

# Notes on Metropolis Hastings

Matthias Seeger  
Department of EECS  
University of California at Berkeley  
485 Soda Hall, Berkeley CA  
*mseeger@cs.berkeley.edu*

August 5, 2005

## 1 The Metropolis Hastings Algorithm

Let  $\pi(\mathbf{x})$  be the density of a distribution we would like to draw samples from. A *Markov Chain Monte Carlo (MCMC)* method does this by running a Markov chain with a transition kernel  $T(\mathbf{x}^*|\mathbf{x})$  (being a conditional probability density) which leaves  $\pi$  invariant, in the sense that if  $\mathbf{x} \sim \pi$ ,  $\mathbf{x}^* \sim T(\cdot|\mathbf{x})$ , then  $\mathbf{x}^* \sim \pi$ . If the chain is also ergodic and nonperiodic (see relevant literature), the marginal distributions of variables down the chain will converge to  $\pi$ .

A simple sufficient condition for invariance is *detailed balance*:

$$T(\mathbf{x}^*|\mathbf{x})\pi(\mathbf{x}) = T(\mathbf{x}|\mathbf{x}^*)\pi(\mathbf{x}^*) \quad \text{for all } \mathbf{x}, \mathbf{x}^*.$$

The Metropolis Hastings procedure is a general way of constructing kernels which fulfil the detailed balance condition. Let  $q(\mathbf{x}^*|\mathbf{x})$  be an arbitrary proposal distribution and define

$$\alpha(\mathbf{x}^*, \mathbf{x}) = \min \left\{ 1, \frac{q(\mathbf{x}|\mathbf{x}^*)\pi(\mathbf{x}^*)}{q(\mathbf{x}^*|\mathbf{x})\pi(\mathbf{x})} \right\}.$$

Then the MH kernel  $T_q(\cdot|\mathbf{x})$  samples  $\mathbf{x}' \sim q(\cdot|\mathbf{x})$ , evaluates  $\alpha = \alpha(\mathbf{x}', \mathbf{x})$  and sets  $\mathbf{x}^* = \mathbf{x}'$  with probability  $\alpha$  (acceptance of the proposal  $\mathbf{x}'$ ),  $\mathbf{x}^* = \mathbf{x}$  otherwise (rejection).

In order to prove detailed balance, let

$$A(\mathbf{x}', \mathbf{x}) = \alpha(\mathbf{x}', \mathbf{x})q(\mathbf{x}'|\mathbf{x})\pi(\mathbf{x}) = \min \{q(\mathbf{x}'|\mathbf{x})\pi(\mathbf{x}), q(\mathbf{x}|\mathbf{x}')\pi(\mathbf{x}')\}.$$

Importantly,  $A$  is symmetric:  $A(\mathbf{x}', \mathbf{x}) = A(\mathbf{x}, \mathbf{x}')$ . Now,

$$T_q(\mathbf{x}^*|\mathbf{x})\pi(\mathbf{x}) = \int A(\mathbf{x}', \mathbf{x})\delta_{\mathbf{x}^*, \mathbf{x}'} + B(\mathbf{x}', \mathbf{x})\delta_{\mathbf{x}^*, \mathbf{x}} d\mathbf{x}' = A(\mathbf{x}^*, \mathbf{x}) + \kappa(\mathbf{x})\delta_{\mathbf{x}^*, \mathbf{x}}$$

with  $B = q\pi - A$ . The r.h.s. is symmetric, which proves detailed balance.

The procedure originally suggested by Metropolis *et.al.* was restricted to symmetric proposal distributions in which case the acceptance probability becomes  $\alpha = \min\{1, \pi(\mathbf{x}^*)/\pi(\mathbf{x})\}$ . The clear advantage of using a symmetric proposal distribution is that  $q(\mathbf{x}^*|\mathbf{x})$  and  $q(\mathbf{x}|\mathbf{x}^*)$  do not have to be computed explicitly. If the proposal is not symmetric, but is still hard to compute, a random MH scheme might work. This is the topic of the next Section.

## 1.1 Random Choice of Proposal Distribution

Suppose there is a family  $q(\cdot|\mathbf{x}, \beta)$  of proposal distributions (the family is indexed by  $\beta$  which is a parameter independent of  $\mathbf{x}$ ) and we would like to choose one of them at random according to  $P(\beta|\mathbf{x})$ . In general we can use the MH kernel with the marginal proposal distribution  $\int q(\cdot|\mathbf{x}, \beta)P(\beta|\mathbf{x})d\beta$ , but this marginal may be hard to compute. If it is true that

$$q(\mathbf{x}'|\mathbf{x}, \beta)P(\beta|\mathbf{x})\pi(\mathbf{x}) = q(\mathbf{x}'|\mathbf{x}, \beta)P(\beta|\mathbf{x}')\pi(\mathbf{x})$$

