Adaptive construction of measure transport with application to Bayesian inference

D. Bigoni[†] (**dabi@mit.edu**), A. Spantini[†], Y.M. Marzouk[†] [†]Massachusetts Institute of Technology

Milan - 16/06/2016

Inference – The Bayesian approach – Seismic tomography



$$\implies |\nabla u(\mathbf{x})| = \frac{1}{v(\mathbf{x})} =$$

 $p(\theta|d) \propto \overbrace{p(d|\theta)}^{ ext{Likelihood}} p(\theta)$





"Knowing the arrival time and the hypocenter of a number of earthquakes, what is the probability distribution of the velocity model of the crust?"

The earthquake image is licensed under Public Domain via Wikimedia Commons. Work in collaboration with A. Zunino

Intractable distributions

1 We have access to densities but sampling is challenging.

• Classical Bayesian inference problems [El Moselhy et al., 2012]

 $p(\theta|d) \propto p(d|\theta)p(\theta)$

• Filtering problems [see Spantini's presentation MS.42]

$$\begin{aligned} X_{t+1} &= \mathcal{F}(t, X_t, w_t) \qquad w_t \sim \mu \\ Y_t &= \mathcal{G}(t, X_t, \gamma_t) \qquad \gamma_t \sim \pi \end{aligned}$$

We have access to samples and we want to characterize their density. [Parno et al., 2015]

•
$$\rho : \mathbb{R}^d \to \mathbb{R}$$
 density of the distribution λ_1

- $\pi: \mathbb{R}^d \to \mathbb{R}$ density of the distribution λ_2
- Let $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ be the map such that either

PF	$\pi = T_{\sharp}\rho = \rho \circ T^{-1} \nabla T^{-1} $
PB	$\rho = T^{\sharp} \pi = \pi \circ T \nabla T $

This means T is such that for $X \sim \lambda_1$, either

$$\begin{array}{ll} {\rm PF} & T(X) = Y \sim \lambda_2 \\ {\rm PB} & T^{-1}(Y) = X \sim \lambda_1 \end{array}$$

$T_{\sharp}\rho$	$T^{\sharp}\pi$

- $\rho: \mathbb{R}^d \rightarrow \mathbb{R}$ density of the distribution λ_1
- $\pi: \mathbb{R}^d \to \mathbb{R}$ density of the distribution λ_2
- Let $T: \mathbb{R}^d \to \mathbb{R}^d$ be the map such that either

$$\begin{array}{ll} \mathbf{PF} & \pi = T_{\sharp}\rho = \rho \circ T^{-1}|\nabla T^{-1}| \\ \mathbf{PB} & \rho = T^{\sharp}\pi = \pi \circ T|\nabla T| \\ \end{array}$$

This means T is such that for $X\sim\lambda_1$, either

- **PF** $T(X) = Y \sim \lambda_2$
- **PB** $T^{-1}(Y) = X \sim \lambda_1$



- $\rho: \mathbb{R}^d \rightarrow \mathbb{R}$ density of the distribution λ_1
- $\pi: \mathbb{R}^d \to \mathbb{R}$ density of the distribution λ_2
- Let $T: \mathbb{R}^d \to \mathbb{R}^d$ be the map such that either

PF	$\pi = T_{\sharp}\rho = \rho \circ T^{-1} \nabla T^{-1} $
PB	$\rho = T^{\sharp}\pi = \pi \circ T \nabla T $

This means T is such that for $X\sim\lambda_1$, either

 $\begin{array}{ll} {\bf PF} & T(X) = Y \sim \lambda_2 \\ {\bf PB} & T^{-1}(Y) = X \sim \lambda_1 \end{array}$



- $\rho: \mathbb{R}^d \rightarrow \mathbb{R}$ density of the distribution λ_1
- $\pi: \mathbb{R}^d \to \mathbb{R}$ density of the distribution λ_2
- Let $T: \mathbb{R}^d \to \mathbb{R}^d$ be the map such that either

$$\begin{array}{ll} \mathbf{PF} & \pi = T_{\sharp}\rho = \rho \circ T^{-1} |\nabla T^{-1}| \\ \mathbf{PB} & \rho = T^{\sharp}\pi = \pi \circ T |\nabla T| \\ \end{array}$$

This means T is such that for $X \sim \lambda_1$, either

 $\begin{array}{ll} \mathsf{PF} & T(X) = Y \sim \lambda_2 \\ \mathsf{PB} & T^{-1}(Y) = X \sim \lambda_1 \end{array}$

Note: T must preserve mass!

