Data assimilation via low-rank couplings

D. Bigoni (**dabi@mit.edu**), R. Baptista, A. Spantini, Y. Marzouk Massachusetts Institute of Technology

> ICIAM Valencia, Spain – 16/07/2019



$$\pi \left(\Theta, \mathbf{Z}_{\Lambda} | \mathbf{y}_{\Xi}\right) \propto \mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda}\right) \pi \left(\Theta, \mathbf{Z}_{\Lambda}\right)$$
$$\mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda}\right) = \prod_{k \in \Xi} \mathcal{L} \left(\mathbf{y}_{k} | \Theta, \mathbf{Z}_{k}\right)$$
$$\pi \left(\Theta, \mathbf{Z}_{\Lambda}\right) = \pi \left(\Theta\right) \pi \left(\mathbf{Z}_{0} | \Theta\right) \prod_{k \in \Lambda} \pi \left(\mathbf{Z}_{k} | \mathbf{Z}_{k-1}, \Theta\right)$$



$$\pi \left(\Theta, \mathbf{Z}_{\Lambda} | \mathbf{y}_{\Xi}\right) \propto \mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda}\right) \pi \left(\Theta, \mathbf{Z}_{\Lambda}\right)$$
$$\mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda}\right) = \prod_{k \in \Xi} \mathcal{L} \left(\mathbf{y}_{k} | \Theta, \mathbf{Z}_{k}\right)$$
$$\pi \left(\Theta, \mathbf{Z}_{\Lambda}\right) = \pi \left(\Theta\right) \pi \left(\mathbf{Z}_{0} | \Theta\right) \prod_{k \in \Lambda} \pi \left(\mathbf{Z}_{k} | \mathbf{Z}_{k-1}, \Theta\right)$$



$$\pi \left(\mathbf{Z}_{k} | \mathbf{y}_{\Xi} \right) \propto \int \mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda} \right) \pi \left(\Theta, \mathbf{Z}_{\Lambda} \right) d\Theta d\mathbf{Z}_{j \neq k}$$
$$\mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda} \right) = \prod_{k \in \Xi} \mathcal{L} \left(\mathbf{y}_{k} | \Theta, \mathbf{Z}_{k} \right)$$
$$\pi \left(\Theta, \mathbf{Z}_{\Lambda} \right) = \pi \left(\Theta \right) \pi \left(\mathbf{Z}_{0} | \Theta \right) \prod_{k \in \Lambda} \pi \left(\mathbf{Z}_{k} | \mathbf{Z}_{k-1}, \Theta \right)$$



$$\pi \left(\mathbf{Z}_{k} | \mathbf{y}_{j \leq k} \right) \propto \int \mathcal{L} \left(\mathbf{y}_{j \leq k} | \Theta, \mathbf{Z}_{j \leq k} \right) \pi \left(\Theta, \mathbf{Z}_{j \leq k} \right) d\Theta \, d\mathbf{Z}_{j < k}$$
$$\mathcal{L} \left(\mathbf{y}_{\Xi} | \Theta, \mathbf{Z}_{\Lambda} \right) = \prod_{k \in \Xi} \mathcal{L} \left(\mathbf{y}_{k} | \Theta, \mathbf{Z}_{k} \right)$$
$$\pi \left(\Theta, \mathbf{Z}_{\Lambda} \right) = \pi \left(\Theta \right) \pi \left(\mathbf{Z}_{0} | \Theta \right) \prod_{k \in \Lambda} \pi \left(\mathbf{Z}_{k} | \mathbf{Z}_{k-1}, \Theta \right)$$

Outline

Transport maps

Data assimilation via low-dimensional couplings

Low-rank structure in lag-1 updates

Results

Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution $\boldsymbol{\nu}_{\rho}$ with density $\rho: \mathbb{R}^d \to \mathbb{R}_{\geq 0}$
- Distribution $\boldsymbol{
 u}_{\pi}$ with density $\pi:\mathbb{R}^d o \mathbb{R}_{\geq 0}$
- For $T: \mathbb{R}^d \to \mathbb{R}^d$ we define
 - $\mathbf{PF} \qquad T_{\sharp}\rho = \rho \circ T^{-1} |\nabla T^{-1}|$
 - **PB** $T^{\sharp}\pi = \pi \circ T |\nabla T|$
- We want T such that
 - **PF** $T_{\sharp}\rho = \pi$ **PB** $T^{\sharp}\pi = \rho$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution $u_{
 ho}$ with density $ho: \mathbb{R}^d o \mathbb{R}_{\geq 0}$
- Distribution u_{π} with density $\pi: \mathbb{R}^d \to \mathbb{R}_{\geq 0}$
- \bullet For $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define
 - $\mathbf{PF} \qquad T_{\sharp}\rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^{\sharp}\pi = \pi \circ T |\nabla T|$

