# In Depth Analysis of Norovirus Data

Andrew Martin

June 6, 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 [2].

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:

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

Other analysis suggests that the second infectious state can be disregarded with equally valid results, reducing it to the $SEIR$ model [3].

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

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.

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

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:

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

The likelihood calculations are made using an amended copy of the final size distribution code from Ross et al[1].

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 $\alpha = 0$. Thee 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. Figure 4 plots the kernel-density curves for the parameters, which is the distribution of candidate values for the parameters. The accepted values correspond to the peaks on the density plots. 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.
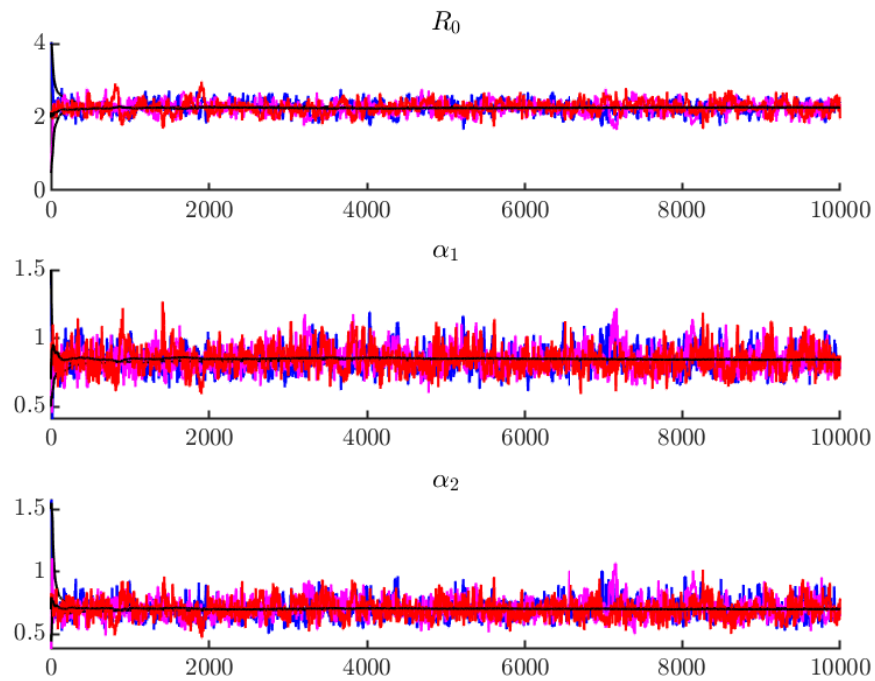
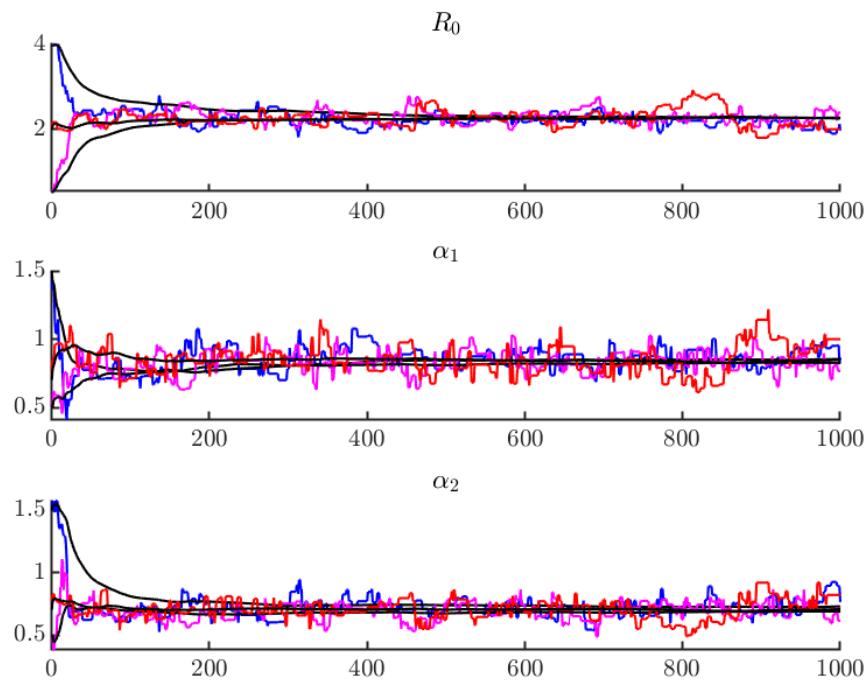Figure 1: Superimposed trace plots for the estimations of the parameters from the modified MH algorithm



