

QUT Digital Repository:  
<http://eprints.qut.edu.au/>



This is the post-print, accepted version of this article. Published as:

McGrory, Clare A. and Pettitt, Anthony N. and Faddy, Malcolm (2009) *A fully Bayesian approach to inference for Coxian phase-type distributions with covariate dependent mean*. *Computational Statistics and Data Analysis*, 53(12). pp. 4311-4321.

© Copyright 2009 Elsevier B.V. All rights reserved.

# A Fully Bayesian Approach to Inference for Coxian Phase-Type Distributions with Covariate Dependent Mean

C.A. McGrory <sup>a,\*</sup>, A.N. Pettitt <sup>a,b</sup> M.J. Faddy <sup>a</sup>

<sup>a</sup>*Queensland University of Technology*

<sup>b</sup>*Lancaster University*

---

## Abstract

Phase-type distributions represent the time to absorption for a finite state Markov chain in continuous time, generalising the exponential distribution and providing a flexible and useful modelling tool. We present a new reversible jump Markov chain Monte Carlo scheme for performing a fully Bayesian analysis of the popular Coxian subclass of phase-type models; the convenient Coxian representation involves fewer parameters than a more general phase-type model. The key novelty of our approach is that we model covariate dependence in the mean whilst using the Coxian phase-type model as a very general residual distribution. Such incorporation of covariates into the model has not previously been attempted in the Bayesian literature. A further novelty is that we also propose a reversible jump scheme for investigating structural changes to the model brought about by the introduction of Erlang phases. Our approach addresses more questions of inference than previous Bayesian treatments of this model and is automatic in nature. We analyse an example dataset comprising lengths of hospital stays of a sample of patients collected from two Australian hospitals to produce a model for a patient's expected length

of stay which incorporates the effects of several covariates. This leads to interesting conclusions about what contributes to length of hospital stay with implications for hospital planning. We compare our results with an alternative classical analysis of these data.

*Key words:* Coxian Phase-type model, Phase-type distribution, Reversible jump Markov chain Monte Carlo, Bayesian analysis, Erlang distribution, Covariate Effects

---

## 1 Introduction

Phase-type models generalise the exponential distribution and are characterised by an underlying finite Markov chain that has one absorbing state. This underlying Markov process passes through a number of transient states, or phases, until eventually being absorbed. Therefore, the phase-type model is the distribution of the time until absorption for a finite Markov process. This is useful in many application areas: phase-type models have been used to analyse hospital length of stay (LoS) data using maximum likelihood-based approaches ([11], [12], [13], [22] and [32]), they have been successfully used in risk analysis ([1], [2]) and queueing theory ([7]), and they can be fitted using the EM algorithm ([3]). There are several subclasses of phase-type distributions; in this paper we focus on Bayesian inference for the highly versatile and popular Coxian subclass of phase-type models. In the Bayesian litera-

---

\* Corresponding author: School of Mathematical Sciences, Queensland University of Technology, GPO Box 2434, Brisbane, Queensland, 4001, Australia, Tel. (61)7-3138-1287, Fax.: (61)7-3138-1508

*Email address:* c.mcgrory@qut.edu.au (C.A. McGrory).

14 ture, phase-type models have been much less well explored. Bayesian Markov  
15 chain Monte Carlo techniques for general phase-type models are explored in  
16 [8], but this is limited to the fixed dimension case. A reversible jump Markov  
17 chain Monte Carlo (RJMCMC) approach which allows the number of transient  
18 phases in the model to vary is taken in [4] and [5]. However, in [4] and [5] an  
19 alternative mixture representation of the Coxian phase-type model, in terms  
20 of a mixed generalised Erlang distribution, is used rather than the matrix ex-  
21 ponential formulation used in this paper. Although these two representations  
22 are mathematically equivalent, in [4] and [5] a number of latent variables were  
23 introduced, which had to be imputed in the RJMCMC scheme. The introduc-  
24 tion of latent variables is not necessary here as we use the matrix exponential  
25 formulation, which has advantages when seeking reasonable acceptance rates  
26 for dimension-changing proposals in RJMCMC algorithms as there are fewer  
27 terms involved in the likelihood. It has also been previously noted that there  
28 is potential for unreliability when mixture type models are fitted using RJM-  
29 CMC (see [21], for example).

30 In this paper we present a novel Bayesian approach in which the Coxian phase-  
31 type model is used as a very general residual distribution. The incorporation  
32 of a covariate dependent mean into the model has not previously been at-  
33 tempted in the Bayesian literature. In the case of regression models, the use  
34 of phase-type distributions allows the error structure of the standard gener-  
35 alised linear model, which is usually a gamma or inverse Gaussian distribution  
36 for positive continuous data, to be more flexible to accommodate, for exam-  
37 ple, long tailedness and a mode near zero simultaneously. This makes these  
38 distributions particularly suited to hospital length of stay (LoS) modelling ap-  
39 plications, such as the one we consider in this paper, where the data typically

40 exhibit these features. In this context, phase-type modelling should result in  
41 more efficient estimation of the covariate dependence than one would obtain  
42 by using a standard exponential family distribution. The phases may or may  
43 not have an interpretation in the context of the application, but our focus here  
44 is on the estimation of the covariate dependence. The aim is to identify factors  
45 leading to increased LoS, which in turn leads to bed occupancy problems, thus  
46 having implications for efficient health-care facility and budget planning. This  
47 is an active research area and various other techniques have been applied to  
48 this problem, examples include the use of classical queuing theory to represent  
49 patient flow through various phases of treatment or centers of care (see [14],  
50 for example) and the use of a stochastic compartmental modelling approach  
51 (see, for instance [30]). Refer to [23] for a useful overview of the directions  
52 that research in this area has taken.

53 In our novel approach we develop an RJMCMC ([17]) analysis of data mod-  
54 elled by a Coxian phase-type distribution. The well-known paper [26] describes  
55 how RJMCMC can be used for mixture model analysis and [27] adapts these  
56 ideas to the hidden Markov model setting. Our Coxian model differs from the  
57 standard Markovian model in that it has additional constraints that must be  
58 taken into consideration in the construction of an appropriate RJMCMC algo-  
59 rithm. The difficulties associated with designing an RJMCMC scheme which  
60 will adequately explore the posterior are well-known, but we have been able to  
61 construct a sampler for this model which traverses the target distribution well.  
62 Our modelling of covariate dependency will also be useful in other applica-  
63 tions. Another contribution of this paper is to devise an RJMCMC scheme for  
64 exploring the inclusion into the phase-type model of an Erlang component,  
65 where specific structure leads to a more peaked mode. Using an RJMCMC

66 scheme we can automatically select the number of transient phases as well as  
67 their associated rate parameters, and estimate the covariate dependence (the  
68 number of covariates is fixed in our scheme, but this could also be estimated  
69 if desired). However, we still have the capability of exploring the important  
70 model features mentioned above. We demonstrate our new RJMCMC ap-  
71 proach with an application to modelling the effects of several covariates on  
72 the length of stay of patients in two Australian hospitals.

