

Bayesian Synthetic Likelihood

Thanasi Bakis, Brian Schetzle

March 16, 2022

- Wood, S. N. (2010), “Statistical Inference for Noisy Nonlinear Ecological Dynamic Systems,” *Nature*.
 - Introduced *synthetic likelihood* technique in a frequentist setting
- L. F. Price, C. C. Drovandi, A. Lee & D. J. Nott (2018), “Bayesian Synthetic Likelihood”, *Journal of Computational and Graphical Statistics*.
 - Extended to *Bayesian synthetic likelihood* technique in a Bayesian setting

Motivating the synthetic likelihood

At a high level, the synthetic likelihood is a replacement for an intractable true likelihood. Synthetic likelihood originates in a *frequentist setting*, when the likelihood function is too irregular to easily maximize (analytically or numerically). Maximizing the synthetic likelihood yields a point estimator for your parameters.

It is trivial to extend this to a Bayesian setup by placing a prior on the parameters, obtaining a posterior distribution of parameters instead of point estimates.

The core requirement of the synthetic likelihood method is that we can still generate samples from the true likelihood.

Upcoming motivating example: Ricker population model

Motivating the synthetic likelihood

In the Ricker model, N_t models the time course of a population, where

$$N_0 = 1$$

$$N_t = rN_{t-1}e^{-N_{t-1}+e_{t-1}}$$

$$e_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

In other words, $N_t|N_{t-1} \sim \text{LogNormal}(\log N_{t-1} + \log r - N_{t-1}, \sigma^2)$.

Additionally, suppose N_t is not actually observed, but a sample from the population $Y_t|N_t \sim \text{Poisson}(\phi N_t)$ is observed.

Our parameters are thus: r , a population growth rate parameter; σ , controlling random noise; and ϕ , a scaling parameter for sampling from the population.

Motivating the synthetic likelihood

To do Bayesian inference, we need:

$$p(r, \sigma, \phi | \mathbf{Y}) \propto p(r, \sigma, \phi) p(\mathbf{Y} | r, \sigma, \phi)$$

.

Can we obtain the likelihood $p(\mathbf{Y} | r, \sigma, \phi)$ required?

Motivating the synthetic likelihood

Let $\theta = (r, \sigma, \phi)$. The joint likelihood with observed and latent variables is easier to work out:

$$\begin{aligned} p(\mathbf{Y}, \mathbf{N} | \theta) &= p(\mathbf{N} | \theta) p(\mathbf{Y} | \mathbf{N}, \theta) \\ &= p(N_1 | \theta) \prod_{t=2}^n p(N_t | N_{t-1}, \theta) \prod_{t=1}^n p(Y_t | N_t, \theta) \\ &= \text{LogNormal}(N_1; \dots) \prod_{t=2}^n \text{LogNormal}(N_t; \dots) \prod_{t=1}^n \text{Poisson}(Y_t; \dots) \end{aligned}$$

Motivating the synthetic likelihood

We would need to marginalize over all N_t :

$$p(\mathbf{Y}|r, \sigma, \phi) = \int_{N_1} \dots \int_{N_n} p(\mathbf{Y}, \mathbf{N}|\theta) dN_1 \dots dN_n$$

This is expensive, in part because the number of integrals grows with the number of data points observed.

Alternatively, instead of having the posterior $p(r, \sigma, \phi|\mathbf{Y})$, we could also try to do inference on the latent population variables: $p(N_1, \dots, N_n, r, \sigma, \phi|\mathbf{Y})$.

However, this would require the number Metropolis-Hastings proposals in each iteration to grow with the number of observed data points too (proposing values of each N_t), which is undesirable.

Bayesian synthetic likelihood

What if, instead of using $p(\mathbf{Y}|\theta)$, we use:

$$p(\mathbf{s}_Y|\theta)$$

where \mathbf{s}_Y is a vector of summary statistics for \mathbf{Y} , eg. the mean, quantiles, etc.

Then, instead of targeting the posterior $p(\theta|\mathbf{Y})$ in our MCMC, we can instead target $p(\theta|\mathbf{s}_Y)$.

(This is similar to obtaining a maximum synthetic likelihood estimator instead of an MLE in the frequentist setting.)

Of course, depending on your choice of statistics, this *synthetic likelihood* may not be tractable either. This method consequently makes a normality assumption:

$$\mathbf{s}_Y|\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta)$$

where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ change with θ .

