

Generalised simulation inference

Concept. To largely generalise the procedure of statistical inference for any model using an algorithm which builds from techniques we developed in the previous chapter. When we say ‘statistical inference’ here; we specifically mean computing the maximum a posteriori (MAP) estimate for any arbitrary stochastic model which has been defined in the *stochadex* simulator. In order for our algorithm to evaluate the MAP of a model, we show that the user must specify the prior distribution over model parameters, and the model must itself be defined within the *stochadex*. In this chapter, we will discuss some concepts which are commonplace within the field of Bayesian inference and provide a few simple examples of how our algorithm might work in various instances. For the mathematically-inclined, this chapter will give a very brief exposition for Bayesian statistical inference methodology — in particular, how it relates to the evaluation of MAP estimates. For the programmers, the software described in this chapter lives in the public Git repository: <https://github.com/umbralcalc/learnadex>.

4.1 Inference formalism

In Bayesian inference, one applies Bayes’ rule to the problem of statistically inferring a model from some dataset. This typically involves the following formula for a posterior distribution

$$\mathcal{P}_{t+1}(z|Y) \propto \mathcal{L}_{t+1}(Y|z)\mathcal{P}(z). \quad (4.1)$$

In the formula above, one relates the prior probability distribution over a parameter set $\mathcal{P}(z)$ and the likelihood $\mathcal{L}_{t+1}(Y|z)$ of some data matrix Y up to timestep $t + 1$ given the parameters z of a model to the posterior probability distribution of parameters given the data $\mathcal{P}_{t+1}(z|Y)$ up to some proportionality constant. All this may sound a bit technical in statistical language, so it can also be helpful to summarise what the formula above states verbally as follows: the initial (prior) state of knowledge about the parameters z we want to learn can be updated by some likelihood function of the data to give a new state of knowledge about the values for z (the ‘posterior’ probability).

From the point of view of statistical inference, if we seek to maximise $\mathcal{P}_{t+1}(z|Y)$ — or its logarithm — in Eq. (4.1) with respect to z , we will obtain what is known as a maximum posteriori

(MAP) estimate of the parameters. In fact, we have already encountered this methodology in the previous chapter when discussing the algorithm which obtains the best fit parameters for the empirical probability reweighting. In this case; while it appears that we optimised the log-likelihood directly as our objective function, one can easily show that this is also technically equivalent obtaining a MAP estimate where one chooses a specific prior $\mathcal{P}(z) \propto 1$ (typically known as a ‘flat prior’).

How might we calculate the posterior in practice with some arbitrary stochastic process model that has been defined in the stochadex? In order to make the comparison to a real dataset, any stochadex model of interest will always need to be able to generate observations which can be directly compared to the data. To formalise this a little; a stochadex model could be represented as a map from z to a set of stochastic measurements $\mathbf{Y}_{t+1}(z), \mathbf{Y}_t(z), \dots$ that are directly comparable to the rows in the real data matrix Y . The values in Y may only represent a noisy or partial measurement of the latent states of the simulation X , so a more complete picture can be provided by the following probabilistic relation

$$P_{t+1}(\mathbf{Y}|z) = \int_{\Omega_{t+1}} dX P_{t+1}(\mathbf{Y}|X) P_{t+1}(X|z), \quad (4.2)$$

where, in practical terms, the measurement probability $P_{t+1}(\mathbf{Y}|X)$ of \mathbf{Y} given X can be represented by sampling from another stochastic process which takes the state of the stochadex simulation as input. Given that we have this capability to compare like-for-like between the data and the simulation; the next problem is to figure out how this comparison between two sequences of vectors can be done in a way which ensures the the statistics of the posterior are ultimately respected.

For an arbitrary simulation model which is defined by the stochadex, the likelihood in Eq. (4.1) is typically not describable as a simple function or distribution. While we could train the probability reweighting we derived in the previous chapter to match the simulation; to do this well would require having an exact formula for the conditional probability, and this is not always easy to derive in the general case. Instead, there is a class of Bayesian inference methods which we shall lean on to help us compute the posterior distribution (and hence the MAP), which are known as ‘Likelihood-Free’ methods [1, 2, 3, 4].

‘Likelihood-Free’ methods work by separating out the components of the posterior which relate to the closeness of rows in \mathbf{Y} to the rows in Y from the components which relate the states X and parameters z of the simulation stochastically to \mathbf{Y} . To achieve this separation, we can make use of chaining conditional probability like this

$$\mathcal{P}_{t+1}(X, z|Y) = \int_{\Upsilon_{t+1}} d\mathbf{Y} \mathcal{P}_{t+1}(\mathbf{Y}|Y) P_{t+1}(X, z|\mathbf{Y}), \quad (4.3)$$

where Υ_{t+1} here corresponds to the domain of the simulated measurements matrix \mathbf{Y} at time $t+1$.

