Fully Adaptive Gaussian Mixture Metropolis-Hastings Algorithm

David Luengo[†], Luca Martino[‡]

[†]Department of Circuits and Systems Engineering, Universidad Politécnica de Madrid. Carretera de Valencia Km. 7, 28031 Madrid, Spain.

[‡]Department of Signal Theory and Communications, Universidad Carlos III de Madrid. Avenida de la Universidad 30, 28911 Leganés, Madrid, Spain.

E-mail: luca.martino@uc3m.es, david.luengo@upm.es

Abstract

Markov Chain Monte Carlo methods are widely used in signal processing and communications for statistical inference and stochastic optimization. In this work, we introduce an efficient adaptive Metropolis-Hastings algorithm to draw samples from generic multi-modal and multi-dimensional target distributions. The proposal density is a mixture of Gaussian densities with all parameters (weights, mean vectors and covariance matrices) updated using all the previously generated samples applying simple recursive rules. Numerical results for the one and two-dimensional cases are provided.

Index Terms

Markov Chain Monte Carlo (MCMC), Gaussian mixtures, adaptive Metropolis-Hastings algorithms

I. INTRODUCTION

Markov Chain Monte Carlo (MCMC) methods [Liu, 2004, Liang et al., 2010] are ubiquitously used for performing inference and solving optimization problems in many scientific fields: statistics, digital communications, machine learning, signal processing, etc. [Robert and Casella, 2004, Wang et al., 2002, Andrieu et al., 2003, Fitzgerald, 2001]. MCMC approaches are able to generate samples virtually from any target distribution (known up to a normalizing constant) by using a simpler proposal distribution. The basic underlying idea of standard MCMC techniques is producing a Markov chain that converges to the target.

The most famous MCMC technique is the Metropolis-Hastings (MH) algorithm [Metropolis et al., 1953, Hastings, 1970, Robert and Casella, 2004]. However, the main drawback of the MH method (and in general of all MCMC methods) is that the correlation among the samples in the Markov chain can be very high when the acceptance rate is low [Liu, 2004, Liang et al., 2010, Martino and Míguez, 2010]. Correlated samples provide less statistical information and the resulting chain can remain trapped almost indefinitely in a local mode, meaning that convergence can be extremely slow. Therefore, since the correlation depends on the discrepancy between the target and proposal distributions, we would like it to be as close to the target as possible.

Several extensions have been proposed in the literature to speed up the convergence and reduce the so called "burn-in" period. Among them, adaptive MH methods (i.e., MH algorithms with adaptive proposal distributions) are particularly interesting [Liang et al., 2010, Martino et al., 2012]. Indeed, MCMC techniques usually need the selection of several parameters by the user before they can be applied to any particular problem. The use of adaptive proposals overcomes this issue, providing *black-box* algorithms with *self-tuning* capabilities. An adaptive MH technique improves the proposal distribution by learning at least some of its parameters from all the previously generated samples. Unfortunately, an important problem with the adaptation of the proposal is that the Markov property is lost and the invariant distribution of the chain could be disturbed. Hence, adaptive MH algorithms must be carefully designed to avoid this issue.

An adaptive Metropolis that uses an adaptive random walk Gaussian proposal, was introduced in [Haario et al., 2001] (we will denote it as AM method). The covariance matrix of the proposal is updated using recursive empirical estimators applied to the samples generated by the chain. The AM algorithm is an example of a partially adaptive MH approach, since it only updates the covariance of the proposal, whereas the mean of the Gaussian jumps to the current state of the chain at each iteration. An attempt of extending the AM algorithm by using a mixture of Gaussians as a proposal and updating all of its parameters (thus obtaining a fully adaptive MH algorithm) can be found in [Giordani and Kohn, 2010]. However, the resulting algorithm is quite complicated, and the proposal is only updated at some iterations.

In this work, we introduce an independent MH technique where the proposal PDF is an adaptive mixture of Gaussians. All the parameters (weights, means and covariance matrices) of the Gaussians in the mixture are updated using empirical estimators with simple recursive formulas (i.e., our method is fully adaptive). After a training period, the proposal is adapted at every iteration. The resulting AGM-MH algorithm can be used to draw samples from arbitrary multi-modal and multi-dimensional targets, always improving the performance w.r.t. a non-adaptive MH scheme using the initial proposal.