Figure 2: First 1000 iterations of the MH algorithm. Clearly the burn in has completed after approximately 500 iterations for the three parameters
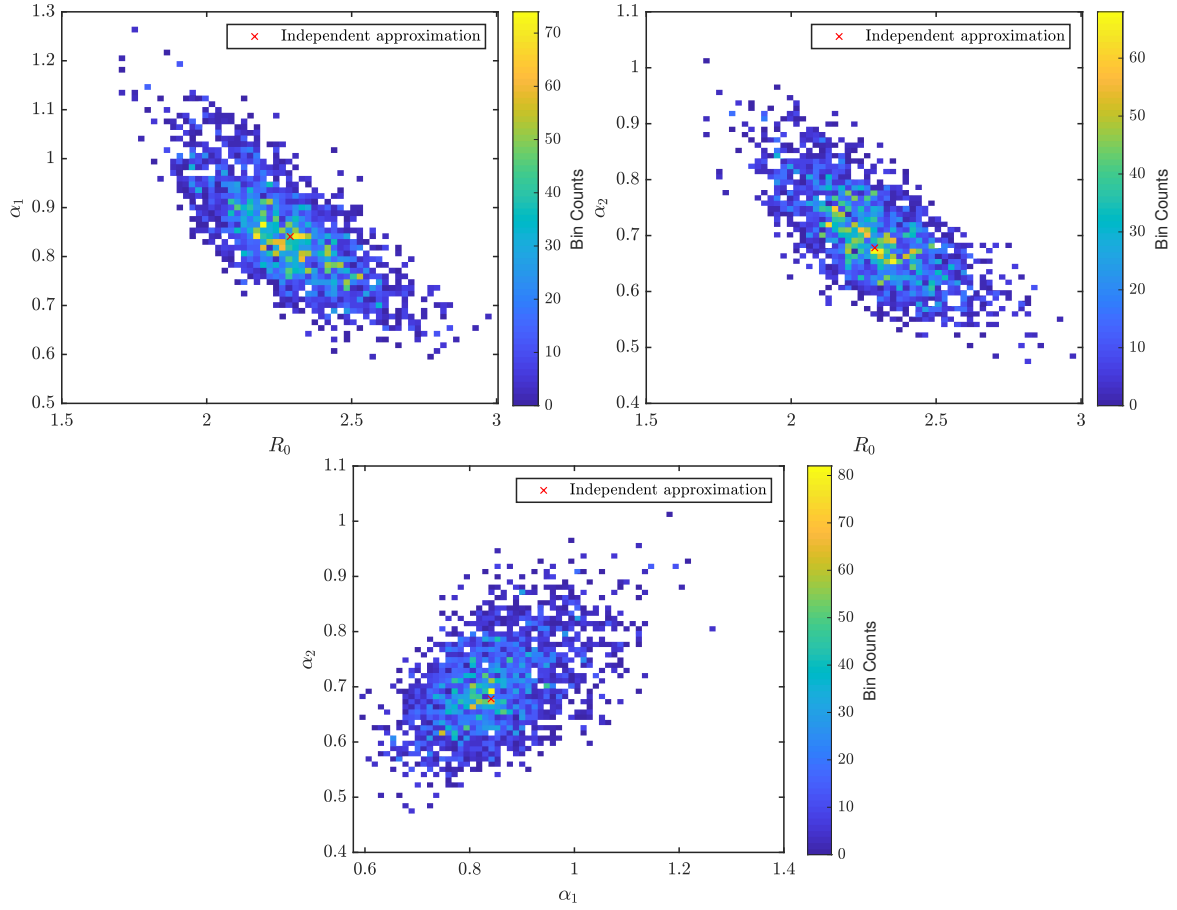
Figure 3: Bin scatter plots of the parameters identifying any correlations and their effect on the value. Top-left: $(R_0, \alpha_1)$, Top-right: $(R_0, \alpha_2)$, Bottom: $(\alpha_1, \alpha_2)$