As we demonstrated in the previous chapter, it’s possible for us to also optimise a probability distribution $\mathcal{P}_{t'}(y|Y) = P_{t'}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}, \dots)$ for each step in time to match the statistics of the measurements in Y as well as possible, given some statistics $\mathcal{M}_{t'} = \mathcal{M}_{t'}(Y)$ and $\mathcal{C}_{t'} = \mathcal{C}_{t'}(Y)$. Assuming the independence of samples (rows) in Y , this distribution can be used to construct the distribution over all of Y through the following product

$$\mathcal{P}_{t+1}(\mathbf{Y}|Y) = \prod_{t'=0}^{t+1} P_{t'}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}, \dots). \quad (4.4)$$

We do not necessarily need to obtain these statistics from the probability reweighting method, but could instead try to fit them via some other objective function. Either way, this represents a lossy *compression* of the data we want to fit the simulation to, and so the best possible fit is desirable; regardless of overfitting. This choice to summarise the data with statistics means we are using what is known as a Bayesian Synthetic Likelihood (BSL) method [2, 3] instead of another class of methods which approximate an objective function directly using a proximity kernel — known as Approximate Bayesian Computation (ABC) methods [1].

Let's consider a few concrete examples of $P_{t'}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}, \dots)$. If the data measurements were well-described by a multivariate normal distribution, then

$$P_{t'}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}, \dots) = \text{MultivariateNormalPDF}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}), \quad (4.5)$$

Similarly, if the data measurements were instead better described by a Poisson distribution, we might disregard the need for a covariance matrix statistic $\mathcal{C}_{t'}$ and instead use

$$P_{t'}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}, \dots) = \text{PoissonPMF}(y; \mathcal{M}_{t'}). \quad (4.6)$$

The more statistically-inclined readers may notice that the probability mass function here would require the integrals in Eq. (4.3) to be replaced with summations over the relevant domains.

Eq. (4.3) demonstrates how one can construct a statistically meaningful way to compare the sequence of real data measurements Y_{t+1}, Y_t, \dots to their modelled equivalents $\mathbf{Y}_{t+1}(z), \mathbf{Y}_t(z), \dots$. But we still haven't shown how to compute $P_{t+1}(X, z|Y)$ for a given simulation, and this can be the most challenging part. To begin with, we can reapply Bayes' rule and the chaining of conditional probability to find

$$P_{t+1}(X, z|Y) \propto P_{t+1}(Y|X)P_{t+1}(X|z)P_t(z|Y'), \quad (4.7)$$

where here Y at $t+1$ is the updated version of itself Y' at t .

The relationship between $P_{t+1}(X|z)$ and previous timesteps can be directly inferred from the probabilistic iteration formula that we introduced in the previous chapter. So we can map probabilities of $X_{0:t+1} = X$ throughout time and learned information about the state of the system can be applied from previous values, given z . But is there a similar relationship we might consider for $P_{t+1}(z|Y)$? Yes there is! The marginalisation

$$P_{t+1}(z|Y) = \int_{\Omega_{t+1}} dX P_{t+1}(X, z|Y) \propto \left[\int_{\Omega_{t+1}} dX P_{t+1}(Y|X)P_{t+1}(X|z) \right] P_t(z|Y'), \quad (4.8)$$

shows how the z updates can occur in an iterative fashion. The reader may also recognize the factor above in brackets as Eq. (4.2). To complete the picture, one can combine the X and z updates into a joint distribution update which takes the following form

$$P_{t+1}(X, z|Y) \propto P_{t+1}(Y|X)P_{(t+1)t}(x|X', z)P_t(X', z|Y'). \quad (4.9)$$

We can also marginalise this distribution over the past state history rows to get a distribution over the latest state row $X_{t+1} = x$ like this

$$P_{t+1}(x, z|Y) = \int_{\Omega_t} dX' P_{t+1}(X, z|Y) \propto \int_{\Omega_t} dX' P_{t+1}(Y|X)P_{(t+1)t}(x|X', z)P_t(X', z|Y'). \quad (4.10)$$

In the next section, we're going to discuss how to translate all of this probabilistic language into some MAP inference algorithms. Before we do this, however, it will be instructive (particularly for 'online' learning algorithms) to consider what happens if the model changes over time and z needs to change in order to better represent the real data. In such situations, we propose to apply the same formula as Eq. (4.10) but instead replace the distribution over (X', z) on the right hand side with its 'past discounted' version¹

$$\int_{\Omega_t} dX' P_t(X', z|Y') \longrightarrow \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' \beta^{t-t'} P_{t'}(x', z|Y'), \quad (4.11)$$

where $0 < \beta < 1$ and we recall the notation which considers distributions over the individual rows x' within the matrix X' in this new version. This time-dependent discount factor could be used to reduce the dependence of the update on data which is much further in the past, and hence will ultimately lead to a more responsive algorithm. This responsiveness would have to be balanced with the tradeoffs associated with discounting potentially valuable data that may offer greater long-term stability. Readers who are familiar with reinforcement learning may be starting to feel in familiar territory here — they will have to wait for the latter parts of the book to see more on discounting though!

