

Problem 1

- a) Integral is approximated by:

$$\int_{\mathbb{R}^n} h(\mathbf{x})g(\mathbf{x})d\mathbf{x} \approx \bar{h}_a = \frac{1}{B} \sum_{i=1}^B h(\mathbf{x}_i)$$

If the samples are independent, the error in approximation measured in terms of the variance is: $\text{Var}\{\bar{h}_a\} = \frac{\text{Var}_g\{h(\mathbf{X})\}}{B}$. Where the subscript indicate that the variance is over the distribution $g(\mathbf{x})$. Thus if $\text{Var}\{h(\mathbf{X})\} < \infty$, and the samples are independent then the limit is zero when $B \rightarrow \infty$.

- b) Integral can be approximated by:

$$\int_{\mathbb{R}^n} h(\mathbf{x})f(\mathbf{x})d\mathbf{x} \approx \bar{h}_b = \frac{1}{B} \sum_{i=1}^B \frac{f(\mathbf{x}_i)}{g(\mathbf{x}_i)} h(\mathbf{x}_i)$$

If the samples are independent, then we can utilize the result from above with the function $\frac{f(\mathbf{x})}{g(\mathbf{x})}h(\mathbf{x})$, thus the error in approximation measured in terms of the variance is:

$$\text{Var}\{\bar{h}_b\} = \frac{\text{Var}_g\left\{\frac{f(\mathbf{X})}{g(\mathbf{X})}h(\mathbf{X})\right\}}{B},$$

where the subscript indicate that the variance is over the distribution $g(\mathbf{x})$. For the convergence: If $\text{Var}_g\left\{\frac{f(\mathbf{X})}{g(\mathbf{X})}h(\mathbf{X})\right\} < \infty$, and the samples are independent then the limit is zero when $B \rightarrow \infty$, which gives convergence. Importantly however for the solution to converge to the right number, we also need that the support of $g(\mathbf{x})$ contains the support of $h(\mathbf{x})f(\mathbf{x})$. Alternatively we could require $E_f\left\{\frac{f(\mathbf{X})}{g(\mathbf{X})}h(\mathbf{X})^2\right\} < \infty$.

Utilizing $q(\mathbf{x})$ rather than $f(\mathbf{x})$ we get:

$$\bar{h}_{b_2} = \frac{\sum_{i=1}^B \frac{q(\mathbf{x}_i)}{g(\mathbf{x}_i)} h(\mathbf{x}_i)}{\sum_{i=1}^B \frac{q(\mathbf{x}_i)}{g(\mathbf{x}_i)}}$$

These methods are denoted by the common name importance sampling. Using un-normalized weights and normalized weights respectively.

- c) In step 2 A1, we sample form the standard normal distribution using the inversion formula. $P(X < x) = P(\Phi^{-1}(u) < x) = P(u < \Phi(x)) = \Phi(x)$. The last equality is a consequence of the normal distribution, the second to last is due to the monotonicity of the transform. In step 2 A2, we sample form an exponential distribution shifted such that the minimum value is x_0 rather than 0.

In formula (3) we use Monte Carlo simulation, such as (1) to compute the integral.

In formula (4) we use importance sampling. Note that we do not need to add the indicator since all samples are larger than x_0 .

- d) Since the jitter is much stronger for algorithm 1 than for algorithm 2. This indicates that the sample variance is larger for this computation. Theoretically, it is much better to have all samples centered at the region where the function is non-zero, since we investigate the part which contributes to the integral much better. Also note that the decay of exponential distribution is slower than the normal distribution. The latter is important for the variance of the ratio. $\text{Var}_g \left\{ \frac{f(\mathbf{X})}{g(\mathbf{X})} h(\mathbf{X}) \right\}$. The jitter is caused by Monte Carlo variability, i.e. the fact that we do not have an infinite number of samples. In order to reduce the variability, we can increase the number of samples. It is also possible to trade some variance in the estimate for bias, by using common random numbers $(u_i, i = 1 \dots, B)$ for all values of x_0 .

Problem 2

- a) The likelihood is:

$$L(p, \lambda) = \prod_{i=1}^n (p \cdot \lambda \cdot \exp -\lambda x + (1 - p) \cdot \lambda^2 x \cdot \exp -\lambda x)$$

In the Newton algorithm, we update the estimate with

$$x_{i+1} = x_i - (H_f)^{-1} \cdot \nabla f$$

Where $f(x)$ is the function, we want to optimize, ∇f is the gradient of this function, and H_f is the Hessian. In the quasi-Newton algorithm we use an approximation to replace the inverse Hessian. This expression is faster to compute. Note that even if the quasi-newton algorithm converges the approximation for the Hessian need not converge to the Hessian.