# 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.
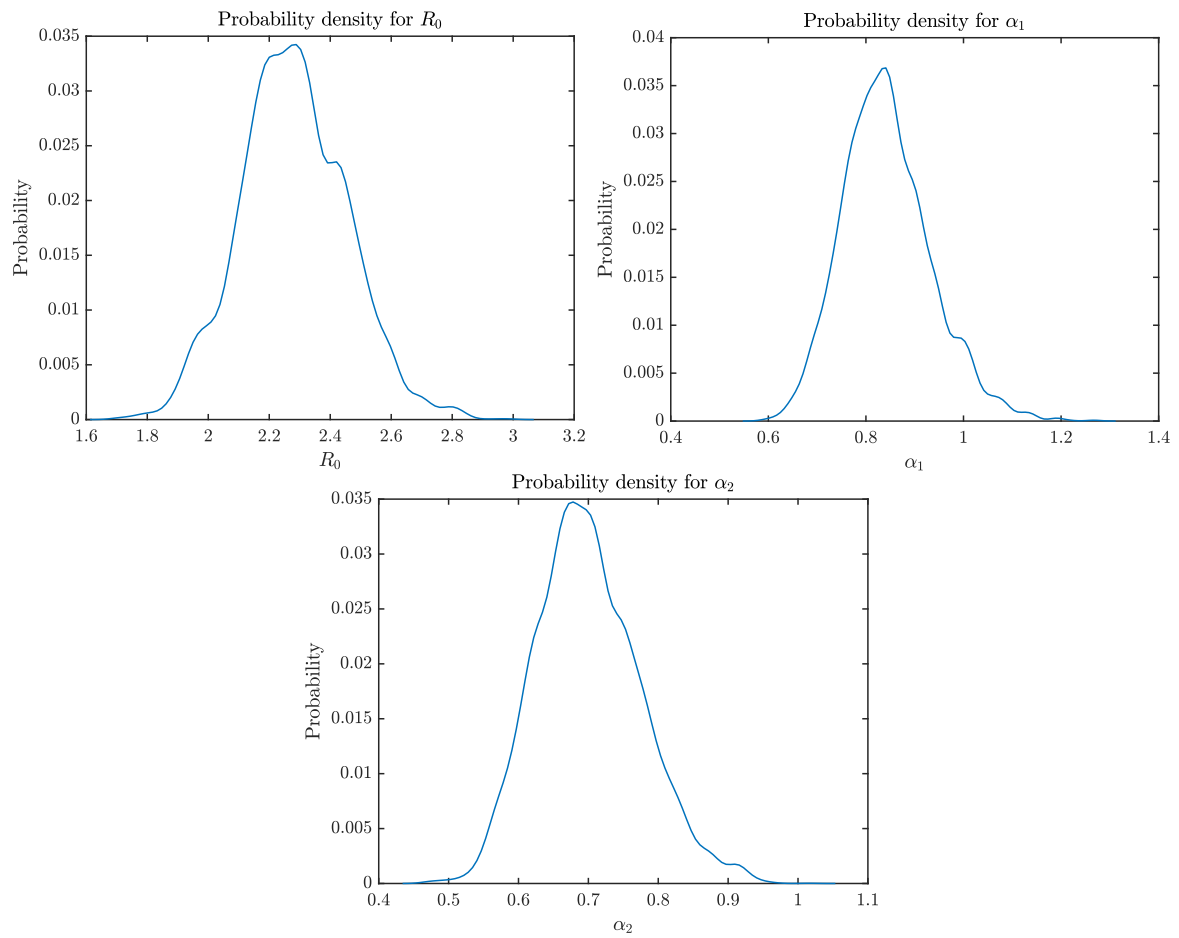
Figure 4: Probability density curves for the parameters - the horizontal axis is the parameter while the vertical is the probability