• We want T such that

PF For $X \sim \boldsymbol{\nu}_{\rho}$, $T(X) \sim \boldsymbol{\nu}_{\pi}$ **PB** For $Y \sim \boldsymbol{\nu}_{\pi}$, $T^{-1}(Y) \sim \boldsymbol{\nu}_{\rho}$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution $\boldsymbol{\nu}_{\rho}$ with density $\rho: \mathbb{R}^d \to \mathbb{R}_{\geq 0}$
- Distribution $\boldsymbol{\nu}_{\pi}$ with density $\pi: \mathbb{R}^d \to \mathbb{R}_{>0}$
- For $T: \mathbb{R}^d \to \mathbb{R}^d$ we define
 - **PF** $T_{\sharp}\rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^{\sharp}\pi = \pi \circ T |\nabla T|$

• We want T such that

PF For $X \sim \boldsymbol{\nu}_o$, $T(X) \sim \boldsymbol{\nu}_{\pi}$ F

PB For
$$Y \sim \boldsymbol{\nu}_{\pi}$$
, $T^{-1}(Y) \sim \boldsymbol{\nu}_{\rho}$

Knothe-Rosenblatt rearrangement

 $\forall \nu_{\rho}, \nu_{\pi}$ Lebesgue absolutely continuous \exists a triangular monotone map s.t. $T_{\sharp}\rho = \pi$



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

Triangular monotone maps

$$\mathcal{T}_{>} = \left\{ T : \mathbb{R}^{d} \to \mathbb{R}^{d} : \overbrace{[T(\mathbf{x})]_{k} = T^{(k)}(x_{1}, \dots, x_{k})}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_{k}} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Triangular monotone maps

$$\mathcal{T}_{>} = \left\{ T : \mathbb{R}^{d} \to \mathbb{R}^{d} : \overbrace{[T(\mathbf{x})]_{k} = T^{(k)}(x_{1}, \dots, x_{k})}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_{k}} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Integrated squared representation – $\varepsilon > 0$

$$T^{(k)}(x_{1:k}) = c_k(x_{1:k-1}) + \int_0^{x_k} \left(h_k(x_{1:k-1}, t)\right)^2 + \varepsilon \, dt$$

Triangular monotone maps



Knothe-Rosenblatt rearrangement

 $\forall \ \boldsymbol{\nu}_{\rho}, \boldsymbol{\nu}_{\pi}$ Lebesgue absolutely continuous \exists a triangular monotone map s.t. $T_{\sharp}\rho = \pi$

> How to find the map $T \in \mathcal{T}_{>}$ such that $T_{\sharp}\rho = \pi$?



Minimize KL-divergence to find optimal map [EIMoselhy et al. 2012]

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \pi} \right]$$

Minimize KL-divergence to find optimal map [EIMoselhy et al. 2012]

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \pi} \right]$$

In particular we will:

- Discretize $\mathbb{E}_{\rho}[f(\mathbf{x})] \approx \sum_{i=1}^{m} f(\mathbf{x}_i) \, \mathbf{w}_i$
- Work on the n-dimensional space $\mathcal{T}^n_>$ of maps $T[\mathbf{a}]$ with parameters $\mathbf{a} \in \mathbb{R}^n$

Minimize KL-divergence to find optimal map [ElMoselhy et al. 2012]

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \pi} \right]$$

In particular we will:

- Discretize $\mathbb{E}_{\rho}[f(\mathbf{x})] \approx \sum_{i=1}^{m} f(\mathbf{x}_i) \, \mathbf{w}_i$
- Work on the *n*-dimensional space $\mathcal{T}^n_>$ of maps $T[\mathbf{a}]$ with parameters $\mathbf{a} \in \mathbb{R}^n$

$$\widetilde{\mathbf{a}} = \operatorname*{arg\,min}_{\mathbf{a} \in \mathbb{R}^n} \sum_{i=1}^m \left[-\log T[\mathbf{a}]^{\sharp} \pi(\mathbf{x}_i) \right] \, \mathbf{w}_i$$

Minimize KL-divergence to find optimal map [EIMoselhy et al. 2012]

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \pi} \right]$$

+ Gradient-based unconstrained optimization if gradients are available

+ We can explore π in parallel

+ We can generate i.i.d. samples from $\hat{T}_{\sharp} \nu_{\rho} = \nu_{\pi}$ in parallel

Minimize KL-divergence to find optimal map [EIMoselhy et al. 2012]

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \pi} \right]$$

+ Gradient-based unconstrained optimization if gradients are available

+ We can explore π in parallel

+ We can generate i.i.d. samples from $\widehat{T}_{\sharp} \nu_{\rho} = \nu_{\pi}$ in parallel

We are working on $\mathcal{T}_{>}^{n} \subset \mathcal{T}_{>}$, so how can we evaluate the quality of the approximation?

