

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data

Vaidotas Šimkus Ben Rhodes Michael Gutmann

School of Informatics
The University of Edinburgh

November 2023



THE UNIVERSITY of EDINBURGH
informatics

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data

- General-purpose method for estimating statistical models from incomplete data.
- Journal of Machine Learning Research, 2023: jmlr.org/papers/v24/21-1373.html.
- Code: github.com/vsimkus/variational-gibbs-inference.
- Demo: nbviewer.org/github/vsimkus/variational-gibbs-inference/blob/main/notebooks/VGI_demo.ipynb.

1. Statistical models and the missing data issue
2. Some problems with direct estimation from incomplete data
3. Variational Gibbs Inference

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

Options:

1. Discard data-points with missing values

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

Options:

1. Discard data-points with missing values \rightarrow loss of information, not sustainable, bias \times

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

Options:

1. Discard data-points with missing values \rightarrow loss of information, not sustainable, bias \times
2. Impute-then-fit

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

Options:

1. Discard data-points with missing values \rightarrow loss of information, not sustainable, bias ✗
2. Impute-then-fit \rightarrow selecting appropriate imputation method, imputation incongeniality ✗

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

Options:

1. Discard data-points with missing values \rightarrow loss of information, not sustainable, bias ✗
2. Impute-then-fit \rightarrow selecting appropriate imputation method, imputation incongeniality ✗
3. Direct fitting by marginalising the missing variables \mathbf{x}_m

- Statistical models $p_{\theta}(\mathbf{x})$ of complex phenomena are a key component of many important tasks in machine learning.
- The models $p_{\theta}(\mathbf{x})$ are specified for fully-observed data \mathbf{x} ,
- And are typically fitted via maximum-likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^i), \quad \text{where } \mathbf{x}^i \in \mathcal{D}.$$

- Real-world data is often incomplete due to: non-response, sensor failure, occlusion, etc.
- What can we do?
- Denote \mathbf{x}_o and \mathbf{x}_m as the observed and missing elements of $\mathbf{x} = \mathbf{x}_o \cup \mathbf{x}_m$.

Options:

1. Discard data-points with missing values \rightarrow loss of information, not sustainable, bias ✗
2. Impute-then-fit \rightarrow selecting appropriate imputation method, imputation incongeniality ✗
3. Direct fitting by marginalising the missing variables \mathbf{x}_m ?

1. Statistical models and the missing data issue
2. Some problems with direct estimation from incomplete data
3. Variational Gibbs Inference

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?
- Expectation-maximisation (EM) (assuming ignorable missingness)

$$\log p_{\theta}(\mathbf{x}_o) = \log \int f(\mathbf{x}_m | \mathbf{x}_o) \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} d\mathbf{x}_m \geq \mathbb{E}_{f(\mathbf{x}_m | \mathbf{x}_o)} \left[\log \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} \right], \quad \text{"ELBO"}$$

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?
- Expectation-maximisation (EM) (assuming ignorable missingness)

$$\log p_{\theta}(\mathbf{x}_o) = \log \int f(\mathbf{x}_m | \mathbf{x}_o) \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} d\mathbf{x}_m \geq \mathbb{E}_{f(\mathbf{x}_m | \mathbf{x}_o)} \left[\log \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} \right], \quad \text{"ELBO"}$$

- **E-step:** Maximise w.r.t. $f(\mathbf{x}_m | \mathbf{x}_o^i)$ for $\forall \mathbf{x}_o^i \in \mathcal{D}$: $f(\mathbf{x}_m | \mathbf{x}_o^i) = p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)$.

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?
- Expectation-maximisation (EM) (assuming ignorable missingness)

$$\log p_{\theta}(\mathbf{x}_o) = \log \int f(\mathbf{x}_m | \mathbf{x}_o) \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} d\mathbf{x}_m \geq \mathbb{E}_{f(\mathbf{x}_m | \mathbf{x}_o)} \left[\log \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} \right], \quad \text{"ELBO"}$$

- **E-step:** Maximise w.r.t. $f(\mathbf{x}_m | \mathbf{x}_o^i)$ for $\forall \mathbf{x}_o^i \in \mathcal{D}$: $f(\mathbf{x}_m | \mathbf{x}_o^i) = p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)$.
- **M-step:** Maximise w.r.t. θ : $\theta^{t+1} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)} [\log p_{\theta}(\mathbf{x}_o^i, \mathbf{x}_m)]$

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?
- Expectation-maximisation (EM) (assuming ignorable missingness)

$$\log p_{\theta}(\mathbf{x}_o) = \log \int f(\mathbf{x}_m | \mathbf{x}_o) \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} d\mathbf{x}_m \geq \mathbb{E}_{f(\mathbf{x}_m | \mathbf{x}_o)} \left[\log \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} \right], \quad \text{"ELBO"}$$

- **E-step:** Maximise w.r.t. $f(\mathbf{x}_m | \mathbf{x}_o^i)$ for $\forall \mathbf{x}_o^i \in \mathcal{D}$: $f(\mathbf{x}_m | \mathbf{x}_o^i) = p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)$.
- **M-step:** Maximise w.r.t. θ : $\theta^{t+1} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)} [\log p_{\theta}(\mathbf{x}_o^i, \mathbf{x}_m)]$
- Monte Carlo EM: Approximate the expectation using Monte Carlo average.

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?
- Expectation-maximisation (EM) (assuming ignorable missingness)

$$\log p_{\theta}(\mathbf{x}_o) = \log \int f(\mathbf{x}_m | \mathbf{x}_o) \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} d\mathbf{x}_m \geq \mathbb{E}_{f(\mathbf{x}_m | \mathbf{x}_o)} \left[\log \frac{p_{\theta}(\mathbf{x}_o, \mathbf{x}_m)}{f(\mathbf{x}_m | \mathbf{x}_o)} \right], \quad \text{"ELBO"}$$

- **E-step:** Maximise w.r.t. $f(\mathbf{x}_m | \mathbf{x}_o^i)$ for $\forall \mathbf{x}_o^i \in \mathcal{D}$: $f(\mathbf{x}_m | \mathbf{x}_o^i) = p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)$.
- **M-step:** Maximise w.r.t. θ : $\theta^{t+1} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)} [\log p_{\theta}(\mathbf{x}_o^i, \mathbf{x}_m)]$
- Monte Carlo EM: Approximate the expectation using Monte Carlo average.
- Then, M-step corresponds to fitting $p_{\theta}(\mathbf{x})$ with completed data.