73 In Section 2 we describe the Coxian subclass of phase-type distributions and  
74 in Section 3 we describe our Bayesian formulation of the model. In Section  
75 4 we present our RJMCMC methodology. In Section 5 we demonstrate the  
76 technique through analysing the hospital LoS data, which leads to conclu-  
77 sions about the effect of several factors on increasing length of stay. Section 6  
78 explores the introduction of Erlang components into the model via RJMCMC  
79 and Section 7 concludes the paper.

## 80 **2 Coxian Phase-Type Distributions**

81 A phase-type distribution describes a Markov process,  $\{X(t); t \geq 0\}$ , say,  
82 where the system moves through some or all of  $K$  transient states, or phases,  
83 before moving to a single absorbing state  $K + 1$ . See [25] for a full description.  
84 The phases are governed by the transition probabilities

$$\begin{aligned} P(X(t + \delta t) = j + 1 | X(t) = j) &= \lambda_j \delta t + o(\delta t), & j = 1, \dots, K - 1 \\ P(X(t + \delta t) = K + 1 | X(t) = j) &= \mu_j \delta t + o(\delta t), & j = 1, \dots, K. \end{aligned}$$

85 Here  $\delta t$  represents a small time increment. The  $\{\lambda_j\}$  are the transition rates  
86 between the transient states and the  $\{\mu_j\}$  describe the transition from any of

87 the transient phases to the absorbing state.

88 In the Coxian phase-type model (see [10]) the system starts in the first phase  
 89 and then moves through the transient phases sequentially before eventually  
 90 being absorbed from any one of them. See Figure 1a for an illustration.

91 The probability density function of the time spent moving through the tran-  
 92 sient states before absorption is  $f(t) = \mathbf{p} \exp\{\mathbf{Q}t\}\mathbf{q}$ , where the infinitesimal  
 93 generator  $\mathbf{Q}$  is given by

$$\mathbf{Q} = \begin{pmatrix} -(\lambda_1 + \mu_1) & \lambda_1 & 0 & \dots & 0 & 0 \\ 0 & -(\lambda_2 + \mu_2) & \lambda_2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & -(\lambda_{K-1} + \mu_{K-1}) & \lambda_{K-1} \\ 0 & 0 & 0 & \dots & 0 & -\mu_K \end{pmatrix},$$

94 and the vectors  $\mathbf{p}$  and  $\mathbf{q}$  take the forms  $\mathbf{p} = (1 \ 0 \ \dots \ 0)$  and  $\mathbf{q} = (\mu_1 \ \mu_2 \ \dots \ \mu_K)^T$ .

95 Here  $\exp\{\cdot\}$  represents the matrix exponential function and we compute this  
 96 using Matlab.

97 The marginal distribution  $[\pi_1(t) \ \dots \ \pi_K(t)] = \mathbf{p} \exp\{\mathbf{Q}t\}$ , describes the proba-  
 98 bility,  $\pi_j(t)$ , that the system is in state  $j$ , where  $j \in [1 : K]$ , at time point  $t$ .

99 The survivor function can be derived from this if it is of interest. The Coxian  
 100 subclass describes any phase-type distribution with a generator matrix  $\mathbf{Q}$  that  
 101 has real eigenvalues and includes the exponential and Erlang distributions.  
 102 Reference [19] describes two algorithms for computing a Coxian representa-

103 tion from a more general phase-type distribution with a generator matrix that  
 104 has real eigenvalues.

105 We can introduce covariate dependency into the model so that the mean  
 106 absorption time is given by the log-linear regression  $\exp\{a + \mathbf{b}^T \mathbf{X}\}$ , where  
 107  $\mathbf{X} = (X_1, \dots, X_c)$  are the covariate values and  $\mathbf{b} = (b_1, \dots, b_c)$  are their co-  
 108 efficients. The expectation of time spent in the system is given by  $\mathbf{E}(T) =$   
 109  $(-1) \mathbf{p} \mathbf{Q}^{-1} (1 \ 1 \ \dots \ 1)^T$  ([25]). Therefore, to incorporate the desired depen-  
 110 dency, we scale the transition rate matrix appropriately as  $\exp\{-\mathbf{b}^T X\} \mathbf{Q}$ ,  
 111 with the intercept term  $a$  given by  $\exp(a) = (-1) \mathbf{p} \mathbf{Q}^{-1} (1 \ 1 \ \dots \ 1)^T$ . In [11]  
 112 covariates are also incorporated in this way in a classical approach, but this  
 113 has not been done in previous Bayesian analyses of these distributions.

### 114 **3 Bayesian Model Formulation for a phase-type Model with an** 115 **Unknown number of Phases**

116 Given observations comprising absorption times  $t_1, \dots, t_n$  from a phase-type  
 117 distribution with  $K$  transient states, and putting  $\boldsymbol{\theta}_K = (\boldsymbol{\lambda}, \boldsymbol{\mu}, \mathbf{b})$ , the likeli-  
 118 hood is given by  $p(t|\boldsymbol{\theta}_K, K) = \prod_{i=1}^n \mathbf{p} \exp\{\mathbf{Q}t_i\} \mathbf{q}$ . The transition rates for the  
 119 transient and the absorbing states are given Gamma prior distributions inde-  
 120 pendent of  $K$ :  $\lambda_j \sim Ga(\alpha_j, \beta_j)$  and  $\mu_j \sim Ga(\gamma_j, \delta_j)$ , where  $Ga(\cdot, \cdot)$  corresponds  
 121 to the Gamma probability density function and  $\{\alpha_j\}$ ,  $\{\beta_j\}$ ,  $\{\gamma_j\}$ , and  $\{\delta_j\}$  are  
 122 hyperparameters. The number of phases  $K$  and the covariates  $\mathbf{b}$  must be as-  
 123 signed prior distributions that are appropriate for the application at hand (we  
 124 specify these in the context of our application later). Our posterior distribution  
 125 has the form



$$\begin{aligned}
p(\boldsymbol{\theta}_K, K | \mathbf{t}) &\propto p(\mathbf{t} | \boldsymbol{\theta}_K, K) p(\boldsymbol{\theta}_K | K) p(K) p(\mathbf{b}) \\
&= \prod_{i=1}^n \mathbf{p} \exp\left(\exp(-\mathbf{b}^T X) \mathbf{Q} t_i\right) \mathbf{q} \prod_{j=1}^{K-1} \frac{1}{\Gamma(\alpha_j)} \frac{1}{\beta_j^{\alpha_j}} \lambda_j^{\alpha_j-1} \exp\left(-\frac{\lambda_j}{\beta_j}\right) \\
&\quad \times \prod_{j=1}^K \frac{1}{\Gamma(\gamma_j)} \frac{1}{\delta_j^{\gamma_j}} \mu_j^{\gamma_j-1} \exp\left(-\frac{\mu_j}{\delta_j}\right) \times p(K) \times p(\mathbf{b}).
\end{aligned} \tag{1}$$

## 126 4 Reversible Jump Markov Chain Monte Carlo Approach

127 RJMCMC techniques (see [17]) allow us to fully explore the available param-  
128 eter space by moving or jumping between models with a varying number of  
129 phases. We devised a RJMCMC scheme for this Coxian model, something for  
130 which we found no guidance in the literature, which enables transdimensional  
131 moves to occur with good acceptance rates. At each iteration of our algorithm  
132 we randomly choose, with equal probability, one of the following three move  
133 types: perform a fixed dimension parameter update, split a phase in two or  
134 combine two existing phases into one, the birth of a new phase or the death  
135 of an existing phase.