for all  $\mathbf{x}, \mathbf{x}', \beta$ , we can also simply sample  $\beta$  along with  $\mathbf{x}'$  and use the conditional acceptance probability  $\alpha(\mathbf{x}', \mathbf{x}, \beta)$  where  $\beta$  is simply plugged in. To prove detailed balance, define  $A(\mathbf{x}', \mathbf{x}, \beta) = \alpha(\mathbf{x}', \mathbf{x}, \beta)q(\mathbf{x}'|\mathbf{x}, \beta)\pi(\mathbf{x})$ , and note that  $A(\mathbf{x}', \mathbf{x}, \beta)P(\beta|\mathbf{x})$  is symmetric for every  $\beta$  by the assumption on  $P(\beta|\mathbf{x})$ . Therefore, the marginal  $\tilde{A}(\mathbf{x}', \mathbf{x}) = \int A(\mathbf{x}', \mathbf{x}, \beta)P(\beta|\mathbf{x})d\beta$  is symmetric, which implies detailed balance. The assumption on  $P(\beta|\mathbf{x})$  is of course satisfied if  $P(\beta|\mathbf{x})$  does not depend on  $\mathbf{x}$ . A more interesting and practically useful case is that  $P(\beta|\mathbf{x}) = P(\beta_2|f(\mathbf{x}, \beta_1))P(\beta_1)$  for  $\beta = (\beta_1, \beta_2)$  and some mapping  $f$  such that whenever a transition from  $\mathbf{x}$  to  $\mathbf{x}'$  is possible under  $q(\cdot|\mathbf{x}, \beta)$ ,  $\beta \sim P(\cdot|\mathbf{x})$ , then  $f(\mathbf{x}, \beta_1) = f(\mathbf{x}', \beta_1)$ .

An example of this setup can be found in [1]. There,  $\mathbf{x}$  consists of a partition of  $n$  observations into groups.  $\mathbf{x}$  is to be updated using split and merge steps. For a group picked as a split candidate, the authors suggest to run Gibbs sampling on the group indicators, based on a “posterior” which is restricted to the data within the group. After a while, the new state is taken as MH proposal. Clearly this is not a symmetric proposal, and if each indicator is updated more than once, then it is hard in general to compute the proposal probabilities. To overcome this problem, the authors use random MH, where all Gibbs updates except for a single final run over the group indicators are lumped into a random proposal. In other words, the family of proposals is indexed by states  $\hat{\mathbf{x}}$  which agree with  $\mathbf{x}$  outside the group, and  $P(\beta|\mathbf{x})$  is implemented using restricted Gibbs sampling. States picked in this way are referred to as *launch states*. Using our notation,  $\beta_1$  is a pair  $k, l$  of observations chosen at random, and a split is proposed if these observations belong to the same group (a merge is proposed otherwise).  $f(\mathbf{x}, \beta_1)$  is the set of all observations assigned to the same group(s) than  $k, l$ , this set is called  $S$  in [1]. If  $k, l$  belong to the same group, the proposal is to split the corresponding group. In this case,  $\beta_2$  is the launch state obtained using restricted Gibbs sampling on the points in  $S \setminus \{k, l\}$ . The final proposal starts from  $\beta_2$  for a single restricted Gibbs run over these points, for which the proposal probabilities can be computed easily.<sup>1</sup> If  $k, l$  belong to different groups, the proposal is to merge them into a single one. Note that even in this case we have to actually sample  $\beta_2$  in order to compute the MH acceptance probability, namely because  $q(\mathbf{x}|\mathbf{x}^*, \beta)$  appears in the MH ratio.

Since the launch state  $\beta_2$  depends on  $\mathbf{x}$  only through  $S = f(\mathbf{x}, \beta_1)$  which is clearly the same as  $S' = f(\mathbf{x}', \beta_1)$  (*i.e.* a split or merge for a given  $\beta_1$  leaves the set  $S$  invariant, as long as the same  $k, l$  are picked), we can use MH without having to compute the marginal proposal distribution (which would be intractable in the case of [1])<sup>2</sup>.

If  $P(\beta|\mathbf{x})$  does not fulfil the symmetry property and the marginal proposal distribution cannot be computed easily, we have the option of extending the state space to include  $\beta$

<sup>1</sup>Since there are no latent variables which have to be summed out.

<sup>2</sup>Unfortunately, Jain and Neal [1] do not mention the argument made here in their paper, and the reference they cite for justifying random MH kernels only deals with the trivial case  $P(\beta|\mathbf{x}) = P(\beta)$ .

along with  $\mathbf{x}$ , with an obvious modification of the MH procedure. In this case,  $\beta$  is an *auxiliary variable* which is dragged along for computational reasons (faster mixing, *etc.*) only.

## 1.2 Gibbs Sampling

Suppose  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ , and that we can sample from  $\pi(\mathbf{x}_j | \mathbf{x}_{3-j})$  easily (the latter is called *full conditional distribution*). Then, the *Gibbs kernels* are MH kernels:

$$q^{(j)}(\mathbf{x}^* | \mathbf{x}) = \pi(\mathbf{x}_j^* | \mathbf{x}_{3-j}) \delta_{\mathbf{x}_{3-j}^*, \mathbf{x}_{3-j}}, \quad j = 1, 2.$$

Note that if  $q^{(j)}$  is used, then  $\alpha(\mathbf{x}^*, \mathbf{x}) = 1$  if  $\mathbf{x}_{3-j}^* = \mathbf{x}_{3-j}$ , so all Gibbs proposals are accepted. The generalization to more than two blocks of  $\mathbf{x}$  is straightforward. In order to ensure ergodicity, one has to cycle through all different  $j$  using some scheme which visits each block infinitely often (this is necessary, but not sufficient).

## References

- [1] S. Jain and M. Neal, R. A split-merge Markov chain Monte Carlo procedure for the Dirichlet process mixture model. Technical Report 2003, Department of Statistics, University of Toronto, July 2000.