

## SC7 Bayes Methods

### Fourth problem sheet (Sections 8.4-9 of lecture notes).

---

#### Section A questions

1. (RJ-MCMC) For  $m \in \{1, 2\}$  and  $x \in (0, 1)$  let

$$\pi_{X,M}(x, m) = \pi_{X|M}(x|m)\pi_M(m)$$

with  $\pi_M(m = 1) = 1/3$ ,  $\pi_M(m = 2) = 2/3$  and

$$\pi_{X|M}(x|m = 1) = \mathbb{I}_{x=1/2}$$

$$\pi_{X|M}(x|m = 2) = 2x.$$

In the joint  $\pi_{X,M}(x, m)$ , we have  $(x, m) \in \Omega^*$  with  $\Omega^* = \{(1/2, 1)\} \cup ((0, 1) \times \{2\})$ .

The marginal distribution for  $X$  is

$$\pi_X(x) = \sum_{m=1}^2 \pi_{X,M}(x, m).$$

Let  $F_X(x) = \Pr(X \leq x)$ ,  $x \in (0, 1)$  be the CDF of  $X \sim \pi_X(\cdot)$ .

- (a) Show that  $F_X(x) = \frac{2}{3}x^2 + \frac{1}{3}\mathbb{I}_{x \geq 1/2}$  and give a simple algorithm realising iid  $X \sim F_X$ .
- (b) Give a RJ-MCMC algorithm targeting  $\pi(x, m)$  and say how you would use it to simulate  $X \sim F_X$ .

*Hint: See code. This gave Figure 1 at the end of the PS.*

2. (Dirichlet process) Let  $H$  be a continuous distribution on  $\Omega = \mathbb{R}^p$ ,  $p \geq 1$  and suppose  $G \sim \Pi(\alpha, H)$  is a DP with  $\alpha > 0$  a real parameter.

- (a) Let  $A \subseteq \Omega$ . Calculate  $\text{var}(G(A))$ . Briefly interpret  $\alpha$  and  $H$  as model “parameters”.
- (b) Suppose for  $i = 1, 2, 3, \dots$ ,  $\theta_i \sim G$  are iid, with  $G \sim \Pi(\alpha, H)$ . Recall (lectures) that marginally  $\theta_1 \sim H$  and  $G|\theta_1 \sim \Pi(\alpha + 1, (\alpha H + \delta_{\theta_1})/(\alpha + 1))$ . Show that for  $n \geq 1$ ,

$$G|\theta_{1:n} \sim DP\left(\alpha + n, \frac{\alpha H + \sum_{i=1}^n \delta_{\theta_i}}{\alpha + n}\right).$$

- (c) Let  $\theta_1^*, \dots, \theta_K^*$  denote the distinct values of  $\theta$  with associated partition  $S = (S_1, \dots, S_K)$ ,  $S_k = \{i : \theta_i = \theta_k^*, i \in [n]\}$  for  $k = 1, \dots, K$ . Show that

$$E(K) = \sum_{i=1}^n \frac{\alpha}{\alpha + i - 1}$$

## Section B questions

3. (Reversible jump MCMC) The skew-normal distribution<sup>1</sup> with density  $Q(y; \mu, \sigma^2, \xi)$  is obtained from the normal by skewing it with a weight  $\xi > 0$ . The skewing is negative for  $0 < \xi < 1$ , positive for  $\xi > 1$  and absent for  $\xi = 1$ , ie  $N(y; \mu, \sigma^2) = Q(y; \mu, \sigma^2, 1)$ .

The Shoshoni data  $y = (y_1, \dots, y_{20})$  give the values of 20 scalar width-to-length ratios of beaded rectangles used by the Shoshoni Indians. They are available here,

[www.statsci.org/data/general/shoshoni.html](http://www.statsci.org/data/general/shoshoni.html).

