

ISSN 1440-771X



**MONASH** University

**Australia**

Department of Econometrics  
and Business Statistics

<http://www.buseco.monash.edu.au/depts/ebs/pubs/wpapers/>

**Bandwidth Selection for Multivariate Kernel Density  
Estimation Using MCMC**

**Xibin Zhang Maxwell L. King Rob J. Hyndman**

**Working Paper 09/04**

# Bandwidth Selection for Multivariate Kernel Density Estimation Using MCMC

Xibin Zhang   Maxwell L. King   Rob J. Hyndman

Department of Econometrics and Business Statistics  
Monash University, Clayton, Victoria, 3800, Australia

Correspondence to: Xibin.Zhang@buseco.monash.edu.au

April 2004

**Abstract:** We provide Markov chain Monte Carlo (MCMC) algorithms for computing the bandwidth matrix for multivariate kernel density estimation. Our approach is based on treating the elements of the bandwidth matrix as parameters to be estimated, which we do by optimizing the likelihood cross-validation criterion. Numerical results show that the resulting bandwidths are superior to all existing methods; for dimensions greater than two, our algorithm is the first practical method for estimating the optimal bandwidth matrix. Moreover, the MCMC algorithm for bandwidth selection for multivariate data has no increased difficulty as the dimension of data increases.

**Key words:** bandwidth selection, cross-validation, multivariate kernel density estimation, sampling algorithms.

**JEL classification:** C14, C11, C51.

## 1 Introduction

Multivariate kernel density estimation is an important technique in multivariate data analysis and it has a wide range of applications (see, for example, Scott 1992; Ait-Sahalia 1996; Donald 1997; Stanton 1997; Ait-Sahalia and Lo 1998). However, its widespread usefulness has been limited by the difficulty in computing an optimal data-driven bandwidth. We remedy this deficiency in this paper.

Let  $\mathbf{X} = (X_1, X_2, \dots, X_d)'$  denote a  $d$ -dimensional random vector with density  $f(\mathbf{x})$  defined on  $\mathbf{R}^d$ , and let  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  be an independent random sample drawn from  $f(\mathbf{x})$ . The general form of the kernel estimator of  $f(\mathbf{x})$  is (Wand and Jones 1993):

$$\hat{f}_H(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K_H(\mathbf{x} - \mathbf{x}_i),$$

where  $K_H(\mathbf{x}) = |H|^{-1/2}K(H^{-1/2}\mathbf{x})$ ,  $K(\cdot)$  is a multivariate kernel function, and  $H$  is a symmetric positive definite  $d \times d$  matrix known as the bandwidth matrix.

The bandwidth matrix can be restricted to a class of positive definite diagonal matrices, and then the corresponding kernel function is known as a product kernel. However, there is much to be gained by choosing a full bandwidth matrix, where the corresponding kernel smoothing is equivalent to pre-rotating the data by an optimal amount and then using a diagonal bandwidth matrix.

It has been widely recognized that the performance of a kernel density estimator is primarily determined by the choice of bandwidth, and only in a minor way by the choice of kernel function. See, for example, Izenman (1991), Scott (1992) and Simonoff (1996).

A large body of literature exists on bandwidth selection for univariate kernel density estimation; see Marron (1987) and Jones, Marron and Sheather (1996) for surveys. However, the literature on bandwidth selection for multivariate data is very limited. Sain, Baggerly and Scott (1994) discussed the performance of bootstrap and cross-validation methods for bandwidth selection in multivariate density estimation. However, they found that the complexity of finding an optimal bandwidth grows prohibitively as the dimension  $d$  of the data increases. Wand and

Jones (1994) presented a less variable cross-validation algorithm using the plug-in method. However, their algorithm requires auxiliary smoothing parameters, and the technology for choosing these auxiliary smoothing parameters is not well developed. In addition, Duong and Hazelton (2003) showed that the full bandwidth matrix selectors suggested by Wand and Jones (1994) fail to produce plug-in bandwidths for some data sets. In response to this problem, Duong and Hazelton (2003) presented an alternative plug-in bandwidth selector for bivariate kernel density estimation. This plug-in method has the advantage that it always produces a finite bandwidth matrix and requires computation of fewer pilot bandwidths. However, it cannot be directly extended to the general multivariate setting.

When data are observed from the multivariate normal density and the product normal kernel is employed, the optimal bandwidth (minimizing the mean integrated squared error) can be approximated by (Bowman and Azzalini 1997)

$$h_i = \sigma_i \left\{ \frac{4}{(d+2)n} \right\}^{1/(d+4)},$$

for  $i = 1, 2, \dots, d$ , where  $\sigma_i$  can be replaced by its sample estimate in practical implementations. We call this the “normal reference rule”. This method is often used in practice, in the absence of any other practical bandwidth selection schemes, despite the fact that most interesting data are non-Gaussian.

Although cross-validation provides a theoretical solution to optimal bandwidth selection, we face computational difficulties in obtaining an optimal bandwidth matrix using conventional optimization methods as the dimension  $d$  increases. Fortunately, the numerical optimization can be approximated through MCMC simulations, where the bandwidth matrix  $H$  is treated as a matrix of parameters, and the posterior density of  $H$  can be obtained through a likelihood cross-validation criterion. One important advantage of our MCMC approach is that it is applicable to data of any dimension, not only to bivariate data. The difficulty of implementing the sampling algorithm does not increase as the dimension of the data increases. In this paper, we present MCMC algorithms for computing the optimal bandwidth matrix for multivariate kernel density estimation through likelihood cross-validation, and sampling algorithms are developed for both diagonal and full bandwidth matrices.

The rest of this paper is organized as follows. Section 2 briefly discusses the likelihood cross-validation criterion and presents MCMC algorithms for both diagonal and full bandwidth matrices. In Section 3, we study the accuracy of the resulting density estimates using known bivariate densities. We find that the MCMC bandwidth selector works better than either the Duong-Hazelton plug-in method or the normal reference rule in the bivariate setting. Section 4 applies the MCMC bandwidth selector to several sets of data drawn from known multivariate densities and finds that our MCMC approach to bandwidth selection works much better than the normal reference rule. Section 5 illustrates the use of our MCMC algorithms for bandwidth selection with an application to some earthquake data; we provide conclusions in Section 6.