Ideally, your statistics truly are normally distributed, but this may not be the case. We will see examples of both cases.

$$\mathbf{s}_{\mathbf{Y}} \sim \mathcal{N}(\boldsymbol{\mu}(\boldsymbol{\theta}), \boldsymbol{\Sigma}(\boldsymbol{\theta}))$$

We still need to be able to evaluate this likelihood, though.

Under the assumption that we can sample \mathbf{Y} from the real likelihood, we can estimate the parameters of the synthetic likelihood via Monte Carlo approximation. . .

Suppose we have a proposed value of $\boldsymbol{\theta}$. Let $\mathbf{Y}_1^*, \dots, \mathbf{Y}_n^* \stackrel{\text{iid}}{\sim} p(\mathbf{Y}|\boldsymbol{\theta})$. In otherwords, given a proposed value of $\boldsymbol{\theta}$, generate n iid datasets from the distribution that hypothetically also generated the observed data \mathbf{Y} .

Bayesian synthetic likelihood

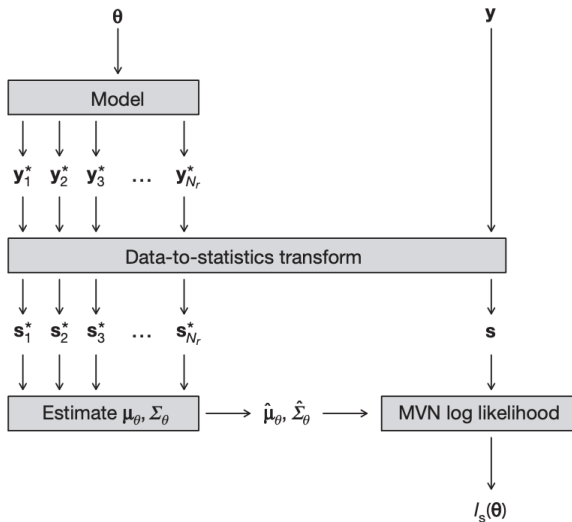
Now, for each simulated dataset \mathbf{Y}_i^* , calculate the corresponding summary statistics $\mathbf{s}_{\mathbf{Y}_i^*}$ in the same manner as $\mathbf{s}_{\mathbf{Y}}$.

Then,

$$\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}} = \frac{1}{n} \sum_{i=1}^n \mathbf{s}_{\mathbf{Y}_i^*}$$
$$\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{s}_{\mathbf{Y}_i^*} - \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}})(\mathbf{s}_{\mathbf{Y}_i^*} - \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}})^{\top}$$

With these estimates, we can evaluate the synthetic likelihood of various proposed values of $\boldsymbol{\theta}$ for our observed summary statistics $\mathbf{s}_{\mathbf{Y}}$.

Bayesian synthetic likelihood



A toy example

Consider the model:

$$\begin{aligned} Y_i | \lambda &\stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda) & i = 1, \dots, 100 \\ \lambda &\sim \text{Gamma}(\alpha = 0.001, \beta = 0.001) \end{aligned}$$

Suppose we generate observations Y_1, \dots, Y_{100} using $\lambda = 30$ and want to conduct inference on λ .

We want to find $p(\lambda | \mathbf{Y}) \propto p(\mathbf{Y} | \lambda)p(\lambda)$ without evaluating $p(\mathbf{Y} | \lambda)$.

A toy example

Note:

$$\begin{aligned} p(\lambda | \mathbf{Y}) &\propto \left[\prod_{i=1}^{100} \frac{\lambda^{Y_i}}{Y_i!} e^{-\lambda} \right] \left[\frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \right] \\ &\propto \lambda^{\alpha + \sum_{i=1}^{100} Y_i - 1} e^{-\lambda(\beta + n)} \\ &\sim \text{Gamma}(\alpha = 0.001 + \sum_{i=1}^{100} Y_i, \beta = 100.001) \end{aligned}$$

So, in this toy example, the posterior distribution is known analytically.

A toy example

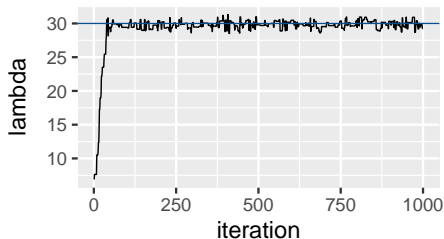
How should we choose a statistic? The paper uses the mean:

$$s_Y = \frac{1}{100} \sum_{i=1}^{100} Y_i$$

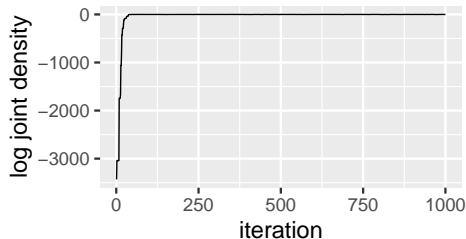
This is the sufficient statistic for the Poisson distribution; all the information contained in the data is also contained in this statistic. Also, by the central limit theorem, the distribution of the mean of a Poisson sample can be adequately approximated by a normal distribution, so synthetic likelihood should perform well in this setting.

