

Weekly Report (4.4-4.10)*

Yanpei Cai[†]

April 11, 2025

During the past week, I mainly used my time to tackle with coding while I spent a little amount of time to write a summary of maths in BOLFI. I have to say most of my time was spent on making up the oversight of my previous works, leading to a reconstruction of codes for the model which has yet completed. The followings are break down to progress and problems identified in each part with more details.

Coding

With available datasets, what can be observed is the percentage of never-smokers (N), smokers (S), and quitters (Q) in the population by the end of each year. With the agent-based model depicting the intricate dynamics of state change process for each agent, through simulations we are able to get the simulated data of each kind by the end of each step (year), which follows the methodology of BOLFI. By implementing BOLFI, what we want are inferences for the parameters which characterise the state change process as well as the social network topology.

The main problem emerged in my codes was actually caused by the oversight on the order of state change for each agent, and the insufficient proficiency of Mesa. When I was summarising the simulated data, the huge observed discrepancies reminded me to double check the model construction in my codes. I noticed that for simplicity of my codes' structure, I initially set each agent to update its state consecutively in the order of its node's **unique_id** instead of set the **SimultaneousActivation** mechanism. This actually causes a natural contagion of the cluster of nodes in the same state to nodes in other states in their neighbourhood. That explains why the contagion can be very fast, as can be seen in the wrong simulated data with moderate parameter values.

To solve this problem, I went back to the codes of **NSQ_Model** and **NSQ_Agent** to set up a more systematic formalism with specific settings. Given the short of time, this has yet completed but I am looking forward to completing the codes by next Monday.

Maths

Handwritten notes is attached after the main parts of this report.

[†]Artificial Intelligence and its Applications Institute, School of Informatics, The University of Edinburgh
Email Address: yanpei.cai@ed.ac.uk

*This is a summary of the works was done in the past week.

Questions Raised for This Week

- **Coding:** (first three tasks are expected to be completed by next Monday and send to Valerio for inspection)
 - **NSQ_Model** and **NSQ_Agent** settings for allowing all agents to update their states simultaneously.
 - **NSQ_Model** and **NSQ_Agent** settings in a more systematic formalism.
 - Data generating mechanism to compile with the time scale of both the UK and the US dataset.
 - **BOLFI** settings to cope with inferences for model parameters given the simulated data obtained in the last task.
- **Maths:** (in the arXiv version paper, S3 Appendix. Parameter estimation from Christakis)
 - Why does the interaction parameter **g** between SMOKER and QUITTER can be replaced by the probability of quitter over t years, as in (20)?
 - How does the expression in (22) are derived with applying (21) to both the interaction term and spontaneous term?
- **General Question (Physics to Complexity Science):**

How do statistical physics and those physical rules further assist our modeling of social systems and real-world social phenomena (e.g., heterogeneity of agents; capturing complexity versus model simplicity, etc.)?

- **General Question (Long-term Plan and Short-term Goals in PhD Journey):**

From my perspective, the short-term goals are coding the agent-based modelling mechanism for our smoking contagion model and implementing BOLFI to obtain better inferences for the model parameters. Can we have a discussion about the long-term plan (since you both seem to have many potential options)?

Plan for Next Week

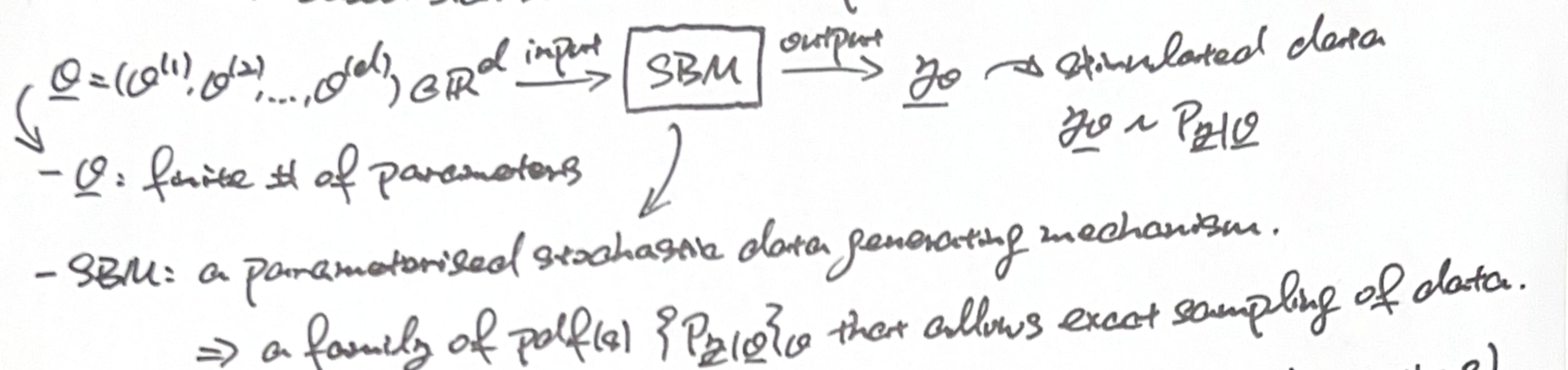
- **Coding:**
 - The first three tasks in **Questions.Coding** by next Monday and get feedback if further problems arise.
 - Testing the **BOLFI** implementation to our smoking contagion model.
IDEA: Firstly, obtain inference for single parameter by keep other model parameters and network generation parameters fixed.
- **Maths:**

Try to answer the questions raised in **Questions.Math** (seek help and get feedback early next week if no progress are made during the weekend).