## 2 MCMC for optimal bandwidth selection

### 2.1 Likelihood cross-validation

The Kullback-Leibler information is a measure of distance between two densities. Our interest is in choosing the approximate density  $\hat{f}_H(\mathbf{x})$  to minimize its distance from the true density  $f(\mathbf{x})$ . In this case, the Kullback-Leibler information is defined as

$$d_{KL}(f, \hat{f}_H) = \int \log \left( \frac{f(\mathbf{x})}{\hat{f}_H(\mathbf{x})} \right) f(\mathbf{x}) d\mathbf{x} = \int \log[f(\mathbf{x})] f(\mathbf{x}) d\mathbf{x} - \int \log[\hat{f}_H(\mathbf{x})] f(\mathbf{x}) d\mathbf{x}, \quad (1)$$

and we want a bandwidth matrix  $H$  that minimizes  $d_{KL}(f, \hat{f}_H)$ , or, equivalently, maximizes  $E \log[\hat{f}_H(\mathbf{x})] = \int \log[\hat{f}_H(\mathbf{x})] f(\mathbf{x}) d\mathbf{x}$ . We estimate the latter quantity by the log pseudo-likelihood function

$$L(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | H) = \sum_{i=1}^n \log \hat{f}_{H,i}(\mathbf{x}_i) \quad (2)$$

where  $\hat{f}_{H,i}$  is the leave-one-out estimator

$$\hat{f}_{H,i}(\mathbf{x}_i) = \frac{1}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n |H|^{-1/2} K \left( H^{-1/2}(\mathbf{x}_i - \mathbf{x}_j) \right).$$

The likelihood cross-validation criterion is to select  $H$  by maximizing  $n^{-1}L(\cdot | H)$ .

## 2.2 The sampling algorithm for bandwidth selection

As the bandwidth matrix is symmetric positive definite, we can obtain its Cholesky decomposition  $H = LL'$ , where  $L$  is a lower triangular matrix. Let  $B = L^{-1}$  which is also lower triangular. Then the kernel estimator of  $f(\mathbf{x})$  is

$$\hat{f}_B(\mathbf{x}) = \frac{1}{n} |B| \sum_{i=1}^n K(B(\mathbf{x} - \mathbf{x}_i)),$$

and the leave-one-out estimator of  $f(\mathbf{x})$  is

$$\hat{f}_{B,i}(\mathbf{x}_i) = \frac{1}{n-1} |B| \sum_{\substack{j=1 \\ j \neq i}}^n K(B(\mathbf{x}_i - \mathbf{x}_j)).$$

We treat the non-zero elements of the bandwidth matrix as parameters, and we obtain the posterior density of the parameters based on the likelihood function given in (2).

We assume that the prior density of each non-zero component of  $B$  is (up to a normalizing constant)

$$\pi(b_{ij} | \lambda) \propto \frac{1}{1 + \lambda b_{ij}^2} \quad (3)$$

for  $j \leq i$  and  $i = 1, 2, \dots, d$ , where  $\lambda$  is a hyperparameter controlling the shape of the prior densities. Uniform priors of bandwidths are generally unsuitable, because the update of each  $b_{ij}$  has negligible effect when  $b_{ij}$  is already very large. However, priors of  $b_{ij}$  defined in (3) can prevent the update of  $b_{ij}$  from getting too large. In a different context, Bauwens and Lubrano (1998) used a similar prior for the degrees of freedom parameter of the  $t$ -distribution. The purpose of such priors is to put low prior probability on the “problematic” region of the parameter space, where the likelihood function is flat. The joint prior of all elements of  $B$  is the product of these marginal priors. Then, using Bayes theorem, the logarithmic posterior of  $B$  is (up to an additive constant)

$$\pi(B | \lambda, \mathbf{x}) \propto \sum_{i=1}^d \sum_{j=1}^i \log \pi(b_{ij} | \lambda) + \sum_{i=1}^n \log \hat{f}_{B,i}(\mathbf{x}_i), \quad (4)$$

and all elements of  $B$  can be sampled through the Metropolis-Hastings algorithm with the

acceptance probability computed through (4).

In the case of a diagonal bandwidth matrix, the algorithm proceeds as for the full bandwidth matrix except that only diagonal elements need to be sampled. When developing the sampling algorithm for a diagonal bandwidth matrix, we actually sample  $H^{1/2}$  instead of  $H^{-1/2}$ .

### 2.3 Transformation of data

The plug-in algorithm for bandwidth selection developed by Duong and Hazelton (2003) uses a simple form for the pilot bandwidths, which is inappropriate when the dispersion of the data differs markedly between the two variates. Hence Duong and Hazelton (2003) suggested that the data be pre-scaled before the plug-in algorithm is implemented.

Given a set of bivariate data denoted by  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , let  $S$  denote the sample variance-covariance matrix with diagonal components  $s_1^2$  and  $s_2^2$ . Duong and Hazelton (2003) defined the sphering and scaling transformations, respectively, by

$$\mathbf{x}_i^* = S^{-1/2}\mathbf{x}_i, \quad \text{and} \quad \mathbf{x}_i^* = S_d^{-1/2}\mathbf{x}_i,$$

for  $i = 1, 2, \dots, n$ , where  $S_d = \text{diagonal}(s_1^2, s_2^2)$ . When the optimal bandwidth matrix, denoted by  $\hat{H}^*$ , for the transformed data is obtained, the optimal bandwidth matrix for the original data can be calculated through the reverse transformation,  $\hat{H} = S^{1/2}\hat{H}^*(S^{1/2})'$  or  $\hat{H} = S_d^{1/2}\hat{H}^*S_d^{1/2}$ .

In contrast, the MCMC algorithm does not require such pre-transformations of data. However, if we choose to make a sphering transformation of the data and use the diagonal bandwidth matrix, the resulting bandwidth estimator for the original data is a full matrix. When the variates are correlated and the diagonal bandwidth matrix is used, the bandwidth matrix estimator obtained through the sphering transformation of original data might produce a better performance than that obtained directly from original data, because the sphering transformation is equivalent to pre-rotating the data (Wand and Jones 1993).

### 3 Numerical studies with bivariate densities

This section examines the performance of the proposed MCMC methods for bandwidth selection via several sets of bivariate data, generated from known densities. As the true density is known in each case, the performance of the bandwidth can be measured by the accuracy of the corresponding kernel density estimate via Kullback-Leibler information.

The Kullback-Leibler information defined in (1) is the mean of  $\log(f(\mathbf{x})/\hat{f}_H(\mathbf{x}))$  under density  $f(\mathbf{x})$ , and so it measures the discrepancy of the estimated density from the true density. If a large number of random vectors, denoted by  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , can be drawn from  $f(\mathbf{x})$ , the Kullback-Leibler information can be estimated by

$$\hat{d}_{KL}(f, \hat{f}_H) = \frac{1}{N} \sum_{i=1}^N \log(f(\mathbf{x}_i)/\hat{f}_H(\mathbf{x}_i)). \quad (5)$$

#### 3.1 True densities

We consider five target densities labelled A, B, C, D and E. Contour plots of these densities are shown in Figure 1.