A toy example

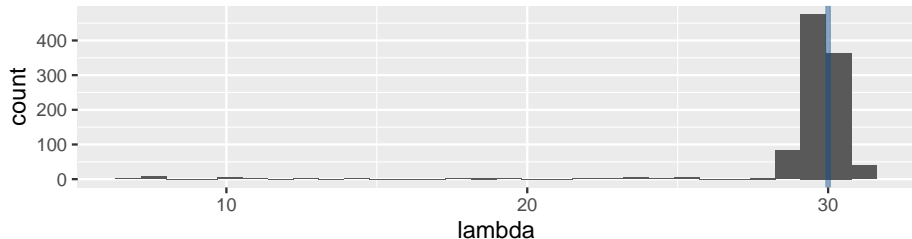
Traceplot (samples)



Traceplot (log joint density)

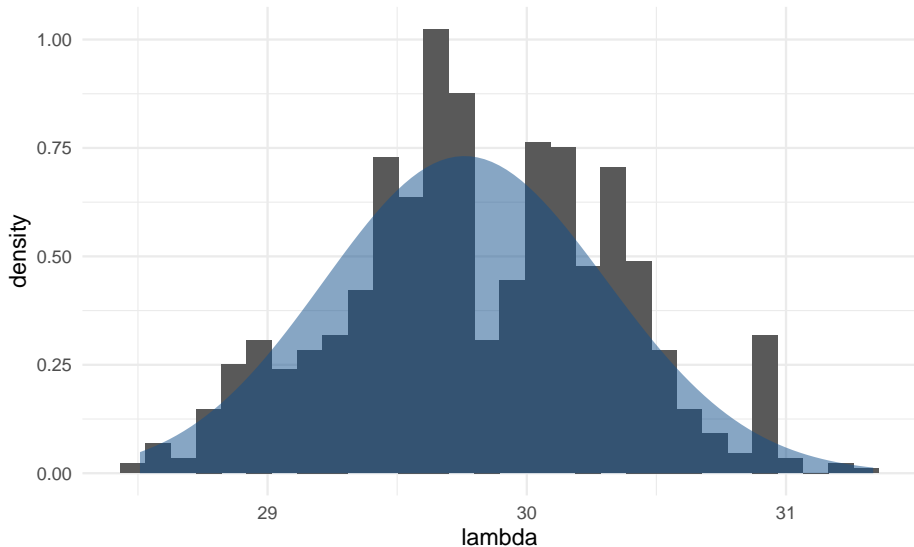


Posterior Histogram



A toy example

Posterior samples (grey) with analytic posterior (blue)



Further exploration

We were interested in trying other statistics to see how well synthetic likelihood performed. This was our own exploration and was not addressed in the paper.

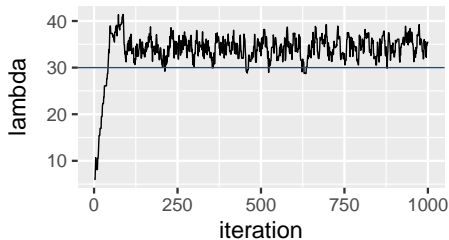
We started with the maximal statistic:

$$s_Y = \max(\mathbf{Y})$$

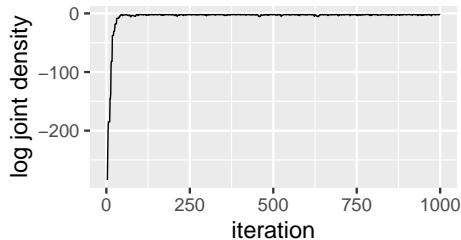
This is not a sufficient statistic for Poisson data and also not approximately normally distributed across many samples. We would expect the synthetic likelihood method to have a harder time identifying the true analytic posterior.

