

# In Depth Analysis of Norovirus Data

Andrew Martin

June 5, 2019

## 1 Introduction

This report gives an in depth explanation and interpretation of a statistical analysis of norovirus. The analysis is based on a dataset containing information on independent outbreaks of norovirus from a series of different hospital wards. The wards use one of three control actions to attempt to manage the rate of infection. A Markov-Chain Monte Carlo method is used to predict the effective reproduction number for norovirus under normal circumstances, and two independent control strategies.

## 2 Background information

In general, epidemics are modelled with the assumption that a single person can transition between a series of states over the course of an infection. An individual will typically be (S)usceptible to the disease initially, until they are (E)xposed to the disease, not yet infectious, but will soon become (I)nfectious, where they expose others, until this individual will (R)ecover from the disease with some level of immunity.

It is possible for there to be other variety of states, and for there to be loops, but this is a simple linear SEIR model for a disease.

$R_0$  denotes the reproduction number, the expected number of secondary infections caused by a single infected individual in an otherwise entirely susceptible population.

$$R_0 = \frac{\beta}{\gamma}$$

Where  $\beta$  is the infection rate, and  $\gamma$  is the recovery rate for the norovirus.  $R_0 > 1$  corresponds to a possibility for the disease to become a large-scale epidemic, while for a value  $R_0 < 1$  the virus should die out.

The Metropolis-Hastings (MH) algorithm is a Markov Chain Monte Carlo (MCMC) method, which uses Markov chains to obtain the distribution of a random parameter vector  $\Theta$ . Assuming the parameter comes from a known (or assumed) prior distribution, the stationary distribution of the CTMC produced will be the distribution of the parameter. Where the prior distribution:  $P(\Theta)$  is the assumed distribution for the parameter(s)  $\Theta$ . Posterior distribution:  $P(D|\Theta)$  is the distribution of data (in this case known) given the parameter  $\Theta$ .

### 3 Assumptions

To obtain a model, for norovirus some basic information on norovirus is considered.

Norovirus is a highly infectious disease. With a time-line of: 12-48 hours after exposure, symptoms appear, after an incubation period. After this incubation period, infection begins and typically lasts 24-48 hours. There is another asymptomatic infectious state for up to 24 hours after this [1].

This suggests some extension of the linear SEIR model, with no loops, and possibly with a secondary asymptomatic infectious state, I.e. the model where an individual can be:

$$\text{Full model: } S \rightarrow E \rightarrow I_1 \rightarrow I_2 \rightarrow R$$

Other analysis suggests that the second infectious state can be disregarded with equally valid results, reducing it to the *SEIR* model [2].

$$\text{Reduced model: } S \rightarrow E \rightarrow I \rightarrow R$$

For the purposes of this report, since the given data is limited to the number of infected and total number of people, there is no possible way to identify a number of exposed, and hence no real value in the inclusion of the exposed state. The information we are seeking (namely the rate of reproduction and effectiveness of control strategies) is not affected by the number of infectious states. Hence the model is fully reduced to the *SIR* model.

$$\boxed{\text{Final model: } S \rightarrow I \rightarrow R}$$

The rate of reproduction  $R_0$  has already been researched by the Chief Medical Officer, with the claim that

$R_0$  for norovirus in typical hospital settings is between 2 and 3 with (approximately) 66% probability

No prior information has been given regarding either of the intervention strategies. With no background on their effectiveness, denoted  $\alpha_1, \alpha_2$ . The expectation is that the intervention strategies would be beneficial in reducing the effective  $R_0$  (I.e.  $\alpha_i < 1$ ) but it is possible that they could have adverse effects ( $\alpha_i > 1$ ). Since no more can be inferred, a reasonable prior

$$\alpha_i \sim U(0, 2)$$

is chosen, noting that a value of 1 would correspond to no change, with a smaller value corresponding to a reduction in infectiousness, and a larger value would cause an increase. Given this prior, if the MH algorithm fails to converge to a value, and appears to increase, this may suggest that  $\alpha_i > 2$  and hence  $\alpha_i$  would be an extremely detrimental strategy and would need to be discarded immediately.

