Optimal Transport for FWI : dealing with oscillatory signals

Aude Allain¹, R. Brossier², Q. Mérigot³, L. Métivier^{1,2}, E. Oudet¹, J.Virieux² Wednesday 3^{rd} May, 2017



¹LJK, CNRS, Univ. Grenoble Alpes, France ²ISTerre, Univ. Grenoble Alpes, France ³CEREMADE, CNRS, Univ. Paris Dauphine, France

05-03-2017, OILWS2 workshop, IPAM, UCLA, USA



Motivation for using optimal transport for FWI

A first implementation: a distance based on the dual form of W_1

Towards a generalization of the approach

Conclusion & perspectives



Motivation for using optimal transport for FWI

A first implementation: a distance based on the dual form of ${\it W}_1$

Towards a generalization of the approach

Conclusion & perspectives



- The formalism of optimal transport goes back to (1781) and is due to Gaspard Monge, French Mathematician (one funder of the Ecole Normale and Ecole Polytechnique)
- Initial purpose: find the more efficient way to move sand piles p_i to fill holes q_j during the building of a bridge

Optimal transport: an optimal assignment problem

- The formalism of optimal transport goes back to (1781) and is due to Gaspard Monge, French Mathematician (one funder of the Ecole Normale and Ecole Polytechnique)
- Initial purpose: find the more efficient way to move sand piles p_i to fill holes q_j during the building of a bridge



The matrix representing this assignment is denoted by

$$\gamma = \begin{pmatrix} 3 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 2 & 0 & 0 \end{pmatrix}$$
(1)





• Kantorovich (1942) relaxation: Optimal transport consists in finding the matrix γ (also known as transport plan) satisfying the **linear** programming problem

$$\min_{\gamma_{ij} \in \Pi(p,q)} \sum_{ij} \gamma_{ij} \|x_i - y_j\|,$$
where $\Pi(p,q) = \left\{\gamma_{ij} \ge 0, \sum_j \gamma_{ij} = p_i, \sum_i \gamma_{ij} = q_j\right\}$
(2)

and $||x_i - y_j||$ denotes a distance between x_i and y_j which is called the ground distance (often the Euclidean distance)



- Assume p_i and q_j are discretization of data d_{cal} and d_{obs}
- The minimal cost in the sense of the optimal transport defines a distance between the distribution d_{cal} and d_{obs}
- This distance is also referred to as the Wasserstein distance $W_p, p \geq 1$

$$W_{p}(d_{cal}, d_{obs}) = \left(\min_{\gamma_{ij} \in \Pi(d_{cal}, d_{obs})} \sum_{ij} \gamma_{ij} \|x_{i} - y_{j}\|^{p}\right)^{1/p}$$
(3)



+ Consider d_{cal} and d_{obs} are two 1D Gaussian distribution shifted by Δt



• Then we have

$$W_p^p(d_{cal}, d_{obs}) = |\Delta t|^p, \ p \ge 1$$
(4)

Convexity of the Wasserstein distance with respect to Δt



- Large-scale to medium-scale velocity perturbations (1000 m to 100 m at the exploration scale) are mainly responsible for **time-shifts**
 - ⇒ from an **inverse problem** point of view, recovering these velocity perturbations require to correctly interpret these **time-shifts**
- This is difficult when the data is compared with a L² distance because of cycle skipping/phase ambiguity





• Simple test: computing the L² misfit between two Ricker with respect to **time-shifts** yields a **multi-modal** misfit function





• Computing the optimal transport misfit between two Ricker with respect to **time-shifts** yields a **convex** misfit function (Engquist and Froese, 2014)





In the general case, computing the Wasserstein distance between d_{cal} and d_{obs} requires two assumptions to be satisfied

• positivity

$$d_{cal} \ge 0, \ d_{obs} \ge 0$$

• mass conservation

$$\sum_{i} (d_{cal})_i = \sum_{j} (d_{obs})_j$$



Consider d_{cal} and d_{obs} represent two discretized seismic traces

- d_{cal} and d_{obs} are oscillatory: positivity assumption breaks down
- the zero frequency of the signals is 0: the mass conservation holds

$$\int d_{cal}(t)dt = \int d_{obs}(t)dt = 0$$
(5)