- Marginalising the missing variables $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally not tractable.
- What can we do if simplifying assumptions cannot be inserted?
- Expectation-maximisation (EM) (assuming ignorable missingness)

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

- **E-step:** Maximise w.r.t. $f(\mathbf{x}_m | \mathbf{x}_o^i)$ for $\forall \mathbf{x}_o^i \in \mathcal{D}$: $f(\mathbf{x}_m | \mathbf{x}_o^i) = p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)$.
- **M-step:** Maximise w.r.t. θ : $\theta^{t+1} = \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o^i)} [\log p_{\theta}(\mathbf{x}_o^i, \mathbf{x}_m)]$
- Monte Carlo EM: Approximate the expectation using Monte Carlo average.
- Then, M-step corresponds to fitting $p_{\theta}(\mathbf{x})$ with completed data.

Issue with Monte Carlo EM

- Conditional sampling of $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$ is generally intractable or inefficient.

Variational inference (VI)

- $\forall \mathbf{x}_o \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ .
- M-step: Sample $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ to approximate the expectation.

Variational approximation to $p_{\theta^t}(\mathbf{x}_m \mid \mathbf{x}_o)$



Variational inference (VI)

- $\forall \mathbf{x}_o \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ .
- M-step: Sample $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ to approximate the expectation.

Advantages of VI

- Choice of $\mathcal{Q}(\phi)$ is in our control.
- Turns inference to optimisation.
- Can fit using SGD.
- Efficient if $|\mathcal{D}|$ is small.

Variational approximation to $p_{\theta^t}(\mathbf{x}_m \mid \mathbf{x}_o)$



Variational inference (VI)

- $\forall \mathbf{x}_o \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ .
- M-step: Sample $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ to approximate the expectation.

Advantages of VI

- Choice of $\mathcal{Q}(\phi)$ is in our control.
- Turns inference to optimisation.
- Can fit using SGD.
- Efficient if $|\mathcal{D}|$ is small.

Disadvantages of VI

- Is inefficient if $|\mathcal{D}|$ is large.

Variational approximation to $p_{\theta^t}(\mathbf{x}_m \mid \mathbf{x}_o)$



Variational inference (VI)

- $\forall \mathbf{x}_o \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ .
- M-step: Sample $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ to approximate the expectation.

Amortised VI

- Parametrise $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ with a *single* neural network $\text{NN}_{\phi}(\mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$.

Advantages of VI

- Choice of $\mathcal{Q}(\phi)$ is in our control.
- Turns inference to optimisation.
- Can fit using SGD.
- Efficient if $|\mathcal{D}|$ is small.

Disadvantages of VI

- Is inefficient if $|\mathcal{D}|$ is large.

Variational approximation to $p_{\theta^t}(\mathbf{x}_m \mid \mathbf{x}_o)$



Variational inference (VI)

- $\forall \mathbf{x}_o \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ .
- M-step: Sample $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ to approximate the expectation.

Amortised VI

- Parametrise $f_{\phi}(\mathbf{x}_m \mid \mathbf{x}_o)$ with a *single* neural network $\text{NN}_{\phi}(\mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$.

Advantages of VI

- Choice of $\mathcal{Q}(\phi)$ is in our control.
- Turns inference to optimisation.
- Can fit using SGD.
- Efficient if $|\mathcal{D}|$ is small.

Disadvantages of VI

- Is inefficient if $|\mathcal{D}|$ is large.

Advantages of amortised VI

- Efficient for large $|\mathcal{D}|$.

Variational approximation to $p_{\theta^t}(\mathbf{x}_m | \mathbf{x}_o)$



Variational inference (VI)

- $\forall \mathbf{x}_o \in \mathcal{D}$ specify a $f_{\phi}(\mathbf{x}_m | \mathbf{x}_o) \in \mathcal{Q}(\phi)$.
- E-step: Maximise the ELBO w.r.t. ϕ .
- M-step: Sample $f_{\phi}(\mathbf{x}_m | \mathbf{x}_o)$ to approximate the expectation.

Amortised VI

- Parametrise $f_{\phi}(\mathbf{x}_m | \mathbf{x}_o)$ with a *single* neural network $\text{NN}_{\phi}(\mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$.

	d_1	d_2	d_3	d_4	$f_{\phi}(\mathbf{x}_m^i \mathbf{x}_o^i)$
\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				\vdots

Advantages of VI

- Choice of $\mathcal{Q}(\phi)$ is in our control.
- Turns inference to optimisation.
- Can fit using SGD.
- Efficient if $|\mathcal{D}|$ is small.

Disadvantages of VI

- Is inefficient if $|\mathcal{D}|$ is large.

Advantages of amortised VI

- Efficient for large $|\mathcal{D}|$.

Disadvantages of amortised VI

- Need one $f_{\phi}(\mathbf{x}_m | \mathbf{x}_o)$ for each pattern of missingness (2^M in total).

1. Statistical models and the missing data issue
2. Some problems with direct estimation from incomplete data
3. Variational Gibbs Inference

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

- General-purpose method for estimating $p_{\theta}(x)$ from incomplete data.
- Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

- General-purpose method for estimating $p_{\theta}(x)$ from incomplete data.
 - Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.
1. Core idea: Turn the 2^M conditional distribution problem into M conditional distributions.

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

- General-purpose method for estimating $p_{\theta}(x)$ from incomplete data.
 - Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.
1. Core idea: Turn the 2^M conditional distribution problem into M conditional distributions.
 2. To make $f_{\phi}^t(x_m | x_o)$ flexible:

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

- General-purpose method for estimating $p_{\theta}(x)$ from incomplete data.
 - Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.
1. Core idea: Turn the 2^M conditional distribution problem into M conditional distributions.
 2. To make $f_{\phi}^t(x_m | x_o)$ flexible:
 - Specify it to be the marginal of a Markov chain with a *learnable* kernel $\kappa_{\phi}(x_m^{\tau+1} | x_o, x_m^{\tau})$.

