

# Learning summary statistics for Bayesian inference with Autoencoders

Carlo Albert OWABC, September 28 2023





## **Problem Statement**

Likelihood functions involving high-dimensional integration/summation over latent variables, because there is either

• a large number of unobserved latent variables:

$$f(\mathbf{y}^{(\text{obs})} \mid \boldsymbol{\theta}) = \int f(\mathbf{y}^{(\text{obs})} \mid \mathbf{x}, \boldsymbol{\theta}) f(\mathbf{x} \mid \boldsymbol{\theta}) d\mathbf{x}$$

• or an unknown partition function (e.g. Ising-type model):

$$f(\mathbf{y}^{(\text{obs})} \mid \boldsymbol{\theta}) = Z^{-1}(\boldsymbol{\theta}) \exp[-H_{\boldsymbol{\theta}}(\mathbf{y}^{(\text{obs})})]$$
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{y}} \exp[-H_{\boldsymbol{\theta}}(\mathbf{y})]$$



## **Approximate Bayesian Computation (ABC)**

- Simulate model outputs y, for many parameter sets  $\theta$ , from  $f(y \ \theta)$ .
- Accept/reject the corresponding parameter sets depending on whether those outputs agree with the observations in terms of
  - Certain features (summary statistics) s(y)
  - A certain tolerance
- E.g. Ising model

$$f(\mathbf{y} \mid \beta, h) = Z^{-1}(\beta, h) \exp\left[\beta\left(\sum_{\langle i,j \rangle} y_i y_j + h \sum_i y_i\right)\right]$$
$$= Z^{-1}(\beta, h) \exp\left[\beta E(\mathbf{y}) + \beta h M(\mathbf{y})\right]$$

 $\mathbf{s}(\mathbf{y}) := (E(\mathbf{y}), M(\mathbf{y}))^T$ 

sufficient statistics!



# Tuning the tolerance: Simulated Annealing ABC (SABC)

Albert et al., Stat. Comput. 2014

- Use metric  $\rho(\mathbf{s}, \mathbf{s}^{(obs)})$  -> interpret as **energy**.
- Simulate from  $f(\mathbf{s}, \boldsymbol{\theta}) e^{-\rho(\mathbf{s}, \mathbf{s}^{(obs)})/T}$  while continuously lowering **temperature** *T*.
- For  $T \to 0$ ,  $f(\mathbf{s}, \boldsymbol{\theta}) e^{-\rho(\mathbf{s}, \mathbf{s}^{(obs)})/T} \to f(\boldsymbol{\theta} \ \mathbf{s}^{(obs)})$ .

Initialisation: sample an ensemble of particles  $\{\mathbf{s}_j, \boldsymbol{\theta}_j\}_{j=1}^N$  from the prior  $f(\mathbf{s}, \boldsymbol{\theta})$ 

- 1. Draw a random particle  $(\mathbf{s}_i, \boldsymbol{\theta}_i)$  from the ensemble
- 2. Make a jump in parameter space  $\theta_j \rightarrow \theta_i^*$
- 3. Simulate a data set  $\mathbf{y}_{i}^{*}$  from the model  $f(\mathbf{y} \mid \boldsymbol{\theta})$  and calculate  $\mathbf{s}_{i}^{*} = \mathbf{s}(\mathbf{y}_{i}^{*})$
- 4. Accept the move with probability  $\min\left(1, \frac{f(\boldsymbol{\theta}_{j}^{*})}{f(\boldsymbol{\theta}_{j})} \exp\left[-\frac{\rho(\mathbf{s}_{j}^{*}, \mathbf{s}^{(\text{obs})}) \rho(\mathbf{s}_{j}, \mathbf{s}^{(\text{obs})})}{T}\right]\right)$
- 5. Lower the temperature T adaptively so as to minimise entropy production.



# Data compression - fighting the curse of dimensionality

- The computational price to pay for decreasing the tolerance grows exponentially with the output dimension.
- Therefore, we must compress the data to a handful of summary statistics.
- Generally, this leads to a loss of information relevant to constrain the parameters.



Albert et al., SciPost Physics 2022

Consider a generic stochastic model:

$$f(\mathbf{y}, \boldsymbol{\theta}) = f(\mathbf{y} \ \boldsymbol{\theta}) f(\boldsymbol{\theta}) \quad \boldsymbol{\theta} \in \mathbb{R}^p, \quad \mathbf{y} \in \mathbb{R}^N,$$

and a map of summary statistics

$$\mathbf{s}: \mathbb{R}^N \longrightarrow \mathbb{R}^q$$
.

