d} |\nabla h|^2 d\nu, \\ h\|_{\dot{H}^{-1}(f)}^2 &\coloneqq \sup \left\{ \int_{\mathbb{R}^d} h\varphi d\nu \ \bigg| \ \|\varphi\|_{\dot{H}^1(f)}^2 \leq 1 \right\} = \int_{\mathbb{R}^d} h(L^{-1}h) d\nu. \end{split}$$

If f= 1, we reconstruct the unweighted $\dot{\mathcal{H}}_{(\mathbb{R}^d)}^{-1}$ seminorm.

Asymptotic Connection

If μ is the probability measure and $d\pi$ is an infinitesimal perturbation that has zero total mass, then

$$W_2(\mu, \mu + d\pi) = \|d\pi\|_{\dot{\mathcal{H}}_{(d\mu)}^{-1}} + o(d\pi).$$
(3)

Non-Asymptotic Connection

If both $f = d\mu$ and $g = d\nu$ are bounded from below and above by constants c_1 and c_2 , we have the following *non-asymptotic* equivalence between W_2 and $\dot{\mathcal{H}}_{(d\mu)}^{-1}$:

$$\frac{1}{\mathfrak{c}_{2}}\|\mu-\nu\|_{\dot{\mathcal{H}}_{(\mathbb{R}^{d})}^{-1}} \leq W_{2}(\mu,\nu) \leq \frac{1}{\mathfrak{c}_{1}}\|\mu-\nu\|_{\dot{\mathcal{H}}_{(\mathbb{R}^{d})}^{-1}}, \quad (4)$$

A linear inverse problem of finding m from noisy data g_{δ}

$$Am = g_{\delta}.$$
 (5)

A is diagonal in the Fourier domain:

$$\widehat{A}(\boldsymbol{\xi}) \sim \langle \boldsymbol{\xi} \rangle^{-\alpha}.$$
 (6)

We seek the solution by minimizing the objective functional

$$\mathcal{O}_{\mathcal{H}^{s}}(m) \equiv \frac{1}{2} \|f(m) - g\|_{\mathcal{H}^{s}}^{2} := \frac{1}{2} \int_{\mathbb{R}^{d}} \langle \boldsymbol{\xi} \rangle^{2s} |\widehat{f}(m)(\boldsymbol{\xi}) - \widehat{g}(\boldsymbol{\xi})|^{2} d\boldsymbol{\xi},$$
(7)

[Engquist-Ren-Y, 2019]

The solution at frequency $\boldsymbol{\xi}$ is therefore

$$\widehat{m}(\boldsymbol{\xi}) = \left(\widehat{A}^*(\boldsymbol{\xi})(\langle \boldsymbol{\xi} \rangle^{2s}\widehat{A})\right)^{-1}\widehat{A}^*(\boldsymbol{\xi})\left(\langle \boldsymbol{\xi} \rangle^{2s}\widehat{g}_{\delta}(\boldsymbol{\xi})\right).$$

We can then perform an inverse Fourier transform to find the solution in physical space. The result is

$$m = (A^* P A)^{-1} A^* P g_{\delta}, \qquad P := (\mathcal{I} - \Delta)^{s/2},$$

where the operator $(\mathcal{I} - \Delta)^{s/2}$ is defined through the relation

$$(\mathcal{I}-\Delta)^{s/2}m=\mathcal{F}^{-1}\Big(\langle \boldsymbol{\xi}\rangle^{s}\widehat{m}\Big),$$

S = 0, S > 0, S < 0.

[Engquist-Ren-Y, 2019]

Differences Between W_2 and \dot{H}^{-1}



Top row: Geodesics in the \dot{H}^{-1} space Bottom row: Geodesics in the W_2 space

[Papadakis-Peyré-Oudet, 2013]



$$m^* = \underset{m}{\operatorname{argmin}} J(m) + R(m)$$

Regularization does not have to be in the form of R(m).

- The choice of the objective function
- The choice of the data e.g., low-frequency data recovers low-wavenumber model
- The choice of numerical discretization
- The optimization algorithm (fixed step size)

Acceleration Methods



Treat Gradient Descent as a Fixed-Point Iteration

$$p_{k+1} = p_k - \eta \frac{\partial J}{\partial p} \Big|_{p=p_k} = G(p_k).$$
$$p^* = G(p^*).$$

Steepest descent \Leftrightarrow Picard Iteration applied to G

Can we do better?