## 4 Method

### 4.1 Data

The dataset, `NorovirusDataA3.txt`, containing information on independent outbreaks across 125 hospital wards is used for all inference in this report. The text file contains comma separated values (CSV) data, where each row denotes an independent ward, and the columns

correspond to: the number of occupied beds, number of people succumbed to the virus, and the treatment action taken universally in the ward  $(T_0, T_1, T_2)$ , respectively.

Denote the treatments as  $\alpha_1, \alpha_2$  respectively, such that  $R_0, \alpha_1 R_0, \alpha_2 R_0$  are the effective reproduction numbers for treatments  $T_0, T_1$ , and  $T_2$  respectively. The assumption of independence between  $\alpha_1, \alpha_2$ , and  $R_0$  is used. This assumption will significantly improve convergence. The assumption does not necessarily have to hold for the posteriors however.

## 4.2 Distributions

Chosen prior distributions (priors):

$$R_0 \sim N'(\mu, \sigma)$$

$$\alpha_1, \alpha_2 \sim U(0, 2)$$

Where  $N'$  is a Normal distribution, concatenated such that all values outside  $[0, 5]$  are dropped. Parameters are  $\mu = 2.5$  and  $\sigma$  is chosen such that  $\approx 66\%$  of the density is contained in  $(2, 3)$ . Assume  $\alpha_i$  are each uniformly distributed between 0, 2. This is chosen since  $\alpha_i R_0 \geq 0$ , and it is expected that  $\alpha_i$  will not be significantly detrimental (so as to effectively double  $R_0$ ). This has been centered around 1, which would correspond to the treatment supplying no change whatsoever.

Since the assumption of uniform priors is used for  $\alpha_1$  and  $\alpha_2$ , they can be omitted from the likelihood calculations barring the check that they sit within the valid range  $(0, 2)$ .

## 4.3 Algorithm

The Metropolis-Hastings (MH) algorithm takes an initial guess  $x_0$ , and then iterates through (variable  $i$ ): Propose a candidate state from the proposal density,  $r$ :

$$x' \sim r(x'|x_i)$$

And accept this state,  $x'$  with probability (based on the data):

$$a(x_i, x') = \frac{p(x')r(x_i|x')}{p(x_i)r(x'|x_i)}$$

If accepted, set  $x_{i+1} = x'$  or if rejected, set  $x_{i+1} = x_i$  and restart the loop.

The modified MH algorithm uses log likelihood for the  $a$  calculation instead:

$$\log(a(x_i, x')) = \log(p(x')) + \log(r(x_i|x')) - \log(p(x_i)) - \log(r(x'|x_i))$$

Using the modified MH algorithm, attempt to obtain parameters  $R_0, \alpha_1, \alpha_2$  where  $\alpha_i$  corresponds to the reduction of  $R_0$  caused by treatment  $i$ . The code used for this is listed in the appendix. This is done in `MATLAB`, by running the script `ROPredict.m`

## 5 Analysis and Results

An effective treatment will have  $\alpha < 1$ , with a theoretically perfect treatment having  $\alpha = 0$ . The assumption that the wards are *independent*, and there is a control set, where no treatment actions are taken are vital to the analysis. For the methods, with initial guesses for  $R_0, \alpha_1, \alpha_2$ :  $(4, 1.5, 1.5)$  (overshoot),  $(0.5, 0.5, 0.5)$  (undershoot) and  $(2, 0.7, 0.7)$  (approximate guess), the MH algorithm is run.

Figure 1 and figure 2 plot the resulting guesses from the MH algorithm. Figure 2 only shows the first 1000 iterations to demonstrate the burn in time, and convergence of solutions. The

plots show the solutions obtained for three different sets of parameter guesses: an overshoot, undershoot and approximate guess - all three of which clearly converge to the same values.