Knothe-Rosenblatt rearrangement

 $\forall \lambda_1,\lambda_2$ absolutely continuous there exists a triangular monotone map s.t. $T(d\lambda_1)=d\lambda_2$



$$T(\mathbf{x}) = \begin{bmatrix} T^{(1)}(x_1) \\ T^{(2)}(x_1, x_2) \\ \vdots \\ T^{(d)}(x_1, \dots, x_d) \end{bmatrix}$$

Practical setting

 $\begin{array}{l} \rho \text{ is the density of } X \sim \mathcal{N}(\mathbf{0},\mathbf{I}) \text{ or some amenable distribution [reference]} \\ \pi \text{ is the density of an intractable distribution [target]} \\ \text{We seek a map } T^* \text{ such that } \quad T^*_{\scriptscriptstyle \rm H}\rho\approx\pi. \end{array}$

Practical setting

 $\begin{array}{l} \rho \text{ is the density of } X \sim \mathcal{N}(\mathbf{0},\mathbf{I}) \text{ or some amenable distribution [reference]} \\ \pi \text{ is the density of an intractable distribution [target]} \\ \text{We seek a map } T^* \text{ such that } \quad T^*_{\text{t}}\rho \approx \pi. \end{array}$

Minimization statement [El Moselhy et al., 2012]

For \mathcal{H}^{Δ} the space of lower-triangular maps, solve

$$T^* = \underset{T \in \mathcal{H}^{\Delta}}{\operatorname{arg\,min}} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \underset{T \in \mathcal{H}^{\Delta}}{\operatorname{arg\,min}} D_{\mathrm{KL}}(\rho \| T^{\sharp}\pi)$$
$$= \underset{T \in \mathcal{H}^{\Delta}}{\operatorname{arg\,min}} \mathbb{E}_{\rho}[\log \rho] - \mathbb{E}_{\rho}\left[\log \pi \circ T + \log |\nabla T|\right]$$

such that $\partial_{x_i} T^{(i)} > 0$.

Practical setting

 $\begin{array}{l} \rho \text{ is the density of } X \sim \mathcal{N}(\mathbf{0},\mathbf{I}) \text{ or some amenable distribution [reference]} \\ \pi \text{ is the density of an intractable distribution [target]} \\ \text{We seek a map } T^* \text{ such that } \quad T^*_{\text{t}}\rho \approx \pi. \end{array}$

Minimization statement [El Moselhy et al., 2012]

For \mathcal{H}^{Δ} the space of lower-triangular maps, solve

$$T^* = \underset{T \in \mathcal{H}^{\Delta}}{\arg\min} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \underset{T \in \mathcal{H}^{\Delta}}{\arg\min} D_{\mathrm{KL}}(\rho \| T^{\sharp}\pi)$$
$$= \underset{T \in \mathcal{H}^{\Delta}}{\arg\min} \mathbb{E}_{\rho}[\log \rho] - \mathbb{E}_{\rho}\left[\log \pi \circ T + \log |\nabla T|\right]$$

such that $\partial_{x_i} T^{(i)} > 0$.

(1) We can use derivative based optimization if $\nabla_{\mathbf{x}} \log \pi$ or $\nabla_{\mathbf{x}}^2 \log \pi$ are provided

- **2** We can explore π in parallel
- **3** We can generate i.i.d. samples from $T^*_{\sharp} \rho \approx \pi$ in parallel
- **4** We can estimate convergence!

Practical setting

 ρ is the density of $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ or some amenable distribution [reference] π is the density of an intractable distribution [target] We seek a map T^* such that $T^*_{\sharp} \rho \approx \pi$.