We want them to be

asymptotically sufficient:  $I(\mathbf{s}(\mathbf{Y}), \mathbf{\Theta}) = I(\mathbf{Y}, \mathbf{\Theta}) + \mathcal{O}(N^{-1})$ (Accuracy)

Asymptotically concentrated:  $H(\mathbf{s}(\mathbf{Y}) \ \mathbf{\Theta}) \sim -\ln N$ (Efficiency)

-> Thermodynamic state variables





Albert et al., SciPost Physics 2022

• If sufficient statistics cannot easily be found, a reasonable choice is to use parameter estimators, e.g.

$$\mathbf{s}(\mathbf{y}) = \hat{\boldsymbol{\theta}}(\mathbf{y}) = \int \boldsymbol{\theta} f(\boldsymbol{\theta} \ \mathbf{y}) d\boldsymbol{\theta}$$

- Fearnhead and Prangle (2012) use a linear regression (built-into SABC software packages)
- Machine Learning lends itself to this task as well! (e.g. Jiang 2017)



Albert et al., SciPost Physics 2022

- For large *i.i.d.* datasets, parameter estimators capture most of the  $\theta$ -related information in the data.
- If the data is *correlated*, this is no longer the case and we might need more summary statistics than parameters
- This is due to entropic effects: For members of the exponential family

$$f(\mathbf{y} \ \boldsymbol{\theta}) = Z^{-1}(\boldsymbol{\theta})c(\mathbf{y})\exp\left(\sum_{\alpha=1}^{q} s_{\alpha}(\mathbf{y})g^{\alpha}(\boldsymbol{\theta})\right)$$

the free energy splits into en energy and an entropy term

$$F_{\theta}(\mathbf{s}) := -\ln f(\mathbf{s} \ \theta) = -\ln \int f(\mathbf{y} \ \theta) d\Omega_{\mathbf{s}}(\mathbf{y}) = U_{\theta}(\mathbf{s}) - S(\mathbf{s})$$

and is not necessarily concentrated around a single point, for fixed  $\theta$ , even if *N* is large; possibly even when  $N \to \infty$  (*phases!*).



Albert et al., SciPost Physics 2022

$$\text{Toy:} \quad y_{n+1} = \alpha f(y_n) + \sigma \epsilon_n \,, \quad \epsilon_n \sim \mathcal{N}(0,1) \quad \text{i.i.d.} \,, \quad f(y) = y^2(y-1)$$





Albert et al., SciPost Physics 2022

Toy: 
$$y_{n+1} = \alpha f(y_n) + \sigma \epsilon_n$$
,  $\epsilon_n \sim \mathcal{N}(0,1)$  i.i.d.,  $f(y) = y^2(y-1)$ 

$$f(\mathbf{y} \ \boldsymbol{\theta}) \propto \sigma^{-N} \exp\left[-\sum_{n} \frac{(y_{n+1} - \alpha f(y_n))^2}{2\sigma^2}\right], \quad \boldsymbol{\theta} = \begin{pmatrix} \alpha \\ \sigma \end{pmatrix}$$

n=1

Sufficient statistics:

$$\hat{\alpha}(\mathbf{y}) = \frac{\sum_{n=1}^{N} y_n f(y_{n-1})}{\sum_{n=1}^{N} (f(y_{n-1}))^2},$$

$$\hat{\sigma}^2(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^{N} (y_n - \hat{\alpha}(\mathbf{y}) f(y_{n-1}))^2,$$
Parameter estimators
$$o(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^{N} (f(y_{n-1}))^2.$$
Order parameter



Albert et al., SciPost Physics 2022



-N

Sufficient statistics:

$$\hat{\alpha}(\mathbf{y}) = \frac{\sum_{n=1}^{N} y_n f(y_{n-1})}{\sum_{n=1}^{N} (f(y_{n-1}))^2},$$
Parameter estimators
$$\hat{\sigma}(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^{N} (y_n - \hat{\alpha}(\mathbf{y}) f(y_{n-1}))^2,$$
Order parameter
$$o(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^{N} (f(y_{n-1}))^2.$$
Order parameter



Albert et al., SciPost Physics 2022



Explicit Noise Conditional Autoencoder ENCA





$$\mathscr{L} = \sum_{j,\alpha} \left( \frac{s_{\alpha}^{(j)} - \theta_{\alpha}}{\theta_{\alpha}} \right)^2 + \sum_{\alpha} \left( \frac{\hat{\theta}_{\alpha} - \theta_{\alpha}}{\theta_{\alpha}} \right)^2$$

Explicit Noise Conditional Autoencoder ENCA



Albert et al., SciPost Physics 2022

Toy:  $y_{n+1} = \alpha f(y_n) + \sigma \epsilon_n$ ,  $\epsilon_n \sim \mathcal{N}(0,1)$  i.i.d.,  $f(y) = y^2(y-1)$ 





Albert et al., SciPost Physics 2022

Toy:  $y_{n+1} = \alpha f(y_n) + \sigma \epsilon_n$ ,  $\epsilon_n \sim \mathcal{N}(0,1)$  i.i.d.,  $f(y) = y^2(y-1)$ 









- 4 features seem to be sufficient to predict streamflow.
- There seems to be **information missing** in the known static catchment attributes.



















## Conclusions

- The **ENCA architecture** can be used to disentangle high-dimensional irrelevant from low-dimensional relevant information.
- It can be used to find near-sufficient and concentrated summary statistics in simulated data from stochastic models.
- It can be used to find catchment features in **observed** streamflow data.



C.A. al., SciPost Physics 2022

Thanks to my collaborators: Simone Ulzega, Antonietta Mira, Firat Oezdemir, Fernando Perez-Cruz, Alberto Bassi.