Figure 3 displays the dependency between the different parameters. Clearly there is some level of dependency. Particularly between  $R_0$  and the different treatments. This is to be expected, as it is a shared parameter, and changes it it would be reflected in  $\alpha_i$ , similarly it would be expected that a decrease in  $R_0$  would correspond to an increase in  $\alpha_i$ , to continue to fit with the data. Importantly, there is no correlation between  $\alpha_1, \alpha_2$ .

Figure 3 also shows that the expected values for  $R_0, \alpha_1, \alpha_2$  taken assuming no covariances closely resemble where they would belong with covariances included.

Since it appears that all the assumptions for the model are reasonably satisfied and there is clear convergence in the MH algorithm, the values obtained should be close to the true values. The values obtained are

est =

2.2883      0.8409      0.6781

I.e.  $R_0 \approx 2.2883$ ,  $\alpha_1 \approx 0.8409$  and  $\alpha_2 \approx 0.6781$ . Since a smaller value of  $\alpha_i$  corresponds to a better means of infection control,  $\alpha_2$  appears to be a better treatment than  $\alpha_1$ .

The value for  $R_0$  appears reasonable, as many of the wards experienced quite large outbreaks of the virus.

The corresponding values for  $\alpha_1 R_0$  and  $\alpha_2 R_0$  obtained are, respectively, 1.9242 and 1.5518. Since both values are still larger than 1 they do not prevent the possibility of a major outbreak, and better alternatives should still be researched. However, the reduction is still significant, and if enforced early in an outbreak, could make a significant difference to the number of infections.

## 6 Conclusion

The effective reproduction numbers,  $R_0$  for the control and two treatment types are found to be, respectively, 2.2883, 1.9242 and 1.5518. The two independent treatments appear to have had positive effects on the reproduction number, with treatment 2,  $T_2$  appearing to be the most effective.

The treatments still do not reduce  $R_0$  below the threshold of  $R_0 = 1$ , meaning the treatments do not remove the possibility of major outbreaks. For now  $T_2$  should still be implemented as it is a significant improvement on taking no action. However, more research should be taken to attempt to find a more effective treatment.

