

Layers of lazy maps for large-scale inference

<https://arxiv.org/abs/1906.00031>

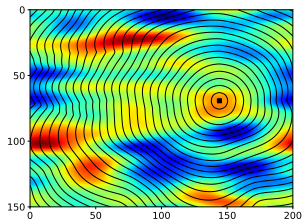
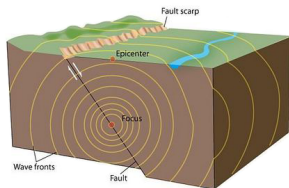
D. Bigoni[†] (**dabi@mit.edu**), O. Zahm[‡], A. Spantini[†], Y.M. Marzouk[†]

[†] Massachusetts Institute of Technology, USA

[‡] INRIA, France

Applied Inverse Problems Conference
Grenoble, France – July 12, 2019

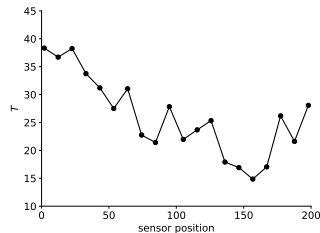
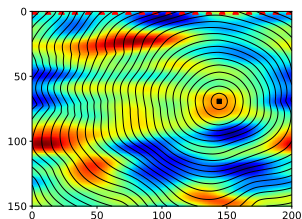
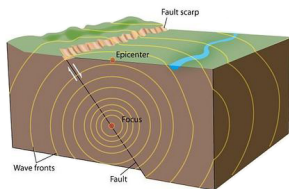
Bayesian inference – an oversimplified example



Mathematical model

$$\overbrace{|\nabla G(\mathbf{x})|}^{\text{travel time}} = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Bayesian inference – an oversimplified example



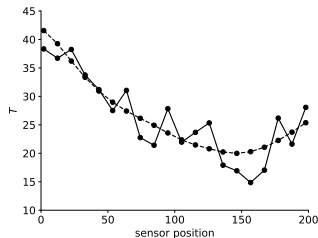
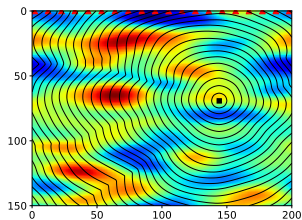
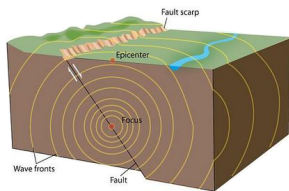
Mathematical model

$$|\nabla \overbrace{G(\mathbf{x})}^{\text{travel time}}| = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Observational model

$$\underbrace{\text{data}}_{\mathbf{d}} = \mathbf{G}(\mathbf{v}) + \underbrace{\varepsilon}_{\text{noise}}$$

Bayesian inference – an oversimplified example



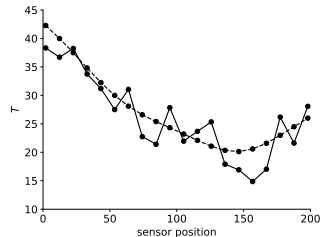
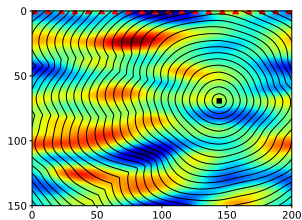
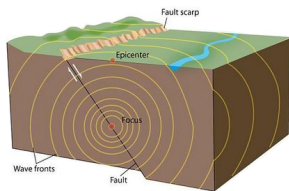
Mathematical model

$$|\nabla \overbrace{G(\mathbf{x})}^{\text{travel time}}| = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Observational model

$$\underbrace{\mathbf{d}}_{\text{data}} = \mathbf{G}(\mathbf{v}) + \underbrace{\varepsilon}_{\text{noise}}$$

Bayesian inference – an oversimplified example



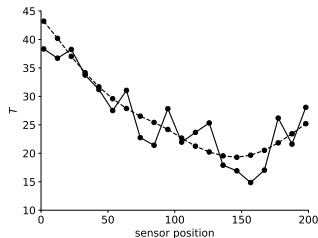
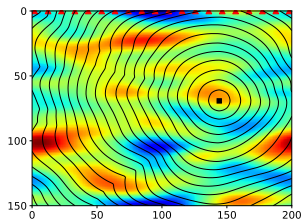
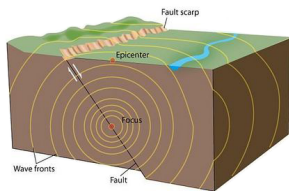
Mathematical model

$$|\nabla \overbrace{G(\mathbf{x})}^{\text{travel time}}| = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Observational model

$$\overbrace{\mathbf{d}}^{\text{data}} = \mathbf{G}(\underbrace{\mathbf{v}}_{\text{velocity field}}) + \underbrace{\varepsilon}_{\text{noise}}$$

Bayesian inference – an oversimplified example



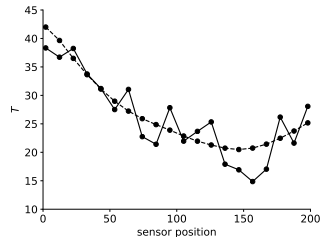
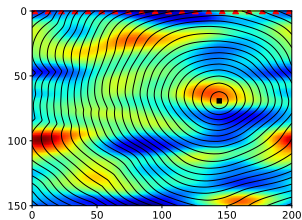
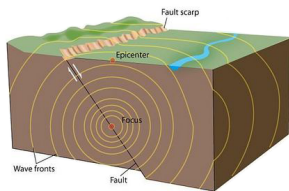
Mathematical model

$$|\nabla \overbrace{G(\mathbf{x})}^{\text{travel time}}| = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Observational model

$$\underbrace{\mathbf{d}}_{\text{data}} = \mathbf{G}(\mathbf{v}) + \underbrace{\varepsilon}_{\text{noise}}$$