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

- General-purpose method for estimating $p_{\theta}(\mathbf{x})$ from incomplete data.
 - Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.
1. Core idea: Turn the 2^M conditional distribution problem into M conditional distributions.
 2. To make $f_{\phi}^t(\mathbf{x}_m | \mathbf{x}_o)$ flexible:
 - Specify it to be the marginal of a Markov chain with a *learnable* kernel $\kappa_{\phi}(\mathbf{x}_m^{\tau+1} | \mathbf{x}_o, \mathbf{x}_m^{\tau})$.
 3. To address the 2^M pattern problem:

Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

- General-purpose method for estimating $p_{\theta}(\mathbf{x})$ from incomplete data.
 - Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.
1. Core idea: Turn the 2^M conditional distribution problem into M conditional distributions.
 2. To make $f_{\phi}^t(\mathbf{x}_m | \mathbf{x}_o)$ flexible:
 - Specify it to be the marginal of a Markov chain with a *learnable* kernel $\kappa_{\phi}(\mathbf{x}_m^{\tau+1} | \mathbf{x}_o, \mathbf{x}_m^{\tau})$.
 3. To address the 2^M pattern problem:
 - We specify the kernel to be Gibbs (updates one dimension of \mathbf{x}_m at a time):

$$\kappa_{\phi}(\mathbf{x}_m^{\tau+1} | \mathbf{x}_m^{\tau}, \mathbf{x}_o) = \mathbb{E}_{\pi(j | \text{idx}(\mathbf{m}))} \left[q_{\phi_j}(x_j | \mathbf{x}_{m \setminus j}^{\tau}, \mathbf{x}_o) \delta(\mathbf{x}_{m \setminus j}^{\tau+1} - \mathbf{x}_{m \setminus j}^{\tau}) \right],$$

where $\pi(j | \text{idx}(\mathbf{m}))$ is the selection probability for the j -th dimension of a Gibbs sampler.

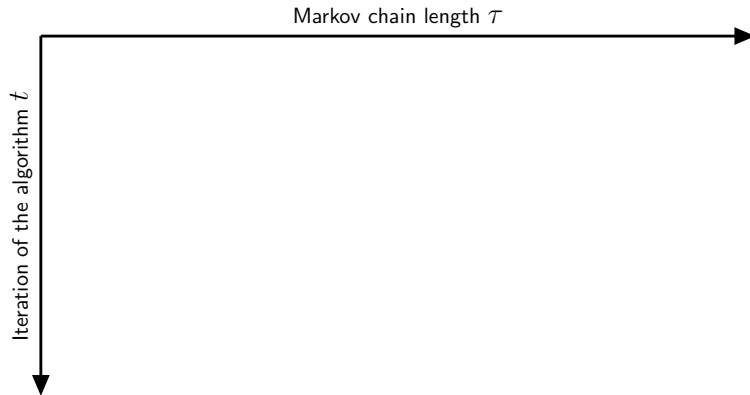
Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data, JMLR, 2023

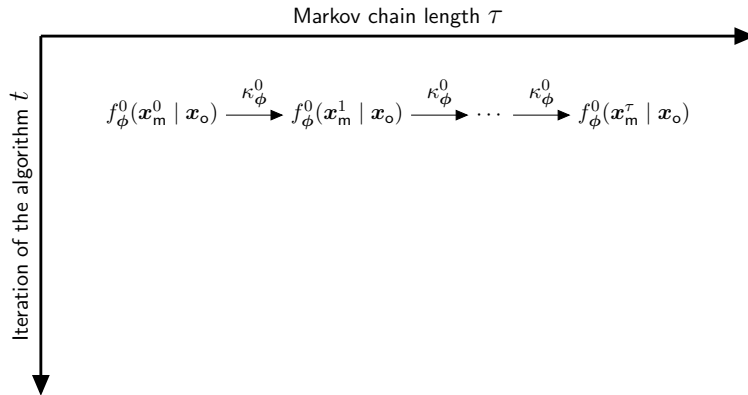
- General-purpose method for estimating $p_{\theta}(\mathbf{x})$ from incomplete data.
 - Efficient for large $|\mathcal{D}|$ and mitigates the need for 2^M conditional distributions.
1. Core idea: Turn the 2^M conditional distribution problem into M conditional distributions.
 2. To make $f_{\phi}^t(\mathbf{x}_m | \mathbf{x}_o)$ flexible:
 - Specify it to be the marginal of a Markov chain with a *learnable* kernel $\kappa_{\phi}(\mathbf{x}_m^{\tau+1} | \mathbf{x}_o, \mathbf{x}_m^{\tau})$.
 3. To address the 2^M pattern problem:
 - We specify the kernel to be Gibbs (updates one dimension of \mathbf{x}_m at a time):

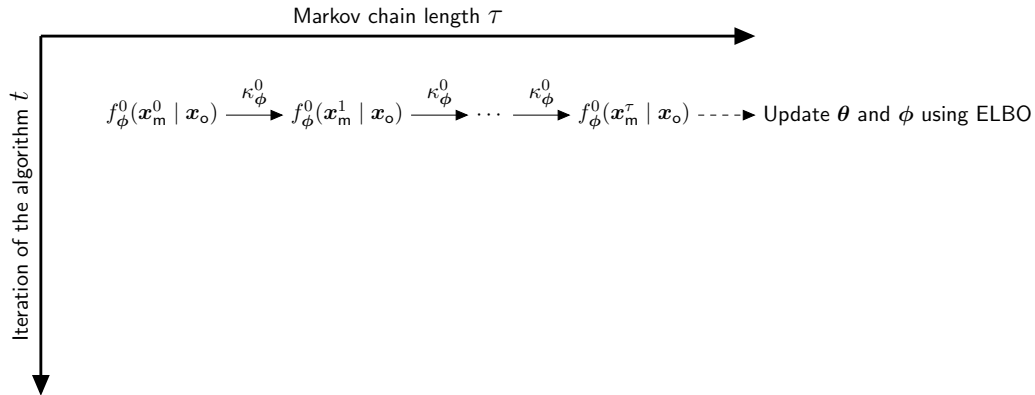
$$\kappa_{\phi}(\mathbf{x}_m^{\tau+1} | \mathbf{x}_m^{\tau}, \mathbf{x}_o) = \mathbb{E}_{\pi(j | \text{idx}(\mathbf{m}))} \left[q_{\phi_j}(x_j | \mathbf{x}_{m \setminus j}^{\tau}, \mathbf{x}_o) \delta(\mathbf{x}_{m \setminus j}^{\tau+1} - \mathbf{x}_{m \setminus j}^{\tau}) \right],$$