Question.

Given : $\underline{\theta}$ (list) \cdot a set of observed data $\underline{y}_0 = (y_0^{(1)}, y_0^{(2)}, \dots, y_0^{(n)})$ of n data points. (possibly dependent).

- a simulator-based statistical model, represented as



Objectives : perform statistical inference of model parameters (given observations).

Obstacles : generative process of the model can be very complex, so that the analytical form of the pdf $P_{Z|\underline{\Theta}}$ is often intractable, leads to inference using the likelihood function $L(\underline{\Theta}) = P_{Z|\underline{\Theta}}(\underline{y}_0|\underline{\Theta})$ is not available.

Solutions \Rightarrow Likelihood-Free Inference (LFI) Methods.

Basic Idea. Since SBM allows data to be generated (via sampling from the pdf $P_{Z|\underline{\Theta}}$ of unknown analytical form), we generate data \underline{y}_0 with multiple values of input $\underline{\Theta}$ and identify the model parameters $\underline{\Theta}$ by finding values which yield simulated data \underline{y}_0 resemble the observed data \underline{y}_0 .

General Manner. constructing approximate likelihood functions with the discrepancy Δ_0 between the observed data \underline{y}_0 and data \underline{y}_0 simulated with parameter values $\underline{\Theta}$.

First Step. On a statistical level, the approximation consists of reducing the observed data \underline{y}_0 to some features (often means summary statistics). * identifying summary statistics helps to reduce dimensionality and to retain sufficient information for the inference of $\underline{\Theta}$.

$$L(\underline{\Theta}) = P_{Z|\underline{\Theta}}(\underline{y}_0|\underline{\Theta}) \Rightarrow L(\underline{\Theta}) = P_{\Phi|\underline{\Theta}}(\underline{\Phi}_0|\underline{\Theta}).$$

$$\underline{\Phi} = (\underline{\Phi}^{(1)}, \underline{\Phi}^{(2)}, \dots, \underline{\Phi}^{(P)}) \in \mathbb{R}^P$$

: summary statistics

* still of unknown analytical form (inherited property from $P_{Z|\underline{\Theta}}$).

$\Rightarrow L(\underline{\Theta})$ still needs to be approximated by some methods.

Examples. - parametric approximation (synthetic likelihood).

- non-parametric approximation (kernel density estimation, i.e. KDE)

Synthetic Likelihood (SL)

Key Ideas. - central limit theorem (CLT) supports the key multivariate normality approximation of the summary statistics, as $\underline{\Phi} \sim N(\underline{\mu}_0, \underline{\Sigma}_0)$, if the # of samples n of \underline{z} is sufficiently large \Rightarrow this suggests unknown.

$$P_{\underline{\Phi}|Q}(\underline{\phi}|Q) \stackrel{(1)}{=} \frac{1}{(2\pi)^{p/2} |\det \underline{\Sigma}_0|^{1/2}} \exp\left(-\frac{1}{2} (\underline{\phi} - \underline{\mu}_0)^T \underline{\Sigma}_0^{-1} (\underline{\phi} - \underline{\mu}_0)\right).$$

\Rightarrow limiting likelihood approximation (with infinite computing power)

$$\tilde{L}_S(Q) = P_{\underline{\Phi}|Q}(\underline{\Phi}|Q).$$

$$\tilde{l}_S(Q) = \log(\tilde{L}_S(Q)) = -\frac{p}{2} \log(2\pi) - \frac{1}{2} \log|\det \underline{\Sigma}_0| - \frac{1}{2} (\underline{\Phi}_0 - \underline{\mu}_0)^T \underline{\Sigma}_0^{-1} (\underline{\Phi}_0 - \underline{\mu}_0)$$

- simulator used to provide estimates of unknown mean $\underline{\mu}_0$ and covariance matrix $\underline{\Sigma}_0$ via a sample average \bar{E}^N over N independently generated summary statistics, allowing the likelihood approximation to be computable.

$$\underline{\Phi}_0 = (\underline{\Phi}_0^{(1)}, \underline{\Phi}_0^{(2)}, \dots, \underline{\Phi}_0^{(N)}). \quad \underline{\Phi}_0^{(i)} \stackrel{(i)}{\text{ind.}} \sim P_{\underline{\Phi}|Q}, \quad \underline{\Phi}_0 = (\underline{\Phi}_0^{(1)(1)}, \underline{\Phi}_0^{(1)(2)}, \dots, \underline{\Phi}_0^{(1)(p)})^T$$

$$\hat{\mu}_0 = \bar{E}^N[\underline{\Phi}_0] = \frac{1}{N} \sum_{i=1}^N \underline{\Phi}_0^{(i)}, \quad \hat{\Sigma}_0 = \bar{E}^N[(\underline{\Phi}_0 - \hat{\mu}_0)(\underline{\Phi}_0 - \hat{\mu}_0)^T].$$

\Rightarrow computable likelihood estimation (with finite computing power)

$$\hat{L}_S(Q) = \frac{1}{(2\pi)^{p/2} |\det \hat{\Sigma}_0|^{1/2}} \exp\left(-\frac{1}{2} (\underline{\Phi}_0 - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (\underline{\Phi}_0 - \hat{\mu}_0)\right).$$

$$\hat{l}_S(Q) = \log(\hat{L}_S(Q)) = -\frac{p}{2} \log(2\pi) - \frac{1}{2} \log|\det \hat{\Sigma}_0| - \frac{1}{2} (\underline{\Phi}_0 - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (\underline{\Phi}_0 - \hat{\mu}_0)$$

Conceptual Diagram