# References

[1]  Andrew J. Black and J.V. Ross. "Computation of epidemic final size distributions". In: *Journal of Theoretical Biology* 367 (2015), pp. 159–165. ISSN: 0022-5193. DOI: `https://doi.org/10.1016/j.jtbi.2014.11.029`. URL: `http://www.sciencedirect.com/science/article/pii/S0022519314006882`.

[2]  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`.

[3]  K. A. M Gaythorpe et al. "Norovirus transmission dynamics: a modelling review". In: 146.2 (2018), pp. 147–158. ISSN: 0950-2688.

# A   Code

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

```matlab
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 %get the dataset once only
15 if ~exist('dataMat','var')
16     %dataMat has format [#people, #infected, #action]
17     dataMat = readmatrix('NorovirusDataA3.txt');
18     %print the first few rows
19     dataMat(1:5,:)
20 end
21
22 %Observation of dataset
23 %how many of them correspond to the disease effectively dying out?
24 amountOfData = length(dataMat(:,1))
25 propDiedOut = sum(dataMat(:,2)<=0.1*dataMat(:,1))/amountOfData
26
27
28
29 %consistency
30 rng(1)
31
32
```

```matlab
33
34  %since we are expecting approximately 66%
35  %getting the right prior
36  %we want the dist to be N'(2.5,sigma)
37  %where sigma is the variance to have approx 66% in (2-3).
38  %N' since we will drop values below 0 and above 5
39  %N' will still be symmetric
40  sigma = -(3-2.5)/norminv((1-0.66)/2);
41  variance = sigma*speye(3);
42
43  varR0 = 0.01;
44  distStruct.ProposalDist = @(x) mvnrnd(x,diag([varR0,0.01,0.01]));
45  distStruct.PriorPDF = @(x) normpdf(x,2.5,sigma);
46  distStruct.LogLikelihoodFunc =@(params) LogLikelihood(params,dataMat);
47  %constrain a_1,a_2 in [0,2]
48  distStruct.ProposalConstraints = @(vars) ...
49      ProposalConstraints(vars,[0,0,0],[5,2,2]);
50  %overshoot, undershoot and approximately close to the true solutions
51  startVars = [4,1.5,1.5;
52              0.5,0.5,0.5;
53              2,0.7,0.7];
54
55  vars = zeros(10000,3,3);
56  for index = 1:3
57      startPoint = startVars(index,:)
58      exitflag = 1;
59      while exitflag~=0
60          [varsTemp,accrate,exitflag]= MetropolisHastingsPassLikelihood...
61              (distStruct,startPoint,10000);
62          %using exitflag change the variance to get a reasonaable acceptance
63          %rate between 0.2, 0.27
64          if exitflag==-1
65              %acceptance was too low
66              %variance is too high
67              varR0 = varR0 *2/3;
68              distStruct.ProposalDist = @(x) mvnrnd(x,diag([varR0,0.01,0.01]));
69              warning('retrying with decreased variance')