Bayesian inference – an oversimplified example



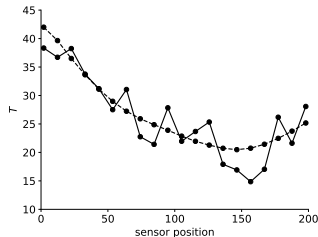
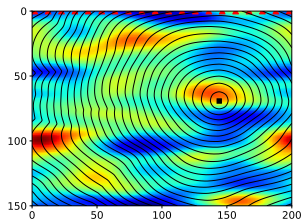
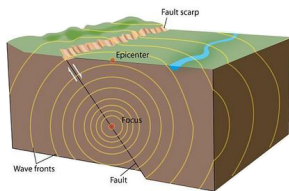
Mathematical model

$$|\nabla \overbrace{G(\mathbf{x})}^{\text{travel time}}| = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Observational model

$$\underbrace{\mathbf{d}}_{\text{data}} = \mathbf{G}(\mathbf{v}) + \underbrace{\varepsilon}_{\text{noise}}$$

Bayesian inference – an oversimplified example



Mathematical model

$$|\nabla \overbrace{G(\mathbf{x})}^{\text{travel time}}| = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

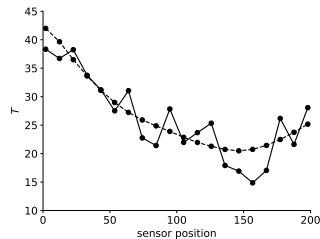
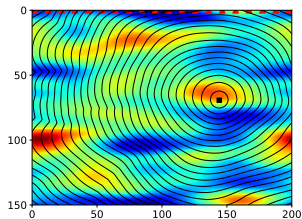
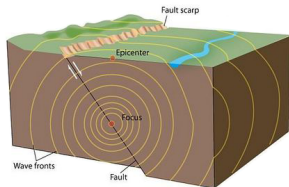
Observational model

$$\overbrace{\mathbf{d}}^{\text{data}} = \mathbf{G}(\mathbf{v}) + \underbrace{\epsilon}_{\text{noise}}$$

Bayesian inference model

$$\underbrace{\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})}_{\text{posterior}} \propto \underbrace{\mathcal{L}_{\mathbf{d}}(\mathbf{v})}_{\text{likelihood}} \underbrace{\pi_{\text{pr}}(\mathbf{v})}_{\text{prior}} = \pi_{\epsilon}(\mathbf{d} - \mathbf{G}(\mathbf{v}))\pi_{\text{pr}}(\mathbf{v})$$

Bayesian inference – an oversimplified example



Bayesian inference model

$$\underbrace{\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})}_{\text{posterior}} \propto \underbrace{\mathcal{L}_{\mathbf{d}}(\mathbf{v})}_{\text{likelihood}} \underbrace{\pi_{\text{pr}}(\mathbf{v})}_{\text{prior}} = \pi_{\epsilon}(\mathbf{d} - \mathbf{G}(\mathbf{v}))\pi_{\text{pr}}(\mathbf{v})$$

Decisions under uncertainty

$$\min_{\delta} \int L(\mathbf{v}, \delta) \pi_{\text{pos}}(\mathbf{v}|\mathbf{d}) d\mathbf{v}$$

Goal: characterize $\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})$, i.e.

- construct approximations

$$\int f(\mathbf{v})\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})d\mathbf{v} \approx \int f(\mathbf{v})\tilde{\pi}_{\text{pos}}(\mathbf{v}|\mathbf{d})d\mathbf{v} \approx \sum_{i=1}^n f(\mathbf{v}^{(i)})\mathbf{w}^{(i)}$$

- control the error between $\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})$ and $\tilde{\pi}_{\text{pos}}(\mathbf{v}|\mathbf{d})$

Difficulties:

- $\mathbf{v} \in \mathbb{R}^d$ where $d \gg 1$
- The model $\mathbf{G}(\mathbf{v})$ is non-linear
- Evaluation of the model $\mathbf{G}(\mathbf{v})$ is expensive

Outline

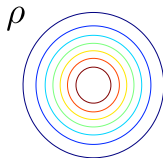
Transport maps

Layers of lazy maps

Results

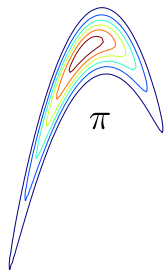
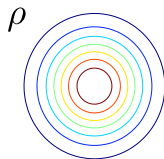
Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$



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

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$

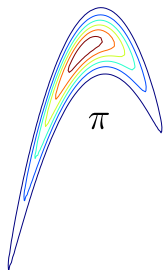
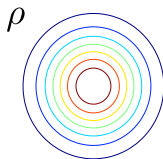


Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

PF $T_\# \rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^\# \pi = \pi \circ T |\nabla T|$



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

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

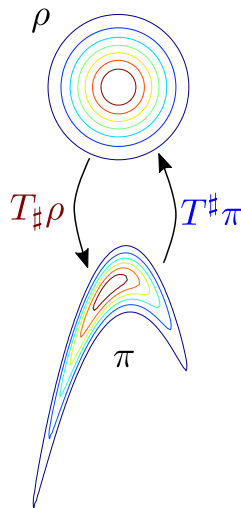
PF $T_\# \rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^\# \pi = \pi \circ T |\nabla T|$