You can see them and an example of the skew-normal in `ProblemSheet4.R`. Consider using Bayesian inference and RJ MCMC to carry out model selection and model averaging over skewed and normal models for the Shoshoni data.

- Suppose the prior probability for normal (model  $m = 1$ ) or skew-normal (model  $m = 2$ ) is  $1/2$ . Write down the joint posterior distribution  $\pi(\theta, m|y)$  for the model index  $m = 1, 2$  and parameters  $\theta = (\mu, \sigma, \xi)$  in as much detail as you can, though without eliciting priors for the parameters.
  - Give a reversible jump MCMC algorithm targeting  $\pi(\theta, m|y)$ . You can omit the fixed dimension updates.
  - Explain how to estimate the Bayes Factor comparing skew-normal and normal models from MCMC output  $\theta^{(t)} = (\mu^{(t)}, \sigma^{(t)}, \xi^{(t)})$  and  $m^{(t)}, t = 1, 2, \dots, T$ . How you would simulate data  $y'$  from the model averaged posterior predictive distribution  $p(y'|y)$ ?
  - (Section C) The code in the R-file `ProblemSheet4.R` implements RJ-MCMC for these data. Use the code to estimate the Bayes factor mentioned above.
4. Let  $\Xi_{[n]}$  be the set of partitions of  $[n] = \{1, \dots, n\}$ . The CRP realises  $S \in \Xi_{[n]}$  with probability

$$P_{\alpha, [n]}(S) = \frac{\Gamma(\alpha)\alpha^K}{\Gamma(\alpha + n)} \prod_{k=1}^K \Gamma(|S_k|).$$

Let  $\mathcal{P}_{[n]}$  be the permutations of  $\{1, \dots, n\}$ .

- For  $\sigma \in \mathcal{P}_n$  let  $P_{\alpha, \sigma}(S)$  be the distribution over partitions we get if the customers arrive in the order  $\sigma = (\sigma_1, \dots, \sigma_n)$  and let  $S(\sigma)$  be the partition obtained by permuting the customer labels in  $S$  according to  $\sigma$ . For example if  $S = (\{1, 2\}, \{3\})$  and  $\sigma = (3, 2, 1)$  then  $S(\sigma) = (\{1\}, \{2, 3\})$  because the new partition is  $\{\{\sigma_1, \sigma_2\}, \{\sigma_3\}\} = \{\{3, 2\}, \{1\}\}$  and recall the convention  $\min(S_k) < \min(S_{k'}) \Leftrightarrow k < k'$ .

Show that  $P_{\alpha, [n]}(S) = P_{\alpha, [n]}(S(\sigma)) = P_{\alpha, \sigma}(S)$  for all  $S \in \Xi_{[n]}$ , so CRP outcomes don't depend on customer arrival order.

---

<sup>1</sup>Fernandez & Steel “*Bayesian Modeling of Skewness and Fat Tails*”, JASA, 1998

- (b) Let  $S \sim P_{\alpha, [n]}$  and  $S^{-i} = (S_1^{-i}, \dots, S_{K^{-i}}^{-i})$  be the partition we get if we realise  $S$  and then remove some  $i \in \{1, \dots, n\}$ . Here  $K^{-i} = K - 1$  if we create an empty cluster when we remove  $i$  and otherwise  $K^{-i} = K$ . For example if  $S = (\{1, 2\}, \{3\})$  then  $K = 2$  and  $S^{-3} = (\{1, 2\})$  so  $K^{-3} = 1$ .

Let  $P_{\alpha, [n] \setminus \{i\}}(S')$  be the probability to realise  $S' \in \Xi_{[n] \setminus \{i\}}$  if  $i$  is removed from the list of customers before  $S'$  is simulated from the CRP. Show that  $S^{-i} \sim P_{\alpha, [n] \setminus \{i\}}(S^{-i})$ .