70          elseif exitflag==1
71              %acceptance was too high
72              %variance is too low
73              varR0 = varR0*2;
74              distStruct.ProposalDist = @(x) mvnrnd(x,diag([varR0,0.01,0.01]));
75              warning('retrying with decreased variance')
76          end
77          %store all the solution sets
78          vars(:,:,index) = varsTemp;
79      end
80
81
82  end
```

```matlab
83  %%
84  cumVars = cumsum(vars,1)./(1:length(vars))';
85  tstr = [" $$R_0$$"," $$\alpha_1$$"," $$\alpha_2$$"];
86  for burnin = 0:1
87      figure
88      for i=1:3
89          hold on
90          subplot(3,1,i)
91          hold on
92          %trace plots & trendlines
93          plot(vars(:,i,1),'b')
94          plot(cumVars(:,i,1),'k')
95          plot(vars(:,i,2),'m')
96          plot(cumVars(:,i,2),'k')
97          plot(vars(:,i,3),'r')
98          plot(cumVars(:,i,3),'k')
99          title(tstr(i))
100
101         if burnin
102             axis([0,1000,-inf,inf])
103         end
104     end
105         saveas(gcf,"MHplot"+num2str(burnin)+".eps","epsc")
106 end
107     axis([0,1000,-inf,inf])
108     xlabel('Iteration')
109
110 %%
111 %%Analysis
112 %just take one set since they all converged
113 %and omit the burnin
114 varsClean = vars(500:end,:,3);
115
116 %density plot
117 %and obtain estimates for R0, a1, a2
118 %axis labels for later
119 labs = [" $$R_0$$"," $$\alpha_1$$"," $$\alpha_2$$"];
120
121 %est is the estimated [R0,a1,a2]
122 est = [0,0,0];
123 for i=1:3
124     figure
125     [prob,val] = ksdensity(varsClean(:,i));
126     prob = prob./sum(prob);
127     plot(val,prob)
128     xlabel(labs(i))
129     ylabel("Probability")
130     title("Probability density for "+labs(i))
131     [~,ind] =max(prob);
132     %assuming there is no dependence
```

```
133        est ( i ) = val ( ind ) ;
134        saveas ( gcf ,"Probdensity"+num2str ( i ) ,"epsc")
135    end
136  %print out R0, a1 , a2
137   est
138  %print out R0a1 , R0a2
139   est ( 2 : 3 ) ∗ est ( 1 )
140
141  %covariance−scatter plots
142   for  i =1:2
143       for  j=i +1:3
144            figure
145            binscatter ( varsClean ( : , i ) , varsClean ( : , j ) ,60 , ' HandleVisibility ' , ' off ' )
146            hold on
147            scatter ( est ( i ) , est ( j ) , ' xr ' )
148            xlabel ( labs ( i ) )
149            ylabel ( labs ( j ) )
150            legend ("Independent approximation")
151            colormap ( gca , ' parula ' )
152            saveas ( gcf ,"BinScatter"+num2str ( i )+num2str ( j ) ,"epsc")
153       end
154   end
```