### 136 4.1 Metropolis-Hastings Fixed Dimension Parameter Update Move

137 We follow standard methods in the literature for updating the rate parameters  
138 and the regression parameters via Metropolis-Hastings.

### 139 4.2 Dimension Changing Reversible Jump Moves

140 We denote the current number of phases in the model by  $K$ , and the proposed  
141 number by  $K^*$ , where  $K^*$  is restricted to be equal to  $K - 1$  or  $K + 1$ . We

142 assume that the maximum potential number of phases is fixed at  $K = K_{max}$ ,  
 143 say. We propose a change of the model parameters from  $\boldsymbol{\theta}_K$  to  $\boldsymbol{\theta}_{K^*}$  through a  
 144 bijective mapping from the parameter space  $(\boldsymbol{\theta}_K, u, v)$  to  $(\boldsymbol{\theta}_{K^*}, u^*, v^*)$ , where  
 145  $u, v, u^*$  and  $v^*$  are auxiliary variables introduced so that dimensionality is the  
 146 same in the current and proposed parameter spaces. These moves are accepted  
 147 with probability  $\min(R, 1)$  where  $R$  is given by

$$\frac{p(t|\boldsymbol{\theta}_{K^*}, K^*)p(\boldsymbol{\theta}_{K^*})p(K^*)}{p(t|\boldsymbol{\theta}_K, K)p(\boldsymbol{\theta}_K)p(K)} \times \frac{J_{K^*,K} p(u^*, v^*|K^*, K, \boldsymbol{\theta}_{K^*})}{J_{K,K^*} p(u, v|K, K^*, \boldsymbol{\theta}_K)} \times \left| \frac{\partial(\boldsymbol{\theta}_{K^*}, u^*, v^*)}{\partial(\boldsymbol{\theta}_K, u, v)} \right|. \quad (2)$$

148 The first term in (2) is the ratio of the likelihood times prior (see (1)) for the  
 149 proposed and current parameter values, and the second is the ratio of proposal  
 150 probabilities. We denote the probability of moving from  $K$  to  $K^*$  phases by  
 151  $J_{K,K^*}$ . The third term is the Jacobian for the transformation.

152 We denote by  $\mu$  and  $\lambda$  the rate parameters associated with the phase of in-  
 153 terest in the model of lower dimension, and by  $\mu_a, \mu_b, \lambda_a$  and  $\lambda_b$  the rates  
 154 associated with the two phases of interest in the model of higher dimension.  
 155 It can be challenging to define a suitable mapping, particularly in the case of  
 156 practically driven applications, as it is often difficult to obtain good mixing  
 157 across model dimensions. To obtain reasonable acceptance rates for proposed  
 158 transdimensional jumps, one requires an appropriately centered proposal dis-  
 159 tribution with well tuned parameters. However, it is not generally obvious  
 160 how best to achieve this ([9] makes some suggestions in this regard). Here,  
 161 we take the approach of constructing our proposal distributions so that the  
 162 proposed parameters are not too distant from the current parameters, and  
 163 we take a matching approach to the construction of our mapping between di-  
 164 mension spaces. We ensure that probability of absorption and the mean time

165 in the phase(s) are matched in the current and proposed dimension spaces  
 166 corresponding to equations (3) and (4), given below.

$$\frac{\mu}{\mu + \lambda} = \frac{\mu_a}{\mu_a + \lambda_a} + \left( \frac{\lambda_a}{\mu_a + \lambda_a} \times \frac{\mu_b}{\mu_b + \lambda_b} \right) \quad (3)$$

$$\frac{1}{\mu + \lambda} = \frac{1}{\mu_a + \lambda_a} + \frac{1}{\mu_b + \lambda_b}. \quad (4)$$

### 167 Split and Combine Moves

168 The design of our split and combine moves does not allow splits or combines of  
 169 the final phase. In all other cases we assume equal probabilities of splitting or  
 170 combining. Figure 1b illustrates these move types graphically. In the combine  
 171 move we have  $(\mu_a, \lambda_a, \mu_b, \lambda_b) \rightarrow (u, v, \mu, \lambda)$ . We put  $u = \mu_a$  and  $v = \lambda_a$ , then  
 172 solving (3) and (4) gives us

$$\mu = \frac{\mu_a \mu_b + \mu_a \lambda_b + \lambda_a \mu_b}{\mu_a + \lambda_a + \mu_b + \lambda_b}$$

$$\lambda = \frac{\lambda_a \lambda_b}{\mu_a + \lambda_a + \mu_b + \lambda_b}.$$

173 In this case, the Jacobian is given by

$$\frac{(\mu_a + \lambda_a)^2 \lambda_a}{(\mu_a + \lambda_a + \mu_b + \lambda_b)^3}.$$

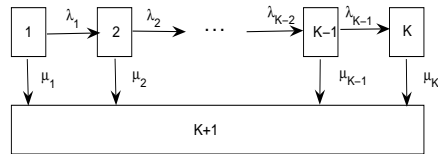
174 Our corresponding split move involves the reverse transition (see Figure 1b).  
 175 In the higher dimension space, we set  $\mu_a$  and  $\lambda_a$  to be equal to the simu-  
 176 lated auxiliary variables  $u$  and  $v$ , respectively, where  $u \sim N_T(2\mu, \sigma^2)$  and  
 177  $v \sim N_T(2\lambda, \sigma^2)$ . Here  $N_T(\cdot, \cdot)$  denotes the Normal density function truncated  
 178 at zero with mean  $\mu$  and suitable tuned variance  $\sigma^2$ . We simulate from the  
 179 truncated distribution since we cannot have negative values for the rate pa-  
 180 rameters. By solving (3) and (4) we obtain

$$\mu_b = \frac{\mu_a^2 \lambda + \mu_a \lambda_a \lambda - \lambda_a \mu \mu_a - \lambda_a^2 \mu}{\lambda_a (-\mu_a - \lambda_a + \mu + \lambda)}$$

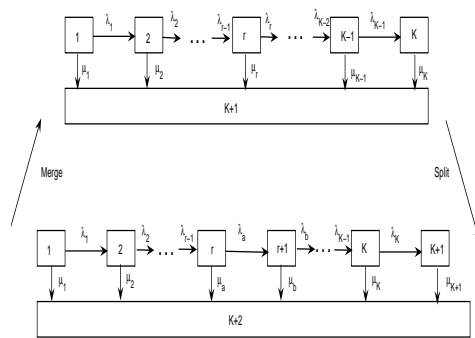
$$\lambda_b = -\frac{(\mu_a + \lambda_a)^2 \lambda}{\lambda_a (-\mu_a - \lambda_a + \mu + \lambda)}.$$

181 If either of  $\mu_b$  or  $\lambda_b$  is negative when calculated the proposal is rejected. The  
 182 Jacobian for the split move is the reciprocal of the corresponding expression  
 183 for the combine move. The acceptance ratio for each of the above moves is  
 184 then given by substituting the appropriate values into (2). These acceptance  
 185 ratios are reciprocals of one another.