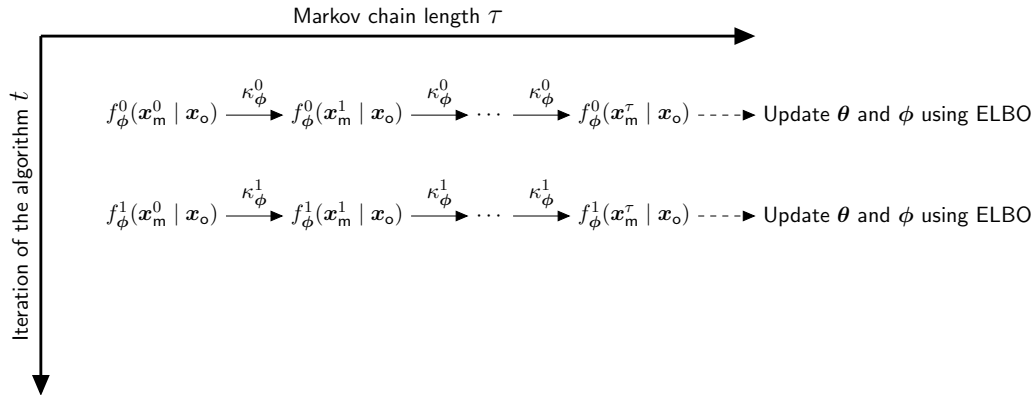
where $\pi(j | \text{idx}(\mathbf{m}))$ is the selection probability for the j -th dimension of a Gibbs sampler.

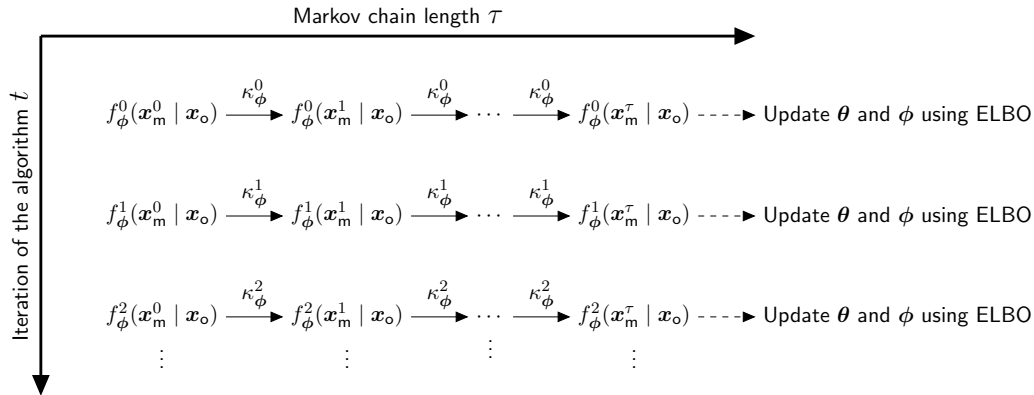
- Hence we have to learn only M variational Gibbs conditional $q_{\phi_j}(x_j | \mathbf{x}_{m \setminus j}, \mathbf{x}_o)$.