Further Exploration

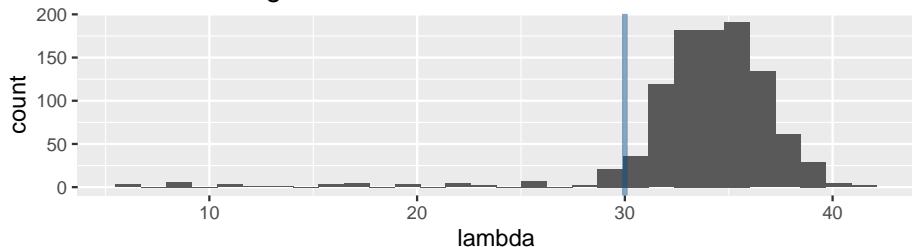
Traceplot (samples)



Traceplot (log joint density)

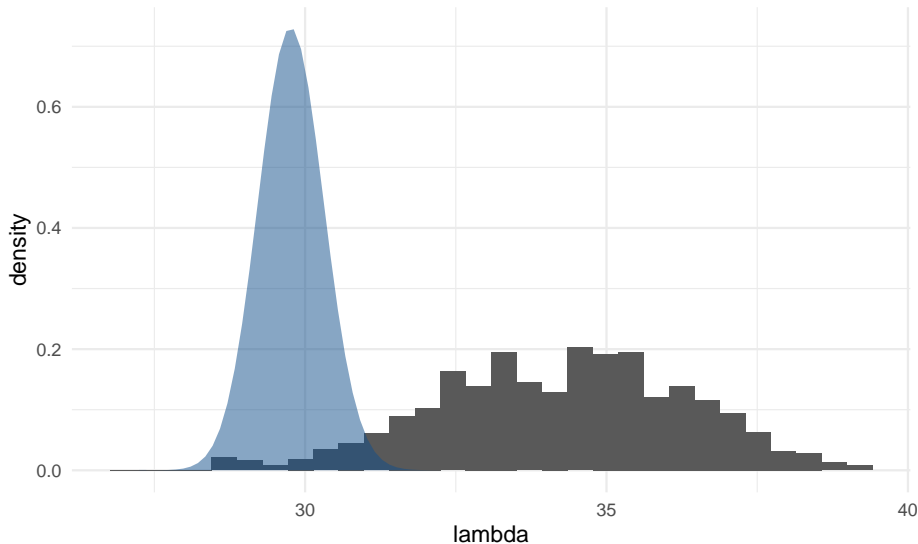


Posterior Histogram



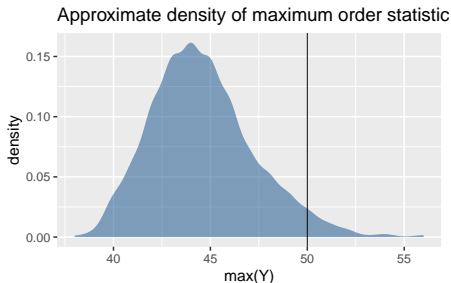
Further exploration

Posterior samples (grey) with analytic posterior (blue)



Further exploration

The synthetic likelihood approximation overestimated the true lambda. This is because the data had an unusually large maximal value.



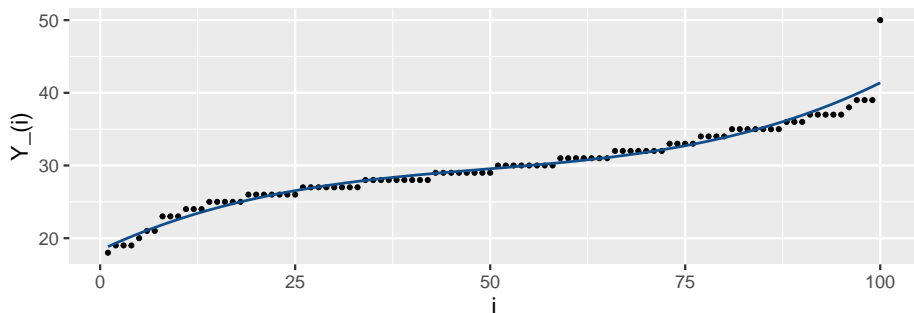
This illustrates how sensitive this method is to the choice of statistic. The approximation will perform optimally on sufficient statistics, but with complicated models, it may not be possible to identify the sufficient statistics. Therefore, statistics must be chosen which capture the important features of the model.

Further exploration

As a final choice of statistics, we used the coefficients of a polynomial regression on the ordered data.