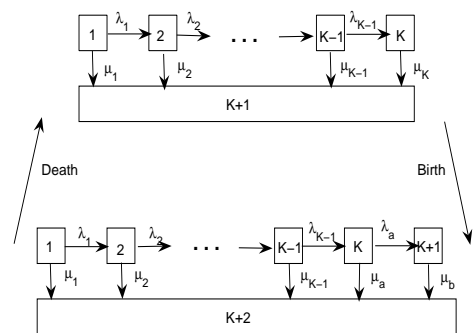
Fig. 1. Diagrammatic representations of (a) the Coxian phase-type model, (b) the effects of our RJMCMC split and merge moves and (c) the effects of our RJMCMC birth and death moves.



(a) Coxian phase-type distribution



(b) Split and Merge Moves



12  
(c) Birth and Death Moves

187 These moves are only applied to the final phase in the current model in a given  
 188 iteration and bring about the birth of a new final phase or the death of the  
 189 existing final phase with equal probability (provided that  $1 < K < K_{max}$ ).  
 190 The death move makes the transition  $(\mu_a, \lambda_a, \mu_b) \rightarrow (u, v, \mu)$ ; see Figure 1c.  
 191 Putting  $u = \mu_a$  and  $v = \lambda_a$  and solving (3) and (4), we obtain

$$\mu = \frac{(\mu_a + \lambda_a)\mu_b}{(\mu_b + \mu_a + \lambda_a)}.$$

192 The Jacobian for the death move is given by

$$\frac{(\mu_a + \lambda_a)^2}{(\mu_b + \mu_a + \lambda_a)^2}.$$

193 In the reverse birth move, we generate  $u$  and  $v$  from the Normal distribution  
 194 truncated at 0 with mean  $\mu$  and variance  $\sigma^2$  and set  $\mu_a = u \sim N_T(\mu, \sigma^2)$  and  
 195  $\lambda_a = v \sim N_T(\mu, \sigma^2)$ . Again,  $\sigma^2$  is chosen to give reasonable rates of acceptance  
 196 for the move. To satisfy (3) and (4), we take  $\mu_b$  to be

$$\mu_b = \frac{(u + v)\mu}{u + v - \mu}.$$

197 If this results in a negative  $\mu_b$ , we reject the proposal. The Jacobian for the  
 198 death move is the reciprocal of that for the corresponding reverse birth move  
 199 described above. We can obtain the acceptance ratio for the birth and the  
 200 death moves by substituting the appropriate quantities into (2) and these are  
 201 of course reciprocals of each other.

## 202 5 Application: Modelling Length of Stay in Hospital

203 The identification of factors that are likely to increase a patient's LoS is a  
204 key goal for hospital planners. By addressing issues that lead to a longer LoS,  
205 health care costs can be reduced. LoS data are characteristically highly right-  
206 skewed making it difficult to fit them with other distributions. See [11] for  
207 a discussion of some of the difficulties associated with modelling LoS data.  
208 Phase-type distributions provide the flexibility that is required to capture the  
209 distributional characteristics of this type of data. The phases may only be  
210 artifacts of the modelling, but could have a physical interpretation in relation  
211 to the context. However, our focus here is on the estimation and interpretation  
212 of the covariate effects.

213 We applied our method to a dataset previously analysed using classical maxi-  
214 mum likelihood techniques ([11]), with our results complementing this analy-  
215 sis. The dataset comprises the lengths of hospital stay of 1901 patients all of  
216 whom were at least 18 years of age. These data were collected from two hospi-  
217 tals in S.E. Queensland, Australia, between October 2002 and January 2003.  
218 Patients were recruited from a range of specialities, but only those whose ad-  
219 missions were considered uncomplicated contributed to the data. The observed  
220 lengths of stay ranged from 0.44 to 170.9 days. The sample mean length of  
221 stay was 7.25 days. Information on ten covariates widely believed to be of  
222 relevance to length of hospital stay was also available for each patient and was  
223 included in our model. Details of the covariate information are given in the  
224 Appendix. (Note that this dataset is a part of a larger dataset collected in a  
225 prospective study and that [15] and [16] provide further details of the method  
226 of collection.)

227 For each patient we also have a predicted length of stay based on the patient's  
228 admission category for an uncomplicated admission. This was obtained from  
229 the Australian Institute of Health and Welfare. The logarithm of the predicted  
230 length of stay,  $x_0$ , was incorporated into our model as an offset variable. In  
231 this way we are modelling a patient's excess length of stay relative to the  
232 prediction and  $\exp\{a + \mathbf{b}^T X\}$  as  $\mathbf{E}(T/x_0)$  where  $T$  is the actual length of stay.

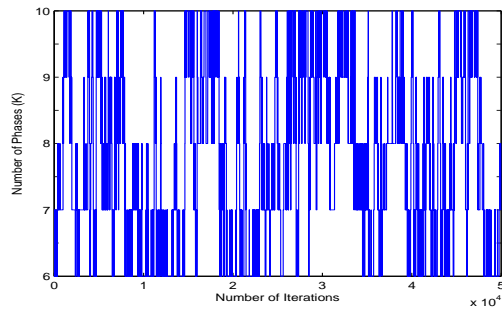
233 We assigned the covariate coefficients  $\mathbf{b}$  uniform priors over the range -10 to  
234 10. We assumed that the maximum potential number of phases in the model  
235 was fixed at  $K_{max} = 10$  and we chose a uniform prior distribution over 1 to  
236 10 for  $K$ . The hyperparameters chosen for the Gamma priors over the rate  
237 parameters were also chosen to be uninformative.