Convergence criterion – Variance diagnostic

$$\widehat{T} = \underset{T \in \mathcal{T}_{>}}{\arg\min} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \underset{T \in \mathcal{T}_{>}}{\arg\min} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \widetilde{\pi}} \right] + \log \int \widetilde{\pi}$$

Optimal
$$\widehat{T} \in \mathcal{T}_{>}$$
 and $\int \widetilde{\pi} = 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(\widehat{T})^{\sharp} \widetilde{\pi}} \right] = 0$
But, optimal $\widetilde{T} \in \mathcal{T}_{>}^{n}$ or $\int \widetilde{\pi} \neq 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(\widetilde{T})^{\sharp} \widetilde{\pi}} \right] \neq 0$

Convergence criterion – Variance diagnostic

$$\widehat{T} = \underset{T \in \mathcal{T}_{>}}{\arg\min} D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \underset{T \in \mathcal{T}_{>}}{\arg\min} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \widetilde{\pi}} \right] + \log \int \widetilde{\pi}$$

Optimal
$$\widehat{T} \in \mathcal{T}_{>}$$
 and $\int \widetilde{\pi} = 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(\widehat{T})^{\sharp} \widetilde{\pi}} \right] = 0$

But, optimal
$$\widetilde{T} \in \mathcal{T}_{>}^{n}$$
 or $\int \widetilde{\pi} \neq 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{\left(\widetilde{T}\right)^{\sharp} \widetilde{\pi}} \right] \neq 0$

$$D_{\mathrm{KL}}(T_{\sharp} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) \approx \frac{1}{2} \mathbb{V} \left[\log \frac{\rho}{T^{\sharp} \tilde{\pi}} \right] \quad \text{as} \quad T \rightarrow \widehat{T}$$

Pros & cons

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp}\pi} \right]$$

+ Gradient-based unconstrained optimization if gradients are available + We can explore π in parallel + We can generate i.i.d. samples from $\widehat{T}_{\sharp}\nu_{\rho} = \nu_{\pi}$ in parallel

+ We can assess convergence!

Pros & cons

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp}\pi} \right]$$

+ Gradient-based unconstrained optimization if gradients are available + We can explore π in parallel + We can generate i.i.d. samples from $\hat{T}_{\sharp}\nu_{\rho} = \nu_{\pi}$ in parallel + We can assess convergence!

> $\widetilde{T}_{\sharp}\rho$ is a **biased** approximation of π . How to reduce such bias?

Use \hat{T} as a preconditioner for ... Markov Chain Monte Carlo

- **1** Generate the Markov chain $\{\mathbf{x}_i\}$ with invariant distribution $\widetilde{T}^{\sharp}\pi$ (use your favorite MCMC method/proposals)
- **2** The Markov chain $\{\widetilde{T}(\mathbf{x}_i)\}$ has invariant distribution π

Use \hat{T} as a preconditioner for ... Markov Chain Monte Carlo

- **1** Generate the Markov chain $\{\mathbf{x}_i\}$ with invariant distribution $\widetilde{T}^{\sharp}\pi$ (use your favorite MCMC method/proposals)
- **2** The Markov chain $\{\widetilde{T}(\mathbf{x}_i)\}$ has invariant distribution π

If
$$\widetilde{T}^{\sharp}\widetilde{\pi} \overset{\propto}{\sim} \rho$$
 then ρ is a **good proposal** for $\widetilde{T}^{\sharp}\widetilde{\pi}$.