- We want T such that

PF $T_\# \rho = \pi$

PB $T^\# \pi = \rho$



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

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

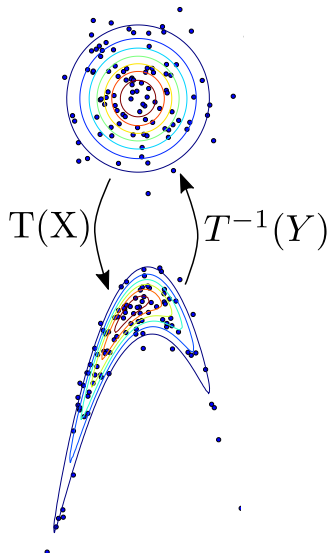
PF $T_\# \rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^\# \pi = \pi \circ T |\nabla T|$

- We want T such that

PF For $X \sim \nu_\rho$, $T(X) \sim \nu_\pi$

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



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

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

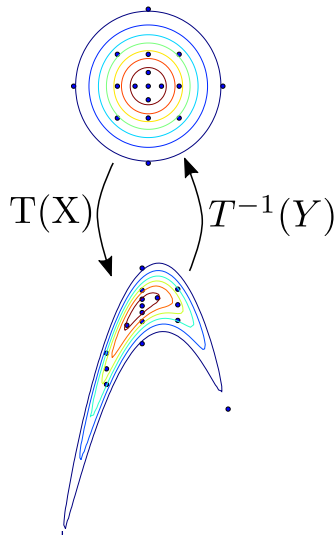
PF $T_\# \rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^\# \pi = \pi \circ T |\nabla T|$

- We want T such that

PF For $X \sim \nu_\rho$, $T(X) \sim \nu_\pi$

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



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

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

PF $T_\# \rho = \rho \circ T^{-1} |\nabla T^{-1}|$

PB $T^\# \pi = \pi \circ T |\nabla T|$

- We want T such that

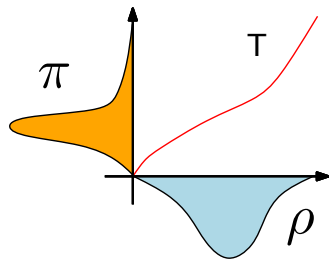
PF For $X \sim \nu_\rho$, $T(X) \sim \nu_\pi$

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

Knothe-Rosenblatt rearrangement

$\forall \nu_\rho, \nu_\pi$ Lebesgue absolutely continuous

\exists a **triangular monotone** map s.t. $T_\# \rho = \pi$



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

Triangular monotone maps

$$\mathcal{T}_{>} = \left\{ T : \mathbb{R}^d \rightarrow \mathbb{R}^d : \overbrace{[T(\mathbf{x})]_k = T^{(k)}(x_1, \dots, x_k)}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_k} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Triangular monotone maps

$$\mathcal{T}_{>} = \left\{ T : \mathbb{R}^d \rightarrow \mathbb{R}^d : \overbrace{[T(\mathbf{x})]_k = T^{(k)}(x_1, \dots, x_k)}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_k} T^{(k)} > 0}^{\text{monotone}} \right\}$$

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

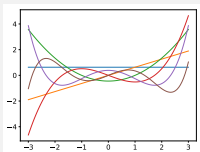
$$\boxed{\mathcal{T}_{>}^n} = \left\{ T : \mathbb{R}^d \rightarrow \mathbb{R}^d : \overbrace{[T(\mathbf{x})]_k = T^{(k)}(x_1, \dots, x_k)}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_k} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Integrated squared representation – $\varepsilon > 0$

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

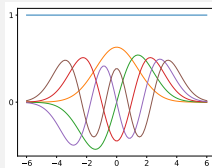
Constant part

$$c_k(x_{1:k-1}) = \sum_{\mathbf{i} \in \mathcal{I}_k} \mathbf{a}_{\mathbf{i}} \Phi_{\mathbf{i}}(x_{1:k-1})$$



Squared part

$$h_k(x_{1:k-1}, t) = \sum_{\mathbf{j} \in \mathcal{J}_k} \mathbf{b}_{\mathbf{j}} \Psi_{\mathbf{j}}(x_{1:k-1}, t)$$

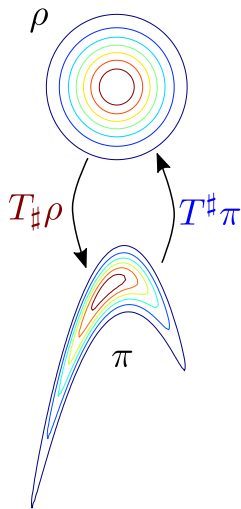


Knothe-Rosenblatt rearrangement

$\forall \nu_\rho, \nu_\pi$ Lebesgue absolutely continuous

\exists a **triangular monotone** map s.t. $T_\# \rho = \pi$

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



Minimize KL-divergence to find optimal map

$$\hat{T} = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_{\#} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#} \pi} \right]$$

Minimize KL-divergence to find optimal map

$$\hat{T} = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_{\#} \nu_{\rho} \| \nu_{\pi}) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#} \pi} \right]$$

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

Minimize KL-divergence to find optimal map

$$\hat{T} = \arg \min_{T \in \mathcal{T}_{>}} D_{\text{KL}}(T_{\#} \nu_{\rho} \| \nu_{\pi}) = \arg \min_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#} \pi} \right]$$

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

$$\hat{T} = \arg \min_{T \in \mathcal{T}_{>}} D_{\text{KL}}(T_{\#} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \arg \min_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#} \tilde{\pi}} \right] + \log \int \tilde{\pi}$$

$$\text{Optimal } \hat{T} \in \mathcal{T}_{>} \text{ and } \int \tilde{\pi} = 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(\hat{T})_{\#} \tilde{\pi}} \right] = 0$$