238 We performed 100 000 iterations of our RJMCMC algorithm and we discarded  
239 the first 50 000 of these iterations to allow for a burn-in period. The algorithm  
240 was tuned so that acceptance rates for fixed dimension updates of the param-  
241 eters  $\boldsymbol{\mu}$ ,  $\boldsymbol{\lambda}$  and  $\mathbf{b}$  were between 30% and 35%. The overall rate of acceptance  
242 for the dimension changing moves was around 7% which although low, is rea-  
243 sonably good for RJMCMC. The trace plot for the number of phases,  $K$ , over  
244 all iterations post burn-in, is given in Figure 2a. The most likely number of  
245 phases was six, having posterior probability of 0.27, followed by the seven-  
246 phase model which had posterior probability of 0.25 (see Figure 2b for the  
247 posterior distribution for  $K$ ). To further examine convergence, we ran our al-  
248 gorithm from two different starting points and thinned the observations to 1  
249 in 250. We then plotted the posterior probability that the number of phases  
250 was six at each iteration point. This plot is shown in Figure 2c. This, together  
251 with the trace plot, suggests that the scheme has converged. The posterior es-  
252 timate of the number of phases as six is in agreement with the classical analysis

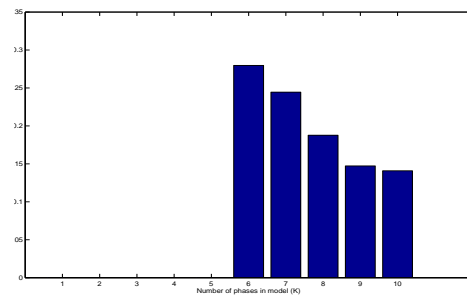


253 [11]. Figure 3 displays the posterior distributions of the parameter estimates  
254 from the six-phase fit. Since our main aim is to model the effects of covariates  
255 on the mean LoS, the actual parameterisation of the Coxian distribution be-  
256 comes irrelevant, as long as the residual variation is adequately described. In  
257 [5] the authors noted that it was necessary to impose an identifiability con-  
258 straint in their RJMCMC analysis of the mixture representation of Coxian  
259 model in order to obtain identifiability of the rate parameter estimates; such  
260 constraints could possibly be considered in our scenario if the rate parameters  
261 were of particular interest in the application. However, it is worth noting that  
262 in our results the posterior distributions for the rate parameters appear to be  
263 unimodal suggesting that identifiability was not a significant problem in our  
264 implementation.

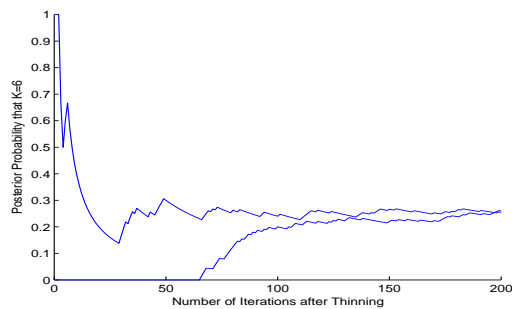
Fig. 2. These plots show the results of 50,000 iterations (after burn-in) of our RJMCMC sampler for the hospital length of stay data. (a) Trace plot of the number of phases ( $K$ ) in the model at each iteration. (b) The posterior distribution of the number of phases ( $K$ ). (c) Plot of the estimated posterior probability that the number of phases in the model is six at each iteration of two different runs of the RJMCMC algorithm.



(a)

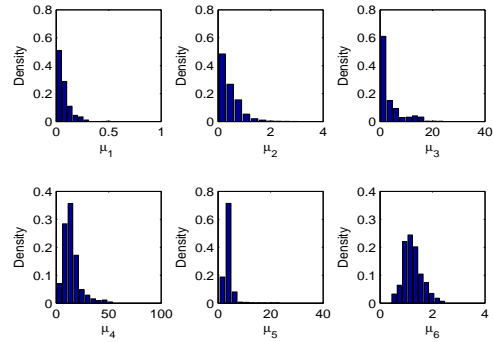


(b)

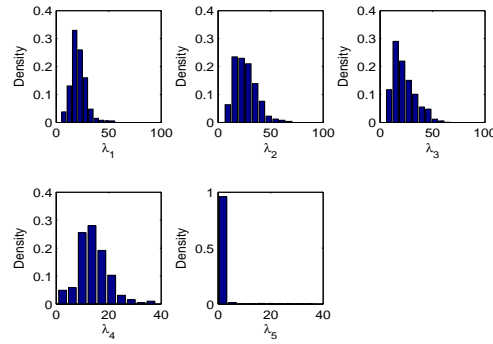


(c)

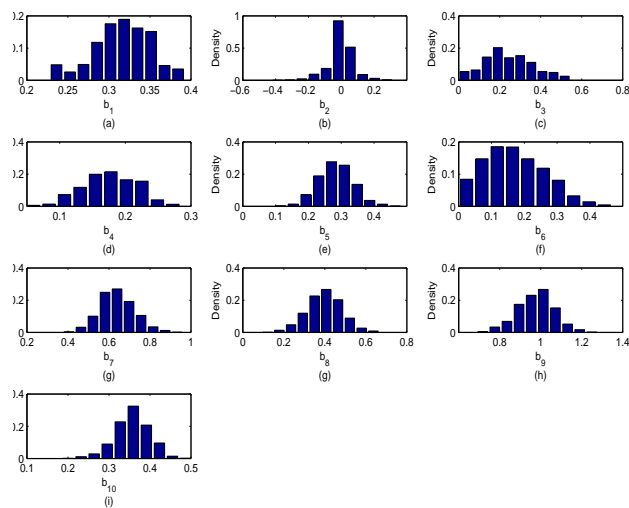
Fig. 3. Posterior distributions of (a) the  $\mu$ 's, (b) the  $\lambda$ 's and (c) the b's (after burn-in) from the six phase model fitted in the RJMCMC analysis of the hospital length of stay data.



(a)



(b)



(c)

265 The estimates of the covariate coefficients  $\mathbf{b}$  are of primary interest in the ap-  
 266 plication and the posterior distributions for these are reasonably symmetrical.  
 267 These can be used to estimate the effect that each of the covariates has on  
 268 increasing the length of the patient's stay beyond the initial prediction made  
 269 upon admission. Since the posterior distributions for the parameters exhibit  
 270 some skewness we used the posterior medians as parameter estimates. The  
 271 posterior medians (posterior standard deviations in brackets) of the intercept  
 272 parameter and covariate coefficients are given by

$$\begin{aligned}
 a &= -1.10(0.04) \\
 \mathbf{b} &= [0.32(0.03) \quad 0.01(0.07) \quad 0.22(0.12) \quad 0.18(0.04) \quad 0.28(0.05) \\
 &\quad 0.16(0.09) \quad 0.64(0.08) \quad 0.39(0.08) \quad 0.98(0.08) \quad 0.35(0.04)].
 \end{aligned}$$

273 Our posterior estimates for the covariate coefficients showed some similarity  
 274 with the maximum likelihood estimates in [11]. Based on our results, we can  
 275 see that contraction of a health care acquired infection (covariate 9) would be  
 276 expected to bring about the greatest increase in length of stay, while faecal  
 277 incontinence (covariate 7) was estimated to be the second most influential  
 278 factor and sex the least influential.

279 Inference about the effect of health care acquired infection is useful to hos-  
 280 pital planners, as health care acquired infections (HAIs) are widely believed  
 281 to place a substantial economic burden upon the health system. Moreover, a  
 282 recent study ([18]) has suggested that HAIs could be prevented in some cases.  
 283 However, [16] has highlighted that despite this consideration there have been  
 284 few published studies on the actual impact that the implementation of infec-  
 285 tion control programs might have in reducing the costs associated with HAIs.  
 286 To estimate what the economic benefits might be, we must first estimate the  
 287 effect of HAIs in real terms. Our estimated coefficient for the HAI covariate

288 was 0.98, with a 95% credible interval of 0.79 to 1.15 for that estimate. Based  
289 on our sample of patients, we would estimate that the contraction of an HAI  
290 would lead to an increased stay of 13.25 days on average, with 95% credible  
291 interval for this estimate of 7.89 days to 15.34 days.