5. Consider the following prior for the cluster labels  $z = (z_1, \dots, z_n)$  of data  $y = (y_1, \dots, y_n)$  in a mixture model with a fixed number  $M$  of components. Let  $w = (w_1, \dots, w_M)$  be a vector of probabilities  $\sum_m w_m = 1$  giving the mixture-component weights.

$$\begin{aligned} w &\sim \text{Dirichlet}(\alpha_1, \dots, \alpha_M), & \text{with } \alpha > 0 \text{ and } \alpha_m = \alpha/M, m = 1, \dots, M \\ z_i &\sim \text{Cat}(w), & \text{iid for } i = 1, \dots, n. \end{aligned}$$

In this model  $z_i \in \{1, \dots, M\}$  is the label of the cluster to which  $y_i$  belongs, and the notation  $z_i \sim \text{Cat}(w)$ ,  $i = 1, \dots, n$  means that for  $m \in \{1, \dots, M\}$  we have  $z_i = m$  with probability  $w_m$ . Suppose the list  $z_1, \dots, z_n$  of cluster labels contains  $K \leq M$  unique distinct values  $m_1, \dots, m_K$ . For  $k = 1, \dots, K$  let  $S_k = \{i : z_i = m_k, i = 1, \dots, n\}$  give the label-grouping determined by  $z$  and let  $S = (S_1, \dots, S_K)$ .

The partition is determined by  $z$ , so that  $S = S(z)$  with  $S \in \Xi_{[n]}$ . There are many  $z$ 's giving the same  $S$ . For example, if  $n = 4$  and  $M = 5$  then  $z = (1, 1, 3, 3)$ ,  $z = (3, 3, 1, 1)$  and  $z = (4, 4, 2, 2)$  determine the same clustering  $S = (\{1, 2\}, \{3, 4\})$ .

- (a) (Section C, but result needed below) Let  $n_k = |S_k|$  for  $k = 1, \dots, K$ . Let  $P_{\alpha, [n]}^M(S)$  be the probability to realise  $S$ . Calculate

$$P_{\alpha, [n]}^M(S) = \sum_{z: S(z)=S} P_{\alpha, [n]}^M(z),$$

where  $P_{\alpha, [n]}^M(z)$  is the probability the process realises  $z = (z_1, \dots, z_n)$ , and show

$$P_{\alpha, [n]}^M(S) = \frac{\Gamma(\alpha)}{\Gamma(\alpha/M)^K} \frac{M!}{(M-K)!} \frac{\prod_{k=1}^K \Gamma(\alpha/M + n_k)}{\Gamma(\alpha + n)}.$$

- (b) Show that, for each  $S \in \Xi_{[n]}$ ,  $\lim_{M \rightarrow \infty} P_{\alpha, [n]}^M(S) = P_{\alpha, [n]}(S)$ , with  $P_{\alpha, [n]}$  from Question (4).

Note:  $x\Gamma(x) = \Gamma(x+1)$  and  $x\Gamma(x) \rightarrow 1$  as  $x \searrow 0$ .

6. A realisation,  $G_M \sim \Pi_M(\alpha, H)$ , of the *multinomial DP* is simulated as follows:

$$\begin{aligned} w &\sim \text{Dirichlet}(\alpha_1, \dots, \alpha_M), & \text{with } \alpha > 0 \text{ and } \alpha_m = \alpha/M, m = 1, \dots, M, \\ \tilde{\theta}_m &\sim H, & \text{iid for } m = 1, \dots, M, \end{aligned}$$

and  $G_M = \sum_{m=1}^M w_m \delta_{\tilde{\theta}_m}$ . Here, for  $m = 1, \dots, M$ ,  $\tilde{\theta}_m \in \mathbb{R}^p$  is a parameter vector of dimension  $p$  and  $H$  is a base distribution with probability density  $h$  on  $\mathbb{R}^p$ .