The rest of the paper is organized as follows. Section II shows the problem formulation. Section III presents the proposed AGM-MH algorithm. Sections IV and V describe efficient parameter update rules and black-box usage. Finally, Sections VI and VII show the results and conclude the paper.

II. PROBLEM FORMULATION

Let us assume that we need to draw samples from a (possibly multi-modal) generic d-dimensional target probability density function (PDF), $p_o(\mathbf{x})$, with support $\mathcal{D} \subseteq \mathbb{R}^d$. The AM algorithm [Haario et al., 2001] uses an adaptive random walk MH with a Gaussian proposal, mean equal to the previous state of the chain ($\mu = \mathbf{x}_{t-1}$), and covariance matrix, \mathbf{C} , estimated from all previous states, i.e., $q(\mathbf{x}|\mathbf{x}_{t-1},\mathbf{C}) \propto \mathcal{N}(\mathbf{x}|\mathbf{x}_{t-1},\mathbf{C})$, where

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \propto \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$
 (1)

denotes a standard multi-variate Gaussian PDF. In order to improve the performance of the AM algorithm, here we consider a mixture of N Gaussians as the proposal PDF, i.e.,

$$q(\mathbf{x}|\mathbf{w}, \boldsymbol{\mu}_{1:N}, \mathbf{C}_{1:N}) = \sum_{i=1}^{N} w_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \mathbf{C}_i),$$
(2)

where $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_i, \mathbf{C}_i)$ is given by (1), $\boldsymbol{\mu}_i = [\mu_{i,1}, ..., \mu_{i,d}]^T$ is the $d \times 1$ mean vector, \mathbf{C}_i is the $d \times d$ positive definite covariance matrix, and $\mathbf{w} = [w_1, ..., w_N]^T$ are the normalized weights (i.e., $w_1 + ... + w_N = 1$). Moreover, we define $\boldsymbol{\mu}_{1:N} = [\boldsymbol{\mu}_1, ..., \boldsymbol{\mu}_N]$ and $\mathbf{C}_{1:N} = [\mathbf{C}_1, ..., \mathbf{C}_N]$. We note that our approach is a fully adaptive MH algorithm, since (unlike the AM algorithm, which only updates the covariance) all the parameters in the mixture are learnt from all the previously generated samples. The resulting algorithm is very simple, since the adaptation is based on empirical estimators that can be implemented efficiently using recursive formulas.

Since the adaptation could disturb the convergence of the generated chain to the target PDF, we consider the possibility of stopping it at an iteration T_{stop} . Hence, for $t > T_{stop}$ our algorithm is a standard MH with an improved proposal PDF w.r.t. the initial choice, thus providing a better performance and guaranteeing convergence. However, the numerical results described in Section VI show that the algorithm seems to maintain the correct ergodicity properties, so we always use $T_{stop} = T_{tot}$. A theoretical convergence proof is under development and will be included in future works. Finally, we note that degeneracy problems can appear during the first iterations in the update of the covariance matrices if we have a poor initialization. In order to avoid this issue, we allow the method to use a few iterations $(t = 1, \dots, T_{train})$ to collect information about the target, assigning the produced state of the chain to the closest Gaussian in the mixture, as in [Haario et al., 2001].

III. AGM-MH ALGORITHM

The proposed AGM-MH algorithm is described below. First of all, notice that, during the first T_{train} time steps the algorithm just assigns the current state \mathbf{x}_t of the chain to a Gaussian among the N in the mixture, according to the minimum Euclidean distance between \mathbf{x}_t and the means $\boldsymbol{\mu}_i^{(t-1)}$, i=1,...,N. Afterwards, the algorithm updates all the parameters in the mixture until $t=T_{stop}$, when adaptation is stopped. In the description of the algorithm, the parameters are updated using a block procedure, but efficient recursive update formulas can be obtained, as shown in Section IV.