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

Convergence criterion – Variance diagnostic

$$\hat{T} = \arg \min_{T \in \mathcal{T}_{>}} D_{\text{KL}}(T_{\#} \boldsymbol{\nu}_{\rho} \| \boldsymbol{\nu}_{\pi}) = \arg \min_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\#} \tilde{\pi}} \right] + \log \int \tilde{\pi}$$

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

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

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

Pros & cons

$$\hat{T} = \arg \min_{T \in \mathcal{T}_{>}} D_{\text{KL}}(T_{\#}\rho \| \pi) = \arg \min_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#}\pi} \right]$$

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

Pros & cons

$$\hat{T} = \arg \min_{T \in \mathcal{T}_{>}} D_{\text{KL}}(T_{\#}\rho \| \pi) = \arg \min_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#}\pi} \right]$$

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

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

Use \hat{T} as a **preconditioner** for ... Importance Sampling

$$\int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int f(\mathbf{x})\frac{\pi(\mathbf{x})}{\hat{T}_{\#}\rho(\mathbf{x})}\hat{T}_{\#}\rho(\mathbf{x})d\mathbf{x} = \int f\circ\hat{T}(\mathbf{x})\frac{\hat{T}^{\#}\pi(\mathbf{x})}{\rho(\mathbf{x})}\rho(\mathbf{x})d\mathbf{x}$$

Use \hat{T} as a **preconditioner** for ... Importance Sampling

$$\int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int f(\mathbf{x})\frac{\pi(\mathbf{x})}{\hat{T}_{\#}\rho(\mathbf{x})}\hat{T}_{\#}\rho(\mathbf{x})d\mathbf{x} = \int f\circ\hat{T}(\mathbf{x})\frac{\hat{T}^{\#}\pi(\mathbf{x})}{\rho(\mathbf{x})}\rho(\mathbf{x})d\mathbf{x}$$

$$\text{If } \hat{T}^{\#}\tilde{\pi} \propto \rho \quad \text{then} \quad \mathbb{V}\left[\frac{\hat{T}^{\#}\tilde{\pi}(\mathbf{x})}{\rho(\mathbf{x})}\right] \approx 0 \quad \text{and}$$

$$\sum_{i=1}^N f\circ\hat{T}(\mathbf{x}_i) \mathbf{w}_i, \quad \mathbf{w}_i = \frac{\tilde{\mathbf{w}}_i}{\sum_{i=1}^N \tilde{\mathbf{w}}_i}, \quad \tilde{\mathbf{w}}_i = \frac{\hat{T}^{\#}\tilde{\pi}(\mathbf{x}_i)}{\rho(\mathbf{x}_i)}$$

will be an **accurate estimator** of $\int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$

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

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

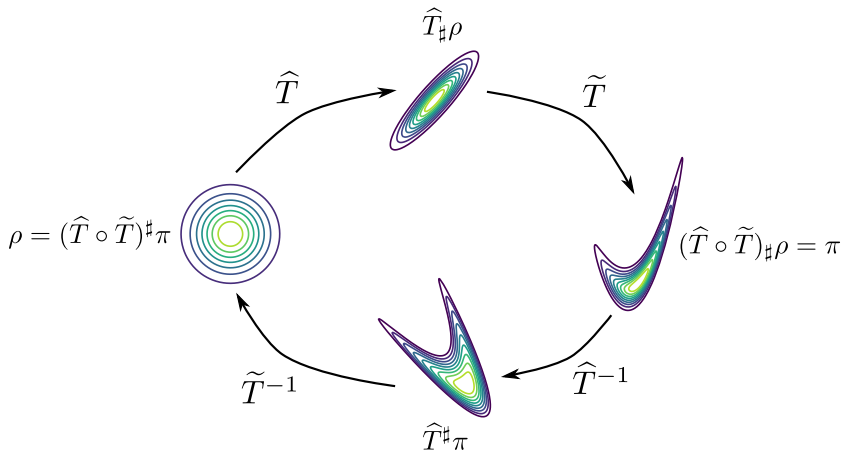
Use \hat{T} as a **preconditioner** for ... Markov Chain Monte Carlo

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

If $\hat{T}^\# \tilde{\pi} \propto \rho$ then ρ is a **good proposal** for $\hat{T}^\# \tilde{\pi}$.

Use \hat{T} as a **preconditioner** for ... Transport Maps!

- 1 Solve $\hat{T} = \arg \min_{T \in \mathcal{T}_{>}^n} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#} \pi} \right]$
- 2 Solve $\tilde{T} = \arg \min_{T \in \mathcal{T}_{>}^n} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#} \hat{T}_{\#} \pi} \right]$



Pros & cons

$$\hat{T} = \arg \min_{T \in \mathcal{T}_{>}} D_{\text{KL}}(T_{\#}\rho \| \pi) = \arg \min_{T \in \mathcal{T}_{>}} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#}\pi} \right]$$

- + **Gradient-based unconstrained optimization** if gradients are available
- + We can **explore π in parallel**
- + We can **generate i.i.d. samples** from $\hat{T}_{\#}\nu_{\rho} = \nu_{\pi}$ **in parallel**
- + We can **assess convergence!**
- + The map can be used as a **preconditioner**

Pros & cons

$$\hat{T} = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_{\#}\rho \parallel \pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\#}\pi} \right]$$

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

Layers of lazy maps

Incrementally construct improving maps
by working on residual distributions.

What is a lazy map?

Few ($r \ll d$) **complex components** and many “**lazy**” identity components:

$$T(\mathbf{x}) = \begin{bmatrix} T^{(1)}(x_1) \\ T^{(2)}(x_1, x_2) \\ \vdots \\ T^{(r)}(x_1, \dots, x_r) \\ x_{r+1} \\ \vdots \\ x_d \end{bmatrix} \in \mathcal{T}_r \subset \mathcal{T}_>$$

Maps of this type are effective if ρ and π agree along $d - r$ coordinates.

