



Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data

Vaidotas Šimkus Benjamin Rhodes Michael U. Gutmann

School of Informatics, The University of Edinburgh



About the paper

- General method for estimating models from incomplete data.
- Needs only D Gibbs conditionals $q_{\phi_j}(x_j | \mathbf{x}_{\setminus j})$ instead of 2^D variational distributions (one for each pattern of missingness).
- Achieves better or competitive performance on normalising flow and VAE estimation from incomplete data.
- Published at Journal of Machine Learning Research, 2023.

Key challenges with deep statistical models

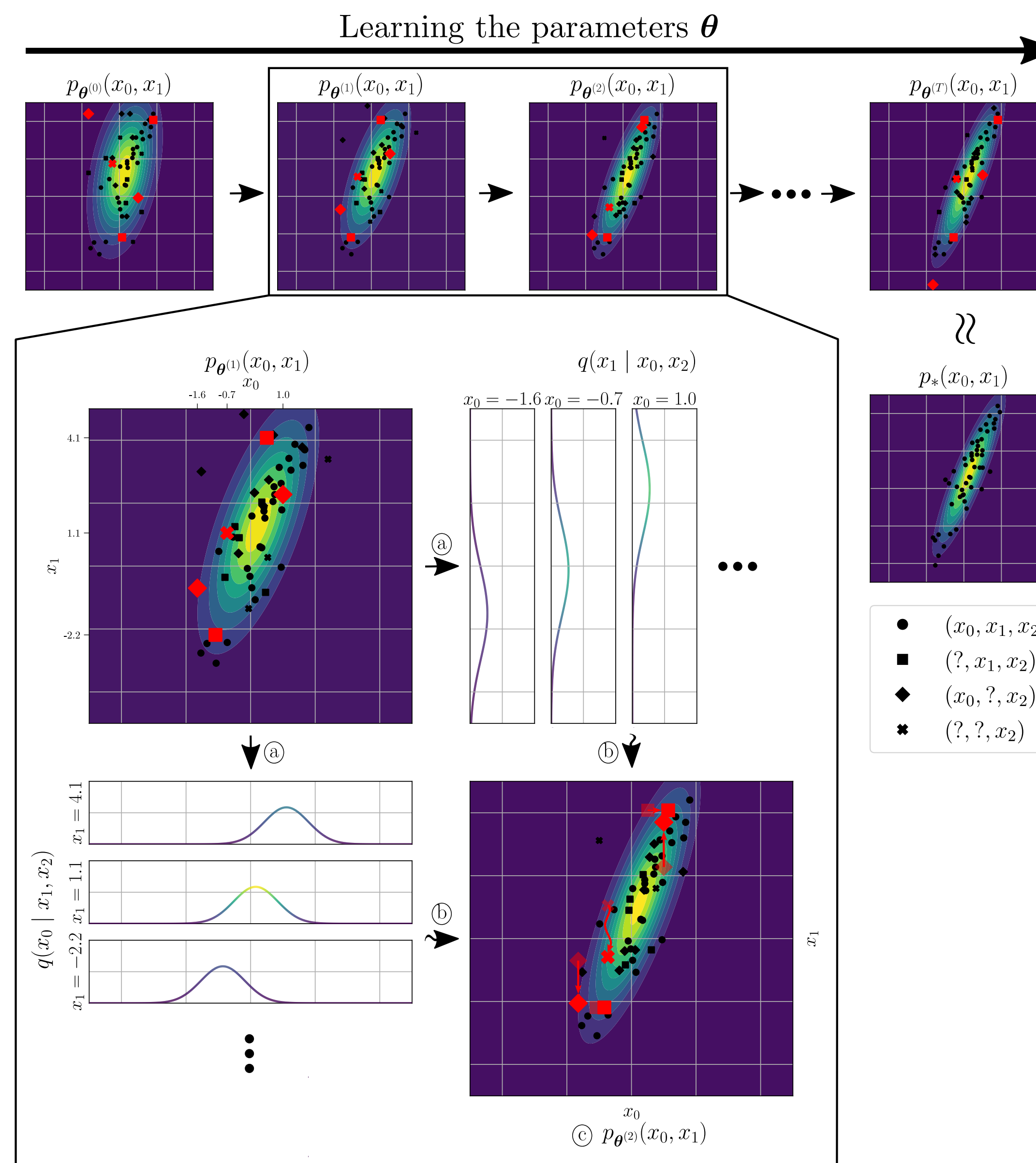
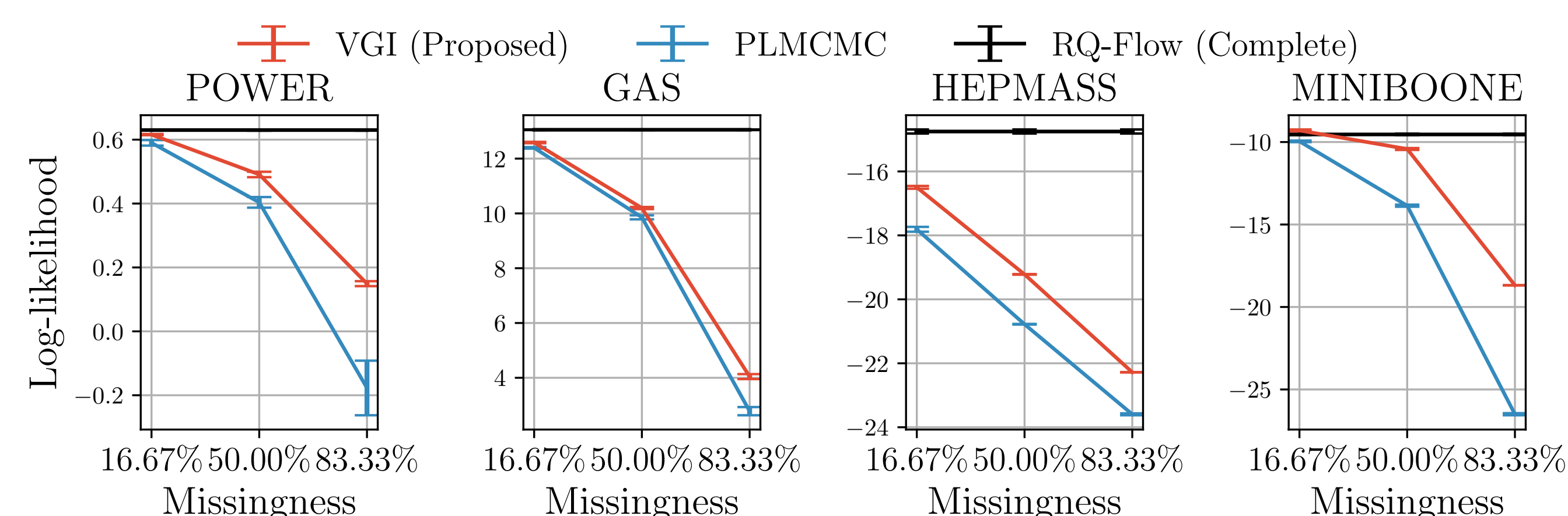
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data $\mathbf{x} \in \mathcal{D}$.
- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}) d\mathbf{x}_{\text{mis}}$ is generally not tractable. (With $\mathbf{x}_{\text{obs}} \cup \mathbf{x}_{\text{mis}} = \mathbf{x}$ and $\mathbf{x}_{\text{obs}} \cap \mathbf{x}_{\text{mis}} = \emptyset$.)
- Conditional sampling of $p_{\theta}(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}})$ is inefficient.