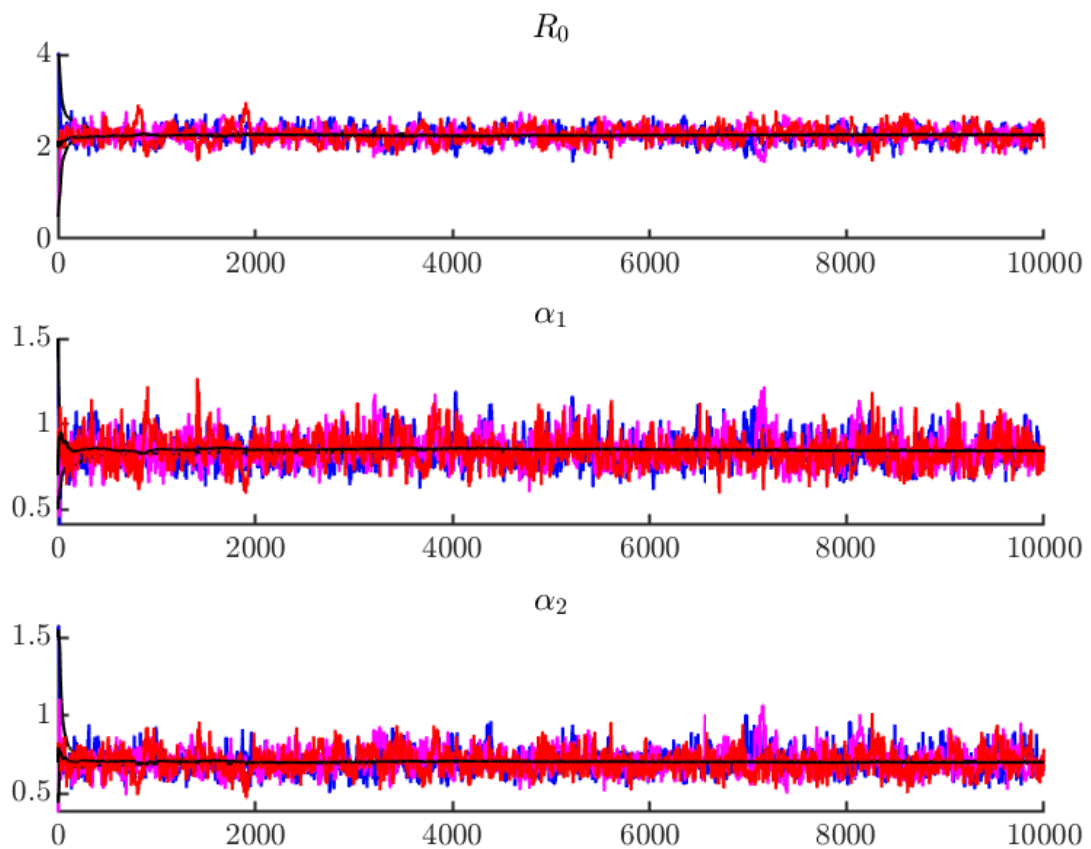


Figure 1: Superimposed trace plots for the estimations of the parameters from the modified MH algorithm

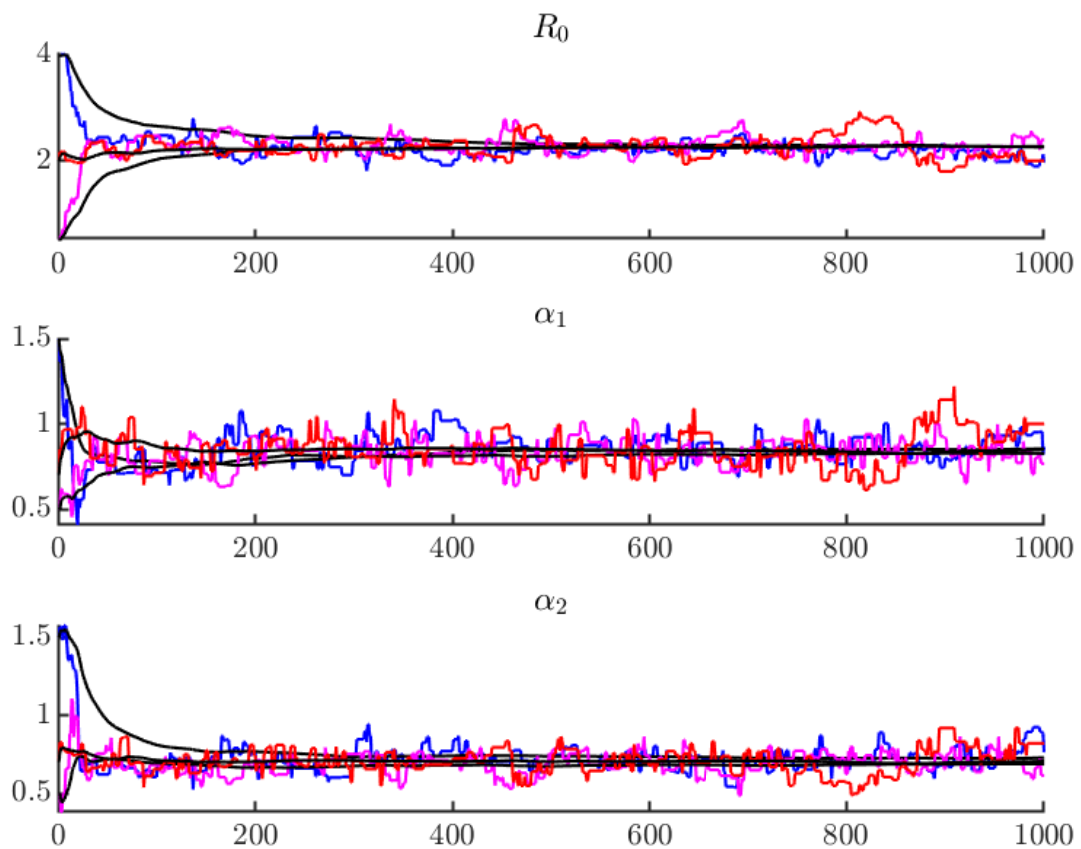


Figure 2: First 1000 iterations of the MH algorithm. Clearly the burn in has completed after approximately 500 iterations for the three parameters