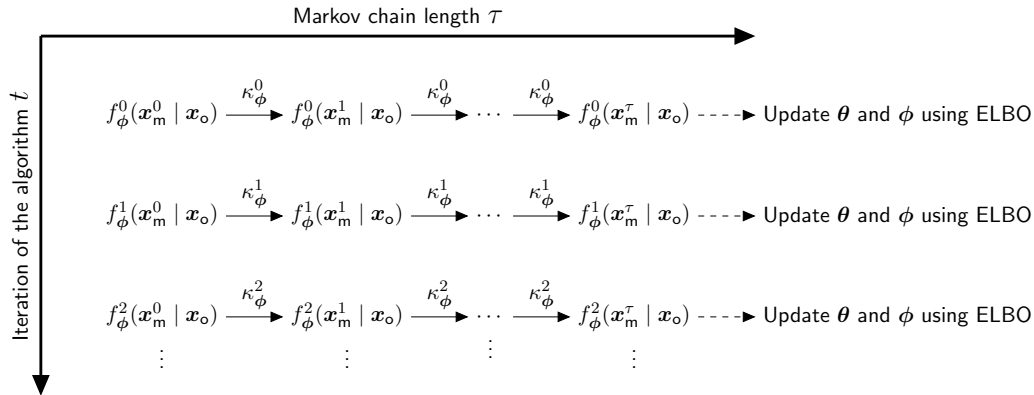




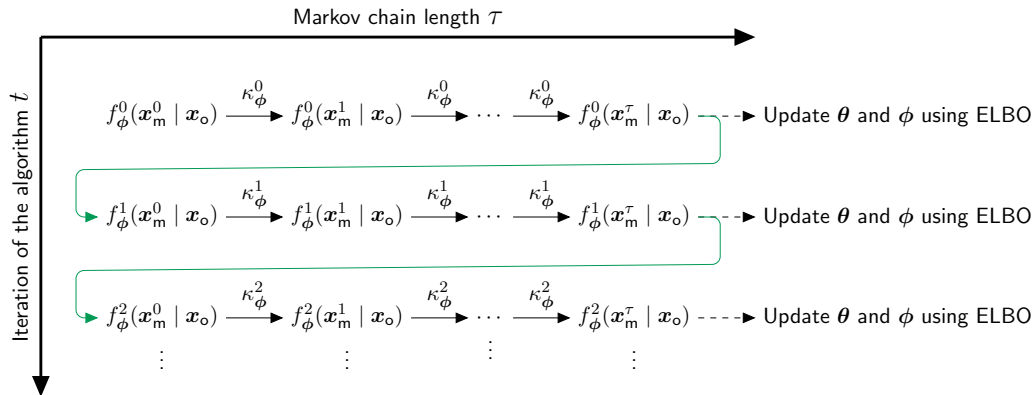




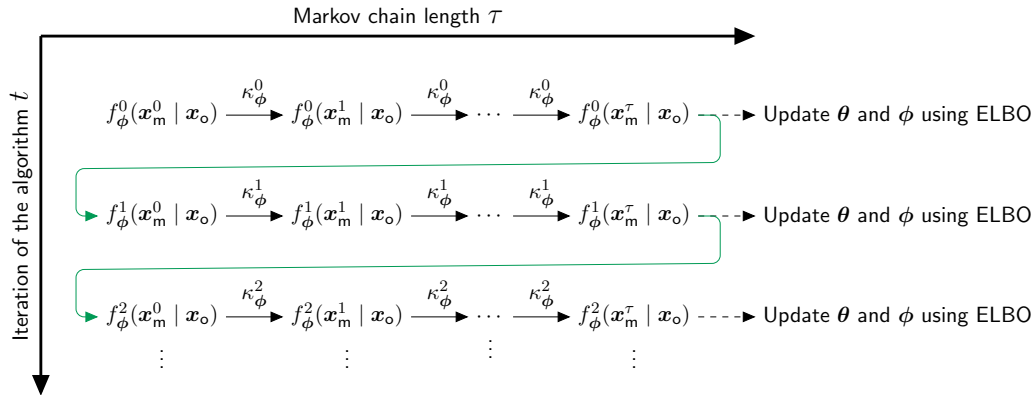
- Sampling long Markov chains at each iteration t of the algorithm is costly.



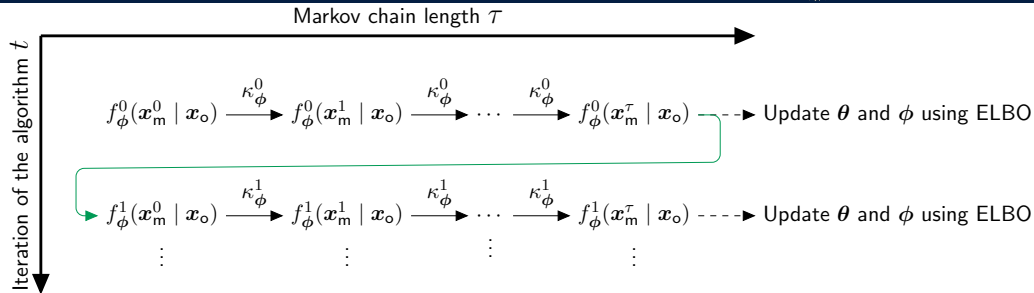
- Sampling long Markov chains at each iteration t of the algorithm is costly.
- Use “persistent” chains: initialise the chains at the last state of the previous iteration.

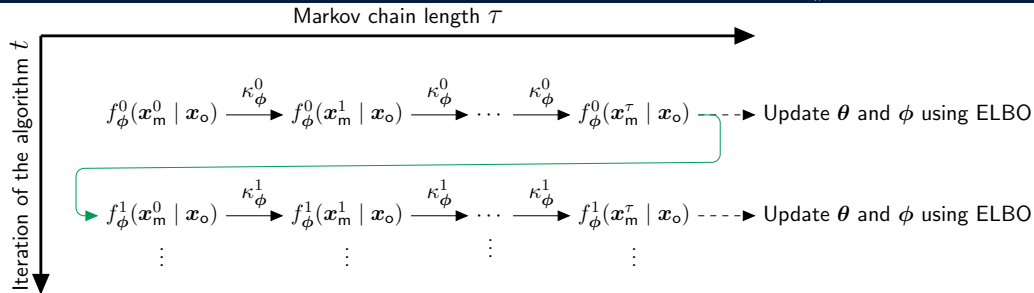


- Sampling long Markov chains at each iteration t of the algorithm is costly.
- Use “persistent” chains: initialise the chains at the last state of the previous iteration.
- Can now use short chains, that is using small τ , at every iteration t .



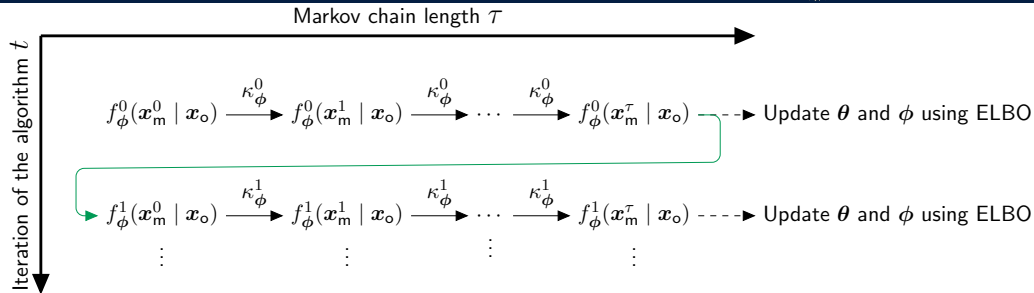
Variational Gibbs Inference: “Cutting” chains





- Computing the marginal density $f_{\phi}^t(\mathbf{x}_m^{\tau} | \mathbf{x}_o)$ of a Markov chain remains intractable:

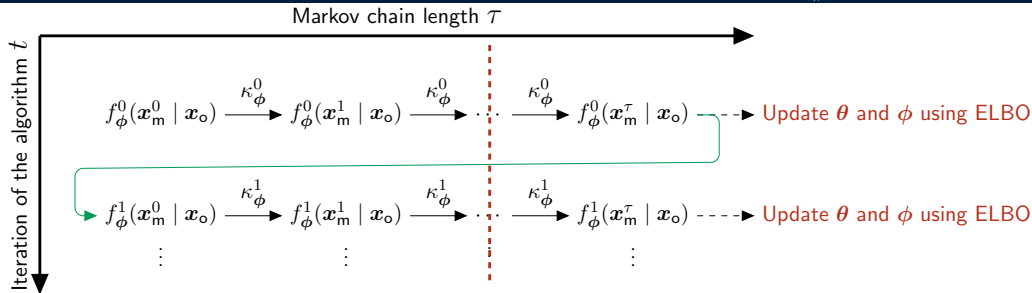
$$f_{\phi}^t(\mathbf{x}_m^{\tau} | \mathbf{x}_o) = \int f_{\phi}^t(\mathbf{x}_m^0 | \mathbf{x}_o) \prod_{h=0}^{\tau-1} \kappa_{\phi}(\mathbf{x}_m^{h+1} | \mathbf{x}_o, \mathbf{x}_m^h) d\mathbf{x}_m^0 \dots d\mathbf{x}_m^{\tau-1}.$$