Linear programming algorithms (simplex) for the computation of the Wasserstein distance require at least ${\cal O}(N^3)$ operations for size N discrete distributions

For acoustic time-domain FWI

- In 2D, the data is discretized with O(10³) time discretization points and O(10²) to receivers $\implies N = O(10^5)$
- In 3D, the data is discretized with O(10³) time discretization points and O(10⁴) to receivers $\implies N = O(10^7)$

Besides the positivity problem, we need an efficient numerical method to solve the optimal transport problem



Motivation for using optimal transport for FWI

A first implementation: a distance based on the dual form of W_1

Towards a generalization of the approach

Conclusion & perspectives

SEISCOPE

To overcome the positivity problem, we focus on the dual formulation of the W_1 distance (Santambrogio, 2015)

where Lip_1 is the space of 1-Lipschitz functions for the ground cost ||x - y||

$$\operatorname{Lip}_{1}\left\{\varphi: x \in X \longrightarrow \mathbb{R}, \ \forall (x_{i}, x_{j}) \in X \times X, \ |\varphi_{i} - \varphi_{j}| \leq ||x_{i} - x_{j}||\right\}$$
(7)

SEISCOPE

To design a fast solver, we focus on the ℓ_1 ground cost

 $\bullet\,$ The discretization of the problem is, for a 3D data cube indexed by

 $1 \leq n \leq N_t$ for time t_n

 $1 \leq i \leq N_x$, $1 \leq j \leq N_y$ for the receiver position x_i, y_j

$$\max_{\varphi_{nij}} \sum_{nij} \varphi_{nij} \left((d_{cal})_{nij} - (d_{obs})_{nij} \right), |\varphi_{nij} - \varphi_{n'i'j'}| < |t_n - t_{n'}| + |x_i - x_{i'}| + |y_j - y_{j'}|,$$
(8)

• The number of unknowns is $N = N_t \times N_x \times N_y$

$O(N^2)$ global linear constraints: too high complexity for efficient numerical algorithms



• Thanks to the sum of absolute values in problem (8), the global system is equivalent to

$$\max_{\varphi_{nij}} \sum_{nij} \varphi_{nij} \left((d_{cal})_{nij} - (d_{obs})_{nij} \right), \tag{9}$$

$$\begin{cases}
|\varphi_{n+1ij} - \varphi_{nij}| < |t_{n+1} - t_n|, \\
|\varphi_{ni+1j} - \varphi_{nij}| < |y_{i+1} - y_i|, \\
|\varphi_{nji+1} - \varphi_{nij}| < |z_{j+1} - z_j|,
\end{cases}$$

O(N) local linear constraints: far better complexity (Métivier et al., 2016c)

SEISCOPE

• The problem (9) is reformulated as the convex non-smooth problem

$$\max_{\varphi} f_1(\varphi) + f_2(A\varphi),$$

$$f_1(\varphi) = \sum_{nij} \varphi_{nij} \left((d_{cal})_{nij} - (d_{obs})_{nij} \right), \quad f_2(\varphi) = -i_K(\varphi),$$
(10)

• The matrix \boldsymbol{A} accounts for the constraints

$$A = \begin{bmatrix} D_x & D_y & D_z & \frac{1}{\lambda} I_N \end{bmatrix}^T \in \mathbb{M}_{P,N}(\mathbb{R})$$
(11)

where D_x, D_y, D_z are the forward finite-differences operators

$$(D_x\varphi)_{ijk} = \frac{\varphi_{i+1,j,k} - \varphi_{ijk}}{h_x}, \quad (D_y\varphi)_{ijk} = \frac{\varphi_{i,j+1,k} - \varphi_{ijk}}{h_y}, \quad (D_z\varphi)_{ijk} = \frac{\varphi_{i,j,k+1} - \varphi_{ijk}}{h_z}$$
(12)

while K is the unit hypercube and i_K the indicator function of K

$$i_K(v) = \begin{vmatrix} 0 & \text{if } v \in K \\ +\infty & \text{if } v \notin K, \end{vmatrix}$$
(13)