Assume there exists a rotation matrix \mathbf{Q} such that

$$\int \pi \circ \mathbf{Q}(\boldsymbol{\xi}_{1:r}, \mathbf{x}_{r+1:d}) \, \mathrm{d}\boldsymbol{\xi}_{1:r} = \int \rho(\boldsymbol{\xi}_{1:r}, \mathbf{x}_{r+1:d}) \, \mathrm{d}\boldsymbol{\xi}_{1:r},$$

Then there exist a lazy map $T \in \mathcal{T}_r$ such that

$$T_{\#}\rho = \mathbf{Q}^{\#}\pi$$

Finding a good rotation \mathbf{Q}

For any distribution ν_η with finite second moment, let

$$(\mathbf{H}_\eta)_{ij} = \int \partial_i \mathfrak{r}(\mathbf{x}) \partial_j \mathfrak{r}(\mathbf{x}) \eta(\mathbf{x}) \, d\mathbf{x} \, , \quad \mathfrak{r} := \log(\pi/\rho).$$

If $\text{rank}(\mathbf{H}_\eta) = r$ and $\nu_\rho = \mathcal{N}(0, \mathbf{I})$, then

there exist a **rotation** \mathbf{Q} and a **lazy map** $T \in \mathcal{T}_r$
such that $T_\# \rho = \mathbf{Q}^\# \pi$

Certified approximation π^* and optimal rotation \mathbf{Q}

[Zahm 2018, Bigoni 2019]

Let the columns of $\mathbf{U} \in \mathbb{R}^{d \times r}$ be the eigenvectors corresponding to the largest r eigenvalues $\{\lambda_i\}_{i=1}^r$ of \mathbf{H}_η and let

$$\pi^*(\mathbf{x}) := f(\mathbf{U}^\top \mathbf{x}) \rho(\mathbf{x}) ,$$

for f given by the conditional expectation

$$f(\mathbf{z}) := \mathbb{E} \left[\pi(\mathbf{X}) / \rho(\mathbf{X}) \middle| \mathbf{U}^\top \mathbf{X} = \mathbf{z} \right] , \quad \mathbf{X} \sim \rho .$$

Then,

$$\mathcal{D}_{\text{KL}}(\pi \| \pi^*) \leq \lambda_{r+1} + \dots + \lambda_d \quad \text{and} \quad \mathbf{Q} = [\mathbf{U} | \mathbf{U}_\perp]$$

In practical problems...

$$(\mathbf{H}_\eta)_{ij} = \int \partial_i \mathfrak{r}(\mathbf{x}) \partial_j \mathfrak{r}(\mathbf{x}) \eta(\mathbf{x}) \, d\mathbf{x} \, , \quad \mathfrak{r} := \log(\pi/\rho).$$

- \mathbf{H}_η will need to be approximated using some quadrature
- \mathbf{H}_η will only be approximately low-rank
- The spectrum of \mathbf{H}_η will depend on the sampling distribution ν_η (the optimal distribution would be ν_π itself)

Construction of one lazy map

- 1: **procedure** LAZYMAP($\pi, \rho, \varepsilon, r_{\max}$)
- 2: Compute $H = \int (\nabla \log \frac{\pi}{\rho})(\nabla \log \frac{\pi}{\rho})^\top d\rho$
- 3: Solve the eigenvalue problem $Hu_i = \lambda_i u_i$
- 4: Let $r = r_{\max} \wedge \min\{r \leq d : \frac{1}{2} \sum_{i>r} \lambda_i \leq \varepsilon\}$
- 5: Assemble $\mathbf{Q} = [\mathbf{U} | \mathbf{U}_\perp]$.
- 6: Find T solution to $\min_{T \in \mathcal{T}_r} \mathcal{D}_{\text{KL}}(\rho || T^\# \mathbf{Q}^\# \pi)$
- 7: **return** $\mathbf{Q} \circ T$
- 8: **end procedure**

Greedy construction of lazy maps

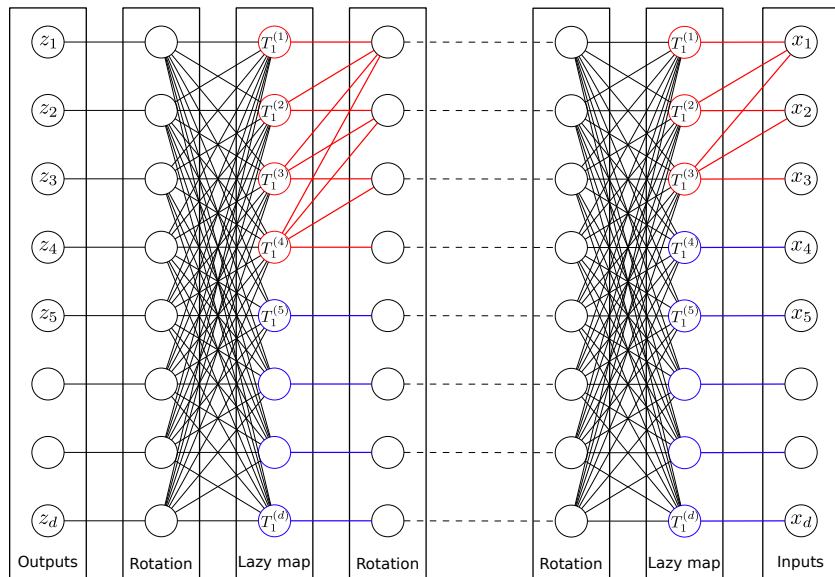
```
1: procedure LAYERSOFLAZYMAPS( $\pi, \rho, \varepsilon, r, \ell_{\max}$ )
2:   Set  $\pi_0 = \pi$  and  $\ell = 0$ 
3:   while  $\ell \leq \ell_{\max}$  and  $\frac{1}{2} \text{Tr}(H_\ell) \geq \varepsilon$  do
4:      $\ell \leftarrow \ell + 1$ 
5:     Compute  $T_\ell = \text{LAZYMAP}(\pi_{\ell-1}, \rho, 0, r)$ 
6:     Update  $\mathfrak{T}_\ell = \mathfrak{T}_{\ell-1} \circ T_\ell$ 
7:     Compute  $\pi_\ell = (\mathfrak{T}_\ell)^\# \pi$ 
8:     Compute  $H_\ell = \int (\nabla \log \frac{\pi_\ell}{\rho})(\nabla \log \frac{\pi_\ell}{\rho})^\top d\rho$ 
9:   end while
10:  return  $\mathfrak{T}_\ell = T_1 \circ \dots \circ T_\ell$ 
11: end procedure
```