**Density A** is bivariate normal with high correlation between the two variates:

$$f_A(\mathbf{x} \mid \mu, \Sigma) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu)\right),$$

with mean  $\mu$  and variance-covariance matrix  $\Sigma$  given by

$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 1 & -0.9 \\ -0.9 & 1 \end{pmatrix}.$$

**Density B** is a mixture of two bivariate normal densities, with high correlation and bimodality:

$$f_B(\mathbf{x} \mid \mu_1, \Sigma_1, \mu_2, \Sigma_2) = \frac{1}{2} f_A(\mathbf{x} \mid \mu_1, \Sigma_1) + \frac{1}{2} f_A(\mathbf{x} \mid \mu_2, \Sigma_2),$$



where

$$\mu_1 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \quad \Sigma_1 = \begin{pmatrix} 1 & -0.9 \\ -0.9 & 1 \end{pmatrix}, \quad \mu_2 = \begin{pmatrix} -1.5 \\ -1.5 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 1 & 0.3 \\ 0.3 & 1 \end{pmatrix}.$$

**Density C** is a bivariate skew-normal density with high correlation:

$$f_C(\mathbf{x} \mid \mu, \Sigma, \alpha) = 2\phi(\mathbf{x} \mid \mu, \Sigma) \Phi(\alpha' W^{-1/2}(\mathbf{x} - \mu)),$$

where  $\phi(\cdot \mid \mu, \Sigma)$  is the bivariate normal density with mean  $\mu$  and variance-covariance matrix  $\Sigma$ ,  $\Phi(\cdot)$  is the cumulative density function of a standard bivariate normal distribution, and  $W$  is a diagonal matrix with diagonal elements the same as those of  $\Sigma$ . This distribution has recently been studied by Azzalini and Dalla Valle (1996), Azzalini and Capitanio (1999, 2003), Jones (2001) and Jones and Faddy (2003) among many others. Here  $\alpha$  is the shape parameter capturing the skewness of the distribution. When  $\alpha = 0$ , the density  $f_C$  becomes the usual normal density. For the purpose of generating a set of data, we use the following parameters,

$$\mu = \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix}, \quad \alpha = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}.$$

**Density D** is a mixture of two bivariate Student  $t$  densities

$$f_D(\mathbf{x} \mid \mu_1, \mu_2, \Sigma, \nu) = \frac{1}{2}t_d(\mathbf{x} \mid \mu_1, \Sigma, \nu) + \frac{1}{2}t_d(\mathbf{x} \mid \mu_2, \Sigma, \nu),$$

where

$$t_d(\mathbf{x} \mid \mu, \Sigma, \nu) = \frac{\Gamma((\nu + d)/2)}{(\nu\pi)^{d/2}\Gamma(\nu/2)|\Sigma|^{1/2}} \left[ 1 + \frac{1}{\nu}(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu) \right]^{-(d+\nu)/2}, \quad (6)$$

has location parameter  $\mu$ , dispersion matrix  $\Sigma$  and degrees of freedom  $\nu$ , and with parameters set to

$$\mu_1 = \begin{pmatrix} -1.5 \\ 0 \end{pmatrix}, \quad \mu_2 = \begin{pmatrix} 1.5 \\ 0 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix},$$

and  $\nu = 5$ . Density D exhibits heavy tail behaviour, high correlation and bimodality.

**Density E** is also a mixture of two bivariate Student  $t$  densities, but has thicker tails than density D:

$$f_E(\mathbf{x} \mid \mu_1, \mu_2, \Sigma, \nu) = \frac{1}{2}t_d(\mathbf{x} \mid \mu_1, \Sigma_1, \nu) + \frac{1}{2}t_d(\mathbf{x} \mid \mu_2, \Sigma_2, \nu),$$

where

$$\mu_1 = \begin{pmatrix} 2.5 \\ 2.5 \end{pmatrix}, \quad \Sigma_1 = \begin{pmatrix} 1 & -0.75 \\ -0.75 & 1 \end{pmatrix}, \quad \mu_2 = \begin{pmatrix} -2.5 \\ -2.5 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 1 & -0.16 \\ -0.16 & 1 \end{pmatrix},$$

and  $\nu = 3$ .

### 3.2 Bandwidth matrix selectors

From each of the proposed bivariate densities, we generate data sets of size  $n = 200, 500$  and  $1000$ , respectively. For each data set, we calculate the bivariate kernel density estimator using the bivariate Gaussian kernel function and bandwidth matrix selected through each of the following selectors.

- M1:** MCMC algorithm for full bandwidth matrix without pre-transformation of data;
- M2:** MCMC algorithm for full bandwidth matrix with scaling transformation of data;
- M3:** MCMC algorithm for full bandwidth matrix with sphering transformation of data;
- M4:** MCMC algorithm for diagonal bandwidth matrix without pre-transformation of data;
- M5:** MCMC algorithm for diagonal bandwidth matrix with scaling transformation of data;
- M6:** MCMC algorithm for diagonal bandwidth matrix with sphering transformation of data;
- P1:** Plug-in selector of full bandwidth matrix with scaling transformation of data;
- P2:** Plug-in selector of full bandwidth matrix with sphering transformation of data;
- P3:** Plug-in selector of diagonal bandwidth matrix with scaling transformation of data;
- P4:** Plug-in selector of diagonal bandwidth matrix with sphering transformation of data;
- A1:** The normal reference rule approach for a diagonal bandwidth.

The plug-in bandwidth selector refers to the algorithm developed by Duong and Hazelton (2003). We have not included the plug-in algorithms of Wand and Jones (1993), because their

algorithm for full bandwidth matrix selection sometimes fails to produce finite bandwidths for some data sets. When their algorithm works, its performance is similar to the plug-in algorithm developed by Duong and Hazelton (2003). See Duong and Hazelton (2003) for a discussion on the advantages of their plug-in algorithm.

### 3.3 MCMC outputs and sensitivity analysis

The hyperparameter of the prior densities defined in (3) is initially set to  $\lambda = 1$  which represents a very flat prior. Given a data set generated from a bivariate density, we sample the diagonal and full bandwidth matrices from their corresponding posterior densities defined in (4) using the random-walk Metropolis-Hastings algorithm, where the proposal density is the multivariate standard normal density, and the tuning parameter is chosen so that the acceptance rate is between 0.2 and 0.3.

The burn-in period is set at 5,000 iterations, and the number of total recorded iterations is 25,000. The initial value of  $B$  is set to the identity matrix. After we obtain the sampled path of  $B$  for each data set, we calculate the ergodic average (or posterior mean) and the batch-mean standard error (see, for example, Roberts 1996), where the number of batches is 50 and there are 500 draws in each batch. The ergodic average acts as an estimator of optimal bandwidth.

We use the batch-mean standard error and the simulation inefficiency factor to check the mixing performance of sampling algorithms (see, for example, Kim, Shephard and Chib 1998; Tse, Zhang and Yu 2004). We use  $f_E(\cdot)$  as an example to illustrate the mixing performance of the sampling algorithm. Table 1 presents a summary of MCMC outputs obtained through methods  $M_1$  and  $M_6$ . Both SIF and the batch-mean standard error show that all the simulated chains have mixed very well. To demonstrate the mixing performance of these samplers graphically, we plot the sampled paths, their associated autocorrelation functions and histograms in Figure 2, which also reveals that these simulated chains have mixed very well. We found similar performance for the other variations on our algorithm, and for other data sets.

We examined the robustness of the results to prior choices, by trying values of  $\lambda = 0.1$  and  $\lambda = 5$  as well as  $\lambda = 1$ . The mixing performance and posterior mean of each sampler was similar in all cases.