- A proximal splitting algorithm can be use to solve this problem: we choose the Simultaneous Direction Method of Multipliers (SDMM) for its good convergence properties (Combettes and Pesquet, 2011)
- The most computational demanding task of the algorithm is the resolution of a linear system involving a matrix which is equivalent to a second-order finite difference discrete of the Laplacian operator with homogeneous Neumann boundary condition
- Fast-solvers can be used to invert this matrix (Poisson's problem), either based on

Fast Fourier Transform (Swarztrauber, 1974): complexity $\mathsf{O}(N\log N)$

Multigrid algorithms (Brandt, 1977; Adams, 1989): complexity in O(N)

A complete description of this numerical strategy is given in (Métivier et al., 2016c) with 2D and 3D examples

SEISCOPE

- We assume the acoustic approximation
- For the L^2 distance, the gradient of the misfit function is

$$\nabla f(x) = \frac{-2}{v_P(x)^3} \int_0^T \partial_{tt} p(x,t) \lambda(x,t) dt$$
(14)

where

- $v_P(x)$: P-wave velocity
- p(x,t): pressure wavefield
- $\lambda(x,t)$: adjoint wavefield backpropagation of the L^2 residuals $d_{obs} d_{cal}$

$$\partial_{tt}\lambda - c^2\Delta\lambda = d_{obs} - d_{cal} \tag{15}$$

SEISCOPE

- We assume the acoustic approximation
- For the KR distance, the gradient of the misfit function is

$$\nabla f(x) = \frac{-2}{v_P(x)^3} \int_0^T \partial_{tt} p(x,t) \lambda(x,t) dt$$
(16)

where

- $v_P(x)$: P-wave velocity
- p(x,t): pressure wavefield
- $\lambda(x,t)$: adjoint wavefield backpropagation of the KR residuals $\overline{\varphi}$

$$\partial_{tt}\lambda - c^2 \Delta \lambda = \overline{\varphi} \tag{17}$$

$$\overline{\varphi_{nij}} = \arg \max_{\varphi_{nij}} \sum_{nij} \varphi_{nij} \left((d_{cal})_{nij} - (d_{obs})_{nij} \right), \quad (18)$$

$$\begin{cases}
|\varphi_{n+1ij} - \varphi_{nij}| < |t_{n+1} - t_n|, \\
|\varphi_{ni+1j} - \varphi_{nij}| < |y_{i+1} - y_i|, \\
|\varphi_{nji+1} - \varphi_{nij}| < |z_{j+1} - z_j|,
\end{cases}$$

For the optimal transport distance, the corresponding adjoint source is equal to the function $\overline{\varphi}$ which is the solution of the constrained maximization problem



• We come back to 1D time-shifted Ricker functions





• Misfit function shape



 L^2 distance function with respect to the time-shift. W_1 distance function with respect to the time shift

We recover a single minimum but lose the convexity





- surface acquisition with 128 sources each 125 m and 168 receivers each 100 m
- acoustic modeling engine to generate synthetic data
- high pass filter of data: no energy below 3 Hz





Inversion using the standard L^2 distance starting from the crude velocity model





Inversion using the W_1 distance starting from the crude velocity model





Inversion using the L^2 distance starting from a crude velocity model





Inversion using the W_1 distance starting from a crude velocity model



The Marmousi case study : example of adjoint-source









(d)1D optimal transport distance result

(e)2D optimal transport distance result



Computation overhead per gradient (FISHPACK and 50 SDMM iterations): 3.8 s (19%)

- L^2 gradient computation time **20,6 s**
- KR gradient computation time 24,4 s

Number of iterations

- L^2 inversion number of iterations: 83
- KR inversion number of iterations: 439



Reconstructing salt structures through a layer stripping approach Exact model



- surface acquisition with 128 sources each 125 m and 161 receivers each 100 m
- acoustic modeling engine to generate synthetic data
- high pass filter of data: no energy below 3 Hz





Reconstructing salt structures through a layer stripping approach Exact model





Reconstructing salt structures through a layer stripping approach Initial model





Reconstructing salt structures through a layer stripping approach Step 1 with L^2 distance





Reconstructing salt structures through a layer stripping approach

Step 1 with optimal transport distance





Reconstructing salt structures through a layer stripping approach Data comparison