292 Our estimate of the effect of the pressure ulcer covariate (covariate 6) is also  
293 worthy of comment since, as [15] points out, many previous studies have sug-  
294 gested that the development of pressure ulcers in hospital has a fairly sig-  
295 nificant effect on lengthening stay. This effect has been estimated as ranging  
296 from a 7 to a 50 day increase in stay for affected patients (references cited in  
297 [15]). However, the authors of [15] suggest that this effect has been overesti-  
298 mated, as they found that the occurrence of pressure ulcers would lead to an  
299 estimated median increase in stay of only 4.31 days (with a 95% confidence  
300 interval of 1.85 to 6.78 for this estimate.) Our results estimate the coefficient  
301 for the pressure ulcer covariate to be 0.16, with a 95% credible interval given  
302 by 0.02 to 0.36. This corresponds to an expected increase in LoS of 1.83 days  
303 on average, with a 95% credible interval for this estimate of 0.20 to 2.61 days.  
304 Therefore, our results also support the view that pressure ulcers may not have  
305 as much of a role in increasing LoS as has previously been suggested.

306 The similarities between our conclusions and those from more classical stud-  
307 ies lends support to the ability of our RJMCMC-based sampling scheme to  
308 obtain useful model estimates in practical applications. With our method the  
309 inference is performed directly, in contrast to the two-tier classical approach  
310 ([11]) of model identification and subsequent maximum likelihood parameter  
311 estimation. It is also worth noting that we reached this solution from start-  
312 ing values that were easily obtained from a simple generalised linear model  
313 fit, rather than multiple iterative searches with different starting values to

314 determine the maximum likelihood solution.

## 315 **6 Reversible Jump Scheme for Initial Erlang Phases**

316 In other analyses of data similar to those here, it has been found that an  
317 adequate model for the data corresponded to having several of the initial values  
318 of  $\mu$  equal to zero with the associated phases having equal values of  $\lambda$ . This  
319 introduces an initial Erlang component leading to a simpler model involving  
320 fewer parameters. To explore the effect this might have on our analysis, we  
321 conceived a move type which we call the birth of an Erlang phase, the reverse  
322 move being the death of an Erlang phase. The essence of this change is to  
323 set the current rate parameter  $\mu_1$  to be equal to zero. If  $\mu_1$  is already zero,  
324 then the move is carried out on  $\mu_2$  and so on. In this way we have developed  
325 an RJMCMC scheme that searches over competing distributional structures.  
326 We describe these moves in specific terms in the following sections. As before,  
327 the acceptance ratio for these moves is obtained by the substitution of the  
328 relevant values into (2).

### 329 *6.1 Birth and Death of the First Erlang Phase*

330 If  $\mu_1$  is currently nonzero, we choose our transformed parameters to satisfy  
331 equation (5) corresponding to matching the mean length of time in the first  
332 phase before and after it becomes an Erlang phase.

$$\frac{1}{\mu_1 + \lambda_1} = \frac{1}{\lambda_a}. \quad (5)$$

333 The birth of the Erlang component involves the transition  $(\mu_1, \lambda_1) \rightarrow (u, \lambda_a)$ .

334 Figure 4a provides an illustration of this move type. Choosing  $\lambda_a$  to satisfy  
335 (5) gives

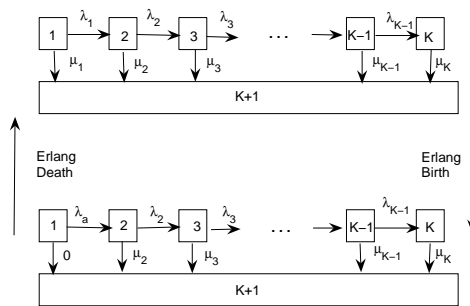
$$\begin{aligned}\lambda_a &= \mu_1 + \lambda_1 \\ u &= \frac{\lambda_1}{\mu_1 + \lambda_1}.\end{aligned}$$

336 The Jacobian for this move is given by

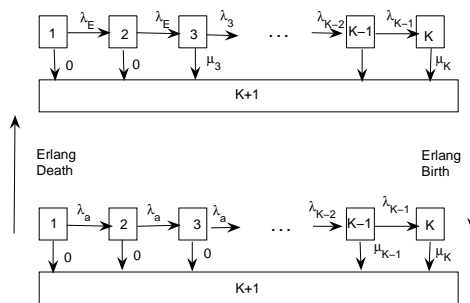
$$\frac{1}{\mu_1 + \lambda_1}.$$

337 The death of the Erlang phase (see Figure 4a), involves the opposite transi-  
338 tion  $(u, \lambda_a) \rightarrow (\mu_1, \lambda_1)$ . We generate our auxiliary variable  $u$  from a uniform  
339 proposal distribution  $u \sim Un(0, 1)$ , where  $Un(\cdot, \cdot)$  represents the uniform dis-  
340 tribution. Then we put  $\mu_1 = u\lambda_a$  and  $\lambda_1 = (1 - u)\lambda_a$ . This choice for  $\mu_1$  and  
341  $\lambda_1$  satisfies (5). The Jacobian is equal to  $\lambda_a$  (the inverse of the Jacobian for  
342 the reverse move).

Fig. 4. Diagrammatic representations of (a) the Erlang birth and death moves when jumping between a general phase model and a one Erlang phase model, and (b) the general Erlang birth and death moves that are performed when there is at least one Erlang phase present in the model.



(a) Birth/death of initial Erlang phase



(b) Example of the birth/death of a general Erlang phase



344 When one or more of the initial  $\mu$ 's have already been set to zero, birth of  
 345 another Erlang component must take into account the equal eigenvalue con-  
 346 straint in the Erlang part of the model. This change will involve two com-  
 347 ponents: the rate parameters for the  $r^{\text{th}}$  phase (the one we are considering  
 348 for incorporation into the Erlang distributed part of the model) and the rate  
 349 parameter for the existing Erlang phase or phases. We denote the latter by  
 350  $\lambda_E$ ; refer to Figure 4b for an illustration with  $r = 3$ . We construct our general  
 351 Erlang birth/death moves so that the mean time in the phases is matched  
 352 before and after the transformation corresponding to equation (6) below.

$$\frac{r-1}{\lambda_E} + \frac{1}{\mu_r + \lambda_r} = \frac{r}{\lambda_a} \quad (6)$$

353 The general birth of an Erlang move increases the number of Erlang phases  
 354 from  $r-1$  to  $r$  and involves the transition  $(\lambda_E, \mu_r, \lambda_r) \rightarrow (\lambda_a, u, v)$ . We put  
 355  $u = \mu_r$  and  $v = \lambda_r$ . Then from (6) we obtain

$$\lambda_a = \frac{r\lambda_E(\mu_r + \lambda_r)}{(r-1)(\mu_r + \lambda_r) + \lambda_E}.$$