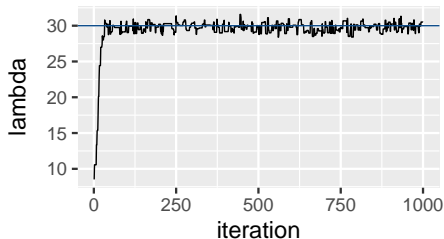
$$Y_{(i)} = \beta_0 + \beta_1 i + \beta_2 i^2 + \beta_3 i^3$$

Plotting $Y_{(i)}$ and the polynomial regression (blue)

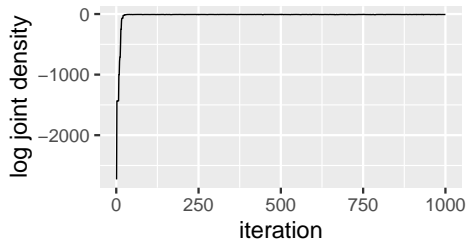


Further exploration

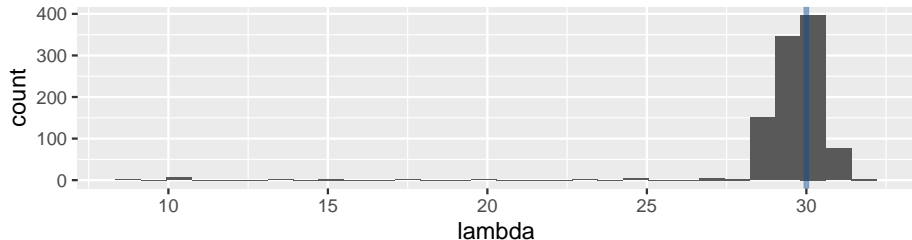
Traceplot (samples)



Traceplot (log joint density)



Posterior Histogram



Ricker population model example

We now return to the Ricker model from the beginning of the presentation to provide a more complex use case.

Some things to consider about the Ricker model vs the toy Poisson model:

- The sufficient statistics are not easily available
- We will need to rely on more complex, non-normal statistics

Ricker population model example

As a refresher, recall we observe a sample Y_t of a population N_t :

$$N_t | N_{t-1} \sim \text{LogNormal}(\log N_{t-1} + \log r - N_{t-1}, \sigma^2)$$
$$Y_t | N_t \sim \text{Poisson}(\phi N_t)$$

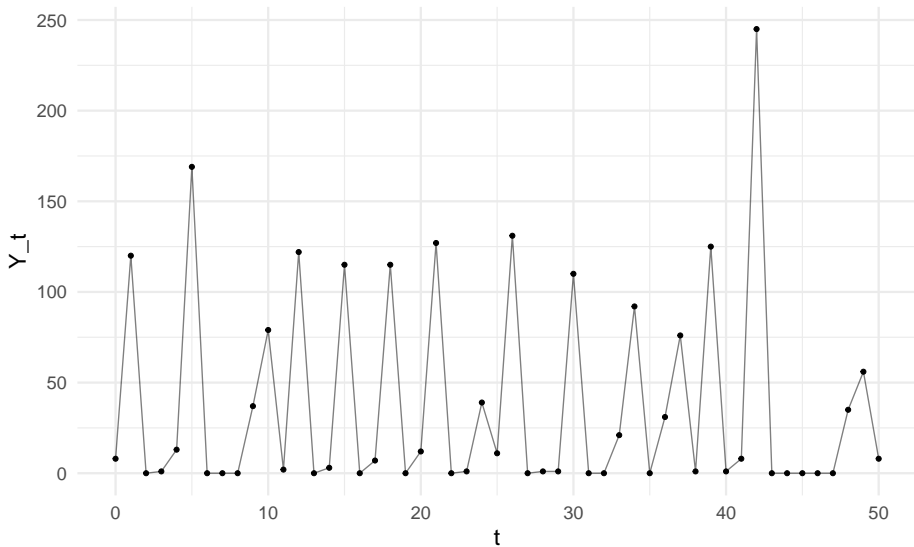
Suppose we observe Y_1, \dots, Y_{50} with $(\log r, \sigma, \phi) = (3.8, 0.3, 10)$.

To use synthetic likelihood, we need to ensure we can sample simulated \mathbf{Y} from the model. Thankfully, sampling is straightforward; using $N_0 = 1$, you can sample N_1, N_2, \dots sequentially, then sample Y_t for each N_t .

We place an uninformative prior over all parameters.