OT for FWI



Reconstructing salt structures through a layer stripping approach Step 2





Reconstructing salt structures through a layer stripping approach Step 3





Reconstructing salt structures through a layer stripping approach Step 4





Reconstructing salt structures through a layer stripping approach Step 5





Reconstructing salt structures through a layer stripping approach Step 6





Motivation for using optimal transport for FWI

A first implementation: a distance based on the dual form of ${\it W}_1$

Towards a generalization of the approach

Conclusion & perspectives

• Addition of a **positive constant mass** to make the signed measures strictly positive

$$\widetilde{W}(d_{cal}, d_{obs}) := W(d_{cal} + \alpha, d_{obs} + \alpha)$$

with $\alpha > |\min(d_{cal}, d_{obs})|$



Drawbacks: the transformation becomes local: no more transportation along the time axis: we loose the convexity with respect to time shifts





- In the literature, we have found 3 ideas based on the separation of positive and negative parts of the signals (Ambrosio et al., 2011; Mainini, 2012)
 - 1. taking the absolute value of the signal
 - 2. transport separately positive and negative part of the signal (Engquist and Froese, 2014)
 - 3. recombine the data using the decomposition in positive and negative part to compare positive measures with mass conservation (Mainini, 2012)







Drawbacks: loss of polarity information: no sensitivity to impedance contrast

Separation of positive and negative parts(Engquist and Froese, 2014)ELSCOPE

$$W_{Engquist}(d_{cal}, d_{obs}) := W(d^+_{cal}, d^+_{obs}) + W(d^-_{cal}, d^-_{obs})$$

with $d_{cal} = d^+_{cal} - d^-_{cal}$ and $d_{obs} = d^+_{obs} - d^-_{obs}$



Drawbacks:

• we lose the mass conservation

$$\int p^+ \neq \int q^+, \quad \int p^- \neq \int q^- \tag{19}$$

• artificial decorrelation between positive and negative part of the signal

Mainin strategy (Mainini, 2012)



$$W_{Mainini}(d_{cal}, d_{obs}) := W(d_{cal}^+ + d_{obs}^-, d_{obs}^+ + d_{cal}^-)$$

with $d_{cal} = d^+_{cal} - d^-_{cal}$ and $d_{obs} = d^+_{obs} - d^-_{obs}$



Advantages:

• positivity and mass conservation

$$\int d_{cal}^{+} - d_{cal}^{-} = \int d_{obs}^{+} - d_{obs}^{-} \Leftrightarrow \int d_{cal}^{+} + d_{obs}^{-} = \int d_{obs}^{+} + d_{cal}^{-}$$
(20)

OT for FWI



• Important: the Mainini strategy for the W_1 distance is equivalent to the dual strategy we have employed so far

$$\begin{split} W_{Mainini}^{1}(d_{cal}, d_{obs}) &= W^{1}(d_{cal}^{+} + d_{obs}^{-}, d_{cal}^{-} + d_{obs}^{+}) \\ &= \max_{\varphi \in \text{Lip}_{1}} \sum_{i=1}^{N} \varphi_{i} \left(d_{cal,i}^{+} + d_{obs,i}^{-} - (d_{cal,i}^{-} + d_{obs,i}^{+}) \right) \\ &= \max_{\varphi \in \text{Lip}_{1}} \sum_{i=1}^{N} \varphi_{i} \left(\underbrace{d_{cal,i}^{+} - d_{cal,i}^{-}}_{d_{cal}} - \underbrace{(d_{obs}^{+} - d_{obs,i}^{-})}_{d_{obs}} \right) \\ &= \max_{\varphi \in \text{Lip}_{1}} \sum_{i=1}^{N} \varphi_{i} \left(d_{cal} - d_{obs} \right) \\ &= W^{1}(d_{cal}, d_{obs}) \end{split}$$



- We have obtained encouraging results with \boldsymbol{W}^1 and a ℓ_1 ground cost
- What about W^2 using the Mainini decomposition ? \rightarrow more convex / smoother w.r.t time shift?