### 3.4 Accuracy of MCMC bandwidth selectors

In order to estimate the Kullback-Leibler information, we generate  $N = 100,000$  bivariate random vectors from the true density and calculate the estimated Kullback-Leibler information defined by (5), which is employed to measure the distance between the bivariate kernel density estimator and the true density. Table 2 presents the estimated Kullback-Leibler information for each density and each bandwidth selector. The simulation study reveals the following evidence.

- For all data sets, the MCMC approach to optimal bandwidth selection performs better than the plug-in approach in computing an optimal diagonal bandwidth matrix, as well as in computing an optimal full bandwidth matrix.
- For all data sets, the MCMC method for both the diagonal and full bandwidth matrices performs better than the normal reference rule method discussed in Scott (1992).
- The scaling transformation adds nothing to the performance of the MCMC method for both diagonal and full bandwidth matrices.
- The sphering transformation of data is only helpful to the MCMC algorithm for a diagonal bandwidth matrix when the two variates are correlated, such as for densities A, C and E. For uncorrelated data, and for the full bandwidth matrix, sphering can degrade performance. This is also supported by Wand and Jones (1993).
- The MCMC algorithm for a diagonal bandwidth matrix applied after sphering does not perform quite as well as the full bandwidth approach. However, the simplicity of using a diagonal bandwidth matrix makes this an attractive approach, especially with high dimensional data.

We have seen that, when applied to bivariate data sets, the MCMC approach to bandwidth selection performs better than the plug-in method developed by Duong and Hazelton (2003) and the normal reference rule discussed in Scott (1992). Furthermore, the MCMC approach has the great advantage that there are no added difficulties in the methodology when the dimension of data increases.

## 4 Numerical studies with multivariate densities

In this section, we examine the accuracy of the MCMC approach in the general multivariate setting. Our examples use  $d = 5$ .

### 4.1 True densities and bandwidth selectors

We consider five target densities labelled F, G, H, I and J, respectively.

**Density F** is a multivariate normal density with location parameter  $\mu$  and variance-covariance matrix defined as

$$\Sigma = \frac{1}{1 - \rho^2} \begin{pmatrix} 1 & \rho & \rho^2 & \dots & \rho^{d-1} \\ \rho & 1 & \rho & \dots & \rho^{d-2} \\ \rho^2 & \rho & 1 & \dots & \rho^{d-3} \\ \dots & & & \dots & \\ \rho^{d-1} & \rho^{d-2} & \rho^{d-3} & \dots & 1 \end{pmatrix},$$

where  $\rho = 0.9$  and  $\mu = (2, 2, 2, 2, 2)'$ .

**Density G** is a mixture of two multivariate normal densities,

$$f_G(\mathbf{x} \mid \mu_1, \mu_2, \Sigma) = \frac{1}{2}f_A(\mathbf{x} \mid \mu_1, \Sigma) + \frac{1}{2}f_A(\mathbf{x} \mid \mu_2, \Sigma),$$

where  $\mu_1 = (2, 2, 2, 2, 2)'$ ,  $\mu_2 = (-1.5, -1.5, -1.5, -1.5, -1.5)'$  and  $\Sigma$  is the identity matrix.

**Density H** is a mixture of two multivariate Student  $t$  densities,

$$f_H(\mathbf{x} \mid \mu_1, \mu_2, \Sigma, \nu) = \frac{1}{2}t_d(\mathbf{x} \mid \mu_1, \Sigma, \nu) + \frac{1}{2}t_d(\mathbf{x} \mid \mu_2, \Sigma, \nu),$$

with  $t_d(\cdot)$  defined in (6),  $\mu_1 = (2, 2, 2, 2, 2)'$ ,  $\mu_2 = (-1.5, -1.5, -1.5, -1.5, -1.5)'$  and  $\Sigma$  is the identity matrix.

**Density I** is the multivariate skew normal density,

$$f_I(\mathbf{x} \mid \mu, \Sigma, \alpha) = 2\phi(\mathbf{x} \mid \mu, \Sigma) \Phi(\alpha'W^{-1/2}(\mathbf{x} - \mu)),$$

where  $\phi(\cdot \mid \mu, \Sigma)$  is the multivariate normal density with location parameter  $\mu$  and variance-covariance matrix  $\Sigma$ ,  $\Phi(\cdot)$  is the cumulative density function of a standard multivariate normal distribution, and  $W$  is a diagonal matrix with diagonal elements the same as those of  $\Sigma$ . To generate a set of data, we define these parameters as  $\mu = (2, 2, 2, 2, 2)'$ , variance-covariance matrix  $\Sigma$  defined in (9) with  $\rho = 0.9$ , and skewness parameter  $\alpha = (-0.5, -0.5, -0.5, -0.5, -0.5)'$ .

**Density J** is the multivariate skew  $t$  density,

$$f_J(\mathbf{x} \mid \mu, \Sigma, \nu, \alpha) = 2t_d(\mathbf{x} \mid \mu, \Sigma, \nu)T_d(\tilde{\mathbf{x}} \mid \nu + d)$$

where  $t_d(\cdot)$  is the multivariate  $t$  density defined in (6),  $T_d(\cdot \mid \nu + d)$  is the cumulative density function of a multivariate  $t$  distribution with mean  $\mathbf{0}$ , identity dispersion matrix and degrees of freedom  $\nu + d$ , and

$$\tilde{\mathbf{x}} = \alpha'W^{-1/2}(\mathbf{x} - \mu) \left( \frac{\nu + d}{(\mathbf{x} - \mu)' \Sigma^{-1}(\mathbf{x} - \mu) + \nu} \right)^{1/2},$$

where  $W$  is a diagonal matrix with diagonal elements the same as those of  $\Sigma$ .

From each of the proposed multivariate densities, we generate data sets of sizes 500, 1000 and 1500, respectively. Then we apply the proposed MCMC algorithms to each generated data set to estimate the optimal bandwidth. As the normal reference rule discussed in Scott (1992) and Bowman and Azzalini (1997) is the only viable alternative, we shall compare the performance of MCMC bandwidth selectors  $M_1$  to  $M_6$  with that of the alternative bandwidth selector  $A_1$ .

The MCMC algorithm and parameter settings are the same as in the bivariate examples.

## 4.2 MCMC outputs and sensitivity analysis

Table 3 shows MCMC output obtained from  $f_F(\cdot)$  with size 1500 to illustrate the mixing performance of the sampling algorithm. Both the batch-mean standard error and SIF show that all the sampled chains have mixed very well. Using the output obtained through  $M_1$ , we plot the sampled paths, their associated autocorrelation functions and histograms in Figure 3, which shows that the simulated chains via this algorithm have mixed well.

The numerical study shows that all algorithms for a diagonal bandwidth matrix have a similar mixing performance, and that all algorithms for a full bandwidth matrix have a similar mixing performance. However, the algorithm for a diagonal bandwidth matrix usually has a better mixing performance than the algorithm for a full bandwidth matrix. Similar results were found with the other data sets. Again, we found the results were insensitive to changes in  $\lambda$ .

