



Enhanced Variational Autoencoder Estimation from Incomplete Data using Mixture Variational Families

Vaidotas Šimkus Michael U. Gutmann

School of Informatics, The University of Edinburgh



Summary

- We show why estimating VAEs from incomplete data is harder compared to the fully-observed case.
- As a result, fitting VAEs from incomplete data requires more flexible variational families, compared to the fully-observed case. But, more flexible families may lack the useful inductive biases!
- We propose two approaches based on variational mixture distributions that can improve VAE estimation from incomplete data while reusing the families from the fully-observed case.

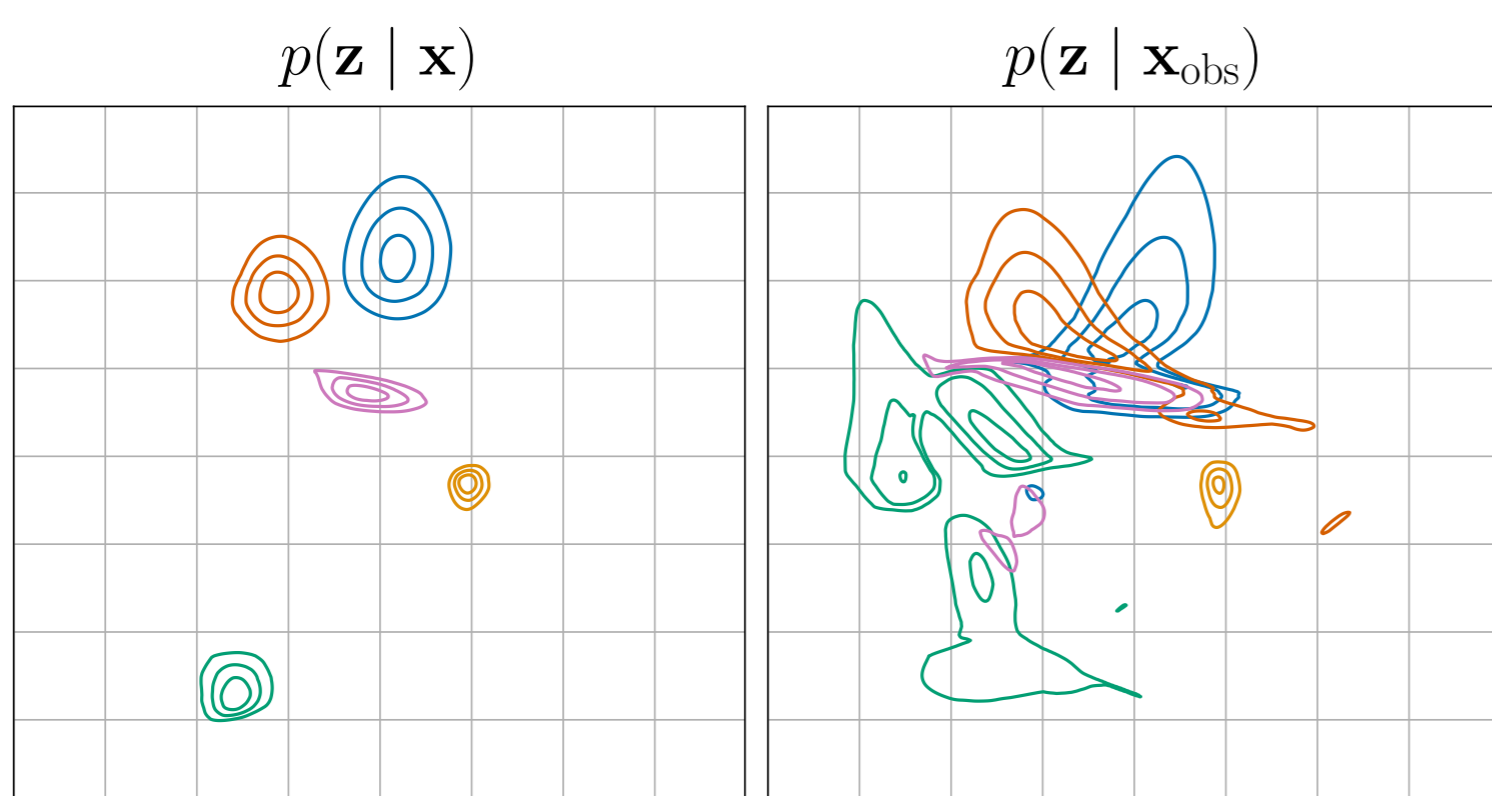
Why is VAE estimation from *incomplete* data hard(er)?

- A VAE model is typically trained by maximising the ELBO:

$$\log p_{\theta}(\mathbf{y}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{y})} \left[\log \frac{p_{\theta}(\mathbf{y}|\mathbf{z})p_{\theta}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{y})} \right] \quad (1)$$

$$= \log p_{\theta}(\mathbf{y}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{y}) \parallel p_{\theta}(\mathbf{z}|\mathbf{y}))$$

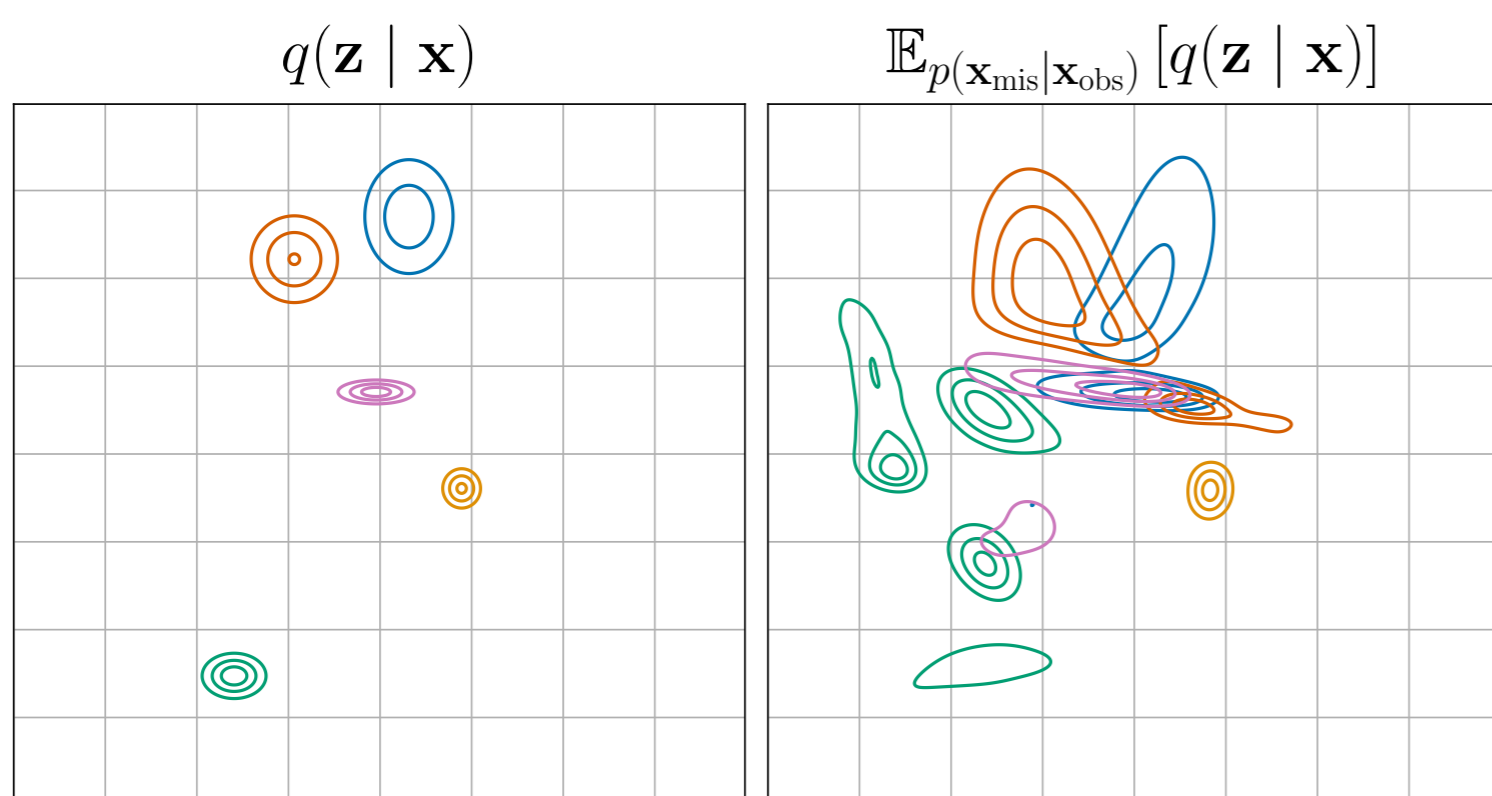
- When data is fully-observed $\mathbf{y} \stackrel{!}{=} \mathbf{x}$ and when incomplete $\mathbf{y} \stackrel{!}{=} \mathbf{x}_{\text{obs}}$, where $\mathbf{x}_{\text{obs}} = \mathbf{x}[\mathbf{m}]$ for a missingness pattern $\mathbf{m} \in \{0, 1\}^D$.
- The objective for fully-observed and incomplete data is similar, so why is it typically harder to train a VAE from incomplete data?
- **Effective training needs to well-minimise the KL term:**
 - When data is fully-observed: $D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}))$,
 - When data is incomplete: $D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}_{\text{obs}}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}_{\text{obs}}))$.



- Need more flexible variational distributions q with incomplete data!
- Potential loss of inductive biases when going from fully-observed to incomplete case, due to modifying the family of q .

Can we re-use variational family from complete case?

- Re-use complete-data approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$ for incomplete-data: $p_{\theta}(\mathbf{z}|\mathbf{x}_{\text{obs}}) = \mathbb{E}_{p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})}[p_{\theta}(\mathbf{z}|\mathbf{x})] \approx \mathbb{E}_{p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})}[q_{\phi}(\mathbf{z}|\mathbf{x})]$.



- Enables re-using the variational family from the fully-observed case.
- As a result we can save some of the useful inductive biases.
- But, the missing-variable distribution $p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})$ in $\mathbb{E}_{p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})}[q_{\phi}(\mathbf{z}|\mathbf{x})]$ is intractable!

Opt. 1: Use finite-mixture variational distribution

- Make the approach more tractable by using finite-mixtures.
- Approximate $\mathbb{E}_{p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})}[q_{\phi}(\mathbf{z}|\mathbf{x})]$ with a finite-mixture:

$$q_{\phi}(\mathbf{z}|\mathbf{x}_{\text{obs}}) = \sum_{k=1}^K q_{\phi}(k|\mathbf{x}_{\text{obs}})q_{\phi}^k(\mathbf{z}|\mathbf{x}_{\text{obs}}).$$
- Where the component distributions $q_{\phi}^k(\mathbf{z}|\mathbf{x}_{\text{obs}})$ re-use the complete-data variational family, and hence they may roughly take the role of $q_{\phi}(\mathbf{z}|\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}})$ for some \mathbf{x}_{mis} .
- Optimise the model using eq. (1) or importance-weighted ELBO.
- Several options available that trade-off variance versus computational cost (see paper).

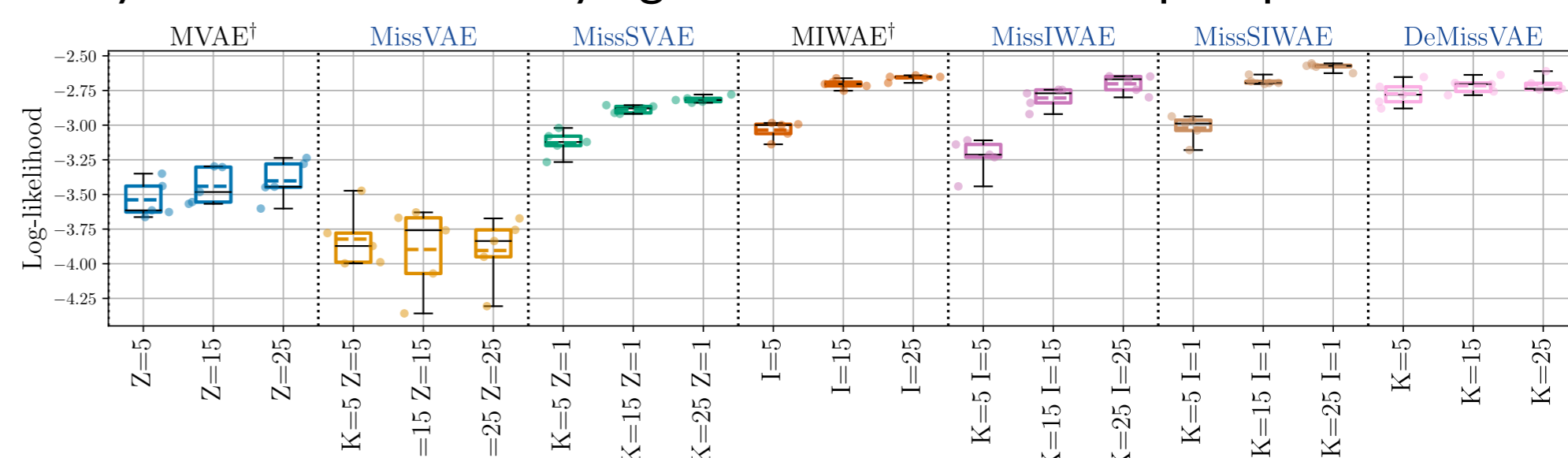
Opt. 2: Use imputation-mixture variational distribution

- Decompose the model estimation task into two: data imputation and model estimation (similar to Monte Carlo EM).
- Approximate $\mathbb{E}_{p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})}[q_{\phi}(\mathbf{z}|\mathbf{x})]$ with an imputation-mixture:

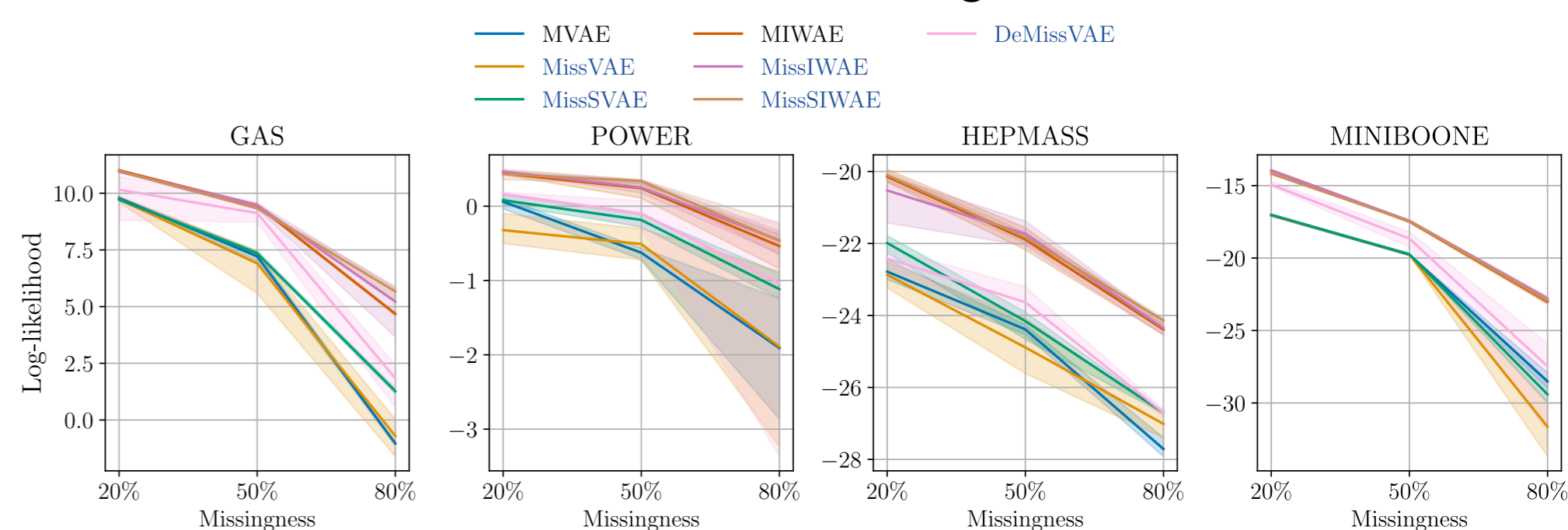
$$q_{\phi, f^t}(\mathbf{z}|\mathbf{x}_{\text{obs}}) = \mathbb{E}_{f^t(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})} \left[q_{\phi}(\mathbf{z}|\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}) \right]$$
- The imputation distribution $f^t(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})$ is represented using *approximate* samples from $p_{\theta}(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})$, using cheap sampling methods such as pseudo-Gibbs, Metropolis-within-Gibbs, or others.
- And, q_{ϕ} is fitted such that $q_{\phi}(\mathbf{z}|\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}) \approx p_{\theta}(\mathbf{z}|\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}})$.
- We thus have a flexible approximation $q_{\phi, f^t}(\mathbf{z}|\mathbf{x}_{\text{obs}}) \approx p_{\theta}(\mathbf{z}|\mathbf{x}_{\text{obs}})$, while retaining the variational family from the complete-data case.
- Method uses separate objectives for θ and ϕ to minimise the potential bias due to approximation errors of the imputation distribution $f^t(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}})$.

Results (proposed methods highlighted in blue)

Toy MoG data with varying number of latent samples per iteration



UCI data with different missingness fractions



MNIST and Omniglot missing 2/4 quadrants at random