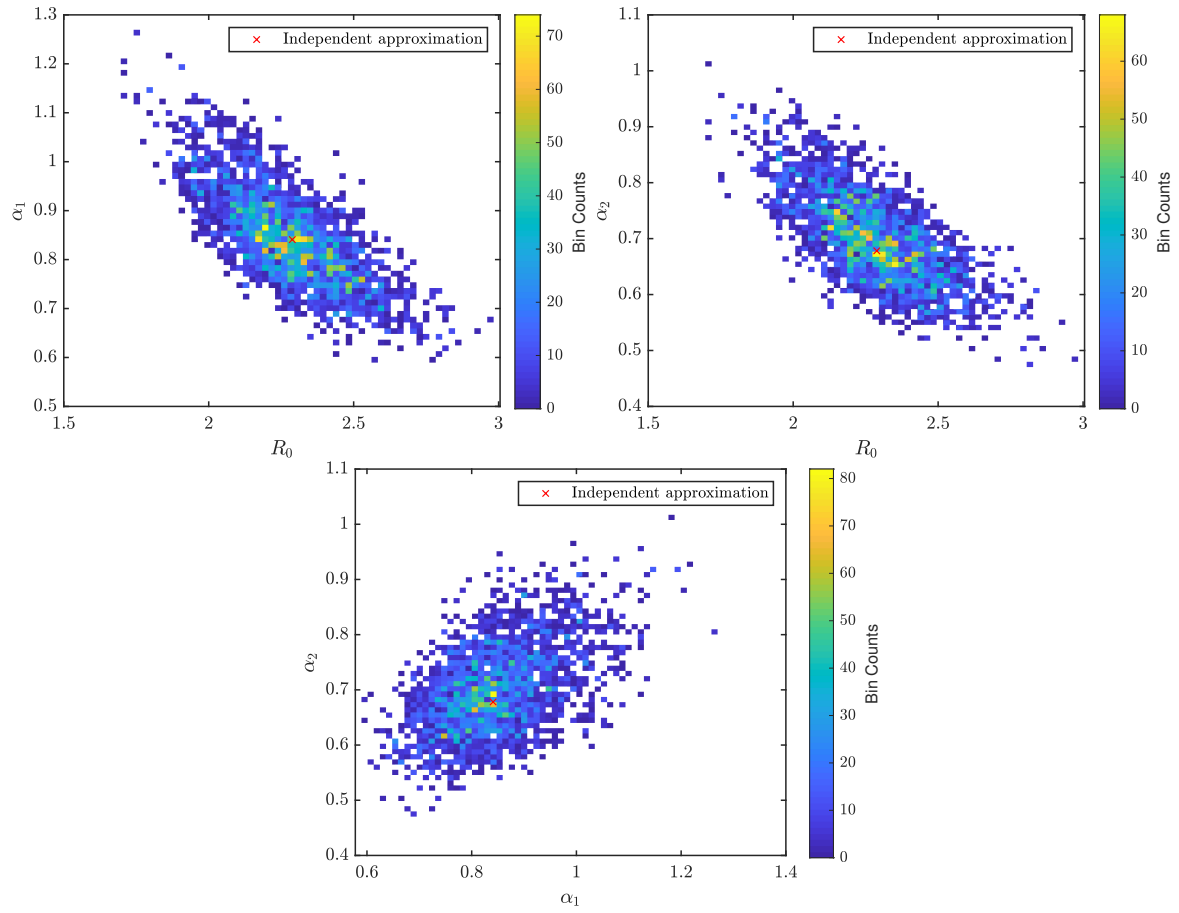


Figure 3: Bin scatter plots of the parameters identifying any correlations and their effect on the value. Top-left:  $(R_0, \alpha_1)$ , Top-right:  $(R_0, \alpha_2)$ , Bottom:  $(\alpha_1, \alpha_2)$