## 4.3 Accuracy of MCMC bandwidth selectors

To estimate the Kullback-Leibler information, we generate  $N = 100,000$  random vectors from the true density and calculate the estimated Kullback-Leibler information defined by (5). Table 4 presents these results for each density and each bandwidth selector.

The simulation study reveals the following evidence. First, all MCMC algorithms perform much better than the normal reference rule. Second, the scaling transformation adds nothing to the performance of the MCMC algorithm for either the diagonal or full matrices. Third, the sphering transformation of data is only useful for the diagonal bandwidth matrix when variables are correlated (such as with densities F, I and J). When there is no correlation, or with the full bandwidth matrix, sphering degrades performance.

## 5 Application to earthquake data

We now apply the methodology to a trivariate data set discussed in Scott (1992). These data represent the epicenters of 510 earthquake tremors that occurred beneath the Mt St Helens vol-

cano in the two months leading up to its eruption in March 1982. The three variables represent latitude, longitude and log-depth below the surface. Scott (1992, plate 8) gives several contours of a kernel density estimate of these data, where the bandwidths appear to have been chosen subjectively. We repeat this plot, but using our optimal bandwidth methodology.

We use the MCMC algorithms  $M_1$  and  $M_5$  to obtain optimal bandwidths, where the hyperparameter  $\lambda = 1$ , the burn-in period consists of 5000 iterations, and the recorded period contains 25000 iterations. Table 5 tabulates a summary of results. Both the batch-mean standard error and SIF show that all sampled chains have mixed very well.

Using the estimated diagonal bandwidth matrix, we compute a kernel density estimator. (The estimate using the full bandwidth matrix was almost identical in this case.) The 98% highest density region (Hyndman 1996) is plotted in Figure 4. The surface was computed using the algorithm of Amenta, Bern and Kamvyselis (1998). Note that the detached shells represent outliers in the data; the large central shell represents the bulk of the epicenters. The figure clearly shows clustering of the epicenters, revealing structure that was not discovered by Scott (1992) using a subjective bandwidth. It would be interesting to identify the clusters with geological features, although this information is not available to us.

## 6 Conclusion

This paper presents MCMC sampling algorithms to estimate the optimal bandwidth for multivariate kernel density estimation via the likelihood cross-validation criterion. This represents the first data-driven bandwidth selection method for density estimation with more than two variables. Our numerical studies show that the algorithms have very good performance in achieving convergence in the simulated Markov chains, and are insensitive to prior choices.

We have found our algorithms to be superior to the bivariate plug-in algorithms of Duong and Hazelton (2003) and the normal reference rule discussed in Scott (1992) and Bowman and Azzalini (1997). Apart from superior performance, the other great advantage of our approach is that it is applicable to high dimensional data with no increased difficulty.



## Acknowledgement

We extend our sincere thanks to Faming Liang for sharing his coding skills and resources, David Scott for providing the earthquake data, Tarn Duong and Martin Hazelton for providing their R library to compute the plug-in bandwidth for bivariate data, and the Victorian Partnership for Advanced Computing for computational support.

## References

- Aït-Sahalia, Y. (1996), "Testing Continuous-Time Models of the Spot Interest Rate," *Review of Financial Studies*, 9, 385-426.
- Aït-Sahalia, Y., and Lo, A.W. (1998), "Nonparametric Estimation of State-Price Densities Implicit in Financial Asset Prices," *The Journal of Finance*, 53, 499-547.
- Amenta, N., Bern, M., and Kamvysselis, M. (1998) "A New Voronoi-based Surface Reconstruction Algorithm", *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, 415-421.
- Azzalini, A., and Capitanio, A. (1999), "Statistical Applications of the Multivariate Skew Normal Distribution," *Journal of the Royal Statistical Society, Series B*, 61, 579-602.
- Azzalini, A., and Capitanio, A. (2003), "Distributions Generated by Perturbation of Symmetry with Emphasis on a Multivariate Skew  $t$ -distribution," *Journal of the Royal Statistical Society, Series B*, 66, 367-389.
- Azzalini, A., and Dalla Valle, A. (1996), "The Multivariate Skew Normal Distribution," *Biometrika*, 83, 715-726.
- Bauwens, L., and Lubrano, M. (1998), "Bayesian Inference on GARCH Models Using the Gibbs Sampler," *Econometrics Journal*, 1, C23-C26.
- Bowman, A.W., and Azzalini, A. (1997), *Applied Smoothing Techniques for Data Analysis*, London: Oxford University Press.
- Donald, S.G. (1997), "Inference Concerning the Number of Factors in a Multivariate Nonparametric Relationship," *Econometrica*, 65, 103-131.

- Duong, T., and Hazelton, M.L. (2003), "Plug-in Bandwidth Selectors for Bivariate Kernel Density Estimation," *Journal of Nonparametric Statistics*, 15, 17-30.
- Hyndman, R.J. (1996), "Computing and Graphing Highest Density Regions," *American Statistician*, 50, 120-126.
- Izenman, A.J. (1991), "Recent Developments in Nonparametric Density Estimation," *Journal of the American Statistical Association*, 86, 205-224.
- Jones, M.C. (2001), "A Skew  $t$  Distribution," in *Probability and Statistical Models with Applications: A Volume in Honor of Theophilos Cacoullos*," eds. C.A. Charalambides, M.V. Koutras, and N. Balakrishnan, London: Chapman & Hall, pp. 269-278.
- Jones, M.C., and Faddy, M.J. (2003), "A Skew Extension of the  $t$ -distribution, with Applications," *Journal of the Royal Statistical Society, Series B*, 66, 159-174.
- Jones, M.C., Marron, J.S., and Sheather, S.J. (1996), "A Brief Survey of Bandwidth Selection for Density Estimation," *Journal of the American Statistical Association*, 91, 401-407.
- Kim, S., Shephard, N., and Chib, S. (1998), "Stochastic Volatility: Likelihood Inference and Comparison with ARCH Models," *Review of Economic Studies*, 65, 361-393.
- Marron, J.S. (1987), "A Comparison of Cross-Validation Techniques in Density Estimation," *Annals of Statistics*, 15, 152-162.
- Roberts, G.O. (1996), "Markov Chain Concepts Related to Sampling Algorithms," in *Markov Chain Monte Carlo in Practice*, eds. W.R. Gilks, S. Richardson, and D.J. Spiegelhalter, London: Chapman & Hall, pp. 45-57.
- Sain, S.R., Baggerly, K.A., and Scott, D.W. (1994), "Cross-Validation of Multivariate Densities," *Journal of the American Statistical Association*, 89, 807-817.
- Scott, D.W. (1992), *Multivariate Density Estimation: Theory, Practice, and Visualization*, New York: John Wiley.
- Simonoff, J.S. (1996), *Smoothing Methods in Statistics*, New York: Springer-Verlag.
- Stanton, R. (1997), "A Nonparametric Model of Term Structure Dynamics and the Market Price of Interest Rate Risk," *The Journal of Finance*, 52, 1973-2002.
- Tse, Y.K., Zhang, X., and Yu, J. (2004), "Estimation of Hyperbolic Diffusion with Markov Chain Monte Carlo Simulation," *Quantitative Finance*, 4, 158-169.

