Monte Carlo and Empirical Methods for Stochastic Inference (MASM11/FMSN50)

Magnus Wiktorsson Centre for Mathematical Sciences Lund University, Sweden

Lecture 12 MCMC for Bayesian computation II Feb 27, 2020

Plan of today's lecture

Last time: Stochastic modeling and Bayesian inference

- 2 Example: Change point detection
- Interlude: Mixing of MCMC samplers
- 4 Example: Korsbetningen

We are here $\longrightarrow \bullet$

Last time: Stochastic modeling and Bayesian inference

2 Example: Change point detection

Interlude: Mixing of MCMC samplers

4 Example: Korsbetningen

The frequentist approach

The frequentist approach to statistics is characterized as follows.

• Data y is viewed as an observation of a random variable Y with distribution \mathbb{P}_0 which most often is assumed to be a member of an exponential family

$$\mathcal{P} = \{ \mathbb{P}_{\theta}; \theta \in \Theta \}.$$

- Estimates $\widehat{\theta}(y)$ are realizations of random variables.
- $\bullet\,$ Confidence intervals are calculated to cover the true value in, say, $95\%\,$ of the cases.
- Hypothesis testing is done by rejecting a hypothesis \mathcal{H}_0 if $\mathbb{P}(\text{data} \| \mathcal{H}_0)$ is small.

The Bayesian approach

Briefly, the Bayesian approach to statistics is as follows.

- The parameter θ is viewed as a random variable, and inference is based completely on the posterior distribution $f(\theta|y)$.
- It is possible to incorporate prior information in terms of the prior distribution $f(\theta)$.
- A 95% credible or posterior probability interval contains θ with a probability of 95% given the observations.
- Hypothesis tests are done by studying $\mathbb{P}(\mathcal{H}_0 \| data)$.

Last time: Stochastic modeling and Bayesian inference

Example: Change point detection Interlude: Mixing of MCMC samplers Example: Korsbetningen

The frequentist vs the Bayesian approach ;-)



We are here $\longrightarrow \bullet$

Last time: Stochastic modeling and Bayesian inference

2 Example: Change point detection

Interlude: Mixing of MCMC samplers

4 Example: Korsbetningen

Example: Change point detection

Consider the following model.

- We have measured the waiting times in a system and suspect that the expected waiting time changed during the monitoring period.
- The observations y_i for i = 1, ..., n are assumed to follow exponential distributions with parameter θ_1 for $i \in \{1, ..., n_b\}$ and parameter θ_2 for $i \in \{n_b + 1, ..., n\}$.
- Further, we put a Gamma prior on θ_k , $\theta_k \sim \Gamma(a,b)$, with a=40 and b=4, and a uniform prior on $n_{\rm b}$



Example: Change point detection (cont.)

Thus, we have unknown parameters $(\theta_1, \theta_2, n_b)$ and data $Y = (y_1, \dots, y_n)$. The posterior becomes

Example: Change point detection (cont.)

The conditional distributions of $heta_1$ and $heta_2$ are easily calculated according to

$$\theta_1 | n_{\mathbf{b}}, \theta_2, y_1, \dots, y_n \sim \Gamma\left(n_{\mathbf{b}} + a, b + \sum_{i=1}^{n_{\mathbf{b}}} y_i\right),$$

$$\theta_2 | n_{\mathbf{b}}, \theta_1, y_1, \dots, y_n \sim \Gamma\left(n - n_{\mathbf{b}} + a, b + \sum_{i=n_{\mathbf{b}}+1}^n y_i\right).$$

The conditional distribution of $n_{\rm b}$ is however more complicated:

$$f(n_{\mathbf{b}}|\theta_1, \theta_2, y_1, \dots, y_n)$$

$$\propto \theta_1^{n_{\mathbf{b}}} \exp\left(-\theta_1 \sum_{i=1}^{n_{\mathbf{b}}} y_i\right) \theta_2^{-n_{\mathbf{b}}} \exp\left(-\theta_2 \sum_{i=n_{\mathbf{b}}+1}^n y_i\right)$$

Example: Change point detection (cont.)

Thus, we sample the posterior of $(\theta_1, \theta_2, n_b)$ using a Gibbs sampler with a MH step for n_b , yielding a hybrid sampler. The MH step goes as follows. Given n_b , we propose a candidate n_b^* uniformly on the integers $\{n_b - R, \ldots, n_b, \ldots, n_b + R\}$, for some R. This forms a symmetric proposal on $\{1, \ldots, n\}$.