Minimization statement [El Moselhy et al., 2012]

For \mathcal{H}^{Δ} the space of lower-triangular maps, solve

$$T^* = \underset{T \in \mathcal{H}^{\Delta}}{\operatorname{arg\,min}} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \underset{T \in \mathcal{H}^{\Delta}}{\operatorname{arg\,min}} D_{\mathrm{KL}}(\rho \| T^{\sharp}\pi)$$
$$= \underset{T \in \mathcal{H}^{\Delta}}{\operatorname{arg\,min}} \mathbb{E}_{\rho}[\log \rho] - \mathbb{E}_{\rho}\left[\log \pi \circ T + \log |\nabla T|\right]$$

such that $\partial_{x_i} T^{(i)} > 0$. Note: Approximation of *d* functions up to *d*-dimensional!

Source of low-dimensional structure Adaptivity Conditional independency Marginal independency I ow-rank structure Adaptive measure transport 16.9.2016

1 Convergence criteria

2 Refinement criteria

1 Convergence criteria

At optimality $\Rightarrow D_{\mathrm{KL}}(T^*_{\sharp}\rho \| \pi) = 0$ $\mathbb{V}\left[\log \frac{\rho}{T^{*\sharp}\pi}\right] \rightarrow \frac{1}{2}D_{\mathrm{KL}}(T^*_{\sharp}\rho \| \pi) \rightarrow 0$

2 Refinement criteria

- **1** Convergence criteria Variance diagnostic $\mathbb{V}\left[\log \frac{\rho}{T^{\sharp}\pi}\right]$
- **2** Refinement criteria

• Minimize
$$\mathcal{J}^{\Delta}(T) = D_{KL}(T_{\sharp}\rho \| \pi)$$
,
over \mathcal{H}^{Δ} (dim $\mathcal{H}^{\Delta} = \infty$)



- Minimize $\mathcal{J}^{\Delta}(T) = D_{KL}(T_{\sharp}\rho \| \pi)$, over \mathcal{H}^{Δ} (dim $\mathcal{H}^{\Delta} = \infty$)
- For $\mathcal{H}_0^{\Delta} \subset \mathcal{H}^{\Delta}$ (dim $\mathcal{H}_0^{\Delta} = n_0$) $T_0^* = \arg \min_{T \in \mathcal{H}_0^{\Delta}} \mathcal{J}^{\Delta}(T)$



- Minimize $\mathcal{J}^{\Delta}(T) = D_{KL}(T_{\sharp}\rho \| \pi)$, over \mathcal{H}^{Δ} (dim $\mathcal{H}^{\Delta} = \infty$)
- For $\mathcal{H}_0^{\Delta} \subset \mathcal{H}^{\Delta}$ (dim $\mathcal{H}_0^{\Delta} = n_0$) $T_0^* = \arg \min_{T \in \mathcal{H}_0^{\Delta}} \mathcal{J}^{\Delta}(T)$
- Equivalent to: $\mathbf{a}_0^* = \arg\min_{\mathbf{a} \in \mathbb{R}^{n_0}} \mathcal{J}^{\Delta}(T[\mathbf{a}])$



- Minimize $\mathcal{J}^{\Delta}(T) = D_{KL}(T_{\sharp}\rho \| \pi)$, over \mathcal{H}^{Δ} (dim $\mathcal{H}^{\Delta} = \infty$)
- For $\mathcal{H}_0^{\Delta} \subset \mathcal{H}^{\Delta}$ (dim $\mathcal{H}_0^{\Delta} = n_0$) $T_0^* = \arg \min_{T \in \mathcal{H}_0^{\Delta}} \mathcal{J}^{\Delta}(T)$
- Equivalent to: $\mathbf{a}_0^* = \arg\min_{\mathbf{a} \in \mathbb{R}^{n_0}} \mathcal{J}^{\Delta}(T[\mathbf{a}])$
- At optimality: $\nabla_{\mathbf{a}} \mathcal{J}^{\Delta}(T[\mathbf{a}_0^*]) = 0$