Greedy construction of lazy maps

```
1: procedure LAYERSOFLAZYMAPS( $\pi, \rho, \varepsilon, r, \ell_{\max}$ )
2:   Set  $\pi_0 = \pi$  and  $\ell = 0$ 
3:   while  $\ell \leq \ell_{\max}$  and  $\frac{1}{2} \text{Tr}(H_\ell) \geq \varepsilon$  do
4:      $\ell \leftarrow \ell + 1$ 
5:     Compute  $T_\ell = \text{LAZYMAP}(\pi_{\ell-1}, \rho, 0, r)$ 
6:     Update  $\mathfrak{T}_\ell = \mathfrak{T}_{\ell-1} \circ T_\ell$ 
7:     Compute  $\pi_\ell = (\mathfrak{T}_\ell)^\# \pi$ 
8:     Compute  $H_\ell = \int (\nabla \log \frac{\pi_\ell}{\rho})(\nabla \log \frac{\pi_\ell}{\rho})^\top d\rho$ 
9:   end while
10:  return  $\mathfrak{T}_\ell = T_1 \circ \dots \circ T_\ell$ 
11: end procedure
```

\mathfrak{T} progressively “Gaussianizes” π .

Composition of layers of lazy transport maps



In practice...

Log-Gaussian Cox process with sparse observations ($d = 64^2$)

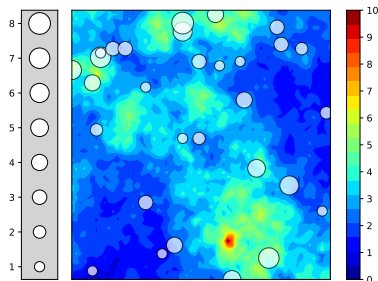
Statistical model

Observables: $\mathbf{Y} := (Y_i)_{i=1}^{30}$, $Y_i \sim \text{Poisson}(\Lambda_i/d)$

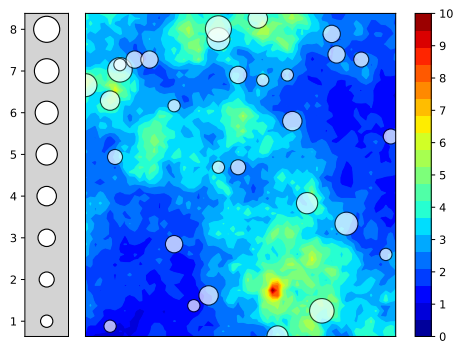
Latent field: $(\Lambda_i)_{i=1}^d \sim \log \mathcal{N}(\mu, \text{cov}(\mathbf{z}, \mathbf{z}'))$

$$\text{cov}(\mathbf{z}, \mathbf{z}') = \sigma^2 \exp(-\|\mathbf{z}_i - \mathbf{z}_j\|_2 / (64\beta))$$

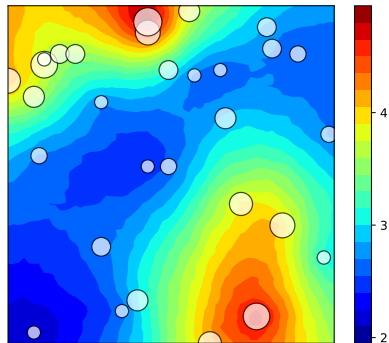
Posterior: $\pi(\Lambda | \mathbf{Y} = \mathbf{y}^*) \propto \pi(\mathbf{Y} = \mathbf{y}^* | \Lambda) \pi(\Lambda)$



Log-Gaussian Cox process with sparse observations ($d = 64^2$)

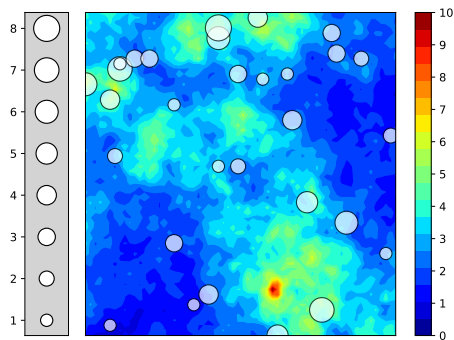


Field Λ^* and observations \mathbf{y}^*

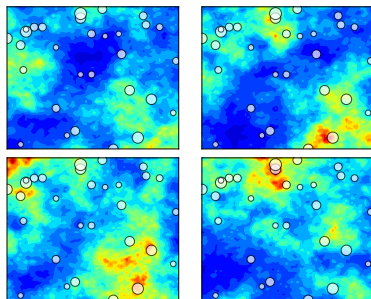


$\mathbb{E}[\Lambda | \mathbf{y}^*]$

Log-Gaussian Cox process with sparse observations ($d = 64^2$)