1) Initialization:

- a) Time instants: Set t = 0. Choose also the values $\mathbf{x}_0 \in \mathcal{D}$, $T_{train} < T_{tot}$ and $T_{train} < T_{stop} < T_{tot}$. Let be T_{tot} the number of total iteration of the chain.
- b) *Proposal:* Choose the number of Gaussians N, as well the initial settings for $\boldsymbol{\mu}_{1:N}^{(0)} = [\boldsymbol{\mu}_1^{(0)},...,\boldsymbol{\mu}_N^{(0)}]$ and $\mathbf{C}_{1:N}^{(0)} = [\mathbf{C}_1^{(0)},...,\mathbf{C}_N^{(0)}]$. Set $\mathbf{w}^{(0)} = \frac{1}{N}\mathbf{1}$.
- c) Auxiliary parameters: Define $\mathbf{S}_i^{(0)} \triangleq [\mathbf{s}_i^{(1)} = \boldsymbol{\mu}_i^{(0)}]$, and $m_i = 1$ represents the number of columns of $\mathbf{S}_i^{(0)}$, with i = 1, ..., N. Let ϵ a small constant value and \mathbf{I}_d a identity matrix.

2) MH steps:

a) Sample x' from a mixture of Gaussian PDFs,

$$\mathbf{x}' \sim q_t(\mathbf{x}'|\mathbf{w}^{(t)}, \boldsymbol{\mu}_{1:N}^{(t)}, \mathbf{C}_{1:N}^{(t)}).$$

b) Accept $\mathbf{x}_{t+1} = \mathbf{x}'$ with probability

$$\alpha = \min \left[1, \frac{p(\mathbf{x}')q(\mathbf{x}_t|\mathbf{w}^{(t)}, \boldsymbol{\mu}_{1:N}^{(t)}, \mathbf{C}_{1:N}^{(t)})}{p(\mathbf{x}_t)q(\mathbf{x}'|\mathbf{w}^{(t)}, \boldsymbol{\mu}_{1:N}^{(t)}, \mathbf{C}_{1:N}^{(t)})} \right], \tag{3}$$

otherwise set $\mathbf{x}_{t+1} = \mathbf{x}_t$.

3) If $t < T_{stop}$, update parameters of the proposal:

a) Find the closest Gaussian to x_{t+1} (w.r.t. Euclidean distance), i.e., find the index

$$j = \arg\min_{i} |\boldsymbol{\mu}_{i}^{(t)} - \mathbf{x}_{t+1}|^{2}. \tag{4}$$

b) Set $m_j = m_j + 1$ and update (adding a new column) the j-th auxiliary matrix

$$\mathbf{S}_{j}^{(t+1)} = [\mathbf{S}_{j}^{(t)}, \mathbf{s}_{j}^{(m_{j})} = \mathbf{x}_{t+1}], \tag{5}$$

whereas $\mathbf{S}_i^{(t+1)} = \mathbf{S}_i^{(t)}$, for all $i \neq j$.

c) If $t > T_{train}$: update the parameters of j-th Gaussian,

$$\mu_j^{(t+1)} = \frac{1}{m_j} \sum_{i=1}^{m_j} \mathbf{s}_j^{(i)},\tag{6}$$

and

$$\mathbf{C}_{j}^{(t+1)} = \frac{\tilde{\mathbf{S}}_{j}^{(t+1)} \cdot \left[\tilde{\mathbf{S}}_{j}^{(t+1)}\right]^{T} + (m_{j} - 1)\epsilon \mathbf{I}_{d}}{m_{j} - 1},$$
(7)

whrere $\tilde{\mathbf{S}}_{j}^{(t+1)} = \mathbf{S}_{j}^{(t+1)} - \boldsymbol{\mu}_{j}^{(t+1)} \otimes \mathbf{1}_{m_{j}}^{\top}$, with \otimes denoting the Kronecker product [Van Loan, 2000]. Morev,er set $\boldsymbol{\mu}_{i}^{(t+1)} = \boldsymbol{\mu}_{i}^{(t)}$, $\mathbf{C}_{i}^{(t+1)} = \mathbf{C}_{i}^{(t)}$, $\forall i \neq j$. Since m_{j} has been incremented, then update also the weights

$$w_i^{(t+1)} = \frac{m_i}{\sum_{k=1}^N m_k}, \quad i = 1, ..., N,$$
 (8)

so that $\mathbf{w}^{(t+1)} = [w_1^{(t+1)}, ..., w_N^{(t+1)}]^T$.

4) If $t < T_{tot}$ repeat from step 2.