- Minimize $\mathcal{J}^{\Delta}(T) = D_{KL}(T_{\sharp}\rho \| \pi)$, over \mathcal{H}^{Δ} (dim $\mathcal{H}^{\Delta} = \infty$)
- For $\mathcal{H}_0^{\Delta} \subset \mathcal{H}^{\Delta}$ (dim $\mathcal{H}_0^{\Delta} = n_0$) $T_0^* = \arg \min_{T \in \mathcal{H}_0^{\Delta}} \mathcal{J}^{\Delta}(T)$
- Equivalent to: $\mathbf{a}_0^* = \arg\min_{\mathbf{a} \in \mathbb{R}^{n_0}} \mathcal{J}^{\Delta}(T[\mathbf{a}])$
- At optimality: $\nabla_{\mathbf{a}} \mathcal{J}^{\Delta}(T[\mathbf{a}_0^*]) = 0$

 $\mathcal{J}^{\triangle}(T_{i+1}) < \mathcal{J}^{\triangle}(T_i)$



• The first variation $\nabla \mathcal{J}^{\Delta}(T[\mathbf{a}_0^*]) \neq 0$, which means There exists $\varepsilon > 0$ such that

$$\mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}] - \varepsilon \nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right)\right) < \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right)$$



• The first variation $\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right): \mathbb{R}^{d} \to \mathbb{R}^{d}$ is a map!

$$\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right) = (\nabla_{\mathbf{x}}T)^{-1}\left(\nabla_{\mathbf{x}}\log\frac{\rho}{T[\mathbf{a}_{0}^{*}]^{\sharp}\pi}\right)$$

• The first variation $\nabla \mathcal{J}^{\Delta}(T[\mathbf{a}_0^*]): \mathbb{R}^d \to \mathbb{R}^d$ is a map!

$$\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right) = (\nabla_{\mathbf{x}}T)^{-1}\left(\nabla_{\mathbf{x}}\log\frac{\rho}{T[\mathbf{a}_{0}^{*}]^{\sharp}\pi}\right)$$

• No new evaluation of $\nabla_{\mathbf{x}} \log \pi$ is required at the sample points $\{\mathbf{x}_k\}$

• The first variation $\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right): \mathbb{R}^{d} \to \mathbb{R}^{d}$ is a map!

$$\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right) = (\nabla_{\mathbf{x}}T)^{-1}\left(\nabla_{\mathbf{x}}\log\frac{\rho}{T[\mathbf{a}_{0}^{*}]^{\sharp}\pi}\right)$$

• No new evaluation of $\nabla_{\mathbf{x}} \log \pi$ is required at the sample points $\{\mathbf{x}_k\}$

• For
$$\mathcal{H}^{\Delta} \supset \mathcal{H}_{1}^{\Delta} \supset \mathcal{H}_{0}^{\Delta}$$
 (dim $\mathcal{H}_{1}^{\Delta} = n_{1}$), solve

$$\mathbf{b}_{1}^{*} = \underset{\mathbf{b} \in \mathbb{R}^{n_{1}}}{\operatorname{arg\,min}} \left\| U[\mathbf{b}] - \nabla \mathcal{J}^{\Delta} \left(T[\mathbf{a}_{0}^{*}] \right) \right\|_{L^{2}_{\rho}}$$

• The first variation $\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right): \mathbb{R}^{d} \to \mathbb{R}^{d}$ is a map!

$$\nabla \mathcal{J}^{\Delta}\left(T[\mathbf{a}_{0}^{*}]\right) = (\nabla_{\mathbf{x}}T)^{-1}\left(\nabla_{\mathbf{x}}\log\frac{\rho}{T[\mathbf{a}_{0}^{*}]^{\sharp}\pi}\right)$$

• No new evaluation of $abla_{\mathbf{x}} \log \pi$ is required at the sample points $\{\mathbf{x}_k\}$

• For
$$\mathcal{H}^{\Delta} \supset \mathcal{H}_{1}^{\Delta} \supset \mathcal{H}_{0}^{\Delta}$$
 (dim $\mathcal{H}_{1}^{\Delta} = n_{1}$), solve

$$\mathbf{b}_{1}^{*} = \operatorname*{arg\,min}_{\mathbf{b}\in\mathbb{R}^{n_{1}}} \left\| U[\mathbf{b}] - \nabla \mathcal{J}^{\Delta} \left(T[\mathbf{a}_{0}^{*}] \right) \right\|_{L^{2}_{\rho}}$$