- (a) For  $i = 1, \dots, n$ , let  $\theta_i = \tilde{\theta}_{z_i}$  with

$$z_i \sim \text{Cat}(w), \quad \text{iid for } i = 1, \dots, n.$$

Show that  $\Pr\{\theta_i \in A | w, \tilde{\theta}\} = G_M(A)$  for  $A \subseteq \mathbb{R}^p$  and  $i = 1, \dots, n$ .

- (b) Let  $\theta_1^*, \dots, \theta_K^*$  denote the distinct values of  $\theta$  with associated partition  $S = (S_1, \dots, S_K)$ ,  $S_k = \{i : \theta_i = \theta_k^*, i \in [n]\}$  for  $k = 1, \dots, K$ . Give the joint distribution  $\pi_M(\theta^*, S)$ .
- (c) Consider the following process.

Step 1 Simulate  $\psi_1 \sim H$

Step 2 Independently for  $i = 1, \dots, n-1$ , and sequentially, simulate

$$\psi_{i+1} \sim \frac{\alpha(1 - K_i/M)H + \sum_{k=1}^{K_i} (n_{i,k} + \alpha/M) \delta_{\psi_k^*}}{\alpha + i}.$$

where  $K_i$  is the number of distinct  $\psi$ -values  $\psi_1^*, \dots, \psi_{K_i}^*$  at the time of the  $i+1$ 'st arrival and  $n_{i,k}$  is the number of times  $\psi_k^*$  appears in the list  $(\psi_1, \dots, \psi_i)$ . Show that  $\psi = (\psi_1, \dots, \psi_n)$  above has the same distribution as  $\theta = (\theta_1, \dots, \theta_n)$  in Question 6a. *Hint: set it up as a variant of a CRP realising  $\psi^*, C$  with  $\psi^*$  the unique values in  $\psi$  and  $C$  the corresponding partition of  $\psi$  and repeat the calculation we did in lectures for  $P_{\alpha, [n]}(S)$  to get  $P(C) = P_{\alpha, [n]}^M(C)$ .*

- (d) (Section C) Let  $\phi_i \sim G$  iid for  $i = 1, \dots, n$  with  $G \sim \Pi(\alpha, H)$  and  $\phi = (\phi_1, \dots, \phi_n)$ . Let  $\phi = \theta(\phi^*, S)$  with  $\theta$  the usual invertible mapping between the two representations. Let  $\psi_i \sim G_M$  iid for  $i = 1, \dots, n$  with  $G_M \sim \Pi_M(\alpha, H)$  and  $\psi = (\psi_1, \dots, \psi_n)$ . Let  $\psi = \theta(\psi^*, C)$  be corresponding unique values and partition representation (ie as in the hint for Question 6c). Show that  $\psi \rightarrow \phi$  in distribution as  $M \rightarrow \infty$  at fixed  $n$ . *Hint show that  $\Pr\{(\psi^*, C) \in A^*\} \rightarrow \Pr\{(\phi^*, S) \in A^*\}$  for all sets  $A^*$ .*

## Section C questions

7. The observation model for data  $y$  is  $y_i \sim f(\cdot | \theta_i)$ , iid for  $i = 1, \dots, n$  with parameter vector  $\theta = (\theta_1, \dots, \theta_n)$  determined from the multinomial Dirichlet process model via a realisation of  $\theta^*$  and  $S$  as in Question 6.
- (a) Write down the posterior  $\pi_M(S, \theta^* | y)$  for  $S, \theta^* | y$  in terms of the model elements.
- (b) Why might we prefer a prior derived from a multinomial Dirichlet process over a prior derived from a Dirichlet process?

- (c) Show that the pairs  $(\theta_i, y_i)_{i=1}^n$  are exchangeable (as pairs, *ie* preserving the association between  $\theta_i$  and  $y_i$ ). Give the  $S, \theta^*$ -update of a Gibbs sampler targeting  $\pi_M(S, \theta^* | y)$ .
8. Mining disasters were common in the period 1850 – 1950. Let  $L = 1850$  and  $U = 1950$  and for  $i = 1, 2, \dots, n$ , let  $y_i \in (L, U)$  be the date of the  $i$ 'th event. Let  $y = (y_1, \dots, y_n)$ .