Observe that the proposal PDF is only updated for $T_{train} < t < T_{stop}$. Moreover, note that the matrix $\mathbf{S}_i^{(t)}$ (and $\tilde{\mathbf{S}}_i^{(t)}$), for some $i \in \{1, ..., N\}$, has dimension $d \times m_i$, so that $\mathbf{C}_i^{(t)}$ has always dimension $d \times d$. The identity matrix $\epsilon \mathbf{I}_d$ is used just to avoid numerical problems (the matrix $\mathbf{C}_i^{(t)}$ must be positive definite), as in [Haario et al., 2001].

IV. EFFICIENT RECURSIVE UPDATE OF THE PARAMETERS

To update the parameters of the selected (j-th) Gaussian PDF in the mixture, we can use recursive expressions. Indeed, Recalling that $m_j = m_j + 1$ is already updated and $\mathbf{s}_j^{(m_j)} = \mathbf{x}_{t+1}$ in step 3b of the algorithm, (6) can be rewritten as

$$\boldsymbol{\mu}_{j}^{(t+1)} = \frac{1}{m_{j}} \mathbf{x}_{t+1} + \frac{m_{j} - 1}{m_{j}} \boldsymbol{\mu}_{j}^{(t)}, \tag{9}$$

and the Eq. (7) becomes

$$\mathbf{C}_{j}^{(t+1)} = \frac{1}{m_{j} - 1} \left[\frac{(\mathbf{x}_{t+1} - \boldsymbol{\mu}_{j}^{(t+1)})(\mathbf{x}_{t+1} - \boldsymbol{\mu}_{j}^{(t+1)})^{T}}{m_{j}} + \epsilon \mathbf{I}_{d} \right] + \frac{m_{j} - 2}{m_{j} - 1} \mathbf{C}_{j}^{(t)}.$$
(10)

Finally, note that

$$\sum_{k=1}^{N} m_k = t + 1 + N,$$

for $T_{train} < t < T_{stop}$, so that

$$w_i^{(t+1)} = \frac{m_i}{t+N+1}, \quad i = 1, ..., N,$$
(11)

In this way, the novel technique becomes computationally efficient in high dimensional problems, as well.

V. AGM-MH AS BLACK-BOX METHOD

The AGM-MH method shows sensitive dependence on the initial conditions. If some prior information about the target is available, it can be used to choose the initial parameters. If no prior informations is available, the AGM-MH can be applied as black-box algorithm in the following way:

- Use a great number N of Gaussians (typically the number must be greater if the dimension d increases).
- Select randomly the means $m{\mu}_{1:N}^{(0)}$ in order to cover as possible the support domain $\mathcal{D}\subseteq\mathbb{R}^d$.
- Choose diagonal covariance matrices $\mathbf{C}_{1:N}^{(0)} = \sigma^2 \mathbf{I}_d$, with a big value of σ^2 in order to explore the space $\mathcal{D} \subseteq \mathbb{R}^d$ in the train period, $t \leq T_{train}$.
- The parameter T_{train} should be chosen bigger for greater dimensions d. Numerical results suggest that $T_{train} = 100d$ could be a suitable choice. In general, for more complicated target distributions a greater T_{train} could be needed.

The use of a huge number of Gaussians does not generate computational problems since in the adaptation of the AGM-MH the weights of the irrelevant vanish quickly to zero. Therefore, the computational cost is controlled by the adaptation so that only the Gaussians located close to high probability regions survive. The useless Gaussian PDFs located faraway to the modes of the target are virtually discarded (the weights becomes quickly zero). Finally it is important to remark that, in general, the proposal is refined from the initial setting and the performance is improved.

VI. SIMULATIONS

A. Example 1

In this toy example, we apply the AGM-MH method to draw from a univariate bimodal target PDF defined as

$$p_o(x) \propto p(x) = \exp\left\{-\frac{(x^2 - 4)^2}{4}\right\}$$

= $\exp\left\{-\frac{x^4 - 4x^2 + 16}{4}\right\}$, (12)

that, clearly, has two modes at $x=\pm 2$. We set N=2, number of Gaussian PDFs in the proposal PDF, $w_i^{(0)}=0.5$, and $(\sigma_i^2)^{(0)}=10$ with i=1,2. The two initial means $\mu_1^{(0)}\sim \mathcal{U}([-4,0]), \ \mu_1^{(0)}\sim \mathcal{U}([0,4])$ are chosen uniformly in [-4,0] and [0,4], respectively. We draw $T_{tot}=5000$ iterations of the chain, and set $T_{train}=200$, $T_{stop}=T_{tot}$ (i.e., the adaptation is never stopped). The initial state is randomly choose as $x_0\sim N(x;0,1)$. We use *all* the generated samples to estimate the mean of the target (that is