Why not use expectation maximisation?

$$\log p_{\theta}(\mathbf{x}_{\text{obs}}) \geq \mathbb{E}_{f(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}})} \left[\log \frac{p_{\theta}(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}})}{f(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}})} \right], \quad \text{“ELBO”}$$

- E-step: Maximise the ELBO w.r.t. $f(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i)$ for $\forall \mathbf{x}_{\text{obs}}^i \in \mathcal{D}$:
 $f(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i) = p_{\theta^t}(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i)$.
- M-step: Maximise the ELBO w.r.t. θ :
 $\theta^{t+1} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{p_{\theta^t}(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i)} \left[\log p_{\theta}(\mathbf{x}_{\text{obs}}^i, \mathbf{x}_{\text{mis}}) \right]$
- Requires sampling $p_{\theta^t}(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i)$ for $\forall \mathbf{x}_{\text{obs}}^i \in \mathcal{D} \rightarrow$ inefficient!

Results: Normalising flows



Issues with (amortised) variational inference

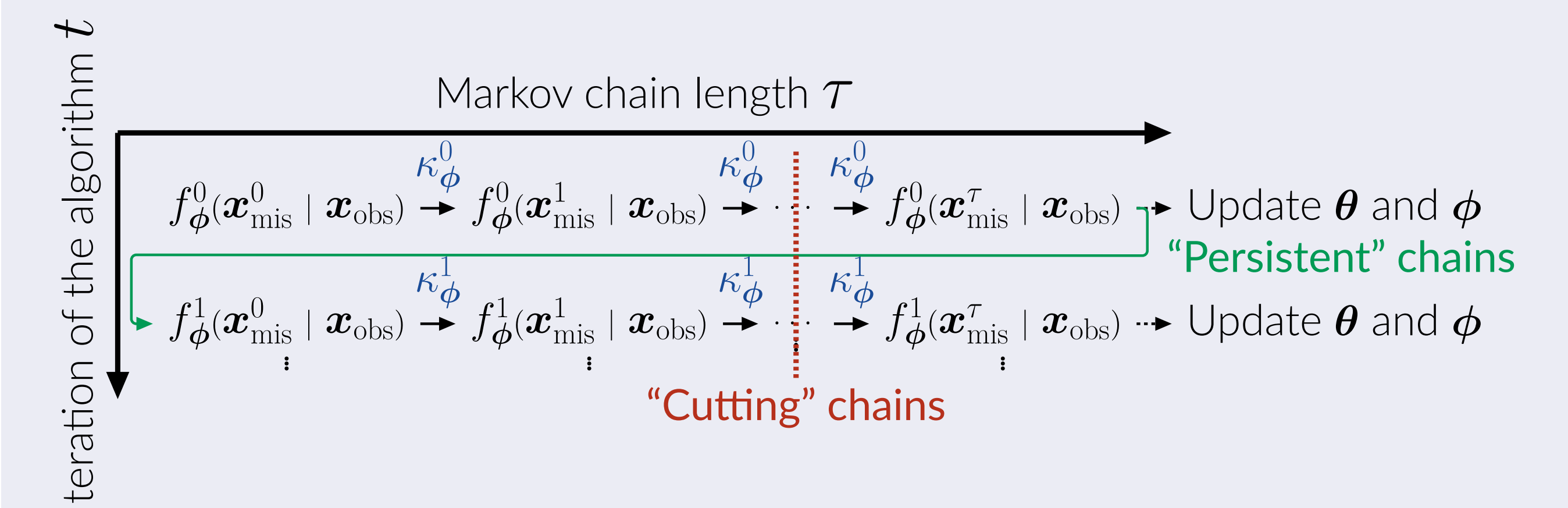
- $\forall \mathbf{x}_{\text{obs}}^i \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ . (Inference \rightarrow optimisation.)
- Is inefficient if $|\mathcal{D}|$ is large.
- Amortised VI: Parametrise all $f_{\phi}(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}}^i)$ with a single neural network $\text{NN}_{\phi}(\mathbf{x}_{\text{obs}}^i)$ shared for $\forall \mathbf{x}_{\text{obs}}^i \in \mathcal{D}$?
- Makes VI efficient for latent-variable models with large $|\mathcal{D}|$.
- But can we use amortised VI for estimation with missing data?

	d_1	d_2	d_3	d_4	$f_{\phi}(\mathbf{x}_{\text{mis}}^i \mathbf{x}_{\text{obs}}^i)$
\mathbf{x}^1	x_1^1	?	x_3^1	x_4^1	$f_{\phi}(x_2^1 x_1^1, x_3^1, x_4^1)$
\mathbf{x}^2	?	x_2^2	x_3^2	?	$f_{\phi}(x_1^2, x_4^2 x_2^2, x_3^2)$
\mathbf{x}^3	?	?	?	x_4^3	$f_{\phi}(x_1^3, x_2^3, x_3^3 x_4^3)$
\vdots					\vdots

- Would require 2^D amortised variational distributions, one for each pattern of missingness \rightarrow inefficient!

Variational Gibbs Inference

- To make $f^t(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}})$ flexible: Specify it to be the marginal of a Markov chain with a learnable kernel $\kappa_{\phi}(\mathbf{x}_{\text{mis}}^{\tau+1} | \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}^{\tau})$.
- To tackle the 2^D scaling problem: Specify the kernel to be Gibbs (updates one dimension of \mathbf{x}_{mis} at a time):
$$\kappa_{\phi}(\mathbf{x}_{\text{mis}}^{\tau+1} | \mathbf{x}_{\text{mis}}^{\tau}, \mathbf{x}_{\text{obs}}) = \mathbb{E}_{\pi(j)} \left[q_{\phi_j}(x_j | \mathbf{x}_{\text{mis}^{\setminus j}}^{\tau}, \mathbf{x}_{\text{obs}}) \delta_{\mathbf{x}_{\text{mis}}^{\tau+1}}(\mathbf{x}_{\text{mis}^{\setminus j}}^{\tau+1}) \right]$$
- To make inference efficient use “persistent” chains: initialise the chains at the last state of the previous iteration.
- To make learning κ_{ϕ} efficient avoid computing the marginal of the Markov chain $f^t(\mathbf{x}_{\text{mis}} | \mathbf{x}_{\text{obs}})$: “cut” the chains just before the last transition and optimise over the last step of the chain.



Results: Variational autoencoder