Model the event times  $y$  as the arrival times of a Poisson process of piecewise constant rate  $\lambda(t)$  per year. For  $m \geq 1$  let  $\theta_0 = L$  and  $\theta_m = U$  and for  $i = 1, \dots, m - 1$  let  $\theta_i \in (L, U)$  be the sorted change-point times at which  $\lambda(t)$  jumps up or down. The number of change-points is  $m - 1$  so if  $m = 1$  then there are no change points and the rate  $\lambda(t)$  is constant for  $t \in (L, U)$ . For  $i = 1, \dots, m$  let  $\lambda_i \geq 0$  give the disaster rate over the interval  $(\theta_{i-1}, \theta_i]$ . The rate function  $\lambda(t) = \lambda(t; \theta, \lambda)$  for  $y$  is

$$\lambda(t) = \sum_{i=1}^m \lambda_i \mathbb{I}_{\theta_{i-1} < t \leq \theta_i} \quad L < t < U.$$

The data and a realisation of  $\lambda(t)$  with  $m = 4$  are shown in Figure 2 below.

Let  $\theta = (\theta_1, \dots, \theta_{m-1})$  and  $\lambda = (\lambda_1, \dots, \lambda_m)$ . Model the change-point times  $\theta$  as arrivals in a Poisson process of unknown rate  $\rho$  per year. The number of intervals  $m$  is unknown. Prior densities  $\pi_R(\rho)$ ,  $\rho \in [0, \infty)$  and  $\pi_\Lambda(\lambda | m) = \prod_{i=1}^m \pi_\Lambda(\lambda_i)$ ,  $\lambda \in [0, \infty)^m$  are given.

- (a) i. Write down the prior  $\pi(\theta, \lambda, m, \rho)$  in as much detail as you can. Specify its parameter space,  $(\theta, \lambda, m, \rho) \in \Omega$  say.
- ii. Write down the posterior  $\pi(\lambda, \theta, m, \rho | y)$  in terms of the model elements.
- (b) In a reversible jump MCMC algorithm targeting  $\pi(\lambda, \theta, m, \rho | y)$ , birth and death updates are chosen with probabilities  $p_{m, m+1}$  and  $p_{m, m-1}$  respectively. A birth proposal  $(\lambda, \theta, m, \rho) \rightarrow (\lambda', \theta', m', \rho)$  with  $m' = m + 1$  is generated as follows: choose an interval  $i \sim U\{1, \dots, m\}$  uniformly; simulate a split point  $\theta^* \sim U(\theta_{i-1}, \theta_i)$ ; simulate two new values  $\lambda_{i,1}, \lambda_{i,2} \sim \text{Exp}(1)$  independently. In the candidate state

$$\begin{aligned} \lambda' &= (\lambda_1, \dots, \lambda_{i-1}, \lambda_{i,1}, \lambda_{i,2}, \lambda_{i+1}, \dots, \lambda_m) \\ \theta' &= (\theta_1, \dots, \theta_{i-1}, \theta^*, \theta_i, \dots, \theta_{m-1}). \end{aligned}$$

Give a matching death proposal  $(\lambda', \theta', m', \rho) \rightarrow (\lambda, \theta, m, \rho)$  and the acceptance probability for the birth proposal. No simplification of expressions is required.