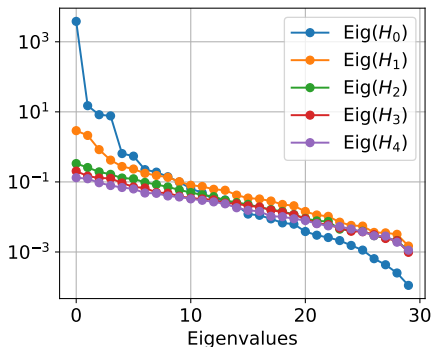
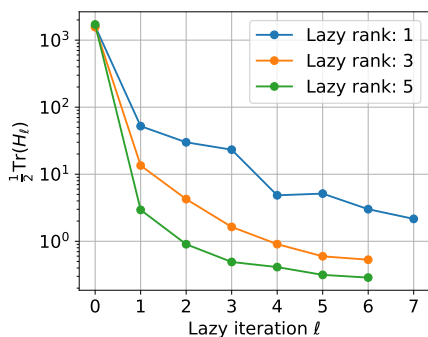


Field Λ^* and observations y^*



Realizations of $\Lambda \sim \pi_{\Lambda|y^*}(\lambda)$

Log-Gaussian Cox process with sparse observations ($d = 64^2$)



Metropolis-Hastings with independent proposals of $\mathcal{T}^\# \pi$

A/R	ESS (worst)
72.6%	26.6%

Elliptic problem with unknown coefficients ($d = 2601$)

Forward model

$G : \kappa \mapsto u$

$$\begin{cases} -\nabla \cdot (\kappa(\mathbf{x}, \omega) \nabla u(\mathbf{x}, \omega)) = 0 & \text{in } \Gamma \times \Omega \\ u(\mathbf{x}, \omega) = 0 & \text{on } \mathbf{x}_1 = 0 \\ u(\mathbf{x}, \omega) = 1 & \text{on } \mathbf{x}_1 = 1 \\ -\frac{\partial u}{\partial n}(\mathbf{x}) = 0 & \text{on } \mathbf{x}_2 \in \{0, 1\} \end{cases}$$

$$\kappa(\mathbf{x}, \omega) = \exp(g(\mathbf{x}, \omega)) , \quad g(\mathbf{x}, \omega) \sim \mathcal{N}(\mathbf{0}, C_g(\mathbf{x}, \mathbf{x}'))$$

$$C_g(\mathbf{x}, \mathbf{x}') = \exp(-|\mathbf{x} - \mathbf{x}'|)$$

Elliptic problem with unknown coefficients ($d = 2601$)

Forward model
 $\mathbf{G} : \kappa \mapsto u$

$$\begin{cases} -\nabla \cdot (\kappa(\mathbf{x}, \omega) \nabla u(\mathbf{x}, \omega)) = 0 & \text{in } \Gamma \times \Omega \\ u(\mathbf{x}, \omega) = 0 & \text{on } \mathbf{x}_1 = 0 \\ u(\mathbf{x}, \omega) = 1 & \text{on } \mathbf{x}_1 = 1 \\ -\frac{\partial u}{\partial n}(\mathbf{x}) = 0 & \text{on } \mathbf{x}_2 \in \{0, 1\} \end{cases}$$

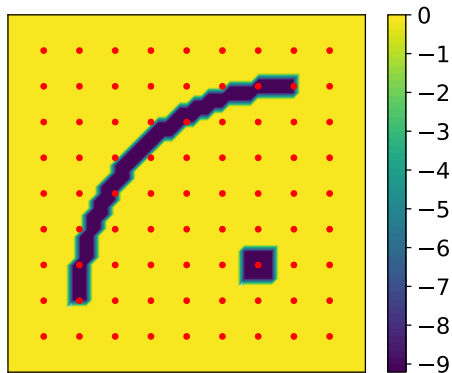
$$\kappa(\mathbf{x}, \omega) = \exp(g(\mathbf{x}, \omega)) , \quad g(\mathbf{x}, \omega) \sim \mathcal{N}(\mathbf{0}, C_g(\mathbf{x}, \mathbf{x}'))$$

$$C_g(\mathbf{x}, \mathbf{x}') = \exp(-|\mathbf{x} - \mathbf{x}'|)$$

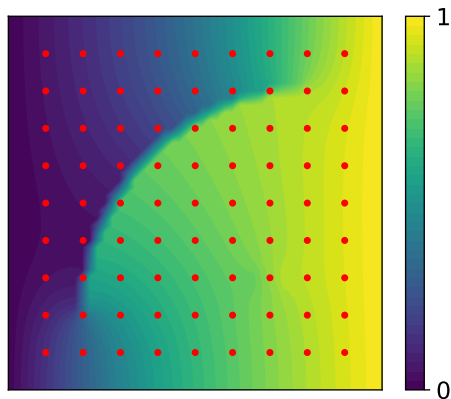
Bayesian inverse problem

$$\underbrace{\pi_{\text{pos}}(\kappa | \mathbf{d})}_{\text{posterior}} \propto \underbrace{\mathcal{L}_{\mathbf{d}}(\kappa)}_{\text{likelihood}} \underbrace{\pi_{\text{pr}}(\kappa)}_{\text{prior}} = \pi_{\varepsilon}(\mathbf{d} - \mathbf{G}(\kappa)) \pi_{\text{pr}}(\kappa)$$

Elliptic problem with unknown coefficients ($d = 2601$)