The acceptance probability for the MH step thus becomes

$$\begin{aligned} \alpha(n_{\mathbf{b}}, n_{\mathbf{b}}^{*}) \\ &= 1 \wedge \frac{\theta_{1}^{n_{\mathbf{b}}^{*}} \theta_{2}^{-n_{\mathbf{b}}^{*}} \exp(-\theta_{1} \sum_{i=1}^{n_{\mathbf{b}}^{*}} y_{i}) \exp(-\theta_{2} \sum_{i=n_{\mathbf{b}}^{*}+1}^{n} y_{i})}{\theta_{1}^{n_{\mathbf{b}}} \theta_{2}^{-n_{\mathbf{b}}} \exp(-\theta_{1} \sum_{i=1}^{n} y_{i}) \exp(-\theta_{2} \sum_{i=n_{\mathbf{b}}+1}^{n} y_{i})}. \end{aligned}$$

Example: Change point detection (cont.)

Running this Gibbs sampler with R = 75 gives an acceptance rate of 33%.



Example: Change point detection (cont.)

The resulting histograms of the parameters are as follows:



Selecting priors

The posterior is computed via Bayes's formula

$$f(\theta|y) = \frac{f(y|\theta)f(\theta)}{\int f(y|\theta')f(\theta') \,\mathrm{d}\theta'} \propto f(y|\theta)f(\theta).$$

In Bayesian modeling there is always an interplay between the prior and the data:

- The posterior is drawn away from the data towards the prior. How far depends on the strength of the prior.
- However, enough data will most likely overwhelm the prior.

Two common prior-types are

- conjugate priors.
- improper (flat) priors.

Conjugate priors

Conjugate priors

- are such that the prior and the posterior belong to the same distribution class for a given likelihood.
- allow theoretical calculations and Gibbs sampling.
- are sometimes criticized since we select priors for ease of calculation. However, often the parameters are flexible enough to allow for a reasonable prior.
- may be hard to derive for complex models.

Conjugate priors

Conjugate priors for θ for some common likelihoods. All parameters except θ are assumed fixed and known and data (y_1, \ldots, y_n) is assumed to be conditionally independent given θ .

Likelihood	Prior	Posterior
$Bin(n,\theta)$	Beta(lpha,eta)	Beta(lpha+y,eta+n-y)
$Ge(\theta)$	Beta(lpha,eta)	$Beta(lpha+n,eta+\sum_{i=1}^n y_i-n)$
$NegBin(n, \theta)$	$Beta(\alpha,\beta)$	Beta(lpha+n,eta+y-n)
Gamma(k, heta)	Gamma(lpha,eta)	$Gamma(lpha+nk,eta+\sum_{i=1}^ny_i)$
$Po(\theta)$	Gamma(lpha,eta)	$Gamma(lpha + \sum_{i=1}^{n} y_i, \beta + n)$
$\mathcal{N}(\mu, \theta^{-1})$	Gamma(lpha,eta)	Gamma $\left(lpha+rac{n}{2},eta+rac{1}{2}\sum_{i=1}^n(y_i-\mu)^2 ight)$
$\mathcal{N}(\theta,\sigma^2)$	$\mathcal{N}(m,s^2)$	$\mathcal{N}\left(rac{m/s^2+nar{y}/\sigma^2}{1/s^2+n/\sigma^2},rac{1}{1/s^2+n/\sigma^2} ight)$

Improper priors

Improper, or flat, priors are used when prior information is deficient.

For instance, if $\theta \in \mathbb{R}$, $f(\theta) \propto 1$ is an improper prior since it is not integrable; however, we allow this as long as the posterior is a well-defined density.

For instance, let y be an observation from $Y \sim \mathcal{N}(\theta, 1)$, where $\theta \in \mathbb{R}$. Since we do not have any prior information concerning θ we put $f(\theta) \propto 1$ for all $\theta \in \mathbb{R}$. After this we proceed, formally, like

$$\begin{split} f(\theta|y) &= \frac{f(y|\theta)f(\theta)}{\int f(y|\theta')f(\theta')\,\mathrm{d}\theta'} = \frac{\mathcal{N}(y;\theta,1)\cdot 1}{\int \mathcal{N}(y;\theta',1)\cdot 1\,\mathrm{d}\theta'} \\ & \stackrel{\mathrm{symm.}}{=} \frac{\mathcal{N}(\theta;y,1)\cdot 1}{\int \mathcal{N}(\theta';y,1)\cdot 1\,\mathrm{d}\theta'} = \mathcal{N}(\theta;y,1). \end{split}$$