0 since p(x) is symmetric). The mean square error (MSE) of the estimation (averaged over 2000 runs) is $\approx 15 \cdot 10^{-4}$. The estimated linear correlation between contiguous samples is ≈ 0.18 . Without using any adaptation, namely using a standard MH, the correlation is ≈ 0.78 . Hence, we can observe as the correlation decreases using the AGM-MH algorithm.

The final averaged locations of the means of the proposal are $\mu_1^{(T_{tot})} \approx -1.88$, $\mu_2^{(T_{tot})} \approx 1.88$. The final weights of the mixture are $w_i^{(T_{tot})} \approx 0.5$ and $\sigma_i^2 \approx 0.16$, i=1,2.

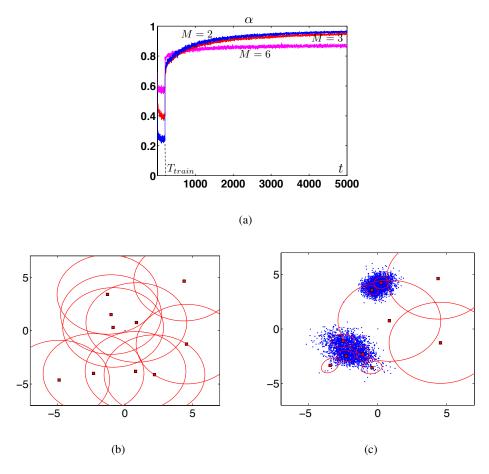


Fig. 1. (a) The averaged values of α as function of the iteration index t for different M=2,3,6. For $t>T_{train}$, the α grows since the proposal is closer to the target. (b) The initial configuration of the means (squares) and the covariance matrices (ellipses). (c) The final configuration at $t=T_{tot}=7000$.

B. Example 2

For the sake of simplicity, we consider again an univariate target density. However, now we consider that target distribution is itself a mixture of Gaussian PDFs. Specifically, the target PDF is formed by M-Gaussians, i.e.,

$$p_o(x) \propto p(x) = \sum_{i=1}^{M} a_i \mathcal{N}(x|\eta_i, \rho_i^2), \tag{13}$$

where the weights are $a_i = 1/M$ and the variances $\rho_i^2 = 4$, i = 1, ..., M. We consider different 3 cases with M = 2, 3, 6. The means are, for each case,

$$\eta_i = -10, 10, \quad \text{for} \quad M = 2,$$

$$\eta_i = -10, 0, 10, \quad \text{for} \quad M = 3,$$

$$\eta_i = -15, -10, -5, 5, 10, 15 \quad \text{for} \quad M = 6,$$

with i=1,...,M. In the proposal we also use N=M Gaussians and each initial mean is chosen uniformly in [-20,20]. All the initial variances are set $(\sigma_j^2)^{(0)}=10$ and the weights $w_j^{(0)}=1/N$, j=1,...,N.

As in the first example, we draw $T_{tot}=5000$ iterations of the chain, set $T_{train}=200$ and $T_{stop}=T_{tot}$ (i.e., the adaptation is never stopped). The initial state of the chain is randomly choose as $x_0 \sim N(x;0,1)$. We use all the generated samples to estimate the normalizing constant of the target. The mean square error (MSE) of the estimation (averaged over 1000 runs) is $1.6 \cdot 10^{-4}$, $1.1 \cdot 10^{-4}$ and $2 \cdot 10^{-5}$ for M=2,3,6 respectively. The resulting correlation is 0.13, 0.14 and 0.16 for M=2,3,6, respectively. With a standard MH (without adaptation) the correlation is 0.81, 0.72 and 0.46, for M=2,3,6. Hence, another time, we remark as the adaptation of the AGM-MH reduces considerably the correlation among the generated samples. Finally, Fig. 1(a) depicts the averaged values of the acceptance function α in Eq. (3) as function of t and for different t. In this case, for $t > T_{train}$, the averaged values of α increase because of the proposal becomes closer to the target PDF owing to the adaptation.