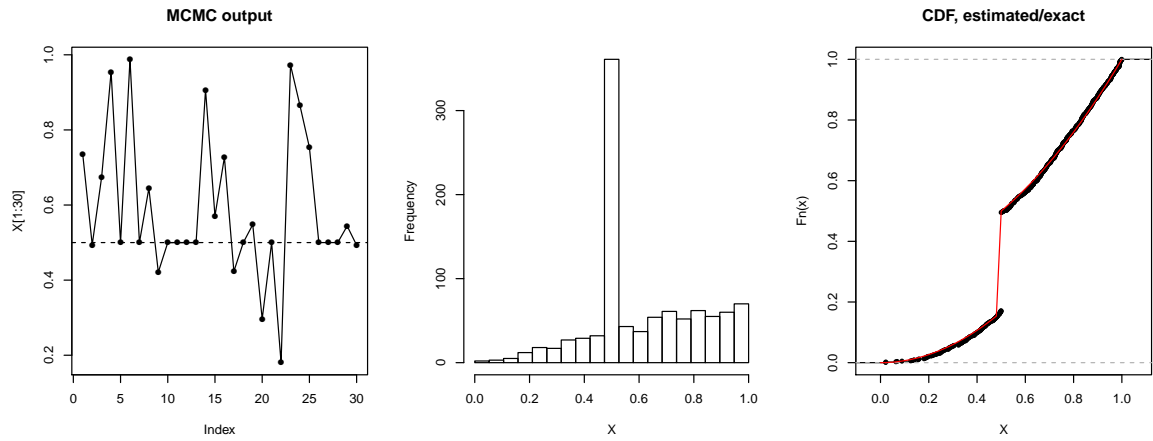


Figure 1: RJ-MCMC targeting  $\pi(x, m)$ : (Left) plot of  $x$ -values realised by the chain (sub-sampled every 10 steps); (Centre) histogram estimate of marginal pdf of  $x$  ( $f_X(x) = \frac{4}{3}x + \frac{1}{3}\delta_{1/2}(x)$ ) showing the atom of probability at  $x = 1/2$ ; (Right) Marginal CDF of  $x$  ( $F_X(x) = \frac{2}{3}x^2 + \frac{1}{3}\mathbb{I}_{x \geq 1/2}$ ).

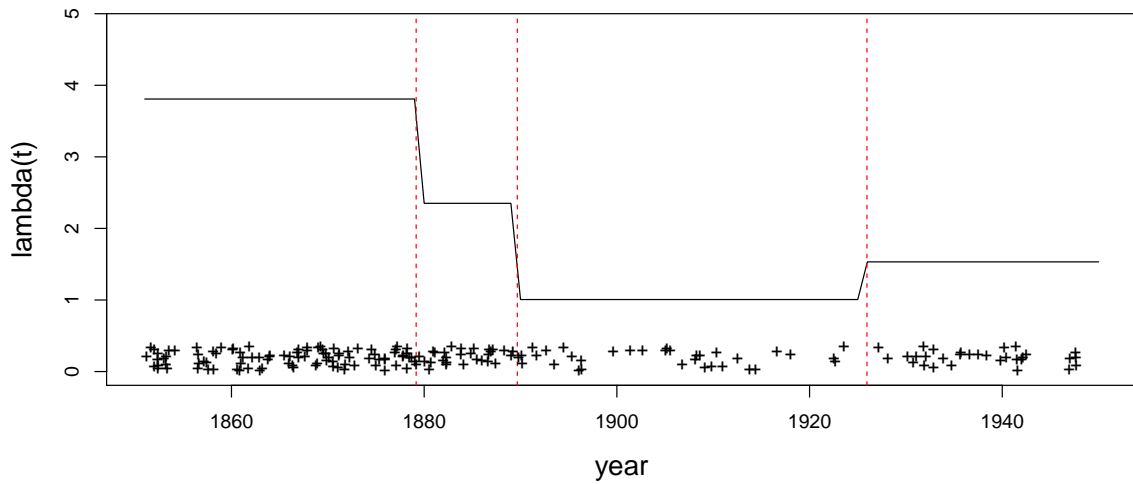


Figure 2: Coal mining disasters: event dates  $y$  (+ signs), change point times ( $\theta$  vertical lines) and  $\lambda(t)$  itself (piecewise constant function of year,  $t$ ).