Wand, M.P., and Jones, M.C. (1993), "Comparison of Smoothing Parameterizations in Bivariate Kernel Density Estimation," *Journal of the American Statistical Association*, 88, 520-528.

Wand, M.P., and Jones, M.C. (1994), "Multivariate Plug-in Bandwidth Selection," *Computational Statistics*, 9, 97-116.

Wand, M.P., and Jones, M.C. (1995), *Kernel Smoothing*, London: Chapman & Hall.

**Table 1:** MCMC results for data generated from  $f_E(\cdot)$ . The first panel is obtained through the algorithm for a diagonal bandwidth matrix ( $M_6$ ), while the second panel is obtained through the algorithm for a full bandwidth matrix ( $M_1$ ).

sample size	bandwidths	mean	standard deviation	batch-mean standard error	SIF	acceptance rate
200	$1/b_{11}$	0.70	0.08	0.0017	10.32	0.224
	$1/b_{22}$	0.75	0.07	0.0015	11.77	
500	$1/b_{11}$	0.68	0.05	0.0011	11.72	0.207
	$1/b_{22}$	0.66	0.05	0.0009	8.73	
1000	$1/b_{11}$	0.69	0.03	0.0006	9.83	0.216
	$1/b_{22}$	0.61	0.03	0.0007	11.65	
200	$b_{11}$	1.18	0.15	0.0035	14.48	0.245
	$b_{21}$	-1.38	0.34	0.0164	57.58	
	$b_{22}$	1.69	0.21	0.0098	51.78	
500	$b_{11}$	1.10	0.08	0.0016	11.41	0.265
	$b_{21}$	-1.58	0.27	0.0137	65.54	
	$b_{22}$	1.91	0.19	0.1920	52.87	
1000	$b_{11}$	1.27	0.07	0.0015	11.68	0.267
	$b_{21}$	-0.79	0.11	0.0028	16.02	
	$b_{22}$	1.61	0.08	0.0016	9.45	

**Table 2:** Estimated Kullback-Leibler information for bivariate densities.

	sample size	Kullback-Leibler information										
		$M_1$	$M_2$	$M_3$	$M_4$	$M_5$	$M_6$	$P_1$	$P_2$	$P_3$	$P_4$	$A_1$
$\hat{E}(\ln f_A) = -2.003$	200	0.046	0.046	0.066	0.101	0.101	0.045	0.117	0.044	0.120	0.203	0.205
	500	0.025	0.024	0.041	0.051	0.051	0.026	0.054	0.030	0.057	0.132	0.134
	1000	0.017	0.017	0.018	0.038	0.038	0.020	0.043	0.019	0.047	0.094	0.096
$\hat{E}(\ln f_B) = -3.099$	200	0.131	0.129	0.158	0.154	0.154	0.228	0.129	0.213	0.153	0.192	0.375
	500	0.074	0.075	0.091	0.094	0.094	0.150	0.075	0.124	0.093	0.112	0.284
	1000	0.042	0.042	0.054	0.058	0.058	0.095	0.040	0.067	0.056	0.067	0.235
$\hat{E}(\ln f_C) = -1.822$	200	0.032	0.032	0.053	0.089	0.089	0.037	0.100	0.050	0.119	0.105	0.114
	500	0.021	0.021	0.037	0.048	0.047	0.022	0.047	0.023	0.055	0.089	0.085
	1000	0.018	0.018	0.040	0.040	0.040	0.021	0.038	0.021	0.043	0.065	0.071
$\hat{E}(\ln f_D) = -3.072$	200	0.299	0.296	0.247	0.394	0.392	0.361	0.357	0.345	0.391	0.325	0.410
	500	0.121	0.121	0.129	0.226	0.226	0.220	0.223	0.197	0.263	0.230	0.327
	1000	0.084	0.084	0.101	0.161	0.161	0.140	0.144	0.135	0.187	0.163	0.255
$\hat{E}(\ln f_E) = -3.850$	200	0.256	0.254	0.281	0.260	0.260	0.258	0.487	0.417	0.488	0.268	0.461
	500	0.219	0.221	0.249	0.240	0.240	0.217	0.333	0.298	0.345	0.240	0.385
	1000	0.149	0.149	0.150	0.178	0.178	0.149	0.260	0.222	0.274	0.173	0.299

**Table 3:** MCMC results for data generated from  $f_F(\cdot)$ . The sample size is 1500.

	bandwidths	mean	standard deviation	batch-mean standard error	SIF	acceptance rate
diagonal matrix	$1/b_{11}$	0.56	0.03	0.0009	21.85	0.250
	$1/b_{22}$	0.58	0.03	0.0009	24.34	
	$1/b_{33}$	0.56	0.03	0.0009	29.25	
	$1/b_{44}$	0.58	0.03	0.0010	36.42	
	$1/b_{55}$	0.58	0.03	0.0009	34.14	
full matrix	$b_{11}$	1.81	0.10	0.0042	41.83	0.272
	$b_{21}$	-0.15	0.15	0.0106	130.54	
	$b_{22}$	1.73	0.09	0.0033	36.26	
	$b_{31}$	0.11	0.18	0.0143	155.34	
	$b_{32}$	-0.15	0.13	0.0076	85.27	
	$b_{33}$	1.80	0.10	0.0031	25.31	
	$b_{41}$	-0.12	0.14	0.0084	93.56	
	$b_{42}$	-0.09	0.14	0.0099	133.07	
	$b_{43}$	-0.02	0.14	0.0083	93.30	
	$b_{44}$	1.74	0.10	0.0041	46.56	
	$b_{51}$	0.00	0.14	0.0084	88.95	
	$b_{52}$	0.07	0.14	0.0098	120.43	
	$b_{53}$	0.05	0.16	0.0114	134.69	
	$b_{54}$	0.18	0.13	0.0087	103.13	
	$b_{55}$	1.78	0.10	0.0042	47.31	