## A.2   MetropolisHastingsPassLikelihood.m

```
1   function [ vars , accRate , exitflag ]= MetropolisHastingsPassLikelihood . . .
2        ( distStruct , startVars , numIterations )
3   %Most generic Log Likelihood MetropolisHastings
4   %%INPUTS
5   %distStruct is a struct containing :
6   %   the function handles
7   %    −PriorPDF − the PDF of the prior to obtain likelihood
8   %    −ProposalDist − the proposed distribution to pull new values
9   %    −ProposalConstraints − the constraints that the params must abide
10  %    −LogLikelihoodFunc − the function to obtain the log likelihood
11  %startVars − the initial values for the vars to estimate
12  %numIterations − the number of times to iterate through the algorithm
13  %%OUTPUTS
14  %vars − the full matrix of variables accepted [ nxd ] where n is
15  %       num iterations and d is the number of variables to predict
16  %accRate − the acceptance rate for the algorithm , a good value is
17  %            in 0.2 − 0.27
18  %exitflag − An enum to log the outcome of the algorithm
19  %            −1, acceptance rate was low ( try decreasing proposed variance )
20  %             0, good acceptance rate
21  %             1, acceptance rate was high ( try increase proposed variance )
22   numAccepted = 0;
23  %extract all the functions
24   PriorPDF = distStruct . PriorPDF ;
25   ProposalDist = distStruct . ProposalDist ;
26   BreaksProposalConstraints = distStruct . ProposalConstraints ;
```

```matlab
27  LogLikelihoodFunc = distStruct.LogLikelihoodFunc;
28  vars = zeros(numIterations,length(startVars));
29  vars(1,:) = startVars;
30
31  for i=2:numIterations
32      %get x'
33      proposal = ProposalDist(vars(i-1,:));
34      %if breaks constraints
35      if BreaksProposalConstraints(proposal)
36          vars(i) = vars(i-1);
37      else
38          %log(a) = ...
39          candidateProbTop = LogLikelihoodFunc(proposal)...
40              + log(PriorPDF(proposal(1)));
41          candidateProbBottom = LogLikelihoodFunc(vars(i-1,:))...
42              + log(PriorPDF(vars(i-1,1)));
43          candidateProb = candidateProbTop - candidateProbBottom;
44          acceptProb= log(rand);
45          %should we accept this state?
46          if acceptProb< candidateProb
47              vars(i,:) = proposal;
48              numAccepted = numAccepted +1;
49          else
50              vars(i,:) = vars(i-1,:);
51          end
52      end
53  end
54  accRate = numAccepted/numIterations;
55
56  %Warn if the acceptance rate is suboptimal
57  %apparently 0.24 is roughly perfect
58  if  accRate < 0.2
59      warning("Bad acceptance rate - too low, accRate = "+num2str(accRate));
60      exitflag = -1;
61  elseif accRate >0.27
62      exitflag = 1;
63      warning("Bad acceptance rate - too high, accRate = "+num2str(accRate));
64  else
65      exitflag = 0;
66
67  end
68
69
70  end
```

## A.3   LogLikelihood.m

```matlab
1  function logLikelihood = LogLikelihood(params,data)
2  %Amended SIR finalsize code from Ross
3  %Amended by Andrew Martin
4  %calculates log likelihood for a given R0,alpha1,alpha2
```

```matlab
5   %Input:
6   %params-   [R0, a1, a2]
7   %data  - [N, I, T_i] where N is the total population, I is number infected
8   %and T_i is the treatment used
9   %Output:
10  %logLikelihood - the logged likelihood of the data given the params
11  logLikelihood = 0;
12  % R0       = params(1);
13  % alpha1   = params(2);
14  % alpha2   = params(3);
15  NVec = data(:,1);
16  %[R0, a1R0, a2R0]
17  paramsModified = [params(1),params(1)*params(2),params(1)*params(3)];
18  %yay matlab uses 1 based indexing
19  R0index = data(:,3)+1;
20
21
22  for iterator=1:length(data)
23      N = NVec(iterator);
24      relevantParam = paramsModified(R0index(iterator));
25      q = zeros(N+1,1);
26      q(2) = 1;
27      %Proportions for each number of infection events
28      %could vectorise, but this is more meaningful
29      for Z2 = 0:N
30          for Z1 = Z2+1:N-1
31              %infection probability (jump prob) - modified from Ross
32              infProb = 1 / ( 1 + ((N-1)/(relevantParam*(N-Z1))));
33              q(Z1+2) = q(Z1+2) + q(Z1+1)*infProb;
34              q(Z1+1) = q(Z1+1)*(1-infProb);
35          end
36      end
37      %sum of the log likelihoods (product of likelihoods)
38      logLikelihood= logLikelihood + log(q(data(iterator,2)+1));
39  end
40
41
42  end
```

## A.4   ProposalConstraints.m

```matlab
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   end
```