In the table we clearly see that there are two modes. One of the modes proposes a probability which is larger than zero. To avoid this issue, we could either use constrained optimization or transform the probability into an unconstrained parameter, e.g. $\theta = \log \left(\frac{1}{p} - 1 \right)$.

- b) One expression for the complete log likelihood is:

$$l(p, \lambda | \mathbf{x}, \mathbf{C}) = \sum_{i=1}^n I(C_i = 1) [\log p + \log \lambda - \lambda x_i] \\ + \sum_{i=1}^n I(C_i = 2) [\log(1 - p) + 2 \log \lambda + \log(x_i) - \lambda x_i]$$

The $Q(\cdot)$ function is the expected value of the complete log likelihood given the observed data and the current estimate of the parameters. In the expression for the

complete log likelihood above we see that the only randomness is in the indicator functions, which are in the expression in a sum. The expectation can be applied to each term in the sum, when we use that $E(I(C_i = 1)|x, p, \lambda) = P(C_i = 1|x, p, \lambda)$, we get the desired result.

$$\begin{aligned} P(C_i = 1|x_i, p, \lambda) &= \frac{P(C_i = 1, x_i|p, \lambda)}{P(x_i|p, \lambda)} \\ &= \frac{P(C_i = 1|p, \lambda)p(x_i|C_i = 1, p, \lambda)}{P(C_i = 1|p, \lambda)p(x_i|C_i = 1, p, \lambda) + P(C_i = 2|p, \lambda)p(x_i|C_i = 2, p, \lambda)} \\ &= \frac{p \cdot \text{Exp}(x_i; \lambda)}{p \cdot \text{Exp}(x_i; \lambda) + (1 - p) \cdot \text{Erlang}(x_i; 2, \lambda)} \end{aligned}$$

Where we have used definition of conditional probability and inserted the expressions.

c) Differentiate $Q(p, \lambda|p^{(t)}, \lambda^{(t)})$, wrt p :

$$\begin{aligned} \frac{\partial Q}{\partial p} &= \sum_{i=1}^n P(C_i = 1|x_i, p^{(t)}, \lambda^{(t)}) \left[\frac{1}{p} \right] - \sum_{i=1}^n \left(1 - P(C_i = 1|x_i, p^{(t)}, \lambda^{(t)}) \right) \left[\frac{1}{1-p} \right] = 0 \\ &\Rightarrow p = \frac{1}{n} \sum_{i=1}^n P(C_i = 1|x_i, p^{(t)}, \lambda^{(t)}) \end{aligned}$$

Differentiate $Q(p, \lambda|p^{(t)}, \lambda^{(t)})$, wrt λ :

$$\begin{aligned} \frac{\partial Q}{\partial \lambda} &= \sum_{i=1}^n P(C_i = 1|x_i, p^{(t)}, \lambda^{(t)}) \left[\frac{1}{\lambda} - x_i \right] + \sum_{i=1}^n \left(P(C_i = 2|x_i, p^{(t)}, \lambda^{(t)}) \right) \left[\frac{2}{\lambda} - x_i \right] = 0 \\ &\Rightarrow \lambda = \frac{\sum_{i=1}^n P(C_i = 1|x_i, p^{(t)}, \lambda^{(t)}) + 2 \cdot \sum_{i=1}^n P(C_i = 2|x_i, p^{(t)}, \lambda^{(t)})}{\sum_{i=1}^n x_i} \end{aligned}$$

In the EM algorithm each iteration will improve the likelihood from 2a.

d) In a non-parametric bootstrap, we resample the dataset with replacement to create a dataset of equal size as the initial. Then we estimate the parameters using the prescribed algorithm. One repeat of this procedure creates one point in the scatterplot in figure 2. We repeat this procedure B times. And uses the set of estimates to represent the sampling uncertainty of the estimator. In the figure we see that the parameter estimates are highly correlated (negative correlation). The lower the value of p the larger the estimate of λ . This is natural since reducing p means that we overlook more samples, which would increase the number

of events in the same time interval. The samples seems to have two modes. One close to $p = 1$ and one centred around $p = 0.75$. There is no obvious reason for this.