To do so, we now need an efficient numerical strategy for general W^p problems Conventional approaches for large scale transport problems:

- Monge-Ampère equation (Philippis and Figalli, 2014)
- Benamou-Brenier formulation (Benamou and Brenier, 2000)
- Entropic regularization (Cuturi, 2013; Benamou et al., 2015)



- We have obtained encouraging results with $\boldsymbol{\mathit{W}}^1$ and a ℓ_1 ground cost
- What about W^2 using the Mainini decomposition ? \rightarrow more convex / smoother w.r.t time shift?

To do so, we now need an efficient numerical strategy for general W^p problems Conventional approaches for large scale transport problems:

- Monge-Ampère equation (Philippis and Figalli, 2014)
- Benamou-Brenier formulation (Benamou and Brenier, 2000)
- Entropic regularization (Cuturi, 2013; Benamou et al., 2015)



• Entropic regularized problem (Cuturi, 2013; Benamou et al., 2015)

$$W_{\varepsilon}^{p}(d_{cal}, d_{obs}) = \min_{\gamma \in \Pi(d_{cal}, d_{obs})} \sum_{i,j=1}^{N} C_{ij}\gamma_{ij} + \underbrace{\varepsilon}_{\text{regularization}} \underbrace{\sum_{i,j=1}^{N} \gamma_{ij}(\log(\gamma_{ij}) - 1)}_{\text{entropy}}$$
(21)

with:

$$\begin{array}{l} \bullet \ \Pi(d_{cal}, d_{obs}) := \left\{ \gamma \in \mathbb{R}_{\geq 0}^{N \times N} \ ; \ \sum_{i=1}^{N} \gamma_{ij} = d_{cal,i} \ , \ \sum_{j=1}^{N} \gamma_{ij} = d_{obs,j} \right\} \\ \bullet \ C_{ij} := \operatorname{dist}(x_i, y_j) \\ \bullet \ \varepsilon > 0 \end{array}$$

• if $\gamma_{ij} = 0$, the convention is that $0 \log(0) = 0$

 An alternate projection algorithm based on the Sinkhorn iteration has been proposed to solve this problem (Bregman, 1967; Sinkhorn and Knopp, 1967; Cuturi, 2013; Benamou et al., 2015)

$$\gamma_{\varepsilon}^{*} = \lim_{n \to +\infty} \gamma^{(n)} = \lim_{n \to +\infty} \operatorname{diag}\left(u^{(n)}\right) K \operatorname{diag}\left(v^{(n)}\right)$$
(22)

$$K_{ij} = e^{-\frac{C_{ij}}{\varepsilon}} \tag{23}$$

• Algorithm

•
$$u_i^{(n+1)} = \frac{d_{cal,i}}{(Kv^{(n)})_i}$$

• $v_j^{(n+1)} = \frac{d_{obs,j}}{(K^T u^{(n+1)})_j}$
• $u^{(n)} \in \mathbb{R}^N, v^{(n)} \in \mathbb{R}^N, K \in \mathbb{R}^{N \times N}$
• $u^{(0)} = v^{(0)} = (1, ..., 1)^T \in \mathbb{R}^N$





The computational complexity depends on the product: Kv

- naive implementation $Kv \to \mathcal{O}(N^2)$
- exploiting the symmetric Toeplitz structure of $K \to \mathcal{O}(N \log(N))$ (FFT acceleration, no approximation)
- exploiting the sparsity of $K \to \mathcal{O}(N)$ (approximation)



• We come back to 1D time-shifted Ricker functions now with the Mainini (2012) approach





Comparison between Mainini cost with W_1 and W_2 computed through the entropic regularization approach (Cuturi, 2013; Benamou et al., 2015)





Motivation for using optimal transport for FWI

A first implementation: a distance based on the dual form of ${\it W}_1$

Towards a generalization of the approach

Conclusion & perspectives



- Optimal transport appears as an interesting approach to mitigate cycle skipping in FWI
- The main difficulty is however to deal with non-positive meausres: the seismic data is oscillatory



- Optimal transport appears as an interesting approach to mitigate cycle skipping in FWI
- The main difficulty is however to deal with non-positive meausres: the seismic data is oscillatory
- A first implementation based on the dual form of the W₁ distance has been proposed in (Métivier et al., 2016a,b,c)
- In this presentation, we have shown that this is actually a special case of a more general approach proposed by Mainini (2012) to extend optimal transport to signed measures