- Computing the marginal density $f_{\phi}^t(\mathbf{x}_m^{\tau} | \mathbf{x}_o)$ of a Markov chain remains intractable:

$$f_{\phi}^t(\mathbf{x}_m^{\tau} | \mathbf{x}_o) = \int f_{\phi}^t(\mathbf{x}_m^0 | \mathbf{x}_o) \prod_{h=0}^{\tau-1} \kappa_{\phi}(\mathbf{x}_m^{h+1} | \mathbf{x}_o, \mathbf{x}_m^h) d\mathbf{x}_m^0 \dots d\mathbf{x}_m^{\tau-1}.$$

- So how can we optimise the parameters ϕ of the kernel κ_{ϕ} ?



- Computing the marginal density $f_{\phi}^t(x_m^{\tau} | x_o)$ of a Markov chain remains intractable:

$$f_{\phi}^t(x_m^{\tau} | x_o) = \int f_{\phi}^t(x_m^0 | x_o) \prod_{h=0}^{\tau-1} \kappa_{\phi}(x_m^{h+1} | x_o, x_m^h) dx_m^0 \dots dx_m^{\tau-1}.$$

- So how can we optimise the parameters ϕ of the kernel κ_{ϕ} ?
- Instead of optimising ϕ over the full length of the Markov chains, we “cut” the chains just before the last transition and optimise over the last step of the chain.

- Objective for learning θ and ϕ :

$$\log p_{\theta}(\mathbf{x}_o) \geq \mathbb{E}_{\pi(j|\text{idx}(\mathbf{m})) f^{t-1}(\mathbf{x}_{\mathbf{m} \setminus j} | \mathbf{x}_o) q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)} \left[\log \frac{p_{\theta}(x_j, \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)}{q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)} \right] + \text{Const.}$$

- Objective for learning θ and ϕ :

$$\log p_{\theta}(\mathbf{x}_o) \geq \mathbb{E}_{\pi(j|\text{idx}(\mathbf{m})) f^{t-1}(\mathbf{x}_{\mathbf{m} \setminus j} | \mathbf{x}_o) q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)} \left[\log \frac{p_{\theta}(x_j, \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)}{q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)} \right] + \text{Const.}$$

- We only need samples from penultimate step of the Markov chain f^{t-1} .

- Objective for learning θ and ϕ :

$$\log p_{\theta}(\mathbf{x}_o) \geq \mathbb{E}_{\pi(j|\text{idx}(\mathbf{m})) f^{t-1}(\mathbf{x}_{\mathbf{m}\setminus j}|\mathbf{x}_o) q_{\phi_j}(x_j|\mathbf{x}_{\mathbf{m}\setminus j}, \mathbf{x}_o)} \left[\log \frac{p_{\theta}(x_j, \mathbf{x}_{\mathbf{m}\setminus j}, \mathbf{x}_o)}{q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m}\setminus j}, \mathbf{x}_o)} \right] + \text{Const.}$$

- We only need samples from penultimate step of the Markov chain f^{t-1} .
- Can optimise w.r.t. θ and ϕ using stochastic gradient ascent.

- Objective for learning θ and ϕ :

$$\log p_{\theta}(\mathbf{x}_o) \geq \mathbb{E}_{\pi(j|\text{idx}(\mathbf{m}))} f^{t-1}(\mathbf{x}_{\mathbf{m} \setminus j} | \mathbf{x}_o) q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o) \left[\log \frac{p_{\theta}(x_j, \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)}{q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)} \right] + \text{Const.}$$

- We only need samples from penultimate step of the Markov chain f^{t-1} .
- Can optimise w.r.t. θ and ϕ using stochastic gradient ascent.
- Maximising the above w.r.t. ϕ corresponds to minimising the KL divergence:

$$\mathbb{E}_{\pi(j|\text{idx}(\mathbf{m}))} f^{t-1}(\mathbf{x}_{\mathbf{m} \setminus j} | \mathbf{x}_o) \left[D_{\text{KL}}(q_{\phi_j}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o) || p_{\theta}(x_j | \mathbf{x}_{\mathbf{m} \setminus j}, \mathbf{x}_o)) \right]$$

- Objective for learning θ and ϕ :

$$\log p_{\theta}(\mathbf{x}_o) \geq \mathbb{E}_{\pi(j|\text{idx}(\mathbf{m}))} f^{t-1}(\mathbf{x}_{m \setminus j} | \mathbf{x}_o) q_{\phi_j}(x_j | \mathbf{x}_{m \setminus j}, \mathbf{x}_o) \left[\log \frac{p_{\theta}(x_j, \mathbf{x}_{m \setminus j}, \mathbf{x}_o)}{q_{\phi_j}(x_j | \mathbf{x}_{m \setminus j}, \mathbf{x}_o)} \right] + \text{Const.}$$

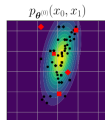
- We only need samples from penultimate step of the Markov chain f^{t-1} .
- Can optimise w.r.t. θ and ϕ using stochastic gradient ascent.
- Maximising the above w.r.t. ϕ corresponds to minimising the KL divergence:

$$\mathbb{E}_{\pi(j|\text{idx}(\mathbf{m}))} f^{t-1}(\mathbf{x}_{m \setminus j} | \mathbf{x}_o) \left[D_{\text{KL}}(q_{\phi_j}(x_j | \mathbf{x}_{m \setminus j}, \mathbf{x}_o) || p_{\theta}(x_j | \mathbf{x}_{m \setminus j}, \mathbf{x}_o)) \right]$$

- The fitted κ_{ϕ} approximates the Gibbs kernel with the stationary distribution $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$.

Algorithm 1 Variational Gibbs inference

1: **Create** K -times imputed data \mathcal{D}_K using f_0



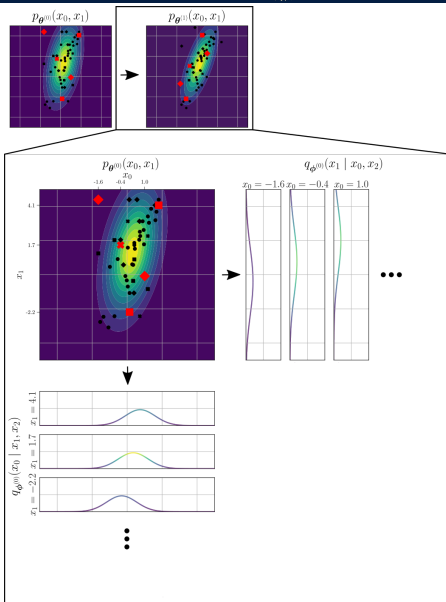
Variational Gibbs Inference: Algorithm



Algorithm 1 Variational Gibbs inference

- 1: **Create** K -times imputed data \mathcal{D}_K using f_0
- 2: **for** t in $[1, \text{max_epochs}]$ **do**
- 3: **Sample** mini-batch \mathcal{B}_K from \mathcal{D}_K

7: **end for**

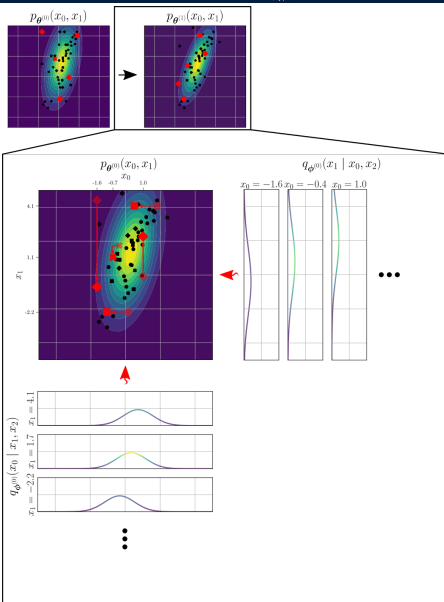


Variational Gibbs Inference: Algorithm



Algorithm 1 Variational Gibbs inference

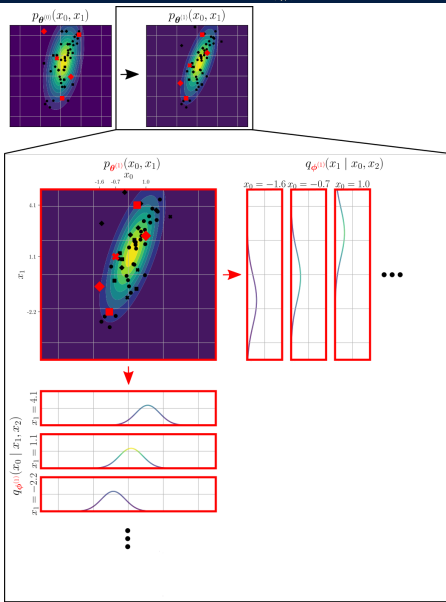
- 1: **Create** K -times imputed data \mathcal{D}_K using f_0
- 2: **for** t in $[1, \text{max_epochs}]$ **do**
- 3: **Sample** mini-batch \mathcal{B}_K from \mathcal{D}_K
- 4: **Update** the imputations in \mathcal{B}_K :
 $\bar{x}_m^{(i,k)} \sim \text{Gibbs}_\tau(x_o^i, \kappa_\phi; x_m^{(i,k)}), \forall x_m^{(i,k)} \in \mathcal{B}_K$
- 5: **Persist** the imputations in \mathcal{B}_K to \mathcal{D}_K
- 7: **end for**



Variational Gibbs Inference: Algorithm

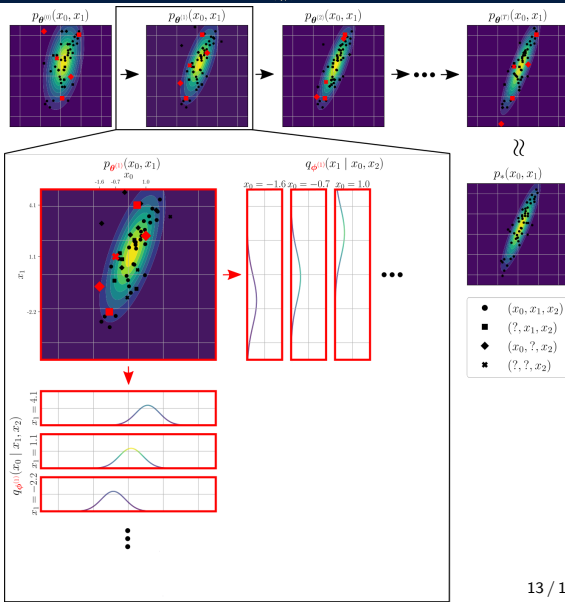
Algorithm 1 Variational Gibbs inference

- 1: **Create** K -times imputed data \mathcal{D}_K using f_0
- 2: **for** t in $[1, \text{max_epochs}]$ **do**
- 3: **Sample** mini-batch \mathcal{B}_K from \mathcal{D}_K
- 4: **Update** the imputations in \mathcal{B}_K :
 $\bar{x}_m^{(i,k)} \sim \text{Gibbs}_\tau(x_o^i, \kappa_\phi; x_m^{(i,k)}), \forall x_m^{(i,k)} \in \mathcal{B}_K$
- 5: **Persist** the imputations in \mathcal{B}_K to \mathcal{D}_K
- 6: **Update** θ and ϕ with SGA.
- 7: **end for**



Algorithm 1 Variational Gibbs inference

- 1: **Create** K -times imputed data \mathcal{D}_K using f_0
- 2: **for** t in $[1, \text{max_epochs}]$ **do**
- 3: **Sample** mini-batch \mathcal{B}_K from \mathcal{D}_K
- 4: **Update** the imputations in \mathcal{B}_K :
 $\bar{x}_m^{(i,k)} \sim \text{Gibbs}_\tau(x_o^i, \kappa_\phi; x_m^{(i,k)}), \forall x_m^{(i,k)} \in \mathcal{B}_K$
- 5: **Persist** the imputations in \mathcal{B}_K to \mathcal{D}_K
- 6: **Update** θ and ϕ with SGA.
- 7: **end for**



- Direct fitting by marginalising the missing variables x_m ?