¹It is important to remark that we consider a mixture of Gaussians as a target PDF just to discuss the performance of the AGM-MH algorithm. More specifically, we desire to show as, in this case, the proposal convergences to real shape of the target, depending on the initial setting. However, clearly, the algorithm can be used to draw from any kind of target distribution.

C. Example 3

In this example, our goal is drawing from a bivariate target PDF using the AGM-MH as a black-box technique. Just for simplicity, we also consider a bivariate mixture of M=2 Gaussians as target distribution, with means $\eta_1=[-2,-2]^T$, $\eta_2=[0,4]^T$ and covariance matrices $\Sigma_1=[0.3\ 0.1;\ 0.1\ 0.3]$, $\Sigma_2=[0.8\ -0.3;\ -0.3\ 0.8]$. The weights are $a_i=0.5,0.5,\ i=1,2$. We set $T_{tot}=7000$, set $T_{train}=200$ and $T_{stop}=T_{tot}$ (i.e., the adaptation is never stopped). First, we use for the proposal N=2 Gaussian PDFs, $w_i^{(0)}=0.5$, $\mathbf{C}_i^{(0)}=10\mathbf{I}_d$ for i=1,2. The means are selected uniformly, $\mu_1\sim\mathcal{U}([-5,5]\times[0,5])$ and $\mu_2\sim\mathcal{U}([-5,5]\times[-5,0])$. In this case, all the parameters of the mixture in the proposal convergences always to the corresponding values of the mixture forming the target PDF.

Moreover, we also consider the case with N=10. We choose the means μ_i of the proposal uniformly in the square $[-5,5] \times [-5,5]$. In this situation, the AGM-MH refines the initial proposal PDF improving the parameters of the Gaussians in the mixture that are in good locations, whereas the weights of the unhelpful Gaussians decrease quickly to zero and their parameters remain invariant, as shown in Fig. 1.

VII. DISCUSSION

We have proposed a novel adaptive independent MH algorithm (AGM-MH) to draw samples from arbitrary multi-modal and multi-dimensional targets. AGM-MH builds on the work of [Haario et al., 2001], extending it by using a Gaussian mixture proposal and updating also the means and the weights of the Gaussians. Compared to a previous extension provided by [Giordani and Kohn, 2010], our approach is more efficient, updates the proposal at every iteration instead of only at a fixed number of iterations.

REFERENCES

- J. S. Liu. Monte Carlo Strategies in Scientific Computing. Springer, 2004.
- F. Liang, C. Liu, and R. Caroll. *Advanced Markov Chain Monte Carlo Methods: Learning from Past Samples*. Wiley Series in Computational Statistics, England, 2010.
- C. P. Robert and G. Casella. *Monte Carlo Statistical Methods*. Springer, 2004.
- Xiaodong Wang, Rong Chen, and Jun S. Liu. Monte Carlo Bayesian signal processing for wireless communications. *Journal of VLSI Signal Processing*, 30:89–105, 2002.
- Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I. Jordan. An introduction to MCMC for machine learning. *Machine Learning*, 50:5–43, 2003.
- W. J. Fitzgerald. Markov chain Monte Carlo methods with applications to signal processing. *Signal Processing*, 81(1):3–18, January 2001.

- N. Metropolis, A. Rosenbluth, M. Rosenbluth, A. Teller, and E. Teller. Equations of state calculations by fast computing machines. *Journal of Chemical Physics*, 21:1087–1091, 1953.
- W. K. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1):97–109, 1970.
- L. Martino and J. Míguez. Generalized rejection sampling schemes and applications in signal processing. Signal Processing, 90(11):2981–2995, November 2010.
- L. Martino, J. Read, and D. Luengo. Improved adaptive rejection Metropolis sampling algorithms. arXiv:1205.5494, May 2012.
- H. Haario, E. Saksman, and J. Tamminen. An adaptive Metropolis algorithm. *Bernoulli*, 7(2):223–242, April 2001.
- P. Giordani and R. Kohn. Adaptive independent Metropolis-Hastings by fast estimation of mixtures of normals. *Journal of Computational and Graphical Statistics*, 19(2):243–259, September 2010.
- Charles F. Van Loan. The ubiquitous Kronecker product. *Journal of Computational and Applied Mathematics*, 123:85–100, 2000.