Input: Given p_0 and $m \ge 1$. *G* is a given fixed-point operator. Set $p_1 = G(p_0)$. **for** k = 0, 1, 2, ... **do Step 1:** Set $m_k = \min(m, k)$ and $F_k = (f_{k-m_k}, ..., f_k)$, where

Step 1. Set $m_k = \min(m, k)$ and $F_k = (j_{k-m_k}, \dots, j_k)$, where $f_i = G(p_i) - p_i$ is the residual. **Step 2:** Find $\alpha^{(k)} = (\alpha_0^{(k)}, \dots, \alpha_{m_k}^{(k)})^T$, $\sum_{i=0}^{m_k} \alpha_i^{(k)} = 1$ to minimize the sum of the weighted residual $||F_k \alpha^{(k)}||$. **Step 3:** Update p_{k+1} according to

$$p_{k+1} = \sum_{i=0}^{m_k} \alpha_i^{(k)} G(p_{k-m_k+i}).$$

end for

[D. G. Anderson, 1965]

Interesting connections with:

- GMRES ($m = \infty$)
- Multi-secant methods
- Krylov-space methods
- Momentum method
- Superlinear w/o approximating inverse Hessian
- Noncontractive (nonlinear) operator convergence

[D. G. Anderson, 1965]

Anderson Acceleration for Seismic Inversion



[Y-Townsend-Appeö,2020],[Y,2020]

Important Components in the deterministic approach



Important Components in the Bayesian approach



Likelihood Function



Likelihood function	Noise model assumption
$\int_D \ \mathcal{G}(u)(x,\cdot) - y(x,\cdot)\ _{L^2(T)}^2 \lambda(\mathrm{d} x)$	$y = \mathcal{G}(u) + \eta, \ \eta \sim N(o, I)$
$\int_D \ \mathcal{G}(u)(x,\cdot)-y(x,\cdot)\ ^2_{\dot{H}^{-1}(T)}\lambda(\mathrm{d} x)$	$y = \mathcal{G}(u) + \eta, \ \eta \sim N(o, -\Delta_T)$
$\int_D W_2^2\left(\widetilde{\mathcal{G}(u)}(x,\cdot),\widetilde{y}(x,\cdot) ight)\lambda(\mathrm{d} x)$	$\widetilde{\mathbf{y}} = \eta \cdot \widetilde{\mathcal{G}(u)}, \ \eta u \sim N(1, \mathcal{L}(u))$
$\int_{D} \left\ \frac{\mathcal{G}(u)(x,\cdot) - y(x,\cdot)}{(y)(x,\cdot)} \right\ _{L^{2}(T)}^{2} \lambda(\mathrm{d} x)$	$\mathbf{y} = \eta \cdot \mathcal{G}(\mathbf{u}), \ 1/\eta \sim N(1, \mathbf{I})$
The W_2 metric can be regarded as asymptotically from the	

state-dependent multiplicative noise data model: measurement error is proportional to the size of the quantity, and the distribution is affected by the model parameter.

[Dunlop-Y,2020]





Given $\nu_0 = N(m_1, C_1) \times N(m_2, C_2)$,

 $\pi_{o} = F^{\sharp}\nu_{o}, \quad F(v,w)(x) = u_{+}\mathbb{1}_{v(x)>o} + w(x)\mathbb{1}_{v(x)\leq o}.$

We don't have to know u^+ a priori.

[Dunlop-Iglesias-Stuart, 2016], [Iglesias-Lu-Stuart, 2016]

Level-Set Prior for Waveform Inversion



[Dunlop-Y,2020]

Why Model Below Reflector Is Hard to Recover by Reflections?





Why Model Below Reflector Is Hard to Recover by Reflections?



[Y-Engquist,2018]

Level-Set Prior for Waveform Inversion



Reconstructions arising from deterministic inversion

[Dunlop-Y,2020]

Level-Set Prior for Waveform Inversion



[Dunlop-Y,2020]

Uncertainty Quantification



Deterministic & Bayesian



For Large-Scale inverse problems

All my collaborators.





Thank you for the attention!