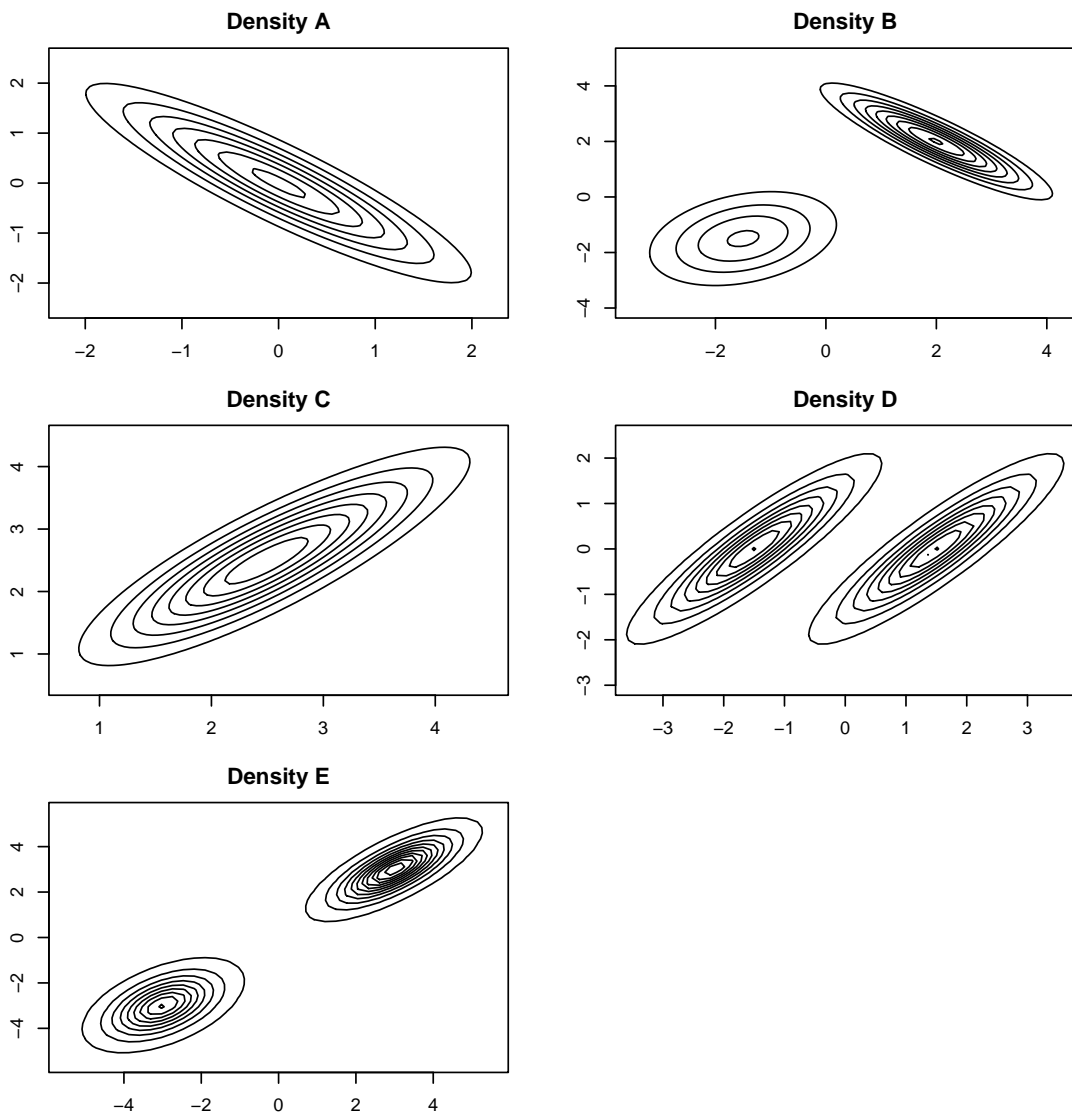
**Table 4:** Estimated Kullback-Leibler information for multivariate densities.

	sample size	Kullback-Leibler information						
		$M_1$	$M_2$	$M_3$	$M_4$	$M_5$	$M_6$	$A_1$
$\hat{E}(\ln f_F) = -7.9283$	500	0.178	0.177	0.539	0.441	0.441	0.186	1.262
	1000	0.127	0.126	0.505	0.304	0.304	0.162	1.235
	1500	0.118	0.117	0.470	0.276	0.276	0.141	1.545
$\hat{E}(\ln f_G) = -7.7934$	500	0.224	0.224	0.548	0.223	0.223	0.381	1.772
	1000	0.148	0.148	0.438	0.144	0.144	0.303	1.604
	1500	0.152	0.151	0.402	0.149	0.149	0.291	1.571
$\hat{E}(\ln f_H) = -9.2232$	500	0.774	0.771	1.147	0.746	0.746	0.915	2.222
	1000	0.687	0.685	1.149	0.677	0.677	0.846	1.862
	1500	0.696	0.696	1.029	0.679	0.680	0.845	1.992
$\hat{E}(\ln f_I) = -7.5123$	500	0.182	0.180	0.668	0.335	0.334	0.206	1.319
	1000	0.141	0.140	0.466	0.272	0.272	0.153	1.112
	1500	0.127	0.126	0.423	0.242	0.242	0.148	1.100
$\hat{E}(\ln f_J) = -7.3760$	500	0.288	0.282	0.725	0.479	0.479	0.247	1.342
	1000	0.142	0.141	0.662	0.331	0.331	0.166	1.204
	1500	0.109	0.109	0.537	0.270	0.270	0.147	1.318

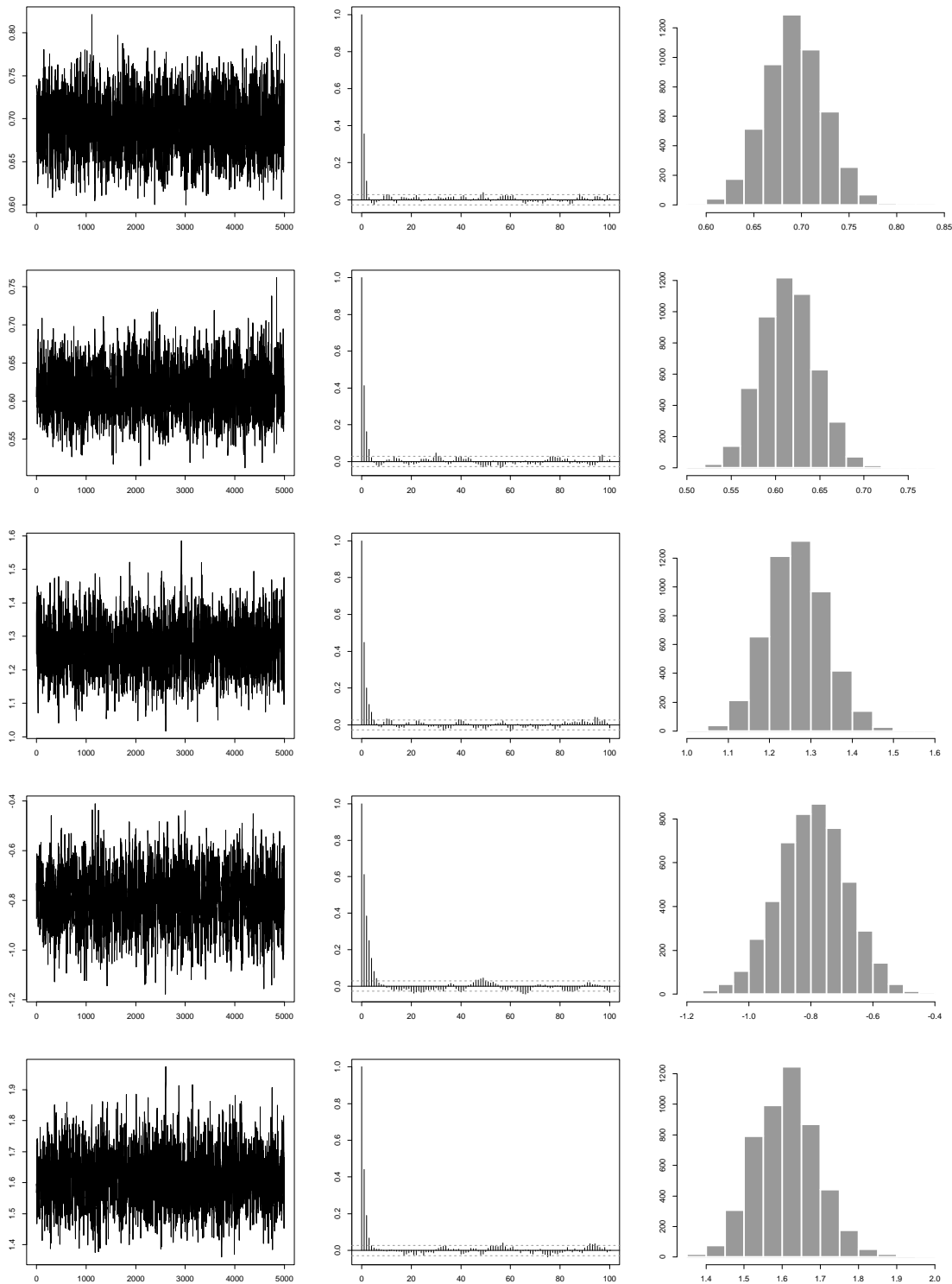
**Table 5:** MCMC results obtained from the Earthquake data.