## References

- [1] South Australian Health Department. *Norovirus infection*. URL: <https://www.sahealth.sa.gov.au/wps/wcm/connect/public+content/sa+health+internet/health+topics/health+conditions+prevention+and+treatment/infectious+diseases/norovirus+infection>.
- [2] K. A. M GAYTHORPE et al. “Norovirus transmission dynamics: a modelling review”. In: 146.2 (2018), pp. 147–158. ISSN: 0950-2688.
- [3] K.A.M. Gaythorpe, C.L. Trotter, and A.J.K. Conlan. “Modelling norovirus transmission and vaccination”. eng. In: *Vaccine* 36.37 (2018), pp. 5565–5571. ISSN: 0264-410X.

## A Code

### A.1 Main Script (R0Predict.m)

```
1 %Pretty plots
2 set(groot, 'DefaultLineLineWidth', 1, ...
3     'DefaultAxesLineWidth', 1, ...
4     'DefaultAxesFontSize', 12, ...
5     'DefaultTextFontSize', 12, ...
6     'DefaultTextInterpreter', 'latex', ...
7     'DefaultLegendInterpreter', 'latex', ...
8     'DefaultColorbarTickLabelInterpreter', 'latex', ...
9     'DefaultAxesTickLabelInterpreter', 'latex');
10
11
12
13 %Speeds up the code for multiple runs
14 if ~exist('dataMat','var')
15     dataMat = readmatrix('NorovirusDataA3.txt');
16     %print the first few rows
17     dataMat(1:5,:)
18 end
19
20 %Observation of dataset
21 %how many of them correspond to the disease effectively dying out?
22 amountOfData = length(dataMat(:,1))
23 propDiedOut = sum(dataMat(:,2) <= 0.1 * dataMat(:,1)) / amountOfData
24
25 %dataMat = [#people, #infected, #action]
26 %consistency
27 rng(1)
28
29
30
31 %since we are expecting approximately 66%
32 %getting the right prior
33 %we want the dist to be N'(2.5,sigma)
34 %where sigma is the variance to have approx 66% in (2-3).
```



```
35 %N' since we will drop values below 0 and above 5
36 %N' will still be symmetric
37 sigma = -(3-2.5)/norminv((1-0.66)/2);
38 variance = sigma* speye(3);
39
40 varR0 = 0.01;
41 distStruct.ProposalDist = @(x) mvnrnd(x, diag([varR0,0.01,0.01]));
42 distStruct.PriorPDF = @(x) normpdf(x,2.5,sigma);
43 distStruct.LogLikelihoodFunc = @(params,data) LogLikelihood(params,data);
44 %constrain a_1,a_2 in [0,2]
45 distStruct.ProposalConstraints = @(vars) ProposalConstraints(vars,[0,0,0],[5
46 %overshoot, undershoot and approximately close
47 startVars = [4,1.5,1.5;
48             0.5,0.5,0.5;
49             2,0.7,0.7];
50
51 vars = zeros(10000,3,3);
52 for index = 1:3
53     startPoint = startVars(index,:);
54     exitflag = 1;
55     while exitflag~=0
56         [varsTemp,accrate,exitflag]= MetropolisHastingsPassLikelihood(distStr
57         if exitflag==-1
58             %acceptance was too low
59             %variance is too high
60             varR0 = varR0 *2/3;
61             distStruct.ProposalDist = @(x) mvnrnd(x,diag([varR0,0.01,0.01]));
62             warning('retrying with decreased variance')
63         else if exitflag==1
64             %acceptance was too high
65             %variance is too low
66             varR0 = varR0*2;
67             distStruct.ProposalDist = @(x) mvnrnd(x,diag([varR0,0.01,0.01]));
68             warning('retrying with decreased variance')
69         end
70     end
71     vars(:, :, index) = varsTemp;
72 end
73
74
75 end
76 %%
77 cumVars = cumsum(vars,1)./(1:length(vars))';
78 tstr = [" $$R_0$$", "$$\alpha_1$$", "$$\alpha_2$$"];
79 for burnin = 0:1
80     figure
81     for i=1:3
82         hold on
83         subplot(3,1,i)
84         hold on
```

```
85     plot(vars(:,i,1), 'b')
86     plot(cumVars(:,i,1), 'k')
87     plot(vars(:,i,2), 'm')
88     plot(cumVars(:,i,2), 'k')
89     plot(vars(:,i,3), 'r')
90     plot(cumVars(:,i,3), 'k')
91     title(tstr(i))
92
93     if burnin
94         axis([0,1000,-inf,inf])
95     end
96 end
97     saveas(gcf,"MHplot"+num2str(burnin)+".eps","eps")
98 end
99     axis([0,1000,-inf,inf])
100    xlabel('Iteration')
101
102 %%
103
104 %%Clean new data
105 %just take one set since they all converged
106 %and omit the burnin
107 varsClean = vars(500:end,: ,3);
108
109 close all
110 %density plot
111 %and obtain estimates for R0, a1, a2
112
113 labs = [" $$R_0$$", "$$\alpha_1$$", "$$\alpha_2$$"];
114
115 est = [0,0,0];
116 for i=1:3
117     figure
118     [prob,val] = ksdensity(varsClean(:,i));
119     prob = prob./sum(prob);
120     plot(val,prob)
121     xlabel(labs(i))
122     ylabel("Probability")
123     title("Probability density for "+labs(i))
124     [~,ind] =max(prob);
125     %assuming there is no dependence
126     est(i) = val(ind);
127     saveas(gcf,"Probdensity"+num2str(i),"eps")
128 end
129 est
130 est(2:3) * est(1)
131
132 for i=1:2
133     for j=i+1:3
134         figure
```

```
135     binscatter(varsClean(:,i),varsClean(:,j),60,'HandleVisibility','off')
136     hold on
137     scatter(est(i),est(j),'xr')
138     xlabel(labs(i))
139     ylabel(labs(j))
140     legend("Independent approximation")
141     colormap(gca,'parula')
142     saveas(gcf,"BinScatter"+num2str(i)+num2str(j),"epsc")
143 end
144 end
```