Ricker population model example

Visualizing observed data over time



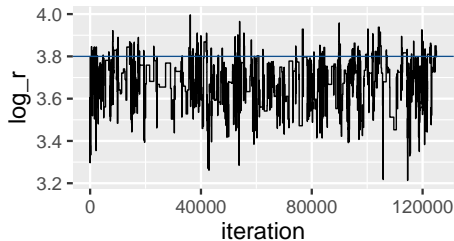
Ricker population model example

This data is noisy and runs over time. What sorts of statistics will we use?

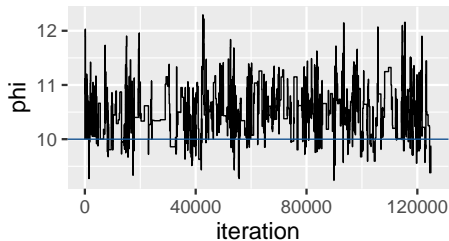
- **Marginal distribution statistics:** to summarize the “shape” of the marginal distribution
 - Mean: \bar{Y}
 - Number of zeros: $\sum_{t=1}^{50} \mathbf{1}_{\{0\}}(Y_t)$
- **Dynamic process statistics:** characterize the relationship between Y_t and Y_{t-1} (and possibly more history)
 - Autoregressive model coefficients: $Y_t^{0.3} \sim Y_{t-1}^{0.3} + Y_{t-1}^{0.6}$
 - Note the exponents were tuned to improve fit
 - Coefficients of regression on ordered differences:
 $(Y_t - Y_{t-1}) \sim Y_t + Y_t^2 + Y_t^3$
- **Time series statistics:** sensitive to the shape and period of fluctuations
 - Coefficients of the autocovariance function, up to lag 5

Ricker population model example

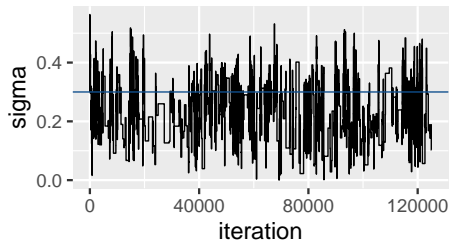
Traceplot (samples)



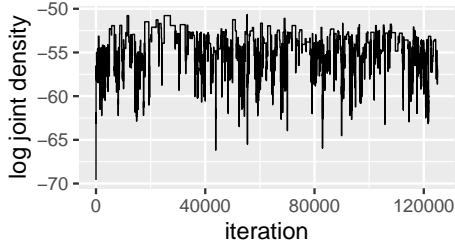
Traceplot (samples)



Traceplot (samples)

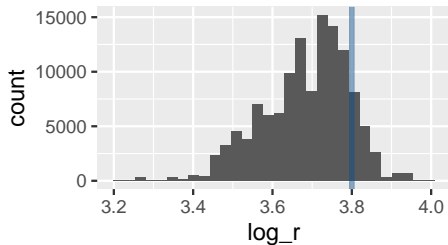


Traceplot (log joint density)

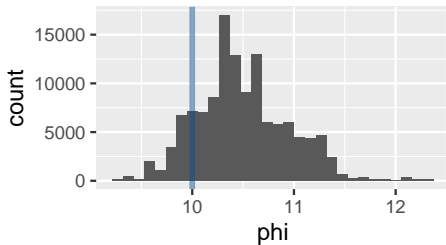


Ricker population model example

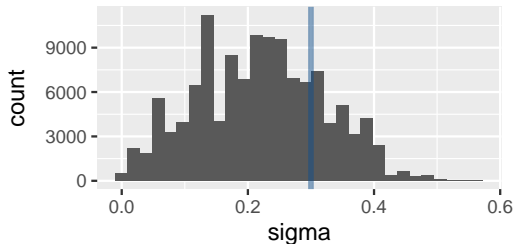
Posterior Histogram



Posterior Histogram



Posterior Histogram



comparison to ABC (or maybe this could go earlier?)

Synthetic Likelihood, Wood (2010), we can get rid of this or move it

- Statistical Problem: Fitting a model to data from chaotic biological systems using maximum likelihood can yield estimates that may be too sensitive to random noise
- Model Choices: This paper is inherently frequentist in nature and uses a parametric Ricker model of population growth over time as an example
- Computational Tools: Wood proposes a multivariate normal approximation to a vector of summary statistics and demonstrates its efficacy using Markov Chain Monte Carlo
- Other Approaches Available: Standard methodology up to that point was to use likelihood based approaches for statistical inference on model parameters but these methods were known to perform poorly on chaotic systems

Paper pros and cons