	bandwidths	mean	standard deviation	batch-mean standard error	SIF	acceptance rate
diagonal matrix	$1/b_{11}$	0.003	0.0001	0.000003	9.07	0.254
	$1/b_{22}$	0.003	0.0001	0.000003	12.60	
	$1/b_{33}$	0.715	0.0383	0.000873	12.96	
full matrix	$b_{11}$	311.65	0.07	0.002	15.80	0.246
	$b_{21}$	101.53	0.10	0.005	62.21	
	$b_{22}$	388.57	0.10	0.003	15.84	
	$b_{31}$	147.45	0.13	0.008	89.38	
	$b_{32}$	97.21	0.16	0.011	118.86	
	$b_{33}$	1.65	0.27	0.012	47.54	

**Figure 1:** *Contour graphs of the proposed bivariate densities.*

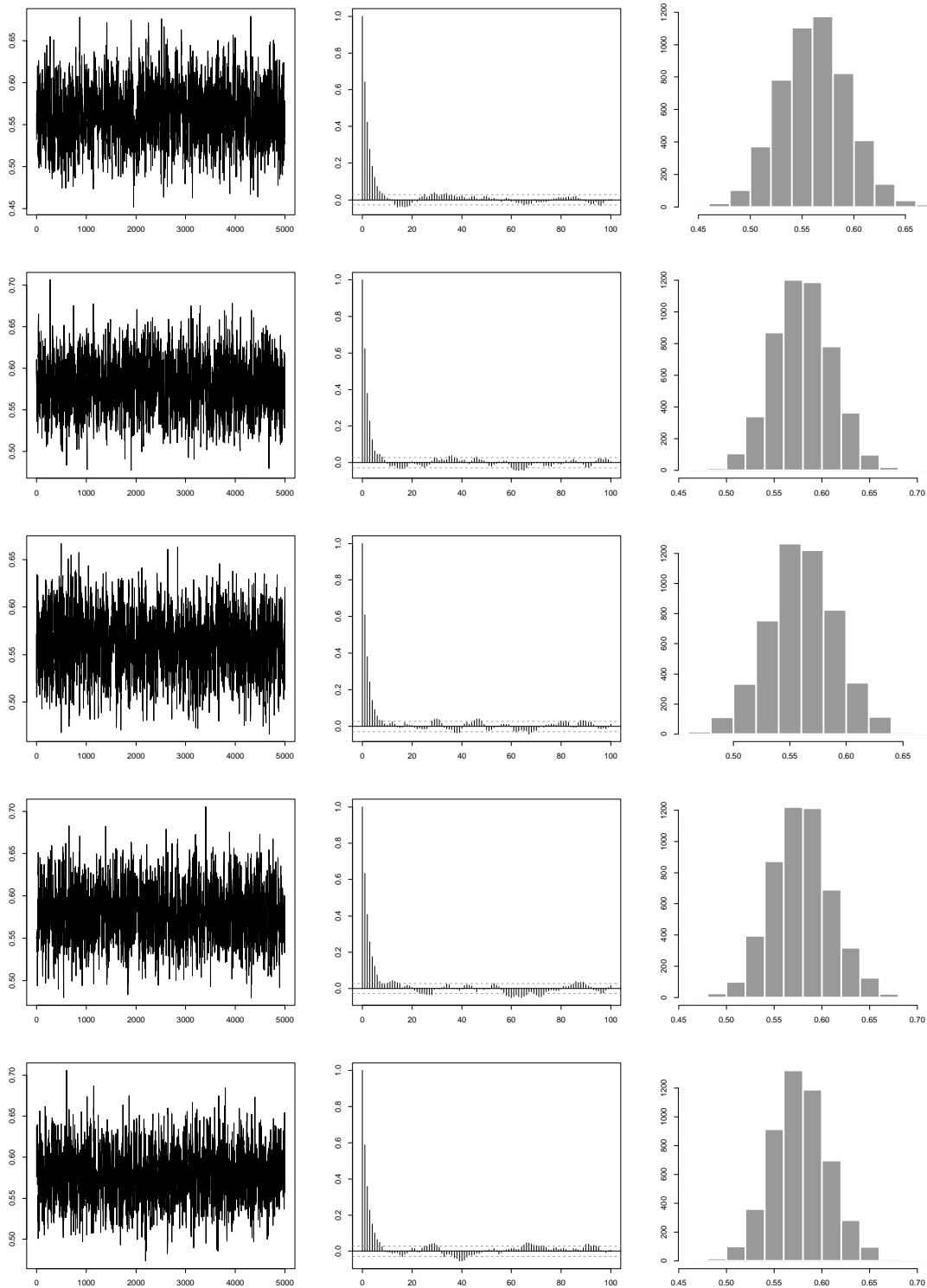


**Figure 2:** MCMC results for the simulated bivariate data. Columns (from left to right) show the sampled paths, their associated autocorrelation functions and histograms. The first two rows represent nonzero elements of  $H^{1/2}$ , while the rest three rows represent elements of  $B$ .





**Figure 3:** MCMC results for the simulated multivariate data. Columns (from left to right) represent the sampled paths, their associated autocorrelation functions and histograms. Rows represent nonzero elements of  $H^{1/2}$ .



**Figure 4:** The 98% highest density region for the earthquake data showing four views looking from north, east, south and west. Negative log-depth is on the vertical axis.