## A.2 MetropolisHastingsPassLikelihood.m

```
1 function [vars,accRate,exitflag]= MetropolisHastingsPassLikelihood(distStruct)
2 %Most generic MetropolisHastings
3 %distStruct has fields
4 %-priorDist - the prior to pull from
5 %-proposalDist - the proposed distribution
6
7 numAccepted = 0;
8 PriorPDF = distStruct.PriorPDF;
9 ProposalDist = distStruct.ProposalDist;
10 BreaksProposalConstraints = distStruct.ProposalConstraints;
11 LogLikelihoodFunc = distStruct.LogLikelihoodFunc;
12 vars = zeros(numIterations,length(startVars));
13 vars(1,:) = startVars;
14
15 for i=2:numIterations
16     proposal = ProposalDist(vars(i-1,:));
17     %if breaks constraints
18     if BreaksProposalConstraints(proposal)
19         vars(i) = vars(i-1);
20     else
21         candidateProbTop = LogLikelihoodFunc(proposal,data) + log(PriorPDF(proposal));
22         candidateProbBottom = LogLikelihoodFunc(vars(i-1,:),data) + log(PriorPDF(vars(i-1,:)));
23         candidateProb = candidateProbTop - candidateProbBottom;
24         acceptProb= log(rand);
25
26         if acceptProb< candidateProb
27             vars(i,:) = proposal;
28             numAccepted = numAccepted +1;
29         else
30             vars(i,:) = vars(i-1,:);
31         end
32     end
33 end
34 accRate = numAccepted/numIterations;
35
36 if accRate < 0.2
37     warning("Bad acceptance rate - too low, accRate = "+num2str(accRate));
38     exitflag = -1;
```

```
39 else if accRate >0.27
40     exitflag = 1;
41     warning("Bad acceptance rate - too high , accRate = "+num2str(accRate));
42 else
43     exitflag = 0;
44 end
45 end
46
47
48 end
```