4.2 Online learning the MAP

Eq. (4.9) tells us how to probabilistically translate the current state of knowledge about (x, z) forward through time in response to the arrival of new data. We also know how to connect the simulated measurements to the real data because Eq. (4.3) essentially gives us an objective function to maximise for each step in time. This is all great in theory; but in practice, this optimisation problem typically has several layers of difficulty to it. Since the model has been defined by its stochastically generated samples of measurements $Y_{t+1}(z), Y_t(z), \dots$, the objective function will manifestly be stochastic too. Another layer of difficulty is that gradients of the objective function are not immediately computable and so navigation around the optimisation domain could be difficult, especially in high-dimensional problems. Lastly, given that the simulation model in the stochadex needs to be running multiple times for each timestep, we need a way of mitigating computational expense.

So how should we proceed? To solve this problem in the general case, Eqs. (4.3) and (4.9) tell us we need to synthesize the following components into a single algorithm:

1. A process $P_{(t+1)t}(x|X', z)$ which iterates the state matrix of the simulation X forward in time.
2. A process $P_{t+1}(Y|X)$ which generates simulated measurements from the simulated states X .
3. A probability distribution $= P_{t'}(y; \mathcal{M}_{t'}, \mathcal{C}_{t'}, \dots)$ which represents the posterior distribution of the simulated measurement vector y given an optimised compression of the real data into summary statistics.

¹In the continuous-time version, this past-discounting factor can depend on the stepsize such that we replace

$$\beta^{t-t'} \longrightarrow \frac{1}{\beta[\delta t(t)]} \prod_{t''=t'}^t \beta[\delta t(t'')].$$

4. Some way of representing samples from the distribution $P_{t+1}(X, z|Y)$ so that their distribution can be updated and will converge towards the posterior over (X, z) .

Before writing this up, should read this paper on efficient amortized inference using neural networks with `BayesFlow` here in particular: [5]. But also, should cite other works to make amortized inference more efficient by using neural networks to learn convenient functions of the Bayes factor in Evidence networks [6].

The algorithm is specifically: 1. if this is a refit step, sample new values for (X, z) for all members of the ensemble from the current (X, z) distribution points and run the iterations for all of these ensemble members from the back of the window all the way up to the current point in time (hence the full matrix X is sampled) 2. take all of the ensemble members a step forward in time 3. approximate the mode by computing the average values of z within the q -th percentile of the sampled probability mass (where q is set by the user and is ideally $< 68\%$) — this idea comes from nested sampling 4. stream in the data for the next point in time and go to 1.

As such an algorithm converges, we can recompute (and hence iteratively improve) the MAP estimate with respect to each iteration of the posterior.

Readers with some machine learning experience may be familiar with the classic exploration vs exploitation tradeoffs. It's clear that these tradeoffs will manifest in our case here when trying to strike a balance between iterating the posterior distribution and optimizing the current posterior with respect to (X, z) to compute the MAP.

Readers of the previous section may also have recognized that Eq. (4.9) contains the same conditional probability $P_{(t+1)t}(x|X', z)$ as the reweighting algorithm. This structure enables us to reuse all of the exposition we provided for the probabilistic reweighting and highlights how the reweighting itself can be used in the algorithm to optimise the posterior.

If we now synthesize both of these observations together, we can see how a stochastic variant of the well-known Expectation-Maximisation Algorithm [7, 8, 9] naturally emerges.

Bibliography

- [1] S. A. Sisson, Y. Fan and M. Beaumont, *Handbook of approximate Bayesian computation*. CRC Press, 2018.
- [2] L. F. Price, C. C. Drovandi, A. Lee and D. J. Nott, *Bayesian synthetic likelihood*, *Journal of Computational and Graphical Statistics* **27** (2018) 1–11.
- [3] S. N. Wood, *Statistical inference for noisy nonlinear ecological dynamic systems*, *Nature* **466** (2010) 1102–1104.
- [4] C. Drovandi and D. T. Frazier, *A comparison of likelihood-free methods with and without summary statistics*, *Statistics and Computing* **32** (2022) 42.
- [5] S. T. Radev, U. K. Mertens, A. Voss, L. Ardizzone and U. Köthe, *Bayesflow: Learning complex stochastic models with invertible neural networks*, *IEEE transactions on neural networks and learning systems* **33** (2020) 1452–1466.
- [6] N. Jeffrey and B. D. Wandelt, *Evidence networks: simple losses for fast, amortized, neural bayesian model comparison*, *arXiv preprint arXiv:2305.11241* (2023) .
- [7] H. O. Hartley, *Maximum likelihood estimation from incomplete data*, *Biometrics* **14** (1958) 174–194.
- [8] A. P. Dempster, N. M. Laird and D. B. Rubin, *Maximum likelihood from incomplete data via the em algorithm*, *Journal of the royal statistical society: series B (methodological)* **39** (1977) 1–22.
- [9] K. P. Murphy, *Machine learning: a probabilistic perspective*. MIT press, 2012.