We are here $\longrightarrow \bullet$

Last time: Stochastic modeling and Bayesian inference

2 Example: Change point detection

Interlude: Mixing of MCMC samplers

4 Example: Korsbetningen

Mixing of MCMC samplers

We recall that the asymptotic variance of au_N is given by

$$\sigma^2 = r(0) + 2\sum_{\ell=1}^{\infty} r(\ell) \quad \text{with} \quad r(\ell) = \lim_{n \to \infty} \mathbb{C}(\phi(X_{n+\ell}), \phi(X_n)).$$

Some remarks:

- Consequently, in order to obtain a low variance of τ_N , the covariance function $r(\ell)$ should decrease rapidly with ℓ .
- For geometrically ergodic chains $r(\ell)$ tends to zero geometrically fast.
- The speed at which $r(\ell)$ tends to zero is typically described using the term mixing.
 - Strong mixing = fast forgetting = rapidly decreasing $r(\ell)$.
 - Bad mixing = slow forgetting = slowly decreasing $r(\ell)$.

Why is good mixing important?

Bad choices of proposal distributions may lead to bad mixing, implying high variance and the need for a large MC sample size to ensure good estimates. This causes problems for the MCMC algorithm in the sense that it may

- fail to converge altogether or
- need a very long time to converge.

Optimal mixing

Let us focus for a while on the MH algorithm.

- When designing a random walk proposal, $X_k^* = X_k + \epsilon$ with $\epsilon \sim \mathcal{N}(0, s\Sigma)$, two things will effect the acceptance rate:
 - \blacksquare How well Σ captures the dependence of the target distribution.
 - 2 How appropriate the scaling s > 0 is.
- One way to obtain a covariance matrix Σ that captures well the dependence structure of the target distribution f(x) is to let

$$\boldsymbol{\Sigma}_{ij} = \frac{2.38}{d} \left(-\frac{\partial^2 \log f(x)}{\partial x_i \partial x_j} \bigg|_{x=x_{\text{mode}}} \right)^{-1}$$

• Rule of thumb: A good acceptance rate is around 30% (23%-44%)!

Mixing—Random walk proposal

Using symmetric normal proposal with three different values for s (small, medium, large, respectively) yields the following trajectories:



Mixing—Random walk proposal

Correlation function for the three chains:



We are here $\longrightarrow \bullet$

Last time: Stochastic modeling and Bayesian inference

- 2 Example: Change point detection
- Interlude: Mixing of MCMC samplers
- 4 Example: Korsbetningen

Example: Korsbetningen



Latin inscription:

Anno Domini MCCCLXI feria III post Jacobi ante portas Visby in manibus Danorum ceciderunt Gutenses, hic sepulti, orate pro eis!

On the third day after saint Jacob, in the year of our lord 1361, the Goths fell outside the gates of Visby at the hands of the Danish. They are buried here. Pray for them!

Example: Korsbetningen—background

The background is the following.

- In 1361 the Danish king Valdemar Atterdag conquered Gotland and captured the rich Hanseatic town of Visby.
- Most of the defenders were killed in the attack and are buried in a field, Korsbetningen, outside of the walls of Visby.
- In 1929–1930 the gravesite (with several graves) was excavated. In grave one e.g. a total of 493 femurs, 237 right and 256 left, were found.
- We want to estimate the number of buried bodies.
- An interesting tv-programme about the event can be found at https://www.youtube.com/watch?v=IGY2OqMXF9w



Example: Korsbetningen-model

We set up the following model.

- Assume that the numbers y_1 and y_2 of left resp. right legs are two observations from a Bin(n,p) distribution.
- Here n is the total number of people buried and p is the probability of finding a leg, left or right, of a person.
- \bullet We put a conjugate $\mbox{Beta}(a,b)\mbox{-prior}$ on p and a $\mathcal{U}(256,2500)$ prior on



Example: Korsbetningen—a hybrid MCMC

We proceed as follows:

• A standard Gibbs step for

$$p|n, y_1, y_2 \sim \mathsf{Beta}(a + y_1 + y_2, b + 2n - (y_1 + y_2)).$$

- MH for n, with a symmetric proposal obtained by drawing, given n, a new candidate n* among the integers {n - R,...,n,...,n + R}.
- The acceptance probability becomes

$$\alpha(n,n^*) = 1 \wedge \frac{(1-p)^{2n^*}(n^*!)^2(n-y_1)!(n-y_2)!}{(1-p)^{2n}(n!)^2(n^*-y_1)!(n^*-y_2)!}.$$

Example: Korsbetningen—a hybrid MCMC