### A.3 LogLikelihood.m

```
1 function logLikelihood = LogLikelihood(params,data)
2 %Amended SIR finalsize code from Josh
3 %Amended by Andrew Martin
4 %calculates log likelihood for a given R0,alpha1,alpha2
5 logLikelihood = 0;
6 % R0      = params(1);
7 % alpha1  = params(2);
8 % alpha2  = params(3);
9 NVec = data(:,1);
10 %[R0, a1R0, a2R0]
11 paramsModified = [params(1),params(1)*params(2),params(1)*params(3)];
12 %yay matlab uses 1 based indexing
13 R0index = data(:,3)+1;
14
15
16 for iterator=1:length(data)
17     N = NVec(iterator);
18     relevantParam = paramsModified(R0index(iterator));
19     q = zeros(N+1,1);
20     q(2) = 1;
21     %Proportions for each number of infection events
22     %could vectorise , but this is more meaningful
23     for Z2 = 0:N
24         for Z1 = Z2+1:N-1
25             %infection probability (jump prob)
26             infProb = 1 / ( 1 + ((N-1)/(relevantParam*(N-Z1))) );
27             q(Z1+2) = q(Z1+2) + q(Z1+1)*infProb;
28             q(Z1+1) = q(Z1+1)*(1-infProb);
29         end
30     end
31     %sum of the log likelihoods (product of likelihoods)
32     logLikelihood= logLikelihood + log(q(data(iterator,2)+1));
33 end
34
35
36 end
```

## A.4 ProposalConstraints.m

```
1 function boolean = ProposalConstraints(vals,min,max)
2 %
3 %OUTPUT:
4 %boolean – true if the boundary constraints are broken
5 boolean = any(vals < min | vals > max);
6
7
8 % %hardcoded way
9 % boolean =vals(1) < 0 || vals(1) > 5 ...
10 %      || vals(2) < min || vals(2) > max ...
11 %      || vals(3) < min || vals(3) > max;
12
13 end
```

## Assignment III

Worth 15% of course assessment; due by 3pm on Friday 7th June, 2019.

Most-relevant lectures: Lectures 21 – 26.

The total marks for this assignment is 45.

Please provide code where appropriate.

### **Report to the Government on the effectiveness of interventions.**

[45 marks]

The CSV file `NorovirusDataA3` contains information regarding independent outbreaks in a set of 125 hospital wards of varying sizes. The first column contains the number of occupied beds in the ward, the second column contains the number of those patients which succumbed to norovirus during the outbreak, and the third column indicates the control action implemented; in the last column, 0 corresponds to standard practices, 1 corresponds to a trial intervention strategy, and 2 corresponds to different trial intervention strategy.

You are to analyse this data to advise the Government on the effectiveness of their interventions, and to advise them on which intervention should be adopted (if any).

You may assume that the interventions work to reduce the effective transmission rate parameter. The Chief Medical Officer has said that  $R_0$  for norovirus in typical hospital settings is between 2 and 3 with (approximately) 66% probability; you may use this expert opinion as prior knowledge.

You are to prepare two reports. One is for the Government. The second is to provide detail on how you have performed the statistical analysis, including the model(s) used, any assumptions you have made and providing evidence that your approach/algorithms are working correctly, for example through the use of simulated data of similar form to the ‘real’ data (e.g., using trace plots from multiple independent chains, and kernel density estimators and box plots). This assignment is deliberately vague – you need to make decisions, but feel free to ask for feedback as you make progress.