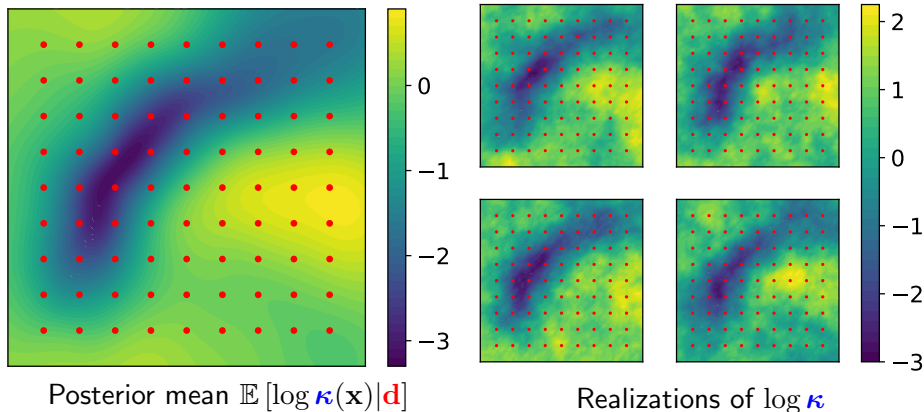


Synthetic field $\log \kappa^*(\mathbf{x})$

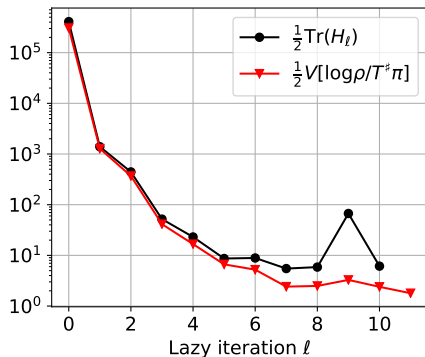


Synthetic solution $\mathbf{G}(\kappa^*)$

Elliptic problem with unknown coefficients ($d = 2601$)



Elliptic problem with unknown coefficients ($d = 2601$)



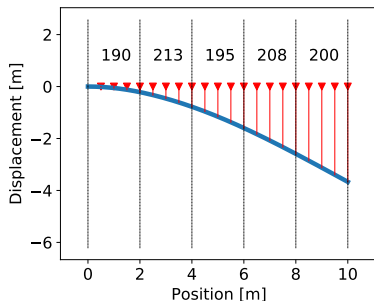
	A/R	ESS		
		worst	best	avg.
π	0.4%	$\sim 0\%$	$\sim 0\%$	$\sim 0\%$
$\mathfrak{T}^\# \pi$	28%	0.2%	1.6%	1.5%

Metropolis-Hastings with pCN
proposal ($\beta = 0.5$)

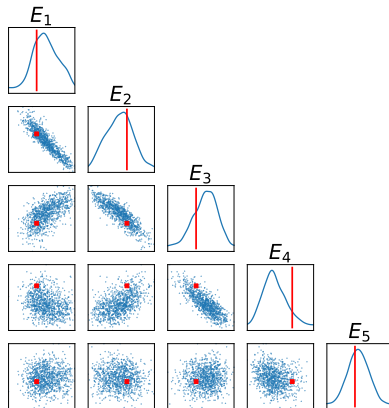
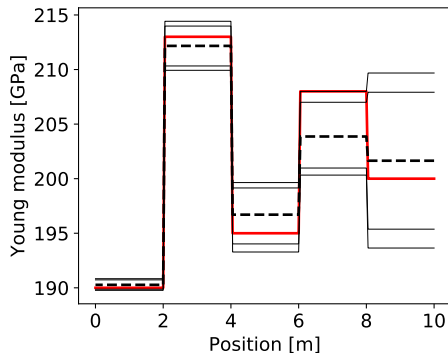
Estimation of the Young's modulus of a cantilever beam ($d = 5$)

Forward model
 $\mathbf{G} : \mathbf{E} \mapsto \mathbf{u}$

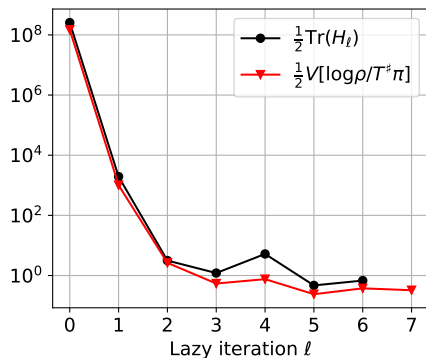
$$\begin{cases} \frac{d}{dx} \left[\frac{E(x)}{2(1+\nu)} \left(\varphi(x) - \frac{d}{dx} u(x) \right) \right] = \frac{q(x)}{\kappa A} , \\ \frac{d}{dx} \left(E(x) I \frac{d}{dx} \varphi(x) \right) = \kappa A \frac{E(x)}{2(1+\nu)} \left(\varphi(x) - \frac{d}{dx} u(x) \right) . \end{cases}$$



Estimation of the Young's modulus of a cantilever beam ($d = 5$)



Estimation of the Young's modulus of a cantilever beam ($d = 5$)



A/R	ESS		
	worst	best	avg.
68.3%	7.0%	38.7%	20.1%

Metropolis-Hastings with
independent proposals

Key contributions

Algorithms for characterizing probability measures
via layers of low-dimensional **deterministic couplings**

Contact: Daniele Bigoni – dabi@mit.edu

Software: <https://transportmaps.mit.edu>

Bigoni et al. “Greedy inference with layers of lazy maps” (arXiv)

Zahm et al. “Certified dimension reduction in nonlinear Bayesian inverse problems” (arXiv)

Bigoni et al. “On the computation of monotone transports” (preprint)

Spantini et al. “Inference via low-dimensional couplings” (JMLR)

Marzouk et al. “Sampling via measure transport: an introduction” (Springer)

Parno et al. “Transport map accelerated Markov chain Monte Carlo” (JUQ)

El Moselhy et al. “Bayesian inference with optimal maps” (JCP)

Thanks to:

