

Missing Data Handling in Amortized Bayesian Inference via Invertible Neural Networks

Zijian Wang

Final Bachelor Thesis Presentation
September 5, 2022

Outline

- 1 Amortized Likelihood-free Bayesian Inference with BayesFlow
- 2 Missing Data Handling
- 3 Experiments and Results
- 4 Summary and Outlook

Problem Setting

- Infer underlying parameters θ from dataset $x_{1:N} = \{x_i\}_{i=1}^N$
- **Bayesian inference** aims at assessing the posterior distribution:

$$p(\theta | x_{1:N}) \propto p(x_{1:N} | \theta) p(\theta)$$

- **Likelihood-free** methods only make use of simulations from the forward model:

$$x_{1:N} \sim p(x_{1:N} | \theta) \iff x_{1:N} = g(\theta, \xi) \text{ with } \xi \sim p(\xi)$$

- **Amortized** inference means to find a global estimator of $p(\theta | x_{1:N})$ that works for any dataset $x_{1:N}$

Learning the Posterior via Normalizing Flows

- Goal: Approximate the true posterior $p(\theta|x) \approx p_\phi(\theta|x)$ for all x
- Idea: Reparametrize p_ϕ in terms of an invertible function $f_\phi: \mathbb{R}^{n_\theta} \rightarrow \mathbb{R}^{n_\theta}$ which performs a normalizing flow between θ and a Gaussian latent variable z :
$$\theta \sim p_\phi(\theta|x) \iff \theta = f_\phi^{-1}(z; x) \text{ with } z \sim \mathcal{N}_{n_\theta}(\mathbf{0}, \mathbf{I})$$
- BayesFlow: Implement f_ϕ as conditional invertible neural network (cINN)

Derivation of Loss Function

- Seek neural network parameters ϕ such that the “distance” between $p_\phi(\boldsymbol{\theta} | \mathbf{x})$ and $p(\boldsymbol{\theta} | \mathbf{x})$ is small
- Formally, we minimize the average KL divergence over all datasets:

$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} \mathbb{E}_{p(\mathbf{x})} [\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}) || p_\phi(\boldsymbol{\theta} | \mathbf{x}))] \\ &= \cdots \\ &= \operatorname{argmin}_{\phi} \left[- \iint p(\boldsymbol{\theta}, \mathbf{x}) \log p_\phi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} d\mathbf{x} \right]\end{aligned}$$

Derivation of Loss Function (cont.)

- Apply the change of variables formula

$$p_{\phi}(\boldsymbol{\theta} | \mathbf{x}) = p_{\mathbf{z}}(f_{\phi}(\boldsymbol{\theta}; \mathbf{x})) |\det \mathbf{J}_{f_{\phi}}|$$

and insert the prescribed density $p_{\mathbf{z}}(\mathbf{z}) \propto \exp(-\frac{1}{2}\|\mathbf{z}\|_2^2)$ to obtain:

$$\hat{\phi} = \operatorname{argmin}_{\phi} \iint p(\boldsymbol{\theta}, \mathbf{x}) \left(\frac{1}{2} \|f_{\phi}(\boldsymbol{\theta}; \mathbf{x})\|_2^2 - \log |\det \mathbf{J}_{f_{\phi}}| \right) d\boldsymbol{\theta} d\mathbf{x}$$

- Monte Carlo estimate (batch size M):

$$\hat{\phi} \approx \operatorname{argmin}_{\phi} \frac{1}{M} \sum_{m=1}^M \left(\frac{1}{2} \|f_{\phi}(\boldsymbol{\theta}^{(m)}; \mathbf{x}^{(m)})\|_2^2 - \log |\det \mathbf{J}_{f_{\phi}}^{(m)}| \right)$$

Summary Statistics Learning

- To handle variable-size data, incorporate a summary network h_ψ that transforms raw data $\mathbf{x}_{1:N}$ into a fixed-size vector $h_\psi(\mathbf{x}_{1:N})$
- Train the summary network jointly with the cINN:

$$\hat{\phi}, \hat{\psi} = \operatorname{argmin}_{\phi, \psi} \mathbb{E}_{p(N, \mathbf{x}_{1:N})} [\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}_{1:N}) || p_\phi(\boldsymbol{\theta} | h_\psi(\mathbf{x}_{1:N})))]$$

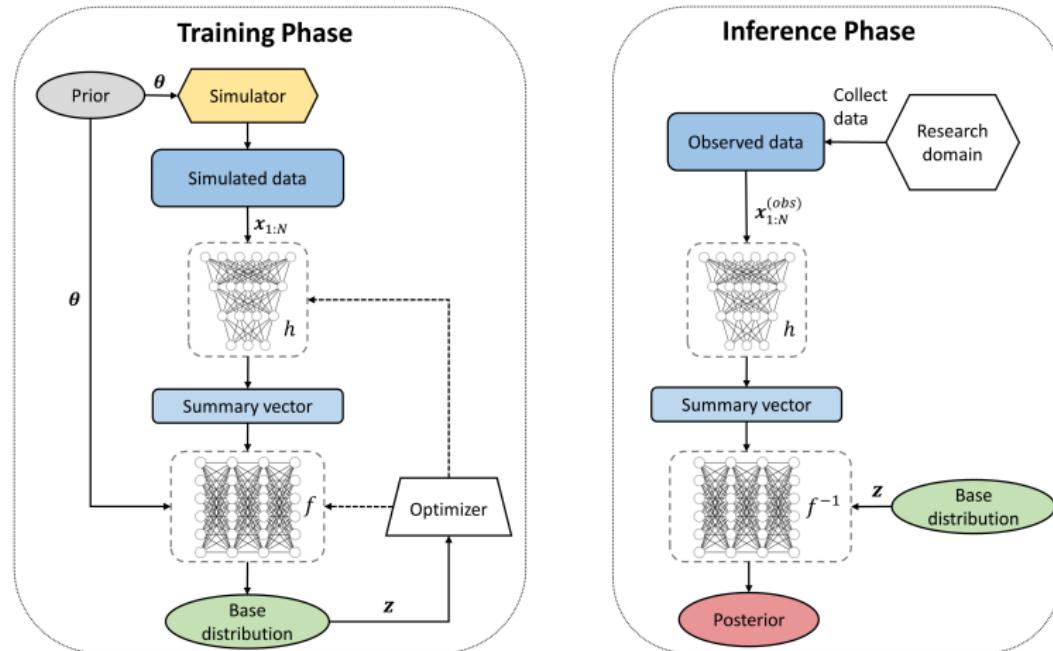
- Monte Carlo estimate of the objective:

$$\hat{\phi}, \hat{\psi} \approx \operatorname{argmin}_{\phi, \psi} \frac{1}{M} \sum_{m=1}^M \left(\frac{1}{2} \left\| f_\phi \left(\boldsymbol{\theta}^{(m)}; h_\psi(\mathbf{x}_{1:N^{(m)}}) \right) \right\|_2^2 - \log \left| \det \mathbf{J}_{f_\phi}^{(m)} \right| \right)$$

Theoretical Exactness

- Assumption: Both the cINN and the summary network are sufficiently deep s.t. perfect convergence is achieved (analytical KL loss $\rightarrow 0$)
- Then it can be proven:
 - The learned posterior $p_{\hat{\phi}}(\theta | h_{\hat{\psi}}(x_{1:N}))$ approximates the true posterior $p(\theta | x_{1:N})$ perfectly for almost every dataset $x_{1:N}$
 - The learned summary statistic $h_{\hat{\psi}}(x_{1:N})$ is “maximally informative”

Amortized Posterior Sampling with BayesFlow



<https://bayesflow.readthedocs.io/en/latest/>

Outline

- 1 Amortized Likelihood-free Bayesian Inference with BayesFlow
- 2 Missing Data Handling
- 3 Experiments and Results
- 4 Summary and Outlook

Missing Data

- Incomplete datasets in which no numerical value is stored for one or multiple observations
 ⇒ Can substantially complicate reliable statistical inference
- Highly relevant e.g. in biomedical sciences:
 - Clinical trials, where probands do not show up at some of the regular screenings
 - COVID-19 pandemic, where infection numbers are not reported on holidays or due to the overload of health departments

Problem Specification

- Time series model producing datasets $\mathbf{x}_{1:N}$ of fixed length N , where the $\mathbf{x}_i \in \mathbb{R}^{n_x}$ are assigned to time points $0 \leq t_1 < \dots < t_N$
- Underlying complete dataset $\mathbf{x}_{1:N}$ is only partially known:

$$\mathcal{S} := \{1 \leq i \leq N : \text{No value available for } \mathbf{x}_i\}, \quad \mathcal{T} := \{1, \dots, N\} \setminus \mathcal{S}$$

- Two possible target distributions:
 - Posterior $p(\boldsymbol{\theta} | \mathbf{x}_{\mathcal{T}})$ cond. on the available data $\mathbf{x}_{\mathcal{T}} := \{\mathbf{x}_i\}_{i \in \mathcal{T}}$
 - Posterior $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$ cond. on the complete original dataset $\mathbf{x}_{1:N}$

\implies Focus on the estimation of $p(\boldsymbol{\theta} | \mathbf{x}_{\mathcal{T}})$

Encoding Missing Values

- Goal: Teach neural network to detect missing values and ignore them for the inference
- Investigated approaches to encode missing data:
 - **Augment by 0/1:** Augment the data matrix $[x_1, \dots, x_N]$ by a binary row, where 0 indicates absence and 1 indicates presence of a data point. Further, insert a constant value $C \in \mathbb{R}$ for all missing data points.
 - **Insert C:** Simply insert a constant $C \in \mathbb{R}$ for all missing data points.
 - **Time labels:** Augment the data matrix containing only the available observations by a row of time points associated to these observations.

Encoding Missing Values - Examples

- Incomplete dataset:

$$[0.9 \quad \text{N/A} \quad 3.1 \quad 4.0 \quad \text{N/A}]$$

- Augment by 0/1 with fill-in value -1 :

$$\Rightarrow \begin{bmatrix} 0.9 & -1.0 & 3.1 & 4.0 & -1.0 \\ 1 & 0 & 1 & 1 & 0 \end{bmatrix}$$

- Insert -1 : $\Rightarrow [0.9 \quad -1.0 \quad 3.1 \quad 4.0 \quad -1.0]$

- Time labels: $\Rightarrow \begin{bmatrix} 0.9 & 3.1 & 4.0 \\ 0.0 & 1.0 & 1.5 \end{bmatrix}$

Integration into BayesFlow

- Choose maximum number $N_{\emptyset}^{\max} < N$ of missing data points
- Train network on datasets with artificially induced missing values that are encoded via one of the three approaches
- For our test models:

Number of missing data points $N_{\emptyset} \in \{0, \dots, N_{\emptyset}^{\max}\}$ and set of missing time indices $\mathcal{S} \subseteq \{1, \dots, N\}$ are sampled uniformly

Outline

- 1 Amortized Likelihood-free Bayesian Inference with BayesFlow
- 2 Missing Data Handling
- 3 Experiments and Results
- 4 Summary and Outlook

Performance on Complete Data - Brief Overview

- Nearly perfect approximation of the true posterior can be achieved.
- For time series models, LSTM networks with approximately as many hidden units as the total number of entries in the data matrix are a reasonable choice of summary network.
- Hard-bounded posteriors resulting from uniform priors are difficult to learn.

Conversion Reaction Model

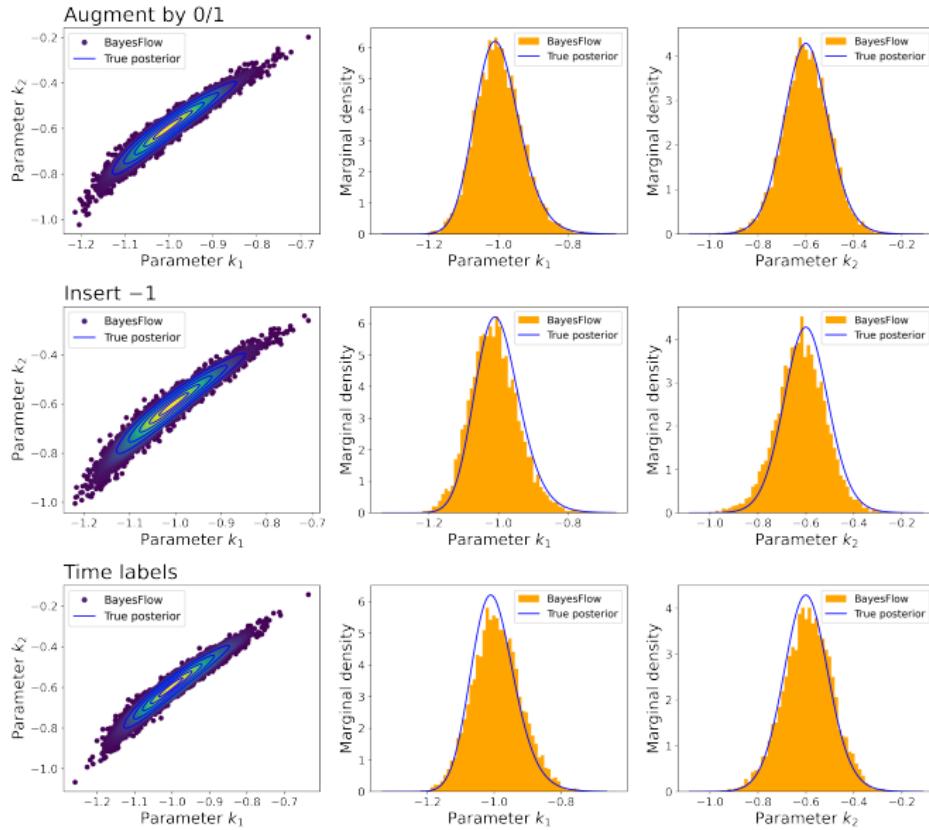
- Forward model:

$$\begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \end{pmatrix} = \begin{pmatrix} -c_1 x_1 + c_2 x_2 \\ c_1 x_1 - c_2 x_2 \end{pmatrix}, \quad \begin{pmatrix} x_1(0) \\ x_2(0) \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$
$$y_t = x_2(t) + \varepsilon_t \text{ with } \varepsilon_t \sim \mathcal{N}(0, 0.015^2)$$

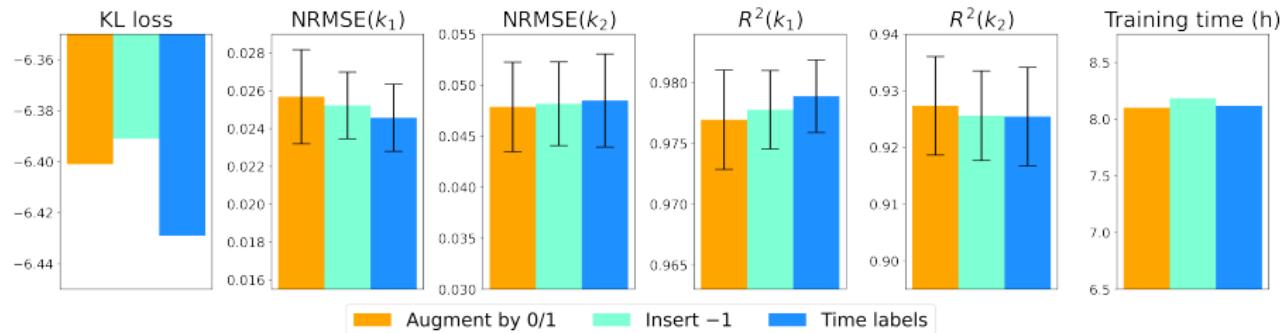
- $N = 11$ observations at $t = 0, 1, \dots, 10$ and $N_{\emptyset}^{\max} = 6$
- Normal prior on log-scale rate parameters $k_j = \log_{10}(c_j)$:

$$k_1, k_2 \sim \mathcal{N}(-0.75, 0.25^2)$$

Comparing the Approaches to Encode Missing Data



Comparing the Approaches to Encode Missing Data (cont.)



⇒ Similarly good performance of all 3 approaches!

Increased Robustness via Augmentation by 0/1

- “Augment by 0/1” and “Insert C ” require a fill-in value C
- CR model (positive trajectories):
 - Natural to choose an unambiguous fill-in value $C = -1$
 - Both networks correctly learn to interpret -1 as missing values
- More complex dynamics: Choice of fill-in value is non-trivial
- In the following: Test ambiguous value $C = 0.5$ on reduced CR model ($N = 3$ observations at $t = 0, 5, 10$ and $N_{\emptyset}^{\max} = 2$)

Increased Robustness via Augmentation by 0/1 (cont.)

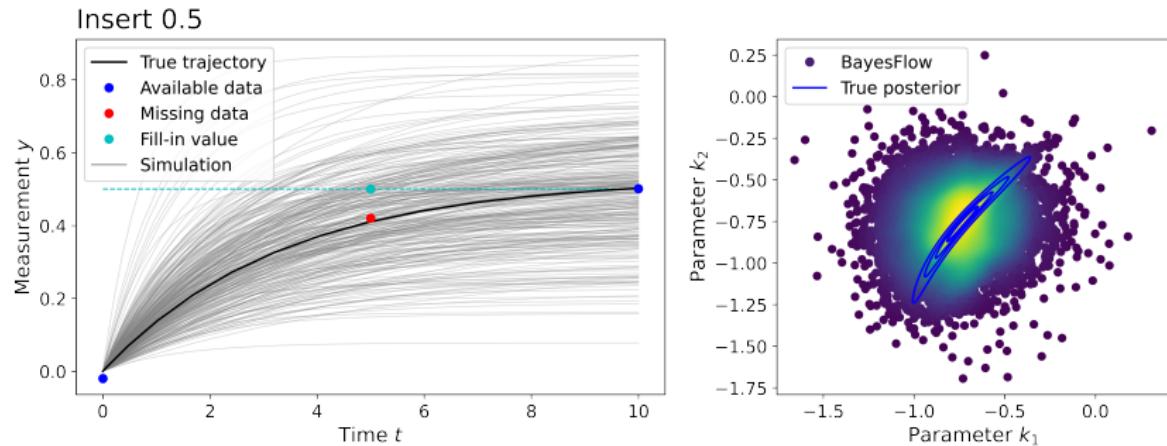


Figure: Test data: $\begin{bmatrix} -0.02 & 0.5 & 0.501 \end{bmatrix}$ (only the second observation is missing). However, the network misinterprets the signal 0.501 as another missing value.

Increased Robustness via Augmentation by 0/1 (cont.)

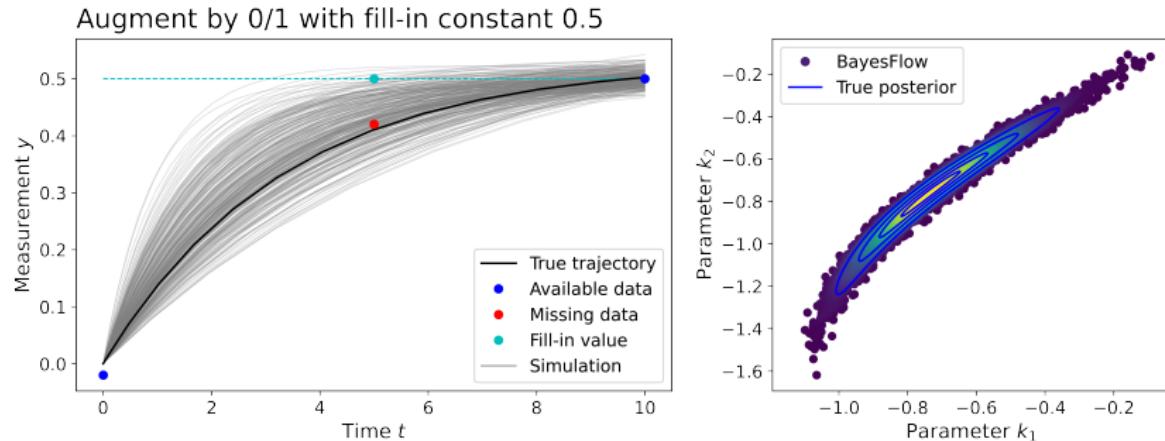


Figure: Test data: $\begin{bmatrix} -0.02 & 0.5 & 0.5 \\ 1 & 0 & 1 \end{bmatrix}$. Based on the augmentation, the network correctly identifies the value 0.5 in the second observation as a missing value and in the third observation as a signal.

Oscillation Model

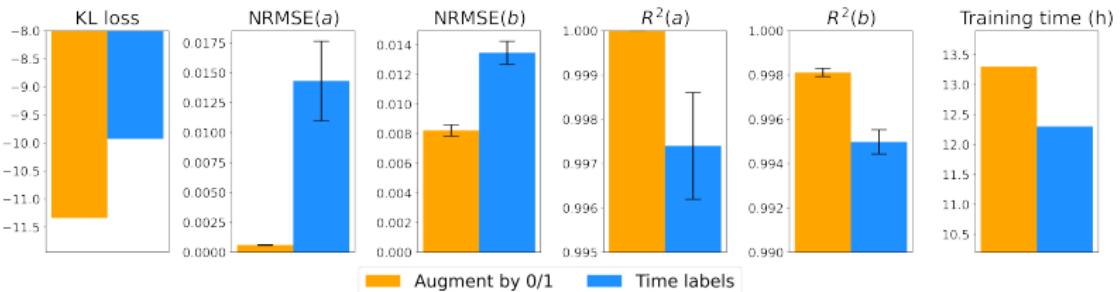
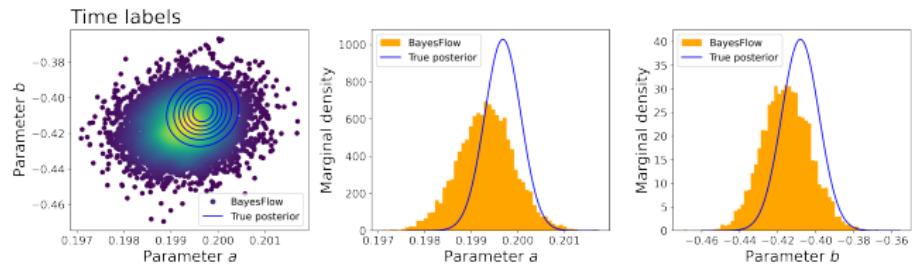
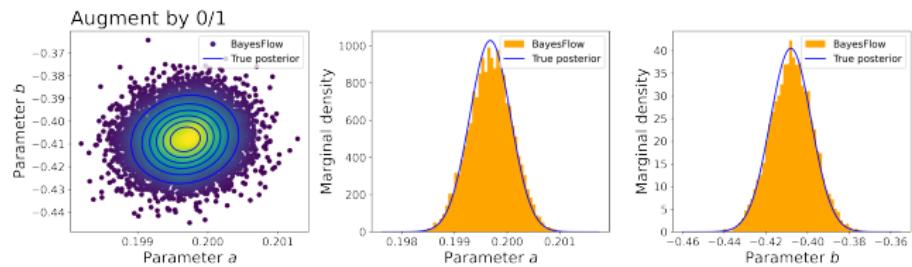
- Forward model:

$$x(t) = \sin(2\pi at) + b$$

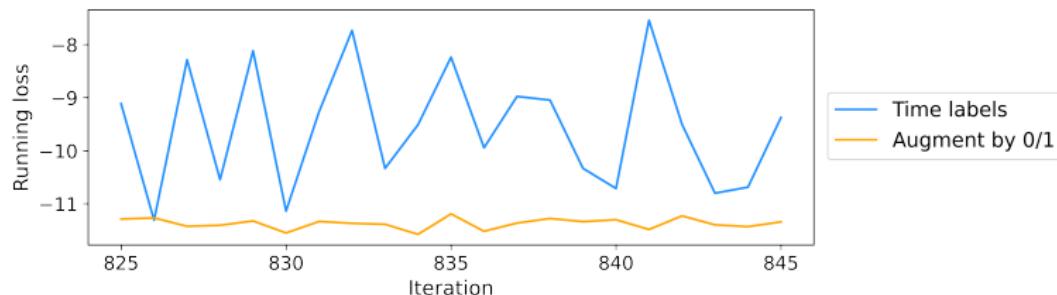
$$y_t = x(t) + \varepsilon_t \text{ with } \varepsilon_t \sim \mathcal{N}(0, 0.05^2)$$

- Prior for frequency parameter: $a \sim \mathcal{U}(0.1, 1)$
- Prior for shift parameter: $b \sim \mathcal{N}(0, 0.25^2)$
- $N = 41$ observations at $t_k = \frac{1}{4}k$ for $k = 0, 1, \dots, 40$ and $N_\emptyset^{\max} = 21$
- Models producing oscillatory data are in general hard to fit (highly irregular objective function landscape)

Poor Convergence of the “Time labels” Approach



Poor Convergence of the “Time labels” Approach (cont.)



- Keeping N_\emptyset fixed for each batch gives an inexact loss approximation and causes the running loss to jump
- Effect is especially strong for the more challenging loss function of the oscillation model
⇒ Convergence issues

SIR Epidemiology Model

- Forward model:

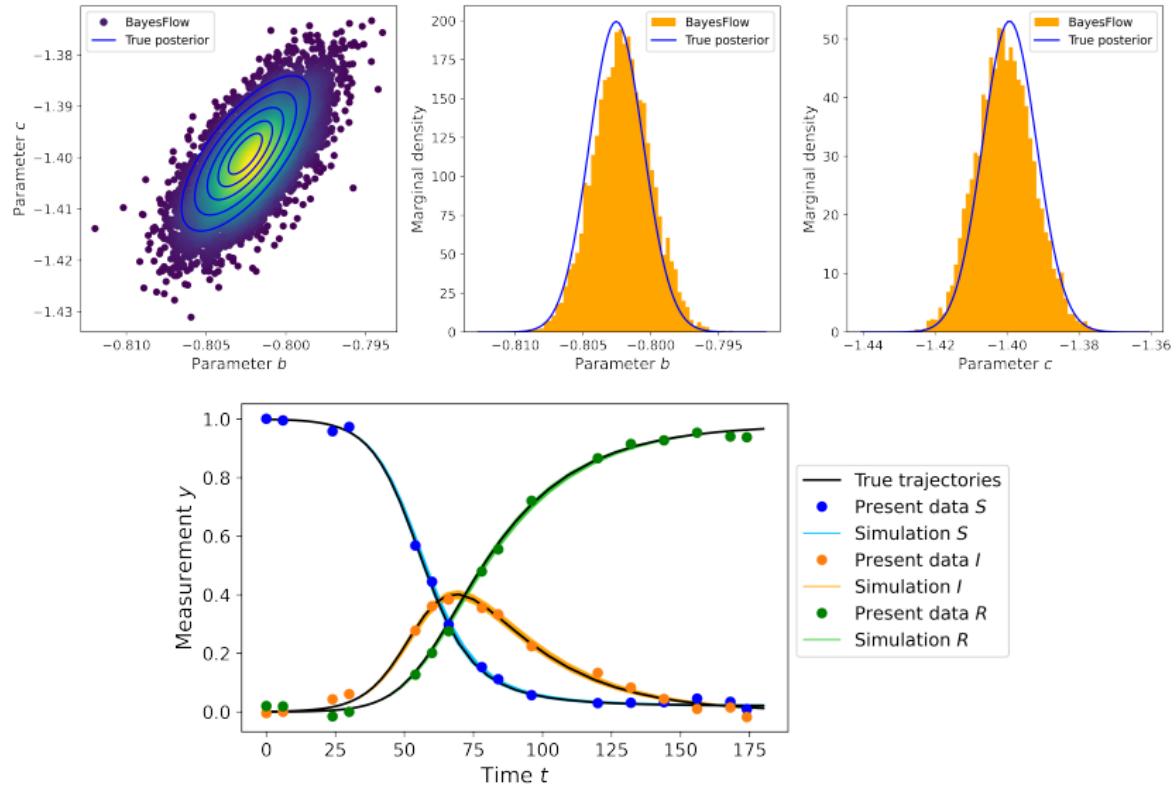
$$\begin{pmatrix} \dot{S} \\ \dot{I} \\ \dot{R} \end{pmatrix} = \begin{pmatrix} -\beta SI \\ \beta SI - \gamma I \\ \gamma I \end{pmatrix}, \quad \begin{pmatrix} S(0) \\ I(0) \\ R(0) \end{pmatrix} = \begin{pmatrix} 0.999 \\ 0.001 \\ 0 \end{pmatrix}$$
$$y_t = \begin{pmatrix} S(t) \\ I(t) \\ R(t) \end{pmatrix} + \begin{pmatrix} \varepsilon_{t,1} \\ \varepsilon_{t,2} \\ \varepsilon_{t,3} \end{pmatrix} \text{ with } \varepsilon_{t,i} \sim \mathcal{N}(0, 0.015^2)$$

- Normal priors on log-scale infection/recovery rate parameters:

$$b = \log_{10}(\beta) \sim \mathcal{N}(-1, 0.3^2), \quad c = \log_{10}(\gamma) \sim \mathcal{N}(-1.5, 0.3^2)$$

- $N = 31$ observations at $t_k = 6k$ for $k = 0, 1, \dots, 30$ and $N_{\emptyset}^{\max} = 15$

Performance of “Augment by 0/1” on SIR Model



Outline

- 1 Amortized Likelihood-free Bayesian Inference with BayesFlow
- 2 Missing Data Handling
- 3 Experiments and Results
- 4 Summary and Outlook

Summary

- Extended the BayesFlow method to perform amortized likelihood-free Bayesian inference on incomplete datasets
 ⇒ Sampling from the posterior conditioned on the available data
- Train network on datasets with artificially induced missing data
- Compared three approaches to encode missing values
 ⇒ Most robust approach: “Augment by 0/1”

Outlook I

- Application to real-world problems from systems biology and epidemiology, e.g. NLME models
- Use parallelization to speed up batch simulation
- Explore more sophisticated LSTM network architecture
- Integrate known patterns behind the data missingness
- If data are missing independently in different observed features, perform augmentation for each feature row

Outlook II

- Recover posterior $p(\theta | \mathbf{x}_{1:N})$ cond. on the complete original data
 \implies Stronger posterior contraction, better point estimates
- Data imputation based on:
 - Spline interpolation
 - Deep learning, e.g. LSTM networks or autoencoders
- Train BayesFlow on imputed data with an added binary row to label imputed/original values
- Continuous weights indicating trustworthiness of imputed values, e.g. from Gaussian process regression

Final Bachelor Thesis Presentation

Thank you for your attention!
I look forward to your questions and feedback!