356 The Jacobian for this move is given by the following expression

$$\frac{r(r-1)(\mu_r + \lambda_r)^2}{((r-1)(\mu_r + \lambda_r) + \lambda_E)^2}.$$

357 The death move involves the reverse transition  $(\lambda_a, u, v) \rightarrow (\lambda_E, \mu_r, \lambda_r)$ . Here  
 358 we put  $\lambda_r = v$  and  $\mu_r = u$ , where  $u \sim N_T(0, \sigma^2)$  and  $v \sim N_T(\lambda_a, \sigma^2)$ . We tune  
 359  $\sigma^2$  to give satisfactory acceptance rates for the move and solve (6) to obtain

$$\lambda_E = \frac{(r-1)\lambda_a(\mu_r + \lambda_r)}{r(\mu_r + \lambda_r) - \lambda_a}.$$

360 The Jacobian for this move is the inverse of the reverse general birth move.

361 We continue to use uninformative Gamma priors for the parameters  $\mu$  and  $\lambda$   
 362 in this scheme. However, when some phases in the model currently correspond  
 363 to an Erlang distribution, the shape and scale parameters of the corresponding  
 364 Gamma prior are multiplied by the current number of Erlang phases in the  
 365 model to give the prior distribution for the Erlang rate parameter. This prior  
 366 was also used in [20].

367 *6.3 Results from Applying Erlang Birth and Death Moves to the Hospital*  
 368 *Length of Stay Data*

369 We ran our Erlang birth/death algorithm using the six phase posterior esti-  
 370 mates from our initial RJMCMC analysis as a starting point. We performed  
 371 10, 000 iterations and discarded the first half of these. We found that the  
 372 most likely number of Erlang phases was two, having posterior probability of  
 373 0.94. The resulting posterior medians (posterior standard deviations given in  
 374 brackets) of the intercept parameter and covariate coefficients were as follows.

$$\begin{aligned} a &= -1.44(0.04) \\ \mathbf{b} &= [0.37(0.03) \quad 0.01(0.01) \quad 0.37(0.12) \quad 0.17(0.04) \quad 0.28(0.05) \\ &\quad 0.15(0.099) \quad 0.63(0.08) \quad 0.40(0.08) \quad 0.96(0.08) \quad 0.37(0.04)]. \end{aligned}$$

375 The regression coefficient posterior medians and standard deviations are very  
 376 similar for the two models except that the coefficient for  $x_3$  has changed from  
 377 0.22 to 0.37 and is more statistically significant. This model is simpler and we

378 have only used nine parameters to describe the phase-type model and ten to  
379 describe the regression part of the model for a dataset of nearly 2000 observa-  
380 tions. We have not reported the posterior distributions of the rate parameters  
381 as they may be subject to some lack of identifiability, but we note that they  
382 have unimodal distributions possibly indicating satisfactory identifiability.

## 383 **7 Conclusions**

384 Our extension of the reversible jump method to Coxian phase-type modelling  
385 with covariate dependent mean provides a fully formal Bayesian method for  
386 fitting such distributions to data and extends previous Bayesian analyses of  
387 this type of model. Our application to hospital LoS data has demonstrated  
388 that our approach can be used to provide valuable statistical inference for real  
389 world problems. In particular, posterior distributions for the number of phases  
390 and the regression parameters are produced, and we have also indicated that  
391 suitable starting values for the RJMCMC algorithm can be easily obtained.  
392 These advantages make this Bayesian approach attractive in practice.

393 We have also devised an RJMCMC method for automatically exploring the  
394 structure of the phase-type model to investigate the inclusion of an initial  
395 Erlang component which, in our case study, gave an improved and simpler  
396 structure for the model. Such modelling can be extended.

397 The phase-type distributions can be interpreted as providing a flexible and par-  
398 tially parametric extension to standard exponential family models, in particu-  
399 lar the gamma density family, while still maintaining a quadratic mean/variance  
400 relationship. Hence in the regression context such models should provide for

401 more robust estimation of regression coefficients. An alternative flexible ap-  
402 proach might be provided by fitting a normal mixture to the logarithm of the  
403 times, but this needs to be investigated. However, such an approach would  
404 not provide the structure of the phase-type model where such a structure may  
405 have a useful interpretation (e.g. hospital LoS) and it is doubtful whether it  
406 could simultaneously capture the mode near zero and the longtailedness of the  
407 data.

408 In our modelling we have not included the case where the covariates are also  
409 selected using the RJMCMC scheme. This would be straightforward to imple-  
410 ment, but it would be best to exercise caution in applications as there could  
411 be confounding between the selection of the number of phases and the covari-  
412 ates. This requires investigation. Other extensions of this theory include the  
413 exploration of cases where we have repeated measures observed on each sub-  
414 ject; this could be achieved through the use of a frailty term. If we wished to  
415 identify the rate parameters then the approach of [5] could straightforwardly  
416 be incorporated into our analysis.

417 Phase-type models are useful in any application where the data exhibit long  
418 tails, and there are many research fields in which this type of data arises in ad-  
419 dition to the applications we have already mentioned. For example, phase-type  
420 models have been used in web site performance optimisation ([31]), wireless  
421 communication system control ([29]), line transect sampling ([28]), gene find-  
422 ing ([24]) and ion channel modelling ([6]).

## 423 **Acknowledgements**

424 The authors' work was supported by Australian Research Council Discovery  
425 and Linkage Grants. Some of the computational resources and services used  
426 in this work were provided by the High Performance Computing and Research  
427 Support Group, Queensland University of Technology, Australia. We wish to  
428 thank the anonymous referees for some helpful comments.

## APPENDIX: Covariate Information Used in Modelling the Hospital Length of Stay Data

Covariate	Description	Range
$x_0$	predicted length of stay in days	1-72
$x_1$	log of age	2.9-4.61
$x_2$	sex (male/female)	binary 0/1
$x_3$	discharge destination (death/survive)	binary 0/1
$x_4$	admission type (emergency/non-emergency)	binary 0/1
$x_5$	anti-coagulant therapy during admission	binary 0/1
$x_6$	pressure ulcer during admission	binary 0/1
$x_7$	faecal incontinence during admission	binary 0/1
$x_8$	gastro-intestinal bleeding during admission	binary 0/1
$x_9$	health care acquired infection	binary 0/1
$x_{10}$	surgical procedure	binary 0/1