Pros & cons

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp}\pi} \right]$$

+ Gradient-based unconstrained optimization if gradients are available

- + We can explore π in parallel
- + We can generate i.i.d. samples from $\hat{T}_{\sharp} \nu_{\rho} = \nu_{\pi}$ in parallel
- + We can assess convergence!
- + The map can be used as a **preconditioner**

Pros & cons

$$\widehat{T} = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} D_{\mathrm{KL}}(T_{\sharp}\rho \| \pi) = \operatorname*{arg\,min}_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp}\pi} \right]$$

+ Gradient-based unconstrained optimization if gradients are available

- + We can explore π in parallel
- + We can generate i.i.d. samples from $\hat{T}_{\sharp} \nu_{\rho} = \nu_{\pi}$ in parallel
- + We can assess convergence!
- + The map can be used as a **preconditioner**
- We need to approximate d functions of up to d variables!

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

Exploit source of low-dimensional structure

1 Smoothness and marginal independence

- Ongoing work...
- 2 Conditional independence
 - Variational filtering/smoothing and parameter estimation

[Spantini, B, Marzouk 2018; Houssineau, Jasra, Singh 2018]

• Ensamble filtering and smoothing

[Spantini, Baptista, Marzouk 2019]

3 Multilevel/multifidelity structure

[Parno, Moselhy, Marzouk 2018; Peherstorfer, Marzouk 2019]

4 Low-rank structure

[B, Zahm, Spantini, Marzouk 2019]

Data assimilation via low-dimensional couplings

How to remove conditional dependencies?





Assume at step $k\in\Lambda$ the map $\mathfrak{M}^1_{k-1}(\mathbf{z})$ is available and

$$\begin{split} (\mathfrak{M}_{k-1}^{1})_{\sharp}\rho(\mathbf{z}_{k}) &= \pi\left(\mathbf{z}_{k}|\mathbf{y}_{j\in\Xi,\;j\leq k}\right) \;.\\ \text{Let } \mathfrak{M}_{k}(\mathbf{z}_{k},\mathbf{z}_{k+1}) &= \left[\begin{array}{c} \mathfrak{M}_{k}^{0}(\mathbf{z}_{k},\mathbf{z}_{k+1})\\ \mathfrak{M}_{k}^{1}(\mathbf{z}_{k+1}) \end{array}\right] \text{ be such that}\\ (\mathfrak{M}_{k})_{\sharp}\rho(\mathbf{z}_{k},\mathbf{z}_{k+1}) &= \left[\begin{array}{c} \mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\\ \mathbf{z}_{k+1} \end{array}\right]^{\sharp}\pi(\mathbf{z}_{k},\mathbf{z}_{k+1}|\mathbf{y}_{j\leq k+1}) \end{split}$$

Then the following holds true [Spantini, B, Marzouk 2018]:

$$\begin{split} & \text{Filtering:} \ (\mathfrak{M}_{k}^{1})_{\sharp}\rho(\mathbf{z}_{k+1}) = \pi \left(\mathbf{z}_{k+1} | \mathbf{y}_{j \in \Xi, \ j \leq k+1} \right) \\ & \text{Full solution:} \ (\mathfrak{T}_{k})_{\sharp}\rho(\mathbf{z}_{0:k+1}) = \pi \left(\mathbf{z}_{0:k+1} | \mathbf{y}_{\Xi \leq k+1} \right) \text{ where,} \end{split}$$

$$\mathfrak{T}_k := \left[egin{array}{c} \mathfrak{M}_0(\mathbf{z}_0, \mathbf{z}_1) \\ \mathbf{z}_2 \\ \mathbf{z}_3 \\ dots \\ \mathbf{z}_{k+1} \end{array}
ight] \circ \left[egin{array}{c} \mathbf{z}_0 \\ \mathfrak{M}_1(\mathbf{z}_1, \mathbf{z}_2) \\ \mathbf{z}_3 \\ dots \\ \mathbf{z}_k
ight] \circ \cdots \circ \left[egin{array}{c} \mathbf{z}_0 \\ dots \\ \mathbf{z}_k
ight] \\ \mathbf{z}_{k-1} \\ \mathfrak{M}_k(\mathbf{z}_k, \mathbf{z}_{k+1}) \end{array}
ight]$$

Find the map

$$\mathfrak{M}_0(oldsymbol{ heta},\mathbf{z}_0,\mathbf{z}_1) = \left[egin{array}{c} \mathfrak{M}_0^0(oldsymbol{ heta}) \ \mathfrak{M}_0^0(oldsymbol{ heta},\mathbf{z}_0,\mathbf{z}_1) \ \mathfrak{M}_0^1(oldsymbol{ heta},\mathbf{z}_1) \end{array}
ight]$$

pushing forward $\mathcal{N}(0, \mathbf{I})$ to the first Markov component of $\pi(\Theta, \mathbf{Z}_{\Lambda} | \mathbf{y}_{\Xi})$:

 $\pi^{0}(\Theta, \mathbf{Z}_{0}, \mathbf{Z}_{1}) := \mathcal{L}(\mathbf{y}_{0}|\Theta, \mathbf{Z}_{0})\mathcal{L}(\mathbf{y}_{1}|\Theta, \mathbf{Z}_{1})\pi(\mathbf{Z}_{1}|\Theta, \mathbf{Z}_{0})\pi(\mathbf{Z}_{0}|\Theta)\pi(\Theta)$







Find the map

$$\mathfrak{M}_{1}(oldsymbol{ heta},\mathbf{z}_{1},\mathbf{z}_{2}) = \left[egin{array}{c} \mathfrak{M}_{1}^{0}(oldsymbol{ heta}) \ \mathfrak{M}_{1}^{0}(oldsymbol{ heta},\mathbf{z}_{1},\mathbf{z}_{2}) \ \mathfrak{M}_{1}^{1}(oldsymbol{ heta},\mathbf{z}_{2}) \end{array}
ight]$$

pushing $\mathcal{N}(0,\mathbf{I})$ to the second Markov component of $T_0^{\sharp}\pi$:

$$\begin{split} \pi^1 \left(\Theta, \mathbf{Z}_1, \mathbf{Z}_2 \right) &:= \rho(\Theta, \mathbf{Z}_1) \, \pi \left(\mathbf{Z}_2 \left| \Theta', \mathbf{Z}_1' \right. \right) \quad \text{(missing observation)} \\ \Theta' &:= \mathfrak{M}_0^{\Theta}(\Theta) \,, \quad \mathbf{Z}_1' := \mathfrak{M}_0^1(\Theta, \mathbf{Z}_1) \end{split}$$





Find the map

$$\mathfrak{M}_{2}(\boldsymbol{\theta},\mathbf{z}_{2},\mathbf{z}_{3}) = \left[\begin{array}{c} \mathfrak{M}_{2}^{\Theta}(\boldsymbol{\theta}) \\ \mathfrak{M}_{2}^{0}(\boldsymbol{\theta},\mathbf{z}_{2},\mathbf{z}_{3}) \\ \mathfrak{M}_{2}^{1}(\boldsymbol{\theta},\mathbf{z}_{3}) \end{array} \right]$$

pushing $\mathcal{N}(0,\mathbf{I})$ to the third Markov component of $T_1^{\sharp}T_0^{\sharp}\pi$:

$$\begin{aligned} \pi^{2}\left(\Theta, \mathbf{Z}_{2}, \mathbf{Z}_{3}\right) &:= \rho(\Theta, \mathbf{Z}_{2}) \,\mathcal{L}\left(\mathbf{y}_{3} \left|\Theta'', \mathbf{Z}_{3}\right.\right) \pi\left(\mathbf{Z}_{3} \left|\Theta'', \mathbf{Z}_{2}''\right.\right) \\ \Theta'' &:= \mathfrak{M}_{0}^{\Theta} \circ \mathfrak{M}_{1}^{\Theta}(\Theta) \;, \quad \mathbf{Z}_{1}'' := \mathfrak{M}_{1}^{1}(\Theta, \mathbf{Z}_{2}) \end{aligned}$$




Removing conditional dependencies



A framework for sequential inference

- Variational Bayesian algorithms for joint state-parameter inference
- Yields filtering, smoothing, full posterior via composition of maps
- Constant cost per assimilation/prediction step
- Evaluation of smoothing/full posterior grows linearly with time
- Map parametrization (basis, sparsity) governs complexity/accuracy trade-off

A framework for sequential inference

- Variational Bayesian algorithms for joint state-parameter inference
- Yields filtering, smoothing, full posterior via composition of maps
- Constant cost per assimilation/prediction step
- Evaluation of smoothing/full posterior grows linearly with time
- Map parametrization (basis, sparsity) governs complexity/accuracy trade-off

We need to characterize maps between $N(0,{\bf I})$ and the lag-1 smoothing

Preconditioning the assimilation step

At step k one needs to solve the problem

$$\mathfrak{M}_{k} = \operatorname*{arg\,min}_{\mathfrak{M}\in\mathcal{T}_{>}} \mathcal{D}_{\mathsf{KL}} \left(\mathfrak{M}_{\sharp}\rho(\mathbf{z}_{k},\mathbf{z}_{k+1}) \middle\| \underbrace{\left[\begin{array}{c} \mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k}) \\ \mathbf{z}_{k+1} \end{array} \right]^{\sharp} \pi(\mathbf{z}_{k},\mathbf{z}_{k+1}|\mathbf{y}_{j\leq k+1})}_{\mathcal{L}(\mathbf{y}_{k+1}|\mathbf{z}_{k+1}) \pi\left(\mathbf{z}_{k+1}|\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\right)\rho(\mathbf{z}_{k})} \right)$$