 $U[\mathbf{b}_1^*]$ informs about:

- active variables to be included in the components $T^{(i)}$
- active coefficients to be included in the parametrization of $T^{(i)}$

Convergence criteria – Variance diagnostic – V [log ^ρ/_{T[‡]π}]
Refinement criteria – First variation – ∇J^Δ (T[a*])



$$\mathcal{J}^{\bigtriangleup}:\mathcal{H}\to\mathbb{R}$$



Figure: Target π

Iteration 1 – Pushforward $T_{\sharp}\rho$



Iteration 2 – Pushforward $T_{\sharp}\rho$



Iteration 3 – Pushforward $T_{\sharp}\rho$





Iteration 1 – Pullback $T^{\sharp}\pi$



Iteration 2 – Pullback $T^{\sharp}\pi$



Iteration 3 – Pullback $T^{\sharp}\pi$



- Latent log-volatilities modeled with an AR(1) process for t = 1, ..., N (N = 30) $X_{t+1} = \mu + \phi(X_t - \mu) + \eta_t$, $\eta_t \sim \mathcal{N}(0, 1)$, $X_1 \sim \mathcal{N}(0, 1/(1 - \phi^2))$
- Observe the mean return for holding the asset at time t

$$Y_t = \varepsilon_t \exp(X_t/2)$$
, $\varepsilon_t \sim \mathcal{N}(0, 1)$

• We want to characterize $\pi \sim \mu, \phi, \mathbf{X}_{1:N} | \mathbf{Y}_{1:N}$



Iteration 1 – Pullback $T^{\sharp}\pi$





Iteration 2 – Pullback $T^{\sharp}\pi$

















































Iteration 9 – Pullback $T^{\sharp}\pi$



Iteration 10 – Pullback $T^{\sharp}\pi$





Iteration 11 – Pullback $T^{\sharp}\pi$





Iteration 12 – Pullback $T^{\sharp}\pi$





Iteration 13 – Pullback $T^{\sharp}\pi$







Conclusions Transport maps potential advantages

- Turns a Bayesian inference problem into an optimization problem
- We can use derivative based optimization if $abla_{\mathbf{x}}\log\pi$ or $abla_{\mathbf{x}}^2\log\pi$ are provided
- We can explore π in parallel
- We can generate i.i.d. samples from $T_{\sharp}^{*}\rho\approx\pi$ in parallel
- We can estimate convergence!

Key contributions

- First variation informed adaptivity
 - \bullet No need for additional evaluation of π
 - Exploits marginal independence present in the problem
- Numerical robust algorithm to perform adaptivity





Parametrization of transport maps

Requirements:

• Lower-triangular structure

$$T : \mathbb{R}^{d} \to \mathbb{R}^{d}$$
$$\mathbf{x} \mapsto \begin{bmatrix} T^{(1)}(x_{1}) \\ T^{(2)}(x_{1}, x_{2}) \\ \vdots \\ T^{(d)}(x_{1}, \dots, x_{d}) \end{bmatrix}$$

Note: $\log |\nabla T| = \sum_{i=1}^{d} \log \partial_{x_i} T^{(i)}$

Parametrization of transport maps

Requirements:

- Lower-triangular structure
- Monotonicity (mass preservation), i.e. $|\nabla T| = \prod_{i=1}^{d} \partial_{x_i} T^{(i)} > 0$

$$T^{(i)}(x_1, \dots, x_i) = c_i(x_1, \dots, x_{i-1}) + \int_0^{x_i} \exp(h_i(x_1, \dots, x_{i-1}, t) dt)$$
$$c_i(x_1, \dots, x_{i-1}) = \sum_{j=1}^N a_j^{(i)} \Phi_j(x_1, \dots, x_{i-1})$$
$$h_i(x_1, \dots, x_{i-1}, t) = \sum_{j=1}^N b_j^{(i)} \Psi_j(x_1, \dots, x_{i-1}, t)$$

($\{\Phi_j\}$ Hermite polynomials, $\{\Psi_j\}$ Hermite functions **(** $\{\Phi_j\}$ and $\{\Psi_j\}$ radial basis functions