429 **References**

- 430 [1] Asmussen, S., 2000. Ruin probabilities, Advanced Series on Statistical Science  
431 and Applied Probability (Vol. 2), Singapore: World Scientific.
- 432 [2] Asmussen, S., and Bladt, M., 1996. Phase-type distributions and risk processes  
433 with state-dependent premiums, Scandinavian Actuarial Journal, 96, 19–36.
- 434 [3] Asmussen, S., Nerman, O., and Olsson, M., 1996. Fitting phase-type  
435 distributions via the EM-algorithm, Scandinavian Journal of Statistics, 23, 419–  
436 441.
- 437 [4] Ausìn, M.C., Lillo, R.E., Ruggeri, F., and Wiper, M.P., 2003. Bayesian  
438 modelling of hospital bed occupancy times using a mixed generalised Erlang  
439 distribution, in Bayesian Statistics 7, eds. J.M. Bernardo, M.J. Bayarri, J.O.  
440 Berger, A.P. David, D. Heckerman, A.F.M. Smith and M. West, Oxford: Oxford  
441 University Press, pp. 443–451.
- 442 [5] Ausìn, M.C., Wiper, M.P., and Lillo, R.E., 2008. Bayesian prediction of the  
443 transient behaviour and busy period in short-and long-tailed  $GI/G/1$  queueing  
444 systems, Computational Statistics and Data Analysis, 52, 1615–1635.
- 445 [6] Ball, F.G., Milne, R.K., and Yeo, G.F. 2000. Stochastic models for systems of  
446 interacting ion channels, IMA Journal of Medicine and Biology, 17, 263–293.
- 447 [7] Bertsimas, D., 1990. An analytic approach to a general class of  $G/G/c$  queueing  
448 systems, Operations Research, 38, 139–155.
- 449 [8] Bladt, M., Gonzalez, A., and Lauritzen, S.L., 2003. The estimation of phase-type  
450 related functionals using Markov chain Monte Carlo methods, Scandinavian  
451 Actuarial Journal, 4, 280–300.
- 452 [9] Brooks, S.P., Guidici, P., and Roberts, G.O., 2003. Efficient construction

- 453 of reversible jump Markov chain Monte Carlo proposal distributions (with  
454 discussion), *Journal of the Royal Statistical Society, Series B*, 57, 473–484.
- 455 [10] Cox, D.R., and Miller, H.D., 1965. *An Introduction to the Theory of Stochastic*  
456 *Processes*. London: Methuen & Co..
- 457 [11] Faddy, M.J., Graves, N., and Pettitt, A.N., 2009. Modeling length of stay in  
458 hospital and other right skewed data: comparison of phase-type, gamma and  
459 log-normal distributions, *Value in Health*, 12, 309–314.
- 460 [12] Faddy, M.J., and McClean, S.I., 1999. Analysing data on lengths of stay of  
461 hospital patients using phase-type distributions, *Applied Stochastic Models in*  
462 *Business and Industry*, 15, 311–317.
- 463 [13] Faddy, M.J., and McClean, S.I., 2005. Markov chain modelling for geriatric  
464 patient care, *Methods of Information in Medicine*, 44, 369–373.
- 465 [14] Gorunescu, F., McClean, S.I., and Millard, P.H., 2002. A queuing model for bed-  
466 occupancy management and planning of hospitals, *Journal of the Operational*  
467 *Research Society*, 53, 19–24.
- 468 [15] Graves, N., Birrell, F. and Whitby, M., 2005. The effect of pressure ulcers on  
469 length of hospital stay. *Infection Control and Hospital Epidemiology*, 26, 293–  
470 297
- 471 [16] Graves, N., Weinhold, D., Tong, E., Birrell, F., Doidge, S., Ramritu, P, Halton,  
472 K. Lairson, D. and Whitby, M., 2007. Effect of healthcare-acquired infection on  
473 length of hospital stay and cost, *Infection Control and Hospital Epidemiology*,  
474 28, 280–292.
- 475 [17] Green, P.J., 1995. Reversible jump Markov chain Monte Carlo computation and  
476 Bayesian model determination, *Biometrika*, 82, 711–732.
- 477 [18] Harbarth, S., Sax, H., and Gastmeier, P., 2003. The preventable proportion of



- 478 nosocomial infections: an overview of published reports, *Journal of Hospital*  
479 *Infection*, 54, 258–266.
- 480 [19] He, Q., and Zhang, H., 2006. Spectral polynomial algorithms for computing  
481 bi-diagonal representation for phase type distributions and matrix-exponential  
482 distributions, *Stochastic Models*, 22, no. 2, 289–317.
- 483 [20] Insua, D.R., Wiper, M. and Ruggeri, F., 1998. Bayesian analysis of  $M/Er/1$   
484 and  $M/H_K/1$  queues, *Queueing Systems*, 30, 289–308.
- 485 [21] Jasra, A., Stephens, D.A. and Holmes, C.C., 2007. Population-based reversible  
486 jump Markov chain Monte Carlo, *Biometrika*, 94, 787–807.
- 487 [22] Marshall, A.H., and McClean, S.I., 2003. Conditional phase-type distributions  
488 for modelling patient length of stay in hospital, *International Transactions in*  
489 *Operational Research*, 10, 565–576.
- 490 [23] Marshall, A., Vasilakis, C., and El-Darzi, E., 2005. Length of stay-based  
491 patient flow models: recent developments and future directions, *Health Care*  
492 *Management Science*, 8, 213–220.
- 493 [24] Munch, K. and Krogh, A., 2006. Automatic generation of gene finders for  
494 eukaryotic species, *BMC Bioinformatics*, 7, 263–275.
- 495 [25] Neuts, M.F. (1981. *Matrix Geometric Solutions in Stochastic Models*,  
496 Baltimore: Johns Hopkins University Press.
- 497 [26] Richardson, S., and Green, P.J., 1997. On Bayesian analysis of mixtures with  
498 an unknown number of components (with discussion) *Journal of the Royal*  
499 *Statistical Society Series B*, 59, 731–792.
- 500 [27] Robert, C.P., Rydén, T., and Titterton, D.M., 2000. Bayesian inference in  
501 hidden Markov models through the reversible jump Markov chain Monte Carlo  
502 method, *Journal of the Royal Statistical Society Series B*, 62, 57–75.

- 503 [28] Skaug, H.J., 2006. Markov Modulated Poisson Processes for Clustered Line  
504 Transect Data, *Environmental and Ecological Statistics*, 13, 199–211.
- 505 [29] Tan, H., Nunez-Queija, R., Gabor, A.F., Boxma, O.J., 2009. Admission control  
506 for differentiated services in future generation CDMA networks, *Performance  
507 Evaluation*, In Press, available online.
- 508 [30] Taylor, G.J., McClean, S.I., and Millard, P.H., 2000. Stochastic models of  
509 geriatric patient bed-occupancy behaviour, *Journal of the Royal Statistical  
510 Society, Series A*, 163, 39–48.
- 511 [31] Van der Weij, W. Bhulai, S., Van der Mei, R. 2009. Dynamic thread assignment  
512 in web server performance optimization, *Performance Evaluation*, 66, 301–310.
- 513 [32] Xie, H., Chausalet, T.J., and Millard, P.H. , 2005. A continuous time Markov  
514 model for the length of stay of elderly people in institutional long-term care,  
515 *Journal of the Royal Statistical Society, Series A*, 168, 51–61.