Preconditioning the assimilation step

At step k one needs to solve the problem

$$\mathfrak{M}_{k} = \operatorname*{arg\,min}_{\mathfrak{M}\in\mathcal{T}_{>}} \mathcal{D}_{\mathsf{KL}} \left(\mathfrak{M}_{\sharp}\rho(\mathbf{z}_{k},\mathbf{z}_{k+1}) \middle\| \underbrace{\left[\begin{array}{c} \mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k}) \\ \mathbf{z}_{k+1} \end{array} \right]^{\sharp} \pi(\mathbf{z}_{k},\mathbf{z}_{k+1}|\mathbf{y}_{j\leq k+1})}_{\mathcal{L}(\mathbf{y}_{k+1}|\mathbf{z}_{k+1}) \pi\left(\mathbf{z}_{k+1}|\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\right)\rho(\mathbf{z}_{k})} \right)$$

Ansatz: the next filtering distribution is similar to the latest.

Preconditioning the assimilation step

At step k one needs to solve the problem

$$\mathfrak{M}_{k} = \operatorname*{arg\,min}_{\mathfrak{M}\in\mathcal{T}_{>}} \mathcal{D}_{\mathsf{KL}}\left(\mathfrak{M}_{\sharp}\rho(\mathbf{z}_{k},\mathbf{z}_{k+1}) \middle\| \underbrace{\left[\begin{array}{c}\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\\\mathbf{z}_{k+1}\end{array}\right]^{\sharp} \pi(\mathbf{z}_{k},\mathbf{z}_{k+1}|\mathbf{y}_{j\leq k+1})}_{\mathcal{L}(\mathbf{y}_{k+1}|\mathbf{z}_{k+1}) \pi\left(\mathbf{z}_{k+1}|\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\right)\rho(\mathbf{z}_{k})}\right)$$

Ansatz: the next filtering distribution is similar to the latest.

Lorenz 63

Dynamics

Assimilation model

$$\begin{cases} \dot{x} = \sigma(y - x) \\ \dot{y} = x(\rho - z) - y \\ \dot{z} = xy - \beta z \end{cases}$$

$$\sigma = 10, \ \beta = 8/3, \ \rho = 28$$



$$\mathbf{u}_{k+1} = \mathbf{u}_k + 0.01 \cdot \mathcal{F}(\mathbf{u}_k) + \mathcal{N}(0, 10^{-1} \cdot \mathbf{I})$$
$$\mathbf{h}_k = \mathcal{O}(\mathbf{u}_k) + \mathcal{N}(0, 2 \cdot \mathbf{I}) , \quad k = 0, 8, 16, \dots$$





























		filtering		smoothing	
		Order 1	Order 3	Order 1	Order 3
at obs.	Avg. RMSE	7.39×10^{-1}	7.46×10^{-1}	4.82×10^{-1}	4.86×10^{-1}
	Med. RMSE	6.49×10^{-1}	6.61×10^{-1}	4.18×10^{-1}	4.23×10^{-1}
	Var. RMSE	1.51×10^{-1}	1.55×10^{-1}	5.62×10^{-2}	5.78×10^{-2}
incl. pred.	Avg. RMSE	8.37×10^{-1}	8.42×10^{-1}	4.67×10^{-1}	4.72×10^{-1}
	Med. RMSE	7.54×10^{-1}	7.50×10^{-1}	4.11×10^{-1}	4.14×10^{-1}
	Var. RMSE	2.23×10^{-1}	2.25×10^{-1}	$5.53 imes10^{-2}$	$5.72 imes 10^{-2}$





		filtering		smoothing	
		Order 1	Order 3	Order 1	Order 3
at obs.	Avg. RMSE	7.39×10^{-1}	7.46×10^{-1}	4.82×10^{-1}	$4.86 imes 10^{-1}$
	Med. RMSE	6.49×10^{-1}	6.61×10^{-1}	4.18×10^{-1}	4.23×10^{-1}
	Var. RMSE	1.51×10^{-1}	1.55×10^{-1}	5.62×10^{-2}	$5.78 imes 10^{-2}$