- SEISCOPE
- Optimal transport appears as an interesting approach to mitigate cycle skipping in FWI
- The main difficulty is however to deal with non-positive meausres: the seismic data is oscillatory
- A first implementation based on the dual form of the W₁ distance has been proposed in (Métivier et al., 2016a,b,c)
- In this presentation, we have shown that this is actually a special case of a more general approach proposed by Mainini (2012) to extend optimal transport to signed measures
- This opens the way to work with more general $W_p, p > 1$ distance, that might be more convex and/or smoother than W_1
- As a preliminary step, we have focused on the entropic regularization approach (Cuturi, 2013; Benamou et al., 2015), an efficient strategy to solve general large scale optimal transport problems



Thank you for your attention

- IDRIS and TGCC, French national computing centers
- CIMENT, Grenoble computing center
- SEISCOPE sponsors : http://seiscope2.osug.fr

Questions?

- Adams, J. C. (1989). MUDPACK: Multigrid portable FORTRAN software for the efficient solution of linear elliptic partial differential equations. *Applied Mathematics and Computation*, 34(2):113–146.
- Ambrosio, L., Mainini, E., and Serfaty, S. (2011). Gradient flow of the Chapman–Rubinstein–Schatzman model for signed vortices. Annales de l'Institut Henri Poincaré (C) Non Linear Analysis, 28(2):217–246.
- Benamou, J.-D. and Brenier, Y. (2000). A computational fluid mechanics solution to the monge-kantorovich mass transfer problem. *Numerische Mathematik*.
- Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. (2015). Iterative Bregman Projections for Regularized Transportation Problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138.
- Brandt, A. (1977). Multi-level adaptive solutions to boundary-value problems. *Mathematics of Computation*, 31:333–390.
- Bregman, L. (1967). The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. USSR computational mathematics and mathematical physics.
- Combettes, P. L. and Pesquet, J.-C. (2011). Proximal splitting methods in signal processing. In Bauschke, H. H., Burachik, R. S., Combettes, P. L., Elser, V., Luke, D. R., and Wolkowicz, H., editors, *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, volume 49 of *Springer Optimization and Its Applications*, pages 185–212. Springer New York.
- Cuturi, M. (2013). Sinkhorn distances: lightspeed computation of optimal transportation distances. Advances in Neural Information Processing Systems.
- Engquist, B. and Froese, B. D. (2014). Application of the wasserstein metric to seismic signals. Communications in Mathematical Science, 12(5):979–988.

Kantorovich, L. (1942). On the transfer of masses. Dokl. Acad. Nauk. USSR, 37:7-8.

- Mainini, E. (2012). A description of transport cost for signed measures. Journal of Mathematical Sciences, 181(6):837–855.
- Métivier, L., Brossier, R., Mérigot, Q., Oudet, E., and Virieux, J. (2016a). Increasing the robustness and applicability of full waveform inversion: an optimal transport distance strategy. *The Leading Edge*, 35(12):1060–1067.
- Métivier, L., Brossier, R., Mérigot, Q., Oudet, E., and Virieux, J. (2016b). Measuring the misfit between seismograms using an optimal transport distance: Application to full waveform inversion. *Geophysical Journal International*, 205:345–377.
- Métivier, L., Brossier, R., Mérigot, Q., Oudet, E., and Virieux, J. (2016c). An optimal transport approach for seismic tomography: Application to 3D full waveform inversion. *Inverse Problems*, 32(11):115008.
- Philippis, G. D. and Figalli, A. (2014). The monge-ampère equation and its link to optimal transportation. BULLETIN (New Series) OF THE AMERICAN MATHEMATICAL SOCIETY.
- Santambrogio, F. (2015). Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling. Progress in Nonlinear Differential Equations and Their Applications. Springer International Publishing.
- Sinkhorn, R. and Knopp, P. (1967). Concerning nonnegative matrices and doubly stochastic matrices. Pacific Journal of Mathematics.
- Swarztrauber, P. N. (1974). A Direct Method for the Discrete Solution of Separable Elliptic Equations. SIAM Journal on Numerical Analysis, 11(6):1136–1150.