Problem 3

a) The conditional probability is computed to proportionality by:

$$\begin{aligned}
p(\boldsymbol{\lambda}, \beta | \alpha, \gamma, \delta, \mathbf{x}, \mathbf{t}) &\propto p(\boldsymbol{\lambda}, \beta, \mathbf{x} | \alpha, \gamma, \delta, \mathbf{t}) \\
&= p(\beta | \gamma, \delta, \alpha, \mathbf{t}) p(\boldsymbol{\lambda} | \beta, \gamma, \delta, \alpha, \mathbf{t}) p(\mathbf{x} | \boldsymbol{\lambda}, \beta, \gamma, \delta, \alpha, \mathbf{t}) \\
&= p(\beta | \gamma, \delta) p(\boldsymbol{\lambda} | \beta, \alpha) p(\mathbf{x} | \boldsymbol{\lambda}, \mathbf{t}) \\
&= p(\beta | \gamma, \delta) \prod_{i=1}^{10} p(\lambda_i | \beta, \alpha) \prod_{i=1}^{10} p(x_i | \lambda_i, t_i) \\
&= \frac{\delta^\gamma \beta^{\gamma-1}}{\Gamma(\gamma)} \exp(-\delta\beta) \prod_{i=1}^{10} \frac{\beta^\alpha \lambda_i^{\alpha-1}}{\Gamma(\alpha)} \exp(-\lambda_i\beta) \prod_{i=1}^{10} \frac{(\lambda_i t_i)^{x_i}}{x_i!} \exp(-\lambda_i t_i)
\end{aligned}$$

Where transition 1 is due to posterior is proportional to the joint distribution (conditioned to given parameters $\alpha, \gamma, \delta, \mathbf{t}$). Second transition is due to common factorization. In transition 3, we impose the dependencies of parameters given in the text. In transition 4, we utilize the prior independence of λ_i , and x_i . In the final step we impose the distributions.

Computation of distributions for the Gibbs-sampler:

$$\begin{aligned}
p(\lambda_i | \alpha, \beta, \gamma, \delta, \mathbf{t}, \mathbf{x}, \boldsymbol{\lambda}_{-i}) &\propto \lambda_i^{\alpha-1} \exp(-\lambda_i\beta) \lambda_i^{x_i} \exp(-\lambda_i t_i) \\
&\propto \lambda_i^{\alpha+x_i-1} \exp(-\lambda_i(\beta + t_i)) \\
&\propto \text{Gamma}(\lambda_i; \alpha + x_i, \beta + t_i)
\end{aligned}$$

$$\begin{aligned}
p(\beta | \alpha, \gamma, \delta, \boldsymbol{\lambda}, \mathbf{x}, \mathbf{t}) &\propto \beta^{\gamma-1} \exp(-\delta\beta) \prod_{i=1}^{10} \beta^\alpha \exp(-\lambda_i\beta) \\
&\propto \beta^{\gamma+10\alpha-1} \exp\left(-\left(\delta + \sum_{i=1}^{10} \lambda_i\right)\beta\right) \\
&\propto \text{Gamma}\left(\beta; \gamma + 10\alpha, \delta + \sum_{i=1}^{10} \lambda_i\right)
\end{aligned}$$

b) In R_H , the $f(x)$, and $f(x^*)$ are the target densities, evaluated in the current sample, x , and the proposed sample, x^* . The function $g(x^*|x)$, is the density for proposing a change to the proposed sample, x^* from the current sample. The function $g(x^*|x)$ is the density for the reverse proposal.

M-H pseudo code:

Initialize $x = x_0$

For i in 1: B

Sample $x^* \sim g(x^*|x_{i-1})$, $u_i \sim \text{Uniform}[0,1]$

Compute R_{MH}

if $u_i < R_{MH}$: set $x_i = x^*$ otherwise set $x_i = x_{i-1}$

The transition kernel in Metropolis Hastings ($x \neq x^*$):

$$P(x, x^*) = g(x^*|x) \max\left\{1, \frac{f(x^*)g(x|x^*)}{f(x)g(x^*|x)}\right\}$$

c) Starting with the integral:

$$\int f(x)P(x, x^*)dx = \int f(x^*)P(x^*, x)dx = f(x^*) \underbrace{\int P(x^*, x)dx}_1 = f(x^*)$$

First transition is just utilizing the detailed balance, next we set what is constant in the integration outside the integral and recognize that the integral of the transition kernel is 1.

For convergence the Markov chain need to be:

- Irreducible – possible to move to any value in a finite number of steps
- Aperiodic – not go into predictable cycles
- Recurrent – always return to sets that has a positive probability