Posterior accuracy ($d = 3000$)	Order 1	Order 3
$\mathbb{V}[\log \frac{\rho}{T^{\sharp}\pi}]$	8.58×10^{-1}	2.11×10^{-1}
Metropolis indep. A/R	33.7%	62.9%
Metropolis indep. ESS (worst)	0.68%	12.54%





























Lorenz 96

Dynamics

Assimilation model

$$\frac{d\mathbf{Z}_{j}}{dt} = (\mathbf{Z}_{j+1} - \mathbf{Z}_{j-2}) \mathbf{Z}_{j-1} - \mathbf{Z}_{j} + F$$
$$j \in \{1, \dots, 40\}, \quad F = 8 \text{ (chaotic)}$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + 0.01 \cdot \mathcal{F}(\mathbf{u}_k) + \mathcal{N}(0, 10^{-1} \cdot \mathbf{I})$$
$$\mathbf{h}_k = \mathcal{O}(\mathbf{u}_k) + \mathcal{N}(0, 0.5 \cdot \mathbf{I}) ,$$
$$k = 0, 10, \dots$$





Map sparsity ansatz





Map sparsity ansatz





Map sparsity ansatz



Filtering





\sim			1
()	rn	er.	
\sim	i u		-

		Filtering	Smoothing ($d = 4 \times 10^4$)
at obs.	Avg. RMSE	8.35×10^{-1}	7.40×10^{-1}
	Med. RMSE	8.11×10^{-1}	7.29×10^{-1}
	Var. RMSE	3.26×10^{-2}	2.42×10^{-2}
incl. pred.	Avg. RMSE	9.20×10^{-1}	7.34×10^{-1}
	Med. RMSE	9.01×10^{-1}	7.25×10^{-1}
	Var. RMSE	4.21×10^{-2}	2.51×10^{-2}
Towards high(er)-dimensions









Work only on directions that depart from the reference



Work only on directions that depart from the reference



How to detect such directions?

Work with lazy maps [B, Zahm, Spantini, Marzouk 2019]

$$\widetilde{\mathfrak{M}}_{k} = \underset{\mathfrak{M}\in\mathcal{T}_{>}}{\operatorname{arg\,min}} \mathcal{D}_{\mathsf{KL}}\left(\mathfrak{M}_{\sharp}\rho(\mathbf{z}_{k},\mathbf{z}_{k+1}) \middle\| \underbrace{\left[\begin{array}{c}\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\\\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k+1})\end{array}\right]^{\sharp}\pi(\mathbf{z}_{k},\mathbf{z}_{k+1}|\mathbf{y}_{j\leq k+1})}_{\widetilde{\pi}(\mathbf{z}_{k},\mathbf{z}_{k+1})}\right)$$

For $\mathfrak{r}({m x}) = \log(\widetilde{\pi}({m x})/
ho({m x}))$ let

$$\mathbf{H} = \int
abla \mathfrak{r}(\boldsymbol{x}) \
abla \mathfrak{r}(\boldsymbol{x})^{ op} \
ho(\boldsymbol{x}) \ \mathrm{d}\boldsymbol{x} \ , \quad ext{and} \quad \mathbf{H} \mathbf{u}_i = \lambda_i \mathbf{u}_i \ .$$

Then, for $\lambda_i \geq \lambda_{i+1}$, $\mathbf{U}_r = [\mathbf{u}_1, \dots, \mathbf{u}_{i+1}]$ and $\mathbf{Q} = \mathbf{U}_r | \mathbf{U}_\perp$, there exist

Work with lazy maps [B, Zahm, Spantini, Marzouk 2019]

$$\widetilde{\mathfrak{M}}_{k} = \underset{\mathfrak{M}\in\mathcal{T}_{>}}{\operatorname{arg\,min}} \mathcal{D}_{\mathsf{KL}}\left(\mathfrak{M}_{\sharp}\rho(\mathbf{z}_{k},\mathbf{z}_{k+1}) \middle\| \underbrace{\left[\begin{array}{c}\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k})\\\mathfrak{M}_{k-1}^{1}(\mathbf{z}_{k+1})\end{array}\right]^{\sharp}\pi(\mathbf{z}_{k},\mathbf{z}_{k+1}|\mathbf{y}_{j\leq k+1})}_{\widetilde{\pi}(\mathbf{z}_{k},\mathbf{z}_{k+1})}\right)$$