- Direct fitting by (approximately) marginalising the missing variables x_m ✓

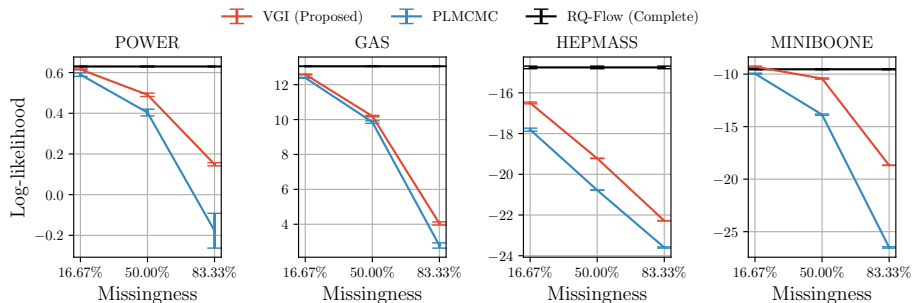
- Direct fitting by (approximately) marginalising the missing variables x_m ✓
- General-purpose method for estimating $p_\theta(x)$ from incomplete data.

- Direct fitting by (approximately) marginalising the missing variables x_m ✓
- General-purpose method for estimating $p_\theta(x)$ from incomplete data.
- Mitigated the need for 2^M conditional distributions to just M by representing the variational distribution via a learnable Gibbs kernel.

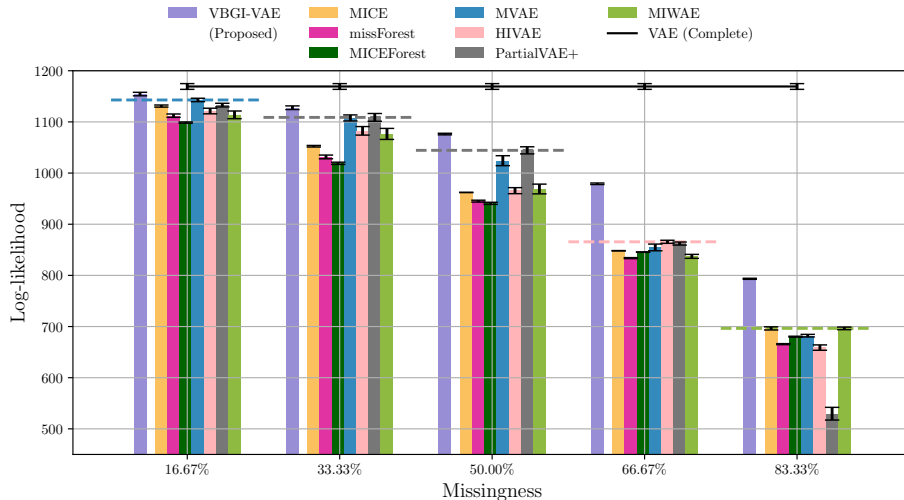
- Direct fitting by (approximately) marginalising the missing variables x_m ✓
- General-purpose method for estimating $p_\theta(x)$ from incomplete data.
- Mitigated the need for 2^M conditional distributions to just M by representing the variational distribution via a learnable Gibbs kernel.
- Used “persistent” chains to efficiently sample imputations using the learnt Gibbs kernel.

- Direct fitting by (approximately) marginalising the missing variables x_m ✓
- General-purpose method for estimating $p_\theta(x)$ from incomplete data.
- Mitigated the need for 2^M conditional distributions to just M by representing the variational distribution via a learnable Gibbs kernel.
- Used “persistent” chains to efficiently sample imputations using the learnt Gibbs kernel.
- “Cut” the Markov chains to make optimisation of ϕ efficient.

Variational Gibbs Inference: Results (Flows)



Variational Gibbs Inference: Results (VAE)



- Statistical models and the missing data issue.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.
 - Marginalisation $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally intractable.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.
 - Marginalisation $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally intractable.
 - EM algorithm requires sampling conditionals $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$, which is expensive.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.
 - Marginalisation $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally intractable.
 - EM algorithm requires sampling conditionals $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$, which is expensive.
 - Standard amortised VI requires 2^M variational distributions, which is inefficient.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.
 - Marginalisation $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally intractable.
 - EM algorithm requires sampling conditionals $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$, which is expensive.
 - Standard amortised VI requires 2^M variational distributions, which is inefficient.
- Variational Gibbs Inference.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.
 - Marginalisation $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally intractable.
 - EM algorithm requires sampling conditionals $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$, which is expensive.
 - Standard amortised VI requires 2^M variational distributions, which is inefficient.
- Variational Gibbs Inference.
 - General purpose method for model estimation from incomplete data.

- Statistical models and the missing data issue.
 - Modern models, such as normalising flows and VAEs, are very flexible.
 - But, they are formulated for complete data.
- Some problems with direct estimation from incomplete data.
 - Marginalisation $\int p_{\theta}(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{x}_m$ is generally intractable.
 - EM algorithm requires sampling conditionals $p_{\theta}(\mathbf{x}_m | \mathbf{x}_o)$ for $\forall \mathbf{x}_o \in \mathcal{D}$, which is expensive.
 - Standard amortised VI requires 2^M variational distributions, which is inefficient.
- Variational Gibbs Inference.
 - General purpose method for model estimation from incomplete data.
 - Achieves good performance on normalising flow and VAE estimation, compared to other methods.

Thank you for listening.
Questions?



Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum Likelihood from Incomplete Data Via the EM Algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22. (Cited on slide 6)



Simkus, V., Rhodes, B., and Gutmann, M. U. (2023). Variational Gibbs Inference for Statistical Model Estimation from Incomplete Data. *Journal of Machine Learning Research*, 24(196):1–72. (Cited on slide 2, 9)



Tieleman, T. (2008). Training restricted Boltzmann machines using approximations to the likelihood gradient. In *International Conference on Machine Learning (ICML)*, pages 1064–1071. (Cited on slide 10)



Wei, G. C. G. and Tanner, M. A. (1990). A Monte Carlo Implementation of the EM Algorithm and the Poor Man's Data Augmentation Algorithms. *Journal of the American Statistical Association*, 85(411):699–704. (Cited on slide 6)