For $\mathfrak{r}(m{x}) = \log(\widetilde{\pi}(m{x})/
ho(m{x}))$ let

$$\mathbf{H} = \int
abla \mathfrak{r}(oldsymbol{x}) \
abla \mathfrak{r}(oldsymbol{x})^{ op} \,
ho(oldsymbol{x}) \ \mathrm{d}oldsymbol{x} \ , \quad ext{and} \quad \mathbf{H} \mathbf{u}_i = \lambda_i \mathbf{u}_i \ .$$

Then, for $\lambda_i \geq \lambda_{i+1}$, $\mathbf{U}_r = [\mathbf{u}_1, \dots, \mathbf{u}_{i+1}]$ and $\mathbf{Q} = \mathbf{U}_r | \mathbf{U}_\perp$, there exist

This would be perfect is **Q** was triangular. But it is NOT!

Build state-dependent subspaces

1 Assemble

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{k,k} & \mathbf{H}_{k,k+1} \\ \mathbf{H}_{k+1,k} & \mathbf{H}_{k+1,k+1} \end{bmatrix}, \qquad \mathbf{H}_{i,j} = \int \left(\nabla_{\mathbf{z}_i} \mathfrak{r} \right) \left(\nabla_{\mathbf{z}_j} \mathfrak{r} \right)^\top \rho(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x}$$

2 Solve two truncated singular value problems:

•
$$[\mathbf{H}_{k,k}|\mathbf{H}_{k,k+1}] = \mathbf{U}_{r_k} \Sigma_{1:r_k} \mathbf{V}_{r_k}$$

• $[\mathbf{H}_{k+1,k}|\mathbf{H}_{k+1,k+1}] = \mathbf{U}_{r_{k+1}} \Sigma_{1:r_{k+1}} \mathbf{V}_{r_{k+1}}$

3 Find a basis \mathbf{U}_r for $\operatorname{span}(\mathbf{U}_{r_k}|\mathbf{U}_{r_{k+1}})$, e.g. using the QR decomposition

$$\mathbf{\mathcal{G}} \text{ Set } \mathbf{Q} = \begin{bmatrix} \mathbf{U}_r | \mathbf{U}_{\perp} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_r | \mathbf{U}_{\perp} \end{bmatrix}$$

$$\mathbf{\mathcal{G}} \text{ Solve } \widetilde{\mathfrak{M}}_k = \arg \min_{\mathfrak{M} \in \mathcal{T}_r^r} \mathcal{D}_{\mathsf{KL}} \left(\mathfrak{M}_{\sharp} \rho(\mathbf{z}_k, \mathbf{z}_{k+1}) \| \mathbf{Q}^{\sharp} \widetilde{\pi}(\mathbf{z}_k, \mathbf{z}_{k+1}) \right)$$

$$\text{ Sparsity of } \widetilde{\mathfrak{M}}_k:$$

Lorenz 96 – model II

$$\begin{aligned} \frac{d\mathbf{X}_n}{dt} &= \sum_{j=-J}^J \sum_{i=-J}^J \left(X_{n-K+j-i} X_{n+K+j} - X_{n-2K-i} X_{n-K-j} \right) / K^2 - X_n + F \\ \mathbf{X} &\in \mathbb{R}^d , \ d = 240 , \ K = 33 , \ J = (K-1)/2 , \ F = 14 \quad \text{(chaotic)} \end{aligned}$$

Assimilation model

$$\mathbf{u}_{k+1} = \mathbf{u}_k + 0.025 \cdot \mathcal{F}(\mathbf{u}_k) + \mathcal{N}(0, 10^{-1} \cdot \mathbf{I})$$
$$\mathbf{h}_i = \mathcal{O}(\mathbf{u}_i) + \mathcal{N}(0, 10^{-1} \cdot \mathbf{I}) , \ i = 0, 4, \dots$$





Using rank-2 lazy maps of order 2 and 50 samples.

Key contributions

Algorithms for filtering, smoothing, and parameters estimation via **deterministic couplings** and **optimization**.

Contact: Daniele Bigoni – dabi@mit.edu

Software: https://transportmaps.mit.edu

Bigoni et al. "Greedy inference with layers of lazy maps" (arXiv) Bigoni et al. "Adaptive construction of measure transports for Bayesian inference" Spantini et al. <u>"Inference via low-dimensional couplings"</u> (JMLR) Marzouk et al. <u>"Sampling via measure transport: an introduction"</u> (Springer) Parno et al. "Transport map accelerated Markov chain Monte Carlo" (JUQ) El Moselhy et al. "Bayesian inference with optimal maps" (JCP)

Thanks to:

