

University of Warwick institutional repository: <http://go.warwick.ac.uk/wrap>

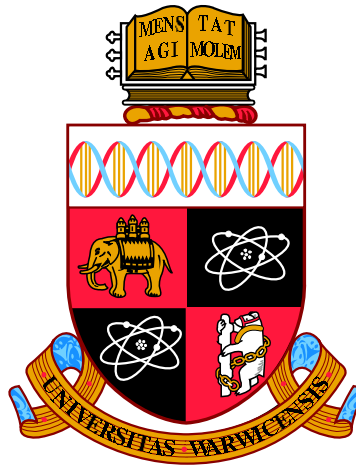
A Thesis Submitted for the Degree of PhD at the University of Warwick

<http://go.warwick.ac.uk/wrap/77458>

This thesis is made available online and is protected by original copyright.

Please scroll down to view the document itself.

Please refer to the repository record for this item for information to help you to cite it. Our policy information is available from the repository home page.



Modelling Health Behaviour

by

Daniel Alexander Sprague

Thesis

Submitted to the University of Warwick

for the degree of

Doctor of Philosophy

Centre for Complexity Science

June 2015

THE UNIVERSITY OF
WARWICK

Contents

List of Tables	v
List of Figures	vii
Acknowledgments	xiv
Declarations	xv
Abstract	xvi
Chapter 1 Introduction	1
1.1 Motivation	1
1.2 The Health Belief Model	3
1.3 Definitions for a general model of health behaviour	3
Chapter 2 Understanding individual-level factors	7
2.1 Introduction	7
2.1.1 Background	8
2.1.2 Previous studies	9
2.1.3 Survey sampling methods	10
2.2 Data	11
2.3 Methods	12
2.3.1 Logistic Regression	12
2.3.2 Clustered bootstrap	14
2.3.3 Model selection	14
2.3.4 Model validation	15
2.4 Analysis	16
2.4.1 FBD mapping	16
2.4.2 Model construction	17
2.4.3 Prediction of priority groups	22

2.5	Results	23
2.6	Discussion	24
2.7	Summary	30
Chapter 3 The role of the health system		32
3.1	Introduction	32
3.1.1	Seeking appropriate healthcare	32
3.1.2	Barriers to access	33
3.1.3	Interventions to improve access	34
3.1.4	Evaluating complex interventions	34
3.1.5	Summary	35
3.2	Methods	35
3.2.1	Approximate string matching	35
3.2.2	Hierarchical models	37
3.2.3	Computational tools	38
3.2.3.1	MCMC	38
3.2.3.2	Hamiltonian Monte Carlo	39
3.2.3.3	MCMC Diagnostics	41
3.2.3.4	Model comparison	42
3.2.3.5	Summarising posterior distributions	42
3.2.4	Spatial modelling	43
3.3	The relationship between health centre quality and usage	44
3.3.1	Data	44
3.3.2	Analysis	46
3.3.3	Results	50
3.3.4	Discussion	57
3.4	The effect of funding mechanism on health facility usage	58
3.4.1	Background	58
3.4.2	Data	59
3.4.3	Analysis	60
3.4.4	Results	63
3.4.5	Discussion	68
3.5	Summary	70
Chapter 4 The spread of health behaviours		72
4.1	Introduction	72
4.1.1	Social Influence	73
4.1.2	Previous models	74

4.1.3	Population models	75
4.1.4	Analysis methods	76
4.1.4.1	Gillespie Algorithm	76
4.1.4.2	Bailey's Method	76
4.2	General complex contagion model	78
4.2.1	Introduction	78
4.2.2	Social influence term	79
4.2.3	Stochastic Model	82
4.2.4	Stochastic results	85
4.2.5	Analysis using Bailey's method	85
4.2.6	Discussion	90
4.3	Deterministic limit	91
4.3.1	Introduction	91
4.3.2	Model	91
4.3.3	Results	92
4.3.4	Discussion	95
4.4	Evidence for complex contagion	95
4.4.1	Introduction	95
4.4.2	Data	96
4.4.3	Statistical model	96
4.4.4	Results	97
4.4.5	Prediction	98
4.4.6	Discussion	102
4.5	Summary	102

Chapter 5 Towards a combined social influence and covariate model104

5.1	Introduction	104
5.2	Coupling behaviour and disease spread	104
5.2.1	Background	104
5.2.2	Model	105
5.2.3	Methods	107
5.2.4	Results	108
5.2.5	Discussion	110
5.3	Towards a combined social and spatial model	114
5.3.1	Background	114
5.3.2	Social model with a spatial component	115
5.3.3	Statistical model with spatial social influence	115

5.3.4	Data	116
5.3.5	Discussion	116
5.4	Summary	117
Chapter 6 Conclusions and Further Work		118
Bibliography		119

List of Tables

2.1	Characteristics of Ugandan mothers of children 0-11 months for each surveyed year.	18
2.2	Characteristics of district-level covariates	18
2.3	Spearman correlation coefficients between pairs of district-level covariates. Population density and health centres per capita, and mean altitude and standard deviation of altitude, showed the highest correlation.	20
2.4	The coefficient for each covariate and interaction term in this predictive model of delivery in a health facility in Uganda. Confidence intervals were calculated using clustered bootstrapping with 1,000 iterations.	25
2.5	The odds ratio for each term included in the model. The odds ratio is given for a unit increase in the (possibly transformed) covariate, for a mother aged 25 and with secondary-level education, surveyed in 2007, and all district-level covariates set to their average. Confidence intervals were calculated using clustered bootstrapping with 1,000 iterations.	26
3.1	Appropriate prescriptions to treat each diagnosis identified in the Sick Child Observation module [1, 2, 3].	46
3.2	Classification of ‘education level’ based on survey response.	46
3.3	Classification of ‘regular partner’ based on survey response.	47
3.4	Gelman-Rubin statistics for the parameters of the best of the three types of model. $\hat{R} \rightarrow 1$ at convergence, so the fact that all of these values are close to 1 is consistent with the models having converged.	52
3.5	Gelman-Rubin convergence diagnostics for each parameter. A value of 1 indicates convergence.	63

3.6	Gelman-Rubin diagnostics for parameters in the two health facility difference models,, consistent with convergence.	66
4.1	Grades of evidence provided by the AIC difference between two models, based on Stylianou et al. [4].	97
4.2	The log-likelihood difference between the simple and complex contagion models. (***) is very strong evidence, (**) is strong evidence, (*) is positive evidence, (.) is no significant evidence for either model, (-) is strong evidence against.	
	† AIC selected models with two peaks.	101

List of Figures

1.1	Illustration of the general model for health behaviour presented in this thesis. The individual-level factors will be investigated in chapter 2, environmental factors in chapter 3, social factors in chapter 4, and disease risk in chapter 5.	5
2.1	Covariates with significant nonlinearity showing untransformed (left) and \log_2 transformed (right).	19
2.2	Covariates that were not improved by the \log_2 transformation.	20
2.3	Scatter plots between pairs of district-level covariates to show the relationship between each pair.	21
2.4	Difference in observed indicator (left) and model residuals (right) against distance for all survey years. Distance between districts is significantly negatively correlated ($p = 0.026$) with the difference in indicator values, but once the model prediction has been subtracted the residuals are no longer significantly correlated ($p = 0.386$).	26
2.5	The ROC curve for this model, showing the true positive rate against false positive rate for various thresholds. The area under this curve is 0.71 [0.70, 0.72], where a perfect classifier would have an AUC of 1.0 and a random classifier would have an AUC of 0.5.	27
2.6	Percentage of mothers of children aged 0-11 months that gave birth in a health facility for A) 2003-04, B) 2006, C) 2009-10, D) 2011. 95% confidence intervals for the indicator are $\pm 11\%$ or lower. Data for 2003 and 2004, and for 2009 and 2010 have been combined for these maps due to the small number of districts surveyed in 2004 and 2010. In the rest of the analysis they are separated.	28

2.7	(A) Priority Map for districts in 2012. Districts are assigned high, mild, low, or unclear priority based on the confidence interval predicted by the model. Kampala, for example, has a predicted indicator confidence interval between 50% and 100%, and is therefore assigned a low priority. (B) Indicator for 2012 as observed in a subsequent LQAS survey. The model was not fitted using this data, and so this map provides an independent test of the predicted confidence interval. All surveyed districts in 2012 have indicator values within the predicted confidence interval, and districts in Eastern Uganda that were predicted to be mild or high priority are observed to have very low values for the indicator.	29
2.8	Predicted 95% confidence range of district indicators for 2012, compared with the observed value. Priority levels are given by the same colour scheme as in Figure 2.7. Only 4 out of 61 (6.6%) of the observed values were outside the predicted range.	29
3.1	Overlap of the token sets for two example ‘names’: “ABC DEF GHI” and “ABC XYZ DEF”. These two names share the tokens “ABC” and “DEF” as a common unique set.	37
3.2	Comparison of two weakly-informative priors for logistic regression.	38
3.3	Comparison of exponential and Gaussian spatial kernels for the same value of the kernel width parameter $\phi = 0.1$	43
3.4	Example semivariograms for an exponential distance kernel with $\phi = 0.1$ (left), and no spatial correlation (right).	45
3.5	Spatial kernel shown in the context of South-West Uganda. The map shows the spatial extent of 75% of the mass of the kernel function evaluated at $\phi = 0.1$ and $\phi = 0.2$, overlaid on a map of the subcounties included in this study. The health facilities and LQAS centroids (where mothers are assumed to be located) are also shown.	49
3.6	Semivariograms for subjective quality (left) and objective quality (right). Neither variable shows a clear semivariance increase with spatial lag, so these variables do not seem to have strong spatial correlation. . .	51

3.7	WAIC values for models of health facility use including subjective and objective quality, for flat models (left) and for hierarchical models (right). The hierarchical models all perform better than the flat models, and most hierarchical models including quality were better than the model that did not. Hierarchical models below $\phi = 0.1$ all performed worse than those with higher values of ϕ , as did those with values higher than $\phi = 0.2$. There therefore appears to be an optimal value of ϕ between 0.1 and 0.2. Finally, hierarchical models which included subjective quality were almost all better than those including objective quality.	52
3.8	Autocorrelations for the subjective (left) and objective (right) model parameters. These plots show that a thinning of 5 to 10 is appropriate.	53
3.9	Autocorrelations for subjective (left) and objective (right) model parameters with a thinning of 10.	53
3.10	Marginal posterior distributions for household-level parameters in the best subjective model. Age is strongly negatively-associated with seeking treatment at a health facility, whereas having a regular partner and secondary education are both weakly positively associated with health facility treatment seeking.	54
3.11	Marginal posterior distributions for health facility-level parameters in the subjective quality model. NGO ownership is strongly positively associated with subjective quality, while health facilities of grade III are likely to have a lower subjective quality. None of the other covariates appeared to predict subjective quality.	55
3.12	Marginal posterior distributions for health facility-level parameters in the objective quality model. NGO ownership is not strongly predictive of objective quality, but health facilities of grade IV are more likely to have higher objective quality as are facilities with lower capacity. .	55
3.13	Marginal posterior distributions for the subjective (left) and objective (right) quality parameters. Both parameters appear to be strongly associated with seeking treatment in a health facility.	56
3.14	Marginal posterior distributions for the standard deviation of the sub-county effect ξ , for the best subjective (left) and objective (right) models.	56

3.15	ROC curves for the model without quality and the best models with subjective and objective quality. The area under curve for the three models are $AUC_{noqual} = 0.73[0.64, 0.80]$, $AUC_{subj} = 0.76[0.66, 0.82]$, and $AUC_{obj} = 0.74[0.64, 0.81]$. All credible intervals overlap significantly, with the subjective model being slightly better than the others as measured by the area under the curve.	57
3.16	Graphical representation of model. N_{hh} is the number of households surveyed in a supervision area (usually 24), N_{sa} is the number of supervision areas per health facility catchment area (always 4 for this data), N_{hc} is the number of health facilities (30 for this data), and t_{before}, t_{after} indicate times before and after the intervention (2012 and 2014) respectively.	62
3.17	Autocorrelation of all household-level covariates (left), and the autocorrelation after a thinning of 8 (right).	64
3.18	Marginal posterior density for coefficients on HFA covariates, showing the parameters for mean price paid by mothers visiting that health facility with sick children, distance to the nearest government health facility of level III or above, and the effect of a health facility being in the treatment group in 2014.	65
3.19	Marginal posterior density for coefficients on household-level covariates, showing parameters for regular partner, age, literacy, year surveyed, and the intercept.	65
3.20	(Left) Selected marginal posterior density for Supervision Area-level random effect. Shown are the largest and smallest effects and a middle-range effect. (Right) Marginal posterior density for the standard deviation of supervision area random effects	66
3.21	ROC curve for this model, with posterior predictive mean and 95% posterior predictive credible interval for the Area Under Curve $AUC = 0.691 [0.685, 0.697]$	67
3.22	Autocorrelation of the beta parameters for the health facility difference models of cost (left) and objective quality (right) at a thinning of 25.	67
3.23	Marginal posterior distributions of the regression parameters for the two health facility difference models of cost (left) and objective quality (right).	68

3.24	Marginal posterior distributions of the standard deviation parameters for the two health facility difference models of cost (left) and objective quality (right).	69
4.1	One realisation of an SIR stochastic process with population size 9, showing the states visited in order (left). Calculation order such that the probability of visiting each state depends only on previously-calculated states (right).	78
4.2	Individual infection rate for simple and complex contagion with $\beta = 1$, $\tau = 5$, and $C = 10$. Simple contagion increases linearly with number of infected contacts, whereas this complex contagion model is a step function with threshold τ	80
4.3	Population-level infection rate for simple contagion and four complex contagions with different parameters. The simple contagion term is again linear, but each complex contagion has a sigmoidal infection function which is steeper as C is increased.	81
4.4	Comparison of Multinomial and Binomial social term for two values of C . The Multinomial term can include multiple types with different thresholds. For the Multinomial model, there is an extra region of low probability in the region where the total infected fraction ($I + J$) is above the lowest threshold, but neither threshold individually has been reached. Parameters are $C = 10$, $\tau_I = 2$, $\tau_J = 4$ (left) and $C = 100$, $\tau_I = 20$, $\tau_J = 40$ (right).	83
4.5	Realisations of the stochastic fad system (left) and the same realisations shifted on the time axis to have the same peak time (right), for different values of τ_j : a) $\tau_j = 5$ b) $\tau_j = 3$. Other parameters are: $N = 500$, $\beta = 1.2$, $\epsilon = 1$, $\tau_i = 2$, $\tau_r = 2$, $I(0) = 5$. All successful fads have very similar shape and timescale	86
4.6	Comparison of the final size distribution calculated using Bailey's Method to the distribution estimated by 1×10^5 simulations using the Gillespie Algorithm. Parameter values are: $N = 100$, $\beta = 0.2$, $\epsilon = 1$, $\tau_i = 1$, $\tau_j = 5$, $\tau_r = 4$, $I(0) = 1$. The two distributions match well, although small variations are noticeable due to Monte Carlo error in the simulation estimate.	87
4.7	Time taken to calculate the final size distribution using this implementation of Bailey's Method against population size N . The time scales as a polynomial between order 2 and order 3.	87

4.8	Probability of more than 50% of the population taking part in the fad for each value of β and ϵ . Other parameter values are: $N = 100$, $\tau_i = 2$, $\tau_j = 5$, $\tau_r = 3$, $I(0) = 2$	88
4.9	Probability of more than 50% of the population taking part in the fad for each value of β and $\frac{I(0)}{N}$. Other parameter values are: $N = 100$, $\epsilon = 1$, $\tau_i = 2$, $\tau_j = 5$, $\tau_r = 3$, $I(0) = 2$	89
4.10	Comparison of the final size distributions for fads with identical parameters before and after one individual has been moved from the S compartment to the I compartment. This small ‘nudge’ leads to a significant chance of the fad infecting the entire population.	89
4.11	Comparison of stochastic and deterministic models for two values of τ_j . Stochastic model is averaged over 100 realisations. Stochastic parameters are $N = 2 \times 10^5$ and $I(0) = 0.01N$. The other parameters are $\beta = 2.7$, $\epsilon = 1$, $\tau_i = 2$, $I(0) = 1$. $\tau_j = 5$ (left) $\tau_j = 3$ (right).	93
4.12	Effect of a) β , b) ϵ , c) τ_i , and d) τ_r on the solution to the ODE system. The two continuous parameters, β and ϵ , both exhibit a sharp threshold between fads that are successful and those that are not. The thresholds τ_i and τ_r control the beginning and end of the fad respectively.	94
4.13	The effect of initial infected fraction on complex contagion (left) and simple contagion (right). For complex contagion, $I(0)$ is very important for determining the success of the fad; for simple contagion $I(0)$ does not affect the success of the fad.	95
4.14	Simple and complex contagion model fits for all fads with sufficient data, ordered by log-likelihood difference. (continued on next page) .	99
4.15	(Continued) Simple and complex contagion model fits for all fads with sufficient data, ordered by log-likelihood difference.	100
4.16	Prediction for the ‘Icebucket Challenge’ fad, using a model fitted to a similar fad. The plot shows the data available when the prediction was made, and the subsequent course of the fad. The model successfully predicts the timing of the peak and the duration of the fad.	102
5.1	Diagram showing the rates between the compartments in the Funk model [5]. Arrows represent possible transitions between the compartments	106

5.2	Timeseries of disease and awareness for the parameter regime with a stable disease-free fixed point, approaching the low-awareness (left) and high-awareness (right) states. In this case the parameter values were: $\beta_u = 1.5, \gamma_u = 1.0, \beta_a = 0.2, \gamma_a = 1.0, \lambda_S = 1.0, \alpha_S = 1.0, \alpha_I = 1.0, \tau_S = 3, \tau_I = 3$. The initial condition for low awareness was $(S_u(0) = 0.98, I_u(0) = 0.01, I_a(0) = 0.0)$, and for high awareness was $(S_u(0) = 0.19, I_u(0) = 0.01, I_a(0) = 0.0)$	109
5.3	Timeseries of disease and awareness for the parameter regime with no stable disease-free fixed point, approaching the low-awareness (left) and high-awareness (right) states. In this case the parameter values were: $\beta_u = 3.0, \gamma_u = 1.0, \beta_a = 0.2, \gamma_a = 1.0, \lambda_S = 1.0, \alpha_S = 1.0, \alpha_I = 1.0, \tau_S = 3, \tau_I = 3$. The initial condition for low awareness was $(S_u(0) = 0.98, I_u(0) = 0.01, I_a(0) = 0.0)$, and for high awareness was $(S_u(0) = 0.19, I_u(0) = 0.01, I_a(0) = 0.0)$	110
5.4	Vector flow diagram cross-section of $(\frac{dS_u}{dt}, \frac{dI_u}{dt}, \frac{dI_a}{dt})$ at $(S_u, I_u, 0)$, for the parameter regime with a stable disease-free fixed point. Stable fixed points are shown as filled circles, and unstable points with open circles.	111
5.5	Vector flow diagram cross-section of $(\frac{dS_u}{dt}, \frac{dI_u}{dt}, \frac{dI_a}{dt})$ at $(S_u, I_u, 0)$, for the parameter regime with no stable disease-free fixed point. Stable fixed points are shown as filled circles, and unstable points with open circles.	112
5.6	Lowest stable endemic disease for parameters β_a and β_u . There is a linear boundary between combinations of β_a and β_u that lead to disease-free states and those that do not.	113
5.7	Critical value of β_a below which there is a disease-free state. This is a linear relation, with slope and intercept dependent on the other parameters. In this case, for $\beta_u > 1.91$ there is no value of β_a that gives a disease-free state.	113

Acknowledgments

I am very grateful to my PhD supervisors Dr Thomas House and Professor Gareth Roberts for giving me their time and energy. Their support has been invaluable, and I have very much enjoyed working with both of them. I would also like to thank Joe Valadez and his team at the Liverpool School of Tropical Medicine for providing data, taking me out to Uganda, and providing public health knowledge.

I have very much enjoyed my time at the Centre for Complexity Science at Warwick, due mainly to the staff and students that I have worked with there. I would also like to thank the WIDER research group for their help and advice.

Finally, I'd like to thank my parents for their advice and support (particularly regarding time management), and Sarah-Emily Mutch for putting up with living with a hermit during the last few months.

This work was supported by the Engineering and Physical Sciences Research Council as part of the University of Warwick Complexity Science Doctoral Training Centre.

Declarations

This thesis is submitted to the University of Warwick in support of my application for the degree of Doctor of Philosophy. It has been composed by myself and has not been submitted in any previous application for any degree.

The work presented (including data generated and data analysis) was carried out by the author except in the cases outlined below:

- Uganda LQAS dataset in Chapter 2 was provided the Liverpool School of Tropical Medicine
- Uganda LQAS and HFA dataset in Chapter 3 was provided by the Liverpool School of Tropical Medicine
- Northern Uganda LQAS and HFA dataset in Chapter 3 was collected by the Liverpool School of Tropical Medicine

The following publications are in the process of submission:

- Sprague, D. A., Jeffery, C., Crossland, N., House, T., Roberts, G. O., Ouma, J., Lwanga, S. K., and Valadez, J. J., Assessing delivery practices of mothers over time and over space in Uganda, 2003-2011. *Tropical Medicine and International Health*, SUBMITTED, 2015.
- Sprague, D. A. and House, T. How behavioural trends spread. *Royal Society Open Science*, SUBMITTED, 2015.

Abstract

Many diseases can be prevented or mitigated through behaviour change, but we lack a quantitative model that can accurately predict these changes and inform policies designed to promote them. Here we introduce a quantitative model of health behaviour that takes into account individual-level barriers, the health system, and spread between individuals. We investigate limits of the model where each of these determining factors is dominant, and use them to predict behaviour from data.

We apply the model to individual-level geographic barriers to mothers giving birth in a health facility, and find evidence that ease-of-access is a major determinant of delivery location. The geographic barriers allow us to explain the observed spatial distribution of this behaviour, and to accurately predict low prevalence regions.

We then apply the model to the role of the health system in determining health facility usage by mothers of sick children. We show that local health facility quality does predict usage, but that this predictive power is significantly less than that gained by including unaccounted-for spatial correlation such as social influence. We also show evidence that results-based funding, rather than traditional input-based funding, increases usage.

We develop a psychologically-motivated ‘complex contagion’ model for social influence and incorporate it into a general model of behaviour spread. We apply this model to short-lived behavioural fads, and show that ‘nudges’ can be very effective in systems with social influence. We successfully fit the model to data for the online spread of real-world behaviour, and use it to predict the peak time and duration of a fad before the peak occurred.

Finally, we discuss ways to incorporate disease state into the model, and to relax the limits used in the rest of the thesis. We consider a model which links health behaviour to disease, and show that complex contagion leads to a feature that is not present in traditional models of disease: the survival of an epidemic depends non-trivially on the initial fraction of the population that is infected. We then introduce two possible models that include both social influence and an inhomogeneous population, and discuss the type of data that might be required to use them predictively.

The model introduced here can be used to understand and predict health behaviours, and we therefore believe that it provides a valuable tool for informing policies to combat disease.

Chapter 1

Introduction

1.1 Motivation

Disease is a far larger cause of death than injury and war, causing around 71% of the 55 million global deaths worldwide in 2013 [6]. Of these deaths, 22% were due to communicable diseases such as HIV/AIDS, tuberculosis, and diarrhoea, and 49% due to non-communicable diseases such as cancers and cardiovascular diseases. Although life expectancy is increasing globally [6], there are still large disparities between the developed and developing worlds and in some countries death rates due to particular diseases are also increasing. This increase is despite the fact that many diseases are known to be prevented or mitigated by changes to an individual's behaviour. Almost all of the major and increasing causes of death specifically noted in the Global Burden of Disease (GBD) Study 2013 [6] can be prevented or mitigated by known behaviour changes: using a condom to prevent HIV/AIDS transmission [7], using a condom and avoiding needle-sharing to reduce the risk of liver cancer-causing Hepatitis C [8], not misusing drugs to avoid drug use disorders, avoiding smoking to reduce the risk of pancreatic cancer [9], and dietary changes to reduce the risk and severity of diabetes [10], chronic kidney disease [11], and atrial fibrillation or flutter [12]. The other major and increasing cause of death specifically noted in the GBD Study, sickle-cell anaemia, is a genetic disease, but taking medication to mitigate this disease could still be considered a health behaviour.

For the purposes of this thesis we consider a change in behaviour to be the action an individual completes, and does not include situations where an individual is willing but unable to complete the change; if clean needles are not available to a needle-sharing individual then we will still consider the individual's behaviour to be sharing needles. It is the individual's action which affects the risk of disease, and so

by understanding these actions we may be able to reverse these increases and further reduce the global burden of disease.

There are four main causal mechanisms for an individual's risk of contracting a disease. To use Type 2 diabetes as an example: genetic factors mean that some individuals are more likely to contract Type 2 diabetes than others [13]; direct environmental risk factors such as the uterine environment can lead to a higher risk of Type 2 diabetes in a child whose mother has contracted the disease [14]; and behaviours such as a healthy diet and physical exercise can help to prevent diabetes [10]. For communicable diseases, the disease state of other individuals in the population is also an important risk factor [15]. We can describe the risk of an individual contracting a disease, taking into account these types of factors, as

$$P(Z_i(t) = 1 | B_i(t), \mathbf{Z}(t), G_i, E_i(t), t),$$

where $Z_i(t)$ is an individual's disease state at time t , G_i are the individual's genetic traits, $E_i(t)$ are risk factors in the individual's environment, $\mathbf{Z}(t)$ is the disease state of the whole population, and $B_i(t)$ is the behaviour of the individual.

Genetic information is increasingly easy to obtain, and its effects on disease can be assessed through, for example, genome-wide association studies [13]. Environmental factors, such as air pollution, are also generally easy to measure directly [16], because they do not have the ethical and practical problems associated with observing human behaviour. The effect of the disease state of others on disease state is the focus of much of infectious disease epidemiology, which has developed many techniques to approach this problem [15]. There is a similar difficulty with understanding the behavioural risk factors $B_i(t)$, in that they are likely to be influenced by the behaviour of others in the population [17], which introduces a feedback loop not present for environmental or genetic factors. To increase the complexity still further, individual behaviour could also depend on the disease state $\mathbf{Z}(t)$ of the rest of the population [18, 19, 20], and on environmental factors such as the ability to find cheap and healthy food in a 'food desert' [21]. This complexity could explain why there is as yet no quantitative general model of health behaviour, and why the burden of diseases preventable by behavioural intervention is still increasing [6].

The aim of this thesis is to develop a general quantitative model of health behaviour that can be fitted to data and used predictively to understand how to encourage good health behaviours. In this chapter we discuss a qualitative model from the psychology literature and use it to motivate our quantitative model. Chapters 2-4 investigate specific limits of the general model, and chapter 5 discusses initial

approaches to combining these limiting cases into a general model.

1.2 The Health Belief Model

There are several qualitative psychological models for an individual's decision to take part in a health behaviour [22], but one of the most widely-accepted is the Health Belief Model (HBM) [23]. The model was originally developed in the 1950s to explain conceptually the failure of individuals in the United States to take part in programs to prevent and detect disease. The HBM defines key concepts that are believed to influence health behaviour:

- Perceived susceptibility: an individual's belief in the likelihood of contracting a disease
- Perceived severity: their belief in the severity of the disease
- Perceived benefits: their belief that the behaviour in question will reduce the risk or severity of the disease
- Perceived barriers: their belief in the physical and psychological costs of the behaviour
- Cues to action: the awareness that an individual has about the behaviour, such as from public health information campaigns, or from social contacts
- Self-efficacy: the individual's confidence in their own ability to take action

Estimates of each of these concepts applied to individuals have been used to predict health behaviours. A major synthesis of evidence for the HBM, conducted by Janz and Becker [23], concluded that susceptibility, severity, benefits, and barriers were statistically significant predictors of preventative health behaviours, sick-role behaviours, and health centre utilisation in most of the studies examined. The HBM has also been used successfully to design interventions to prevent risky sexual behaviours and to promote cancer-screening uptake [22]. The HBM therefore provides a useful qualitative framework for understanding health behaviours, and can help to inform the quantitative and predictive model of health behaviours that is the goal of this thesis.

1.3 Definitions for a general model of health behaviour

Throughout this thesis, health behaviours will be treated as being discrete: each individual may only be in one behavioural state at a time. For our general health

behaviour model, an individual i 's behaviour at time t is given by an integer random variable $B_i(t)$, with each integer referring to a different behavioural state. Possible states could be, for example, $B_i(t) \in \{0, 1\}$, with 1 corresponding to individual i taking part in the behaviour and 0 to not taking part. The behaviour of the entire population at time t is given by the vector $\mathbf{B}(t)$.

Not all individuals are the same: each will have characteristics such as age and level of education. We label these covariates $\Theta(t)$. Individuals may also live in different environments, such as urban or rural, and have different levels of access to their local health system. These environmental factors we label $\mathbf{E}(t)$. For convenience, we combine the individual and environmental covariates into one variable:

$$\mathbf{X}(t) = \{\Theta(t), \mathbf{E}(t)\}.$$

Individuals may be influenced by others in their social contact network, and the effect of this influence will be labelled $f(\mathbf{B}(t))$.

There may be an additional dynamic process occurring at the same time as the behavioural process, $\mathbf{Z}(\mathbf{B}(t), t)$, which depends on the behavioural state. In particular, this could be an infectious disease spreading in the same population as a preventative behaviour which reduces the spread of the disease.

Finally, in some chapters we consider survey data, in which the behaviour of only a small sample of the population is observed. In this case, we define the survey sample $\mathbf{Y}(t) \subset \mathbf{B}(t)$, which is the behaviour of individuals in the sample at time t .

All of the variables defined above are relevant to behaviour, according to the HBM:

- Perceived susceptibility: behaviour depends on the disease state of other individuals in the population, $\mathbf{Z}(\mathbf{B}(t), t)$, which in turn depends on their behavioural states.
- Perceived severity: behaviour depends on both the disease state of others, $\mathbf{Z}(\mathbf{B}(t), t)$, and on covariates such as an individual's level of education given in $\Theta(t)$.
- Perceived benefits: behaviour depends on covariates such as an individual's level of education given in $\Theta(t)$.
- Perceived barriers: some behaviours depend on ease of access or costs of taking part in the behaviour. This could act through the individual's environment, $\mathbf{E}(t)$, such as the local health system.

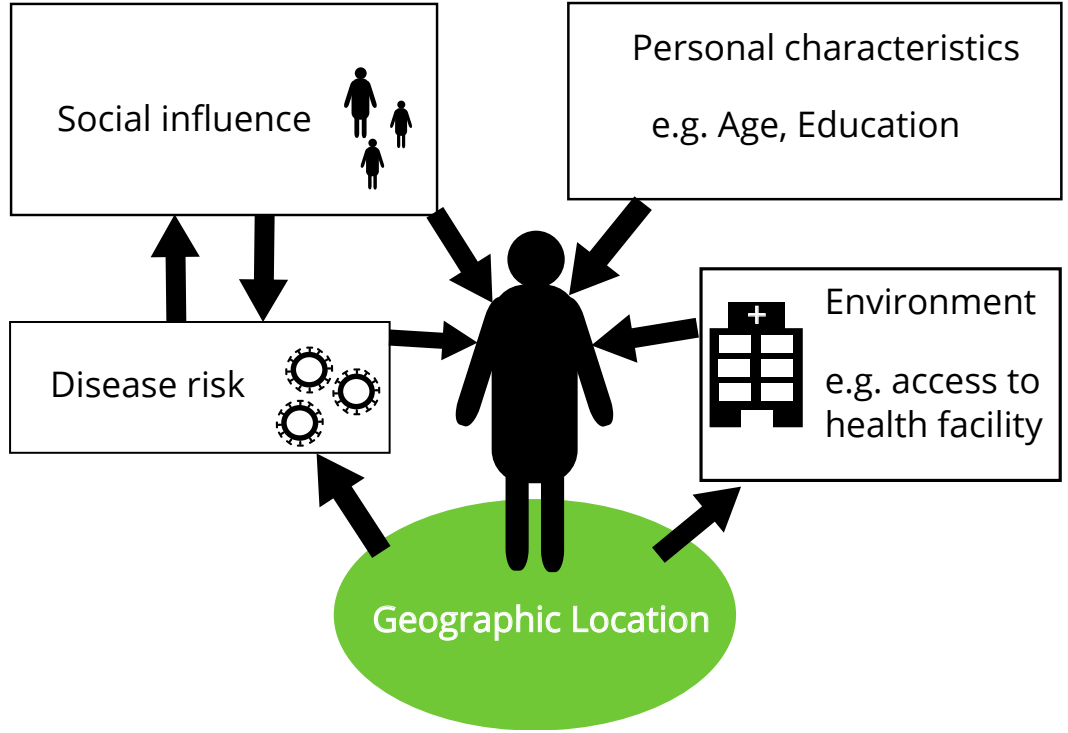


Figure 1.1: Illustration of the general model for health behaviour presented in this thesis. The individual-level factors will be investigated in chapter 2, environmental factors in chapter 3, social factors in chapter 4, and disease risk in chapter 5.

- Cues to action: behaviour depends on the behaviour of social contacts, $f(\mathbf{B}(t))$
- Self-efficacy: may depend on many individual, $\Theta(t)$, and environmental, $\mathbf{E}(t)$, factors.

If these are the important variables for understanding behaviour, then

$$\mathbf{B}(t) \rightarrow \mathbf{B}(t, f(\mathbf{B}(t)), \mathbf{X}(t), \mathbf{Z}(t, \mathbf{B}(t))).$$

The general model is illustrated schematically in Figure 1.1.

We will use two main assumptions to make the model tractable. First, we assume that changes in behaviour are Markovian, which means that the probability of $\mathbf{B}(t)$ transitioning to another state $\mathbf{B}(t + \Delta t)$ depends, at $O(\Delta t)$, only on $\mathbf{B}(t)$, $\mathbf{X}(t)$, and $\mathbf{Z}(\mathbf{B}(t), t)$. This ‘memoryless’ property allows us to model the time each individual spends in a particular behavioural state as following an exponential distribution, with a rate depending on $\mathbf{B}(t)$, $\mathbf{X}(t)$, and $\mathbf{Z}(\mathbf{B}(t), t)$. We refer to the component of the rate that depends on $\mathbf{B}(t)$ as the ‘social influence’ component, since it arises from individuals observing and being influenced by the behaviours of those around

them; the component that depends on $\mathbf{X}(t)$ we refer to as the ‘covariate’ component, which arises from external factors that influence an individual’s behaviour.

The second assumption is that the covariates $\mathbf{X}(t)$ have a linear effect on the log-odds of an individual taking part in a behaviour:

$$\text{logit}(P(B_i(t) = 1)) = \beta^T \boldsymbol{\Theta} + \gamma^T \mathbf{E}. \quad (1.1)$$

The log-odds of a probability p are defined as:

$$\text{logit}(p) = \log \left(\frac{p}{1-p} \right).$$

This assumption is the key one behind the statistical method Logistic Regression, described in chapter 2 and used throughout the thesis.

In the following chapters, we consider three limits:

- Covariate-only, chapter 2 and chapter 3: consider a timescale much slower than that on which social influence operates, so $\mathbf{B}(t) \rightarrow \mathbf{B}(\mathbf{X}(t), \alpha)$ where α is a constant ‘average’ social influence effect on each individual.
- Homogeneous population, chapter 4: consider populations in which the covariates are the same for all individuals and are constant over time, so $\mathbf{B}(t) \rightarrow \mathbf{B}(t, f(\mathbf{B}(t)))$.
- Homogeneous population with disease, chapter 5: consider social influence in an homogeneous population with a separate dynamical process influencing behaviour, so

$$\mathbf{B}(t) \rightarrow \mathbf{B}(t, f(\mathbf{B}(t)), \mathbf{Z}(t, \mathbf{B}(t))).$$

The general model which includes all of the aspects mentioned above is:

$$\mathbf{B}(t) \rightarrow \mathbf{B}(t, f(\mathbf{B}(t)), \mathbf{X}(t), \mathbf{Z}(t, \mathbf{B}(t))).$$

In chapter 5 we discuss the data required to fit such a model, and discuss how this model might be approached based on the insights gained from the other work in this thesis.

Chapter 2

Understanding individual-level factors

2.1 Introduction

Many individual-level factors can influence health behaviour. Age, for example, has been shown to be important in determining whether an ill person will visit a health worker or just stay at home [24]; education may mean the difference between someone using a condom instead of having unprotected sex [25]. An individual's surroundings also affect their behaviour: if the roads in their area are rough then they may be less likely to go to a hospital [26].

Most studies of health behaviour concentrate on a few behaviours and a small number of factors, such as age and education, or wealth, or proximity to health centres [27, 28, 29, 30, 31, 32]; however, they do not combine these into an overall model that takes into account interactions between factors. Due to the difficulty in collecting sufficient data covering a large area, previous studies also tend to ignore the spatial distribution of behaviours. Most data come from expensive surveys, carried out infrequently, and hence cannot be used to investigate time trends. Recently, however, a sampling methodology called Lot Quality Assurance Sampling (LQAS) has been applied to survey collection. It allows smaller samples sizes to be collected at each survey site while maintaining statistical power, and thus the cost of sampling is sufficiently low that surveys may be carried out more frequently and in more areas [33].

In this section I will develop a framework for using LQAS data to visualise the distribution of a health behaviour across a country, understand the demographic and geographic factors influencing that behaviour, and finally predict which areas of the

country are likely to require intervention in future. I will illustrate this framework by applying it to one particular health behaviour: mothers choosing to give birth in a health centre, using LQAS data from Uganda.

2.1.1 Background

There is clear evidence that the prevalence of particular health behaviours varies between socioeconomic groups: Rani et al. [34], for example, show that in India the prevalence of tobacco consumption, a health behaviour strongly associated with mortality [35], was significantly higher in poorer, less educated, or older individuals. The correlation of socioeconomic status and education with an individual's health behaviour has been studied in much more depth than the relationship with the local environment: the limited availability of healthy and affordable food in a 'food desert' [21], for example, acts as a barrier to healthy eating in these areas. Under the HBM framework this would relate to the 'perceived barriers' concept, but in previous empirical HBM studies this component has been measured by self-reporting rather than geographic data. These barriers in the local environment are important to consider because they may prevent an individual from taking part in a behaviour even if they are aware of its health benefits. It is the goal of this chapter to add these environmental aspects to our current understanding of health behaviours, allowing us to better predict which individuals will take part in a health behaviour and which geographical regions have high and lower prevalence of the behaviour.

The biggest impacts that understanding health behaviours can have on reducing disease are in developing countries, where the risk of preventable disease is high and basic positive health behaviours are not being followed due to a combination of availability of health services and individual willingness. More than 10 million children die each year, most from preventable diseases in developing countries [36]. Jones et al. [37] estimate that currently known interventions could prevent around 63% of these deaths, but that these interventions are not being delivered to the mothers and children who need them. Three of the Millennium Development Goals (MDG) have components directly related to health behaviours [38]:

- Goal 4: Reducing the under-5 mortality rate
- Goal 5: Proportion of births attended by a skilled health worker
- Goal 6: Proportion of children under 5 sleeping under insecticide-treated bed nets

The Millennium Development Goals Report 2014 [38] has identified both child mortality and maternal mortality as areas where much more effort is needed to reach set targets. Understanding which children and mothers are least likely to take part in basic positive health behaviours, and which regions have low prevalence of these behaviours, could be an important tool for reaching the MDG targets in future. In this chapter we therefore investigate behaviours relating to maternal and child health in developing countries and, as called for in the MDG Report 2014, use statistical modelling to identify problems in local areas and specific population groups to suggest effective interventions.

The risk of maternal and neonatal mortality is highest around the time of labour [39], and this risk can be significantly reduced if skilled health personnel attend the delivery [40, 41, 42, 43]. In low-income countries, newborns delivered by a skilled birth attendant in a health facility stand a greater chance of survival than newborns delivered elsewhere [40]. Campbell and Graham [44] strongly recommend a policy of encouraging facility-based delivery (FBD) to reduce maternal mortality, due to the higher likelihood of skilled attendance and the provision of emergency obstetric care that would not be available at home. In many developing countries, however, many mothers give birth at home without skilled attendance.

2.1.2 Previous studies

Previous studies of delivery in a health facility have included the following factors: maternal characteristics, index pregnancy characteristics, access, socio-cultural beliefs and past experiences of the mother. Maternal characteristics such as young maternal age, high levels of education and increased autonomy are positively associated with FBD [27, 28, 29, 30, 45]. Parity, or the number of times a woman has given birth, is negatively associated with FBD use [29, 46], and several traits relating directly to the pregnancy in question also affect safe delivery practices: labour onset late at night or short duration of labour can inhibit a mother from accessing a health facility for delivery services [46]. The number of ANC visits a mother has undertaken is positively associated with the use of health facilities for delivery, possibly because ANC staff are more likely to advise in favour of FBD [47].

Ease of physical and financial access to health facilities is positively associated with FBD [29, 47, 48]. Increased distance to facilities or facilities located in difficult terrain, as well as high costs, both formal and informal, are negatively associated with FBD [27, 28, 31, 32]. Most of the studies of FBD mentioned above use data from countrywide cross-section surveys, such as Demographic and Health Survey, at one time point. These studies do not provide information on time trends or on

the spatial distribution of FBD prevalence, because these surveys select clusters of households from across a wide spatial area [49].

2.1.3 Survey sampling methods

Most previous studies investigating facility-based delivery and other health behaviours in developing countries use data from clustered random sampling. A common survey design for estimating level of use of health services is the ‘30 × 7’ method [50], in which 30 clusters are selected across a country and 7 individuals are interviewed in each cluster [51]. A greater numbers of clusters can be chosen to provide better regional coverage: the MEASURE DHS household survey in Uganda in 2011, for example, selected 404 clusters and 25 households from each cluster [45] to obtain confident estimates for the 10 sub-regions in Uganda. This method is practical for estimating health behaviour prevalence across a country or in large spatial regions within it, but to resolve finer spatial differences would require more clusters and hence more expense.

An alternative to cluster sampling which estimates prevalence on finer spatial scales and is cheaper to implement at the local level is LQAS sampling. LQAS is a sampling method that is increasingly being used to monitor health systems across the world [33] because it requires smaller sample sizes than traditional clustered sampling to identify spatial regions that are not reaching a particular target; it is also easy to incorporate into local Monitoring and Evaluation systems and can therefore be carried out more frequently and ensure higher coverage [33]. The LQAS method splits large spatial regions into sub-regions, and calculates the probability that a particular sub-region, or ‘supervision area’ (SA), has met a specified prevalence target, based on a random sample from that supervision area. The sample size n is chosen such that there is acceptably low probability of misclassification [50]; that is, classifying a supervision area that meets the target prevalence as having failed to meet it and vice versa.

All supervision areas within the larger spatial region are surveyed during LQAS, so in addition to providing information on target achievement, the data from these individual supervision areas can be aggregated to provide a stratified random sample of the aggregated area [52]. This is an advantage for this application because some of the supervision areas within a region may be different from each other: some may be urban and others rural, or some may be closer to health centres than others. By obtaining a stratified sample we can be more confident that the aggregated estimate is representative of the whole larger region [53], and obtain an estimate with higher precision than that obtained by clustered sampling [54].

The 95% confidence interval for the aggregate prevalence estimate can be calculated based on the sample size for different values of the true prevalence in a region. The worst-case scenario for the confidence interval occurs for a true prevalence of 50%, and the interval width decreases with sample size. For sample sizes above 95 even this worst-case scenario has a 95% confidence interval of less than $\pm 10\%$, so most LQAS surveys choose supervision area sample sizes such that when aggregated to regions of interest the total sample size is above 95 [52].

The LQAS method's improved spatial information and lower cost for repeated surveys makes it a practical tool for investigating the spatial distribution and temporal trends of health behaviour prevalence in developing countries.

2.2 Data

Uganda, being a Countdown Country, is one of 75 countries where $>95\%$ of maternal and child deaths take place, and these mortality rates are high nationwide. The World Health Organisation and UNICEF recommend that countries analyse data at both national and subnational levels to identify gaps and inequities in health services [55]. They urge improved data collection and analysis methods to decrease the amount of missing data and ensure that timely results are made available to policymakers [55].

It is also important to uncover factors influencing facility-based delivery and take into account potential confounders. The studies mentioned above did not investigate interactions between variables, which are essential to understanding situations where two or more explanatory variables enhance or reduce each other. They generally included at most one geographic covariate, and did not try to explain or predict spatial variations in FBD. In this chapter, we assess variations over time and space in FBD in Uganda, and fit a statistical model to identify factors associated with FBD. We apply this model to 2003-11 data to identify areas expected to have low indicator coverage in 2012, and validate this prediction with the actual data. This approach can therefore inform policy-makers and program managers on the status of FBD, trends and variations occurring over time, and can identify locations needing further investigation.

We analyse a collection of LQAS surveys, in which a set of Ugandan districts were chosen for each year in which a survey took place. Trained district health managers collected data with household surveys conducted in 19-64 districts at 7 time-points during 2003-2011, using standard LQAS methodology [56]. The District Health Management Team divided each district into 4-6 administrative sub-district

strata or supervision areas (SA), and selected 19 mothers (or 24 if 4 SAs) randomly from communities within each SA. Each SA had a similar population to the others in its district, but exact weights for each supervision area were not available. Using the definitions in chapter 1, the behaviour of respondents \mathbf{Y} can be partitioned into $\mathbf{Y} = \{\mathbf{Y}_1^*, \dots, \mathbf{Y}_{N_{sa}}^*\}$, where $\mathbf{Y}_j^* \subset \mathbf{B}_j^*$ is the behaviour of the respondents from one supervision area j , \mathbf{B}_j^* is the behaviour of the whole population in supervision area j , and N_{sa} is the total number of supervision areas in the data.

The SA sample size was selected so that when sub-district data (the SA) are aggregated, the resulting district-level coverage proportion estimates for key indicators are calculated with a 95% confidence interval not exceeding $\pm 10\%$. Communities were selected using Probability Proportionate to Size (PPS) sampling, wherein a comprehensive community population list supplied by each district was the sampling frame to select communities from which the individual samples are taken. PPS sampling ensures that sample communities are selected based upon their proportional representation of the entire population. Individual respondents were then randomly selected from the PPS-selected communities using random sampling techniques [56].

A total of 11,723 randomly selected mothers of children aged 0-11 months were interviewed. Each maternal questionnaire included demographic characteristics and various health-related behaviours. Respondents with missing responses were removed, leaving a total of 11,414 (97%) records with complete information. These data were integrated into a superset and analysed.

For geographic variables we made use of geospatial road and population data from 2009 [57] and 2010 Geographical Information System (GIS) locations of health centres. We calculated the number of health facilities per capita based on the number of health facilities with in-patient beds (Level III and above), since mothers are referred to these higher-level facilities for FBD. A wealth index based on household assets data and calculated in DHS 2011 [45] was used to calculate a mean wealth index for each district. Altitude data was obtained from the US Geological Survey [58].

2.3 Methods

2.3.1 Logistic Regression

In this chapter we will use the model equation (1.1) introduced in chapter 1. We assume that the timescales on which the social effect $f(\mathbf{B}(t))$ changes are much faster than a year, the timescale present in this data. We also assume that mothers are only influenced by others within their supervision area, so social effect for individual

i , $f_i(\mathbf{B}(t))$, is given by $f_i(\mathbf{B}(t)) = f_i(\mathbf{B}_{\text{SA}(i)}^*(t))$, where $\text{SA}(i)$ is the supervision area of individual i . Combined, these two assumptions lead to the following social effect term for each individual:

$$f_i(\mathbf{B}(t), t) = f_i(\mathbf{B}_{\text{SA}(i)}^*) = \alpha_{\text{SA}(i)}.$$

In addition to the social effect, there will also be unknown factors that lead to individuals in the same supervision area having correlated behaviours. It is not possible to separate these factors from the social effect, so we combine both into a general ‘supervision area effect’, ξ .

If we assume that the log-odds of taking part in a behaviour are a linear combination of the individual-level, environmental, and social influence effects, then we are left with the following statistical model:

$$\begin{aligned} \mathbf{Y} &\sim \text{Binomial}(\mathbf{p}) \\ \text{logit}(\mathbf{p}) &= \beta^T \mathbf{\Theta} + \gamma^T \mathbf{E} + \xi. \end{aligned}$$

This model takes the form of a random-intercept logistic regression. The parameters for this model, β_k and γ_k , can be interpreted as the increase in the log-odds of the outcome associated with a unit increase in the covariate X_k . The parameters therefore measure the strength of association between a covariate and the outcome. In this chapter we do not try to estimate the supervision area effect ξ , but we account for it using the clustered bootstrap method described below.

The logistic regression model can also include interactions between covariates by adding multiplicative terms to the linear predictor, like so:

$$\text{logit}(p_i) = \beta_0 1 + X_{i,1}\beta_1 + X_{i,2}\beta_2 + X_{i,1}X_{i,2}\beta_3.$$

Interaction terms in the model represent the possibility that the effect of one covariate depends on the value of another.

In this chapter we only consider behaviours that are binary, so $Y_i \in \{0, 1\}$. The log-likelihood function for binary logistic regression is given by:

$$l(\beta, \gamma | \mathbf{Y}, \mathbf{X}) \propto \sum_{i=1}^N (Y_i \log p_i + (1 - Y_i) \log(1 - p_i)). \quad (2.1)$$

This function is numerically maximized with respect to β to provide the maximum-likelihood estimate for the covariate coefficients. We calculated 95% con-

fidence intervals for the β and γ parameters in a robust way that takes into account the additional correlations due to the supervision area effect ξ , using the clustered bootstrap.

2.3.2 Clustered bootstrap

We calculated 95% confidence intervals for prevalence estimates in each surveyed district, and for logistic regression coefficients, using a clustered bootstrap [59]. The clustered bootstrap procedure is as follows [60]:

1. Original data consists g clusters, $G = (\mathbf{x}_1, \dots, \mathbf{x}_g)$
2. Repeat:
 - (a) Create a new data set G_i^* by sampling g clusters with replacement from G
 - (b) Estimate the parameter of interest, $\hat{\mu}_i$, based on this new sample
3. Calculate the $[2.5, 97.5]^{\text{th}}$ percentiles of $\hat{\mu}$. This corresponds to the 95% confidence interval for the estimate of population parameter μ .

When estimating district prevalence, the original clusters for district d are $\mathbf{Y}_d^\dagger = \{\mathbf{Y}_j^* \forall j \in D(d)\}$, where $D(d)$ are the supervision areas in district d . The parameter of interest is $\hat{\mu}_d = \sum_{Y_i \in \mathbf{Y}_d^\dagger} \mathbf{I}_{\{Y_i=1\}}$. When estimating logistic regression coefficients, the original clusters are $(\mathbf{Y}, \mathbf{X}) = \{(\mathbf{Y}_1^*, \mathbf{X}_1^*), \dots, (\mathbf{Y}_{N_{sa}}^*, \mathbf{X}_1^*)\}$, and the parameters of interest are β and γ . The logistic regression model is fitted during each iteration of the bootstrap, and the values for the β and γ parameters used for $\hat{\mu}_i$.

We used a clustered bootstrap for this model because a simple bootstrap method would ignore the supervision area effect by assuming that the data originate from independent random samples, leading to incorrect inference [61].

2.3.3 Model selection

In this chapter we chose interactions to include in the model using stepwise forward selection [62] based on the Akaike Information Criterion (AIC) [63]. The AIC is a quantity which describes the relative quality of a statistical model, finding a balance between the goodness-of-fit of the model and the number of parameters in the model. It is calculated as follows:

$$\text{AIC} = 2k - 2l^*,$$

where l^* is the maximum value of the log-likelihood function for the statistical model, in this case calculated in equation (2.1). Lower values of the AIC indicate a better model. During stepwise forward selection a base model with no interactions was compared with models including each plausible interaction term, and the interaction model with the lowest AIC was then selected. The process was then repeated with the new selected model as the base model, until no additional interaction term model has a lower AIC.

2.3.4 Model validation

We used an receiver operating characteristic (ROC) curve to summarize the predictive power of our model. The ROC curve is defined as a piecewise linear curve plotting the relationship between the true positive rate (the probability that a true outcome is correctly predicted to be true) and the false positive rate (the probability that a false outcome is predicted to be true) for different classification cutoffs [64]. The true positive and false positive rates for a classification cutoff λ is given by:

$$\begin{aligned}\text{TPR}(\lambda) &= \frac{1}{N} \sum_i^N \mathbf{I}_{\{Y_i=1\}} \mathbf{I}_{\{p_i \geq \lambda\}} \\ \text{FPR}(\lambda) &= \frac{1}{N} \sum_i^N \mathbf{I}_{\{Y_i=0\}} \mathbf{I}_{\{p_i \geq \lambda\}}.\end{aligned}$$

This is repeated for $\lambda_i = p_i \forall i$, where \mathbf{p} is the vector of probabilities predicted by the model.

The predictive power can be summarised by the area under the ROC curve (AUC):

$$\text{AUC} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{I}_{\{p_i^+ - p_j^- > 0\}},$$

where $\{p_1^+, \dots, p_m^+\}$ are the predicted probabilities for the m positive ($Y_i = 1$) datapoints and $\{p_1^-, \dots, p_n^-\}$ are the probabilities predicted for the n negative ($Y_i = 0$) datapoints. The AUC represents the probability that a randomly-selected positive datapoint will have a higher predicted probability than a randomly-selected negative datapoint. An AUC of 1.0 indicates a perfect prediction: all datapoints were correctly classified. An AUC of 0.5 indicates a random test, which allocates positive outcomes at random for half of the datapoints [65]. We calculate the AUC nonparametrically, using the trapezoidal rule to calculate the area under the ROC curve. Both the ROC curve and AUC were calculated using the scikit-learn package in the

Python programming language [66]. We calculated 95% confidence intervals for the AUC by using predictions from the models fitted during the clustered bootstrap process and then taking the [2.5, 97.5]th percentiles.

Although ROC curves provide a measure of within-sample predictive power, they do not measure the accuracy of the model when applied to data that was not in the original sample. We used two-fold Monte Carlo cross-validation [67] to estimate the out-of-sample prediction error for the logistic regression model. This process proceeds as follows:

1. Randomly split data into two halves: a training set and a test set.
2. Fit the model to the training set
3. Use this fitted model to predict outcomes for the test set
4. Store the difference between the true test set outcomes and those predicted by the model
5. Return to step 1.

Once this process has been repeated a sufficient number of times the root-mean-square error between the true and predicted outcomes can be calculated as follows:

$$\text{MSE} = \sqrt{\sum_{i=1}^{N^*} (Y_i^* - \mu_i^*)^2},$$

where \mathbf{Y}^* are the outcome variables of the test dataset and μ^* are the model predictions for the test dataset. The MSE provides an estimate of the average prediction error for unseen data.

2.4 Analysis

The aim of our study is to identify the simultaneous correlates of FBD, in order to provide a framework for classifying districts in different priority levels. Our analysis consists of 3 phases: FBD mapping, model construction, and prediction of priority groups. All analysis was done using the statistical software R version 2.15 [68], and we used the R-package ‘maptools’ [69] to construct the maps.

2.4.1 FBD mapping

We classified mothers as giving birth either at home or in a health facility and plotted on a map the percentage of mothers with FBD for each district surveyed.

We constructed the maps by colouring the districts in a map of Uganda based on the percentage of mothers surveyed in each district who had given birth in a health facility. One map was produced for each cluster of survey years: 2003-04, 2006, 2009-10, and 2011. Survey years were combined so that a similar number of districts were surveyed in each cluster.

One complication for producing these maps was that the number of districts in Uganda increased between 2004 and 2006, and between 2006 and 2009. The indicator for each district was calculated based on the boundaries used at the time of the survey, but to aid comparison between survey years we plotted the indicator for all survey years on the post-2009 district map. For the older districts that were divided in subsequent years, we applied the value of the indicator for the old district to each of the districts in 2009 that it had been divided into. For example, the 2006 district in the Western region of Uganda called ‘Bushenyi’ was subsequently divided into five districts. When plotting a map of the 2006 survey results, we assigned the indicator value of the 2006 ‘Bushenyi’ district to all of the five districts that it became in the 2009 map. We calculated 95% confidence intervals for the prevalence in each district using clustered bootstrapping, as described in section 2.3.2.

2.4.2 Model construction

Using all 2003-2011 data, we fitted a logistic regression model to investigate what and how factors were associated with FBD simultaneously. The individual-level factors included in the model were age, education and the year that the mother was surveyed. Marital status was removed from the model since it did not have a statistically significant effect. The number of mothers in each group of survey years with each level of educations, and with each age group (ages are grouped for conciseness), are given in Table 2.1.

We also included district-level covariates: each mother was assigned a value for the number of health facilities per capita, population density, road density, wealth index, and mean and standard deviation of the altitude of her district. The mean value and range of these covariates across all districts is given in Table 2.2.

These covariates showed significant nonlinearity. We reduced this by applying a base-2 logarithm transformation and checking if the AIC for the fitted model was improved. Four covariates—health facilities per capita, population density, road density, and standard deviation of altitude—were transformed in this way (shown in Figure 2.1). Two covariates—mean altitude and wealth index—were not improved by the transformation and were thus left untransformed (Figure 2.2).

Correlations between covariates are shown in Figure 2.3, and correlation co-

Characteristic	Levels	Sample Size per Survey Year			
		2003-04	2006	2009-10	2011
Education	None	623	180	1542	2506
	Primary	1929	727	729	1110
	Secondary	507	239	319	644
	Post-secondary	63	33	83	180
	Total	3122	1179	2673	4440
Age	<20	724	265	530	843
	20-30	1765	645	1507	2522
	30-40	580	251	567	965
	>40	53	18	69	110
	Total	3122	1179	2673	4440

Table 2.1: Characteristics of Ugandan mothers of children 0-11 months for each surveyed year.

Covariate	Mean value across all districts	Range
Wealth index	0.0	(-2.0 to 3.6)
Health centres per 100,000 capita	4.9	(0.2 to 36.7)
Road density (metres per km ²)	133.8	(0.0 to 358.6)
Standard deviation of altitude (metres)	146	(10 to 956)
Mean altitude (metres)	1221	(701 to 2428)
Population density (per km ²)	279	(4 to 8647)

Table 2.2: Characteristics of district-level covariates

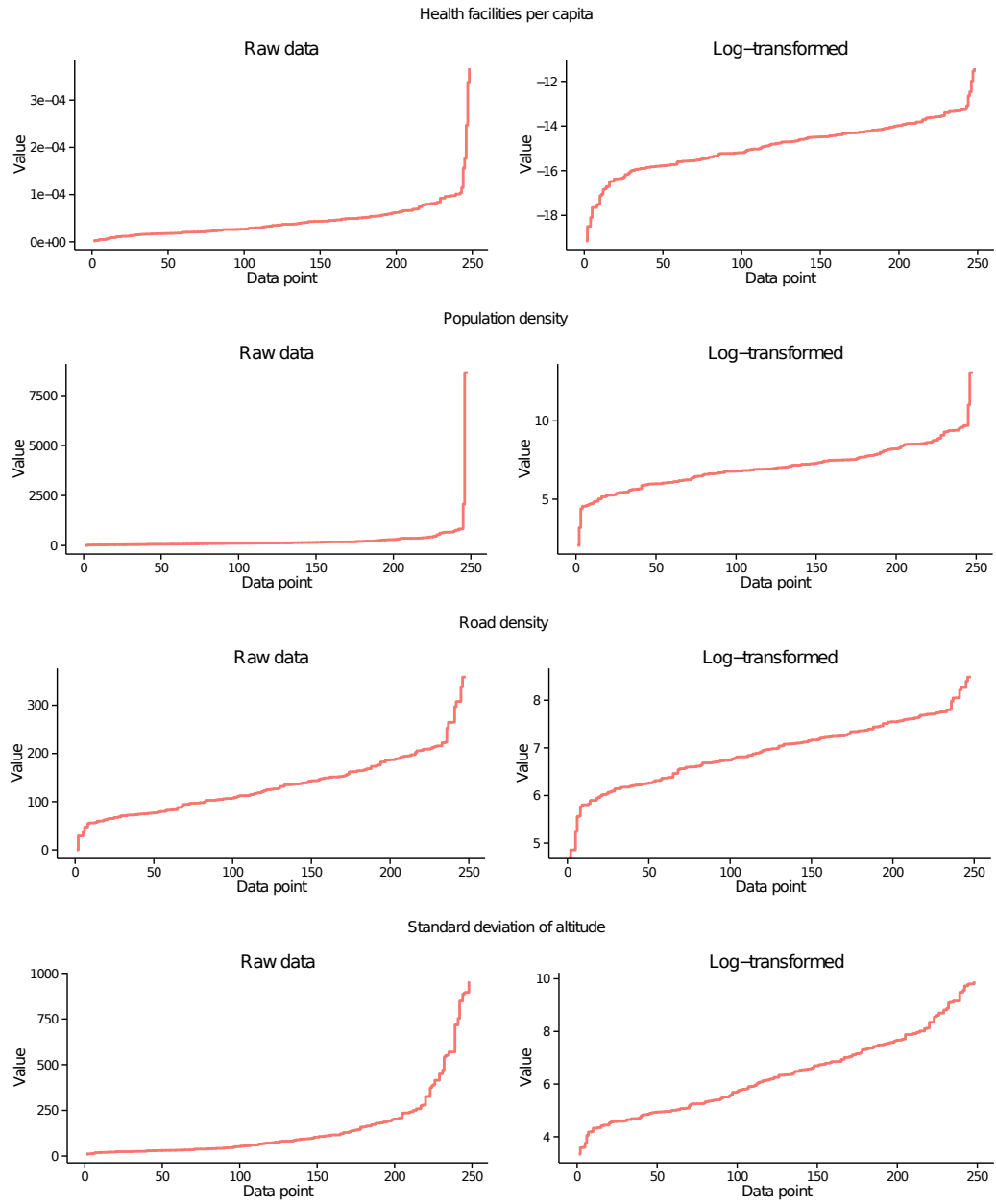


Figure 2.1: Covariates with significant nonlinearity showing untransformed (left) and \log_2 transformed (right).

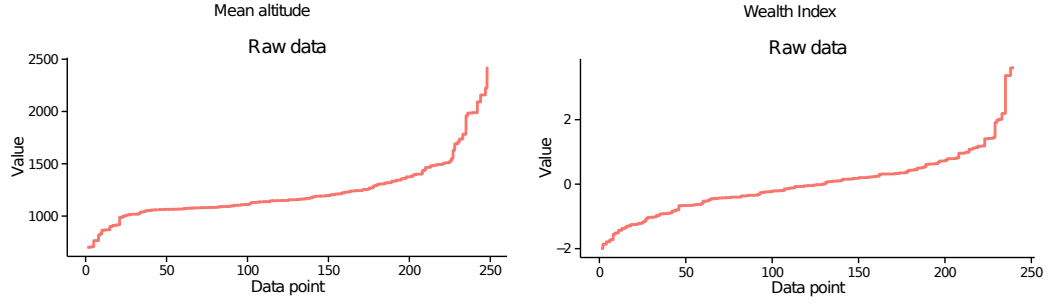


Figure 2.2: Covariates that were not improved by the \log_2 transformation.

	hcpc	roaddens	popdens	wealth	topo_std	topo_mean
hcpc	-	0.33	0.63	0.32	0.16	0.36
roaddens		-	0.39	0.24	-0.29	-0.13
popdens			-	0.21	-0.22	-0.01
wealth				-	0.05	0.21
topo_std					-	0.70
topo_mean						-

Table 2.3: Spearman correlation coefficients between pairs of district-level covariates. Population density and health centres per capita, and mean altitude and standard deviation of altitude, showed the highest correlation.

efficients are given in Table 2.3. The biggest correlations are between the mean altitude in a district and the standard deviation of altitude, and between health centres per capita and population density.

All covariates were included as continuous variables, except for education which was categorical. We used forward selection [62] based on the Akaike Information Criterion (AIC) [63] to include significant interaction terms between the covariates. Our selected model gave an estimate of the odds ratio (OR) for FBD for each covariate. For our model, the OR for a covariate is the ratio between the odds of FBD for two mothers, both of whom, for the covariate being examined, have all other covariates set to their mean values. If the covariate is categorical, such as education level, then the ratio is between each level and the lowest level which, in this example, is ‘no formal education’. If a base-2-log-transformed covariate was used in the model, then the ratio is between the odds calculated for the covariate and double the covariate. For the other continuous covariates, the ratio is between the odds calculated for the covariate and the covariate plus a unit increase. The OR therefore provides an estimate for how strongly each covariate is associated with the odds of FBD.

We validated our selected model by comparing it to a null spatial model, for

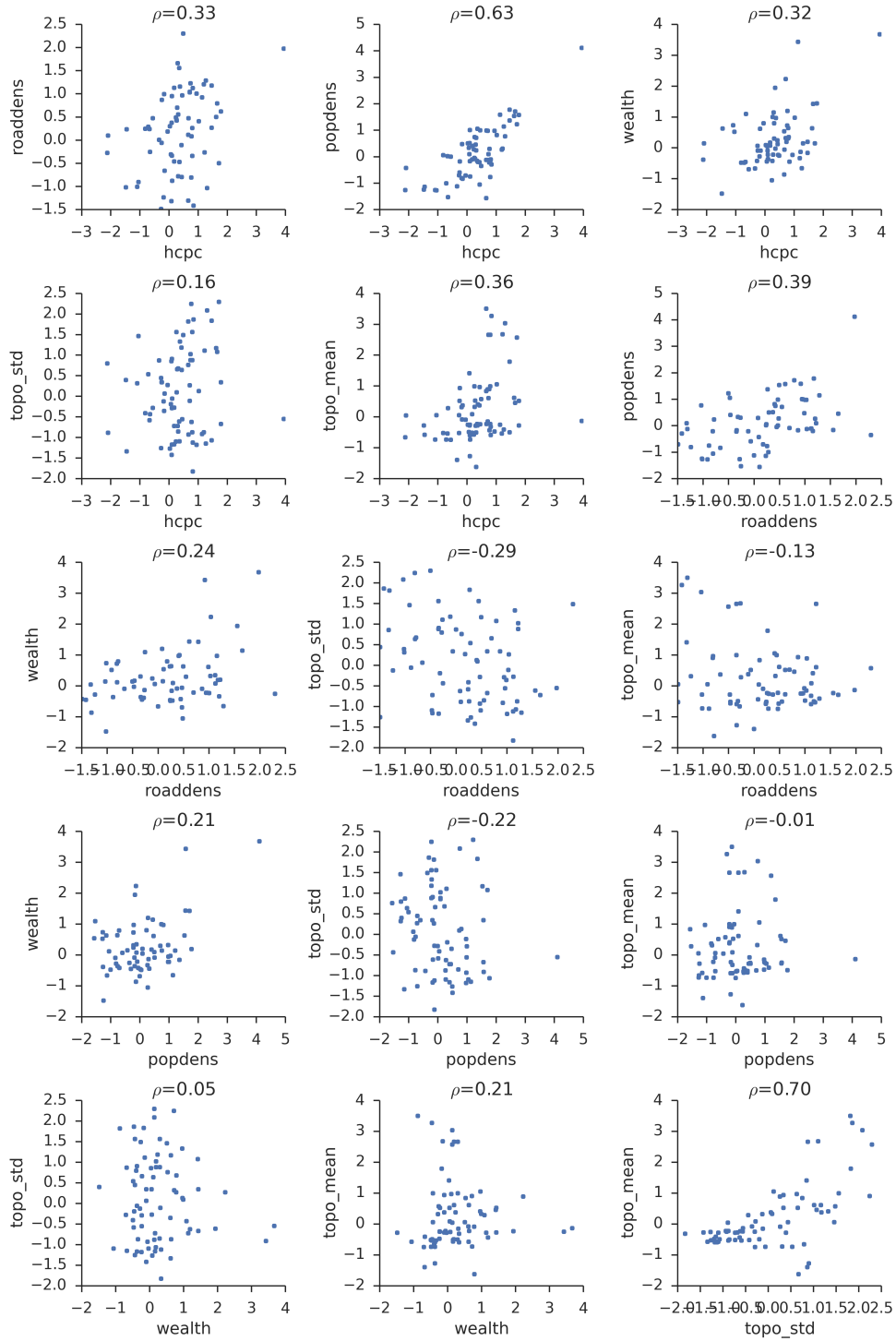


Figure 2.3: Scatter plots between pairs of district-level covariates to show the relationship between each pair.

which the probability of FBD for a mother is predicted to be the average value for her district. This null model represents a situation where each district is totally separate, and the differences between the indicators in each district are due to unknown spatial factors. If this null model has a better AIC then the unknown spatial factors are more important than the covariates we identified. As a second stage of model validation we looked for spatial correlation in the residuals of our selected model: we plotted the difference between the residual for each pair of districts against the distance between the two districts. Significant correlation between the distance between two districts and the difference in their residuals would imply that some spatial patterns in the indicator have not been captured by the model. As a third stage of model validation we constructed an ROC curve and calculated the area under it as a test of the model's accuracy. Finally, as a fourth stage, we used two-fold Monte Carlo cross-validation [67] to estimate the prediction error for unseen data; the model was repeatedly fitted to a randomly chosen half of the 2003-2011 data and then used to predict the FBD values of the other half. For each iteration, we calculated the squared error between the observed and predicted district-level FBD indicator, and took the mean over all 1000 iterations. The resulting mean squared error measured how well the model predicted the indicator for unseen data.

2.4.3 Prediction of priority groups

Finally, we used the model to classify unsurveyed districts into 'priority' groups to flag districts predicted to have particularly low indicator values. Our selected unsurveyed districts are those which were actually surveyed in 2012. Since we do not know the distribution of age and education in these unsurveyed districts, we decided to predict an upper and lower limit for the indicator in each district rather than an average value. We chose the values for age and education most strongly associated with FBD and then the values with the strongest negative association, and then we used the model to predict the probability of FBD for a mother with her age and education set to these values and the survey year set to 2012. To obtain a reasonable upper limit for the indicator in each district, we applied the model to the strongly associated age and education values chosen above. To account for any uncertainty in the model parameters we took the upper part of 95% CI obtained from the model as a conservative estimate of the upper limit. For the lower limit, the same procedure was carried out with the negatively associated values and taking the lower part of the 95% CI. The lower and upper limit form the predicted range of each district. The priority groups were assigned on the basis of these limits. The low-priority group, defined as districts with lower limits between 50-100% FBD

and upper limits between 80-100% FBD, contained districts that were likely to have high indicator values. The mild-priority group, defined as districts with lower limits between 0-30% and upper limits between 60-80%, contained districts likely to have fairly low indicator values. The high-priority group, defined as districts with lower limits between 0-30% and upper limits between 30-60%, contained districts likely to have very low indicator values. All other scenarios were classified as an unclear-priority group. We then validated the projections by checking that the 2012 values lay within their predicted ranges.

2.5 Results

We plotted spatially the percentage of mothers of children aged 0-11 months with FBD for different time points (Figure 2.6). During the time-period over which the data were collected, some districts subdivided; so we plotted the indicator on the 2011 district map to aid comparison across time.

For early survey years (2003 to 2006) most districts have <60% of mothers of infants with FBD, except around the capital, Kampala. For later years (2009 to 2011), >60% of mothers had FBD in most surveyed districts. Although much of northern Uganda was not sampled in the more recent time-points, the one district sampled in both 2003-04 and 2011 increased from 20-30% to 70-80%. This progress is not replicated everywhere. Surveys conducted in eastern Uganda in later years, for example, still reported <50% FBD.

We give the distribution of ages and educational categories for mothers with complete data in Table 2.1. Table 2.4 shows the logistic regression coefficients, and Table 2.5 shows the odds ratio for each covariate in the model for an ‘average’ mother. An ‘average’ mother has secondary-level education, is 25-years old, was surveyed in 2007, and lives in a district with average values of all district-level covariates. The results show that the odds of FBD were significantly lower for each one-year increase in maternal age (OR=0.98 [0.97, 0.99]), whereas they were significantly increased for each additional level of maternal education (Primary: OR=1.59 [1.42, 1.78], Secondary: OR=3.37 [2.88, 3.94], Post-secondary: OR=10.4 [6.28, 18.1]). The odds of FBD were significantly greater in districts with double the health facility density (OR=1.12 [1.02, 1.23]), or road density (OR=1.13 [1, 1.26]), or in districts with a unit increase in the wealth index (OR=1.38 [1.24, 1.53]). Living in the capital, Kampala, was strongly associated with FBD (OR=8.38 [2.24, 23]). Districts with double the standard deviation of altitude (a proxy for the roughness and difficulty of the terrain) were strongly associated with a decrease in the odds of FBD (OR=0.89

[0.84, 0.94]). Finally, there was a significant time trend: mothers surveyed in later years were more likely to have FBDs (OR=1.08 [1.04, 1.13]).

Our model had a significantly lower AIC ($AIC = 13383$) than the null spatial model ($AIC = 13690$). It also performed better under cross-validation: the indicator predicted for each district based on unseen data was within 20% of the observed indicator for 95% of the districts, whereas the null spatial model predicted the indicator of only 75% of districts within 20% of the observed value. Our model is therefore a significantly better explanation of the observed prevalence than a model with separate unknown spatial covariates for each district.

Figure 2.4 shows the spatial correlation of the indicator and of the residuals from our model. The indicator shows significant correlation of the indicator differences for pairs of districts against the distance between them, whereas this correlation is not present in the model residuals. This result suggests that the model is capturing much of the spatial patterns present in the data. The area under the ROC curve, shown in Figure 2.5, was 0.71 [0.70, 0.72], which is significantly higher than that expected for a random classifier, and is similar to good epidemiological models of disease [70] [71].

Using the 2003-2011 fitted logistic regression model, we predicted for all Ugandan districts the reasonable lower and upper limits for FBD during 2012, without requiring LQAS data for that year. We categorise districts into different priority levels using these predicted values. We validate these predicted priority levels by comparing them with indicator values calculated from a subsequent 2012 LQAS survey covering 61 districts (Figure 2.7). The priority map identifies many north-eastern districts as being mild or high priority – classifications that agree with the low values for these districts seen in the observed data. In addition, all but four (6.6%) of the observed indicator values fall within the confidence range predicted by the model.

2.6 Discussion

Our results show the percentage of mothers with FBD increasing over time, but varying among surveyed districts (Figure 2.6). This variation is strongly correlated with geographic and demographic factors. Initiatives meant to increase the uptake of services, including maternal services, have been introduced in Uganda in the past several years [72]. Although this analysis cannot prove causation, the trend over time displays a progressively greater uptake of maternity services in most areas. Despite this pattern, some districts display low FBD, particularly in northeastern Uganda. The logistic regression model provides possible explanations; specifically,

Covariate	Coefficient and 95% CIs §
(Intercept)	0.107 [-0.303,0.503]
Age	-0.0285 [-0.0406,-0.0162] *
Education (primary)	0.408 [0.215,0.607] *
Education (secondary)	1.42 [1.14,1.7] *
Education (post)	2.72 [1.98,3.94] *
Health facility density †	-0.036 [-0.314,0.254]
Road density †	0.0824 [0.002,0.156] *
Population density †	0.297 [0.0743,0.525] *
Living in Kampala	1.9 [0.808,3.14] *
District wealth index	0.307 [0.208,0.415] *
Standard deviation of altitude †	-0.176 [-0.26,-0.0921] *
Mean altitude	0.325 [0.072,0.566] *
Year of survey	0.0777 [0.0164,0.142] *
Interaction Term	
Standard deviation of altitude : Mean altitude	-0.191 [-0.28,-0.0998] *
Health facility density : Year	0.0538 [0.0134,0.0956] *
Road density : Population density	-0.253 [-0.372,-0.145] *
Mean altitude : Year	-0.0696 [-0.105,-0.0329] *
Population density : Year	-0.0854 [-0.119,-0.0506] *
Health facility density : Mean altitude	-0.0962 [-0.333,0.122]
Education (primary) : Year	0.0133 [-0.0183,0.0441]
Education (secondary) : Year	-0.052 [-0.0977,-0.00668] *
Education (post) : Year	-0.119 [-0.281,-0.00558] *
Age : Year	0.00213 [0.000132,0.0041] *
Health facility density : Mean altitude : Year	0.0446 [0.0149,0.0765] *

* A 95%-significant positive or negative effect

† Results for a doubling of this variable, rather than a unit increase

Table 2.4: The coefficient for each covariate and interaction term in this predictive model of delivery in a health facility in Uganda. Confidence intervals were calculated using clustered bootstrapping with 1,000 iterations.

Covariate	Odds ratio and 95% CIs
Age	0.98 [0.974,0.987] *
Education (primary)	1.59 [1.42,1.78] *
Education (secondary)	3.37 [2.88,3.94] *
Education (post)	10.4 [6.28,18.1] *
Health facility density †	1.12 [1.02,1.23] *
Road density †	1.13 [1,1.26] *
Population density †	0.97 [0.892,1.06]
Living in Kampala	8.38 [2.24,23] *
District wealth index	1.38 [1.24,1.53] *
Standard deviation of altitude †	0.89 [0.842,0.941] *
Mean altitude	1 [0.9997,1.001]
Year of survey	1.08 [1.04,1.13] *

* A 95%-significant positive or negative effect

† Results for a doubling of this variable, rather than a unit increase

Table 2.5: The odds ratio for each term included in the model. The odds ratio is given for a unit increase in the (possibly transformed) covariate, for a mother aged 25 and with secondary-level education, surveyed in 2007, and all district-level covariates set to their average. Confidence intervals were calculated using clustered bootstrapping with 1,000 iterations.

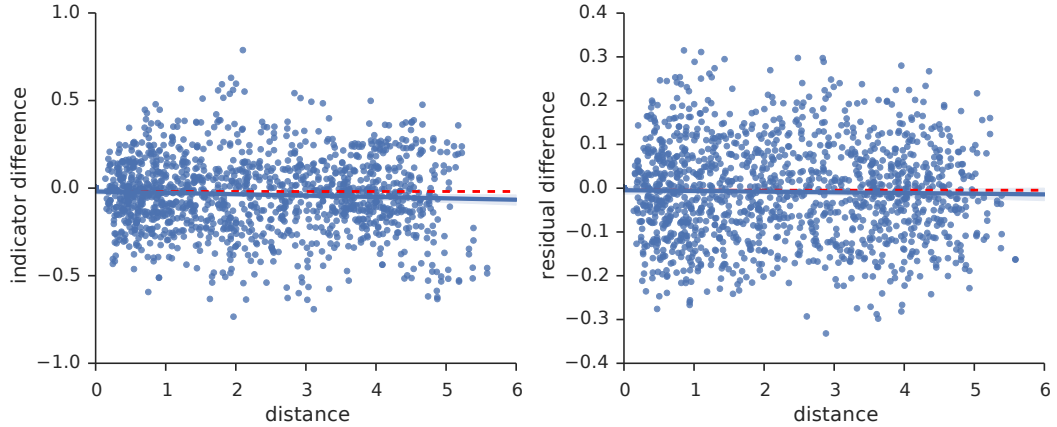


Figure 2.4: Difference in observed indicator (left) and model residuals (right) against distance for all survey years. Distance between districts is significantly negatively correlated ($p = 0.026$) with the difference in indicator values, but once the model prediction has been subtracted the residuals are no longer significantly correlated ($p = 0.386$).

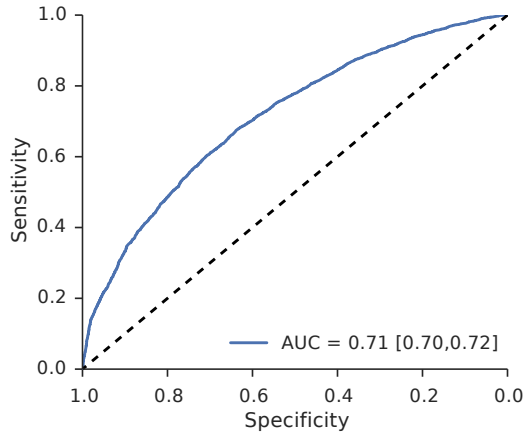


Figure 2.5: The ROC curve for this model, showing the true positive rate against false positive rate for various thresholds. The area under this curve is 0.71 [0.70, 0.72], where a perfect classifier would have an AUC of 1.0 and a random classifier would have an AUC of 0.5.

low uptake is associated with low health facility density, low road density, mountainous terrain, and lack of geographical access to health facilities due to few roads. In other countries, these same factors have been shown to have the same impact on uptake of maternal services [28, 31, 32, 73, 74, 75, 76]. Furthermore, previous research in Uganda shows that “difficult-to-access” areas also suffer acute staffing shortages, high rates of absenteeism, and poor quality of care [77, 78], potentially reducing demand by women in labour.

We also found that age, level of education and wealth status also influence FBD. Younger women and those with higher levels of education are more likely to practice FBD. Our findings are consistent with the 2011 Uganda DHS, which reports that older mothers are less likely to give birth in a health facility [45]. However, this result is inconsistent with a meta-analysis of socio-geographic factors in numerous countries which found age to have no statistical significance in determining FBD; it did find, however, that high parity is negatively associated with FBD [79, 29, 46]. Studies controlling for parity have found either no effect of maternal age on FBD or that increased age was positively related to increased use of delivery services [75]. We were unable to control for parity in our study. As parity is often linked to maternal age, it may be the influence of parity, rather than age, which we have vicariously detected. In our study, increased maternal education was positively associated with FBD, a finding consistent with other studies [73, 74, 79, 80, 81, 45]. We also determined that with increasing wealth index in her district a mother was more likely to use FBD. This finding is consistent with the results of numerous other

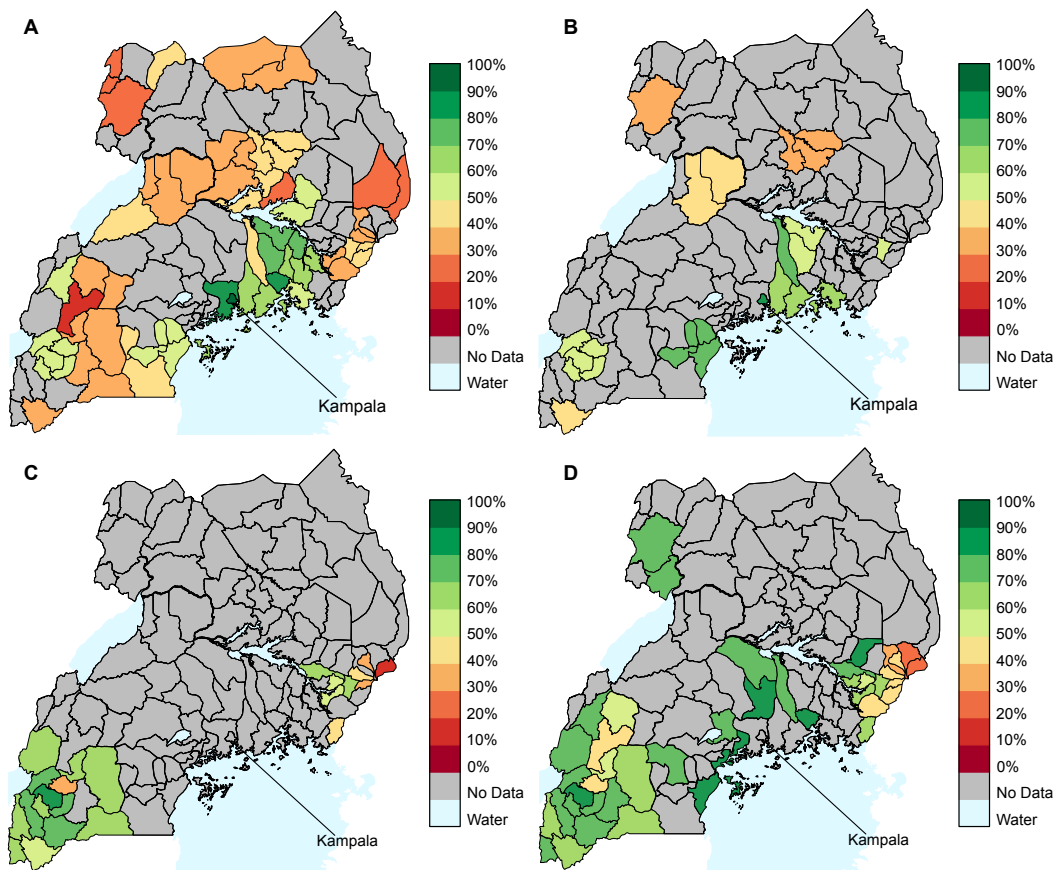


Figure 2.6: Percentage of mothers of children aged 0-11 months that gave birth in a health facility for A) 2003-04, B) 2006, C) 2009-10, D) 2011. 95% confidence intervals for the indicator are $\pm 11\%$ or lower. Data for 2003 and 2004, and for 2009 and 2010 have been combined for these maps due to the small number of districts surveyed in 2004 and 2010. In the rest of the analysis they are separated.

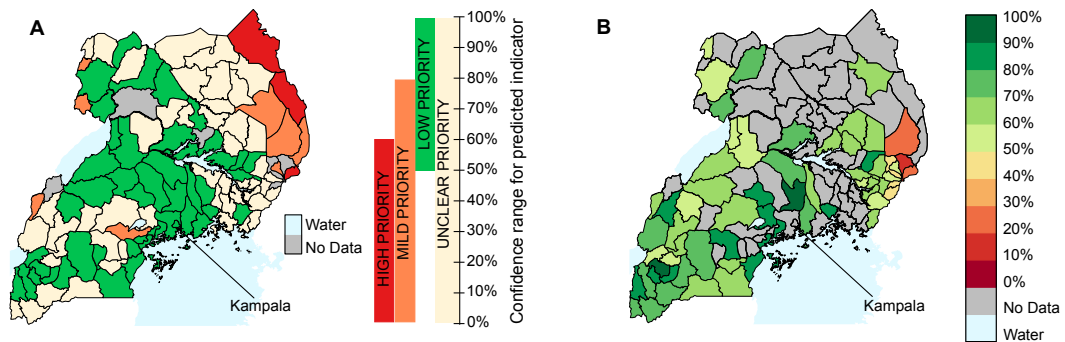


Figure 2.7: (A) Priority Map for districts in 2012. Districts are assigned high, mild, low, or unclear priority based on the confidence interval predicted by the model. Kampala, for example, has a predicted indicator confidence interval between 50% and 100%, and is therefore assigned a low priority. (B) Indicator for 2012 as observed in a subsequent LQAS survey. The model was not fitted using this data, and so this map provides an independent test of the predicted confidence interval. All surveyed districts in 2012 have indicator values within the predicted confidence interval, and districts in Eastern Uganda that were predicted to be mild or high priority are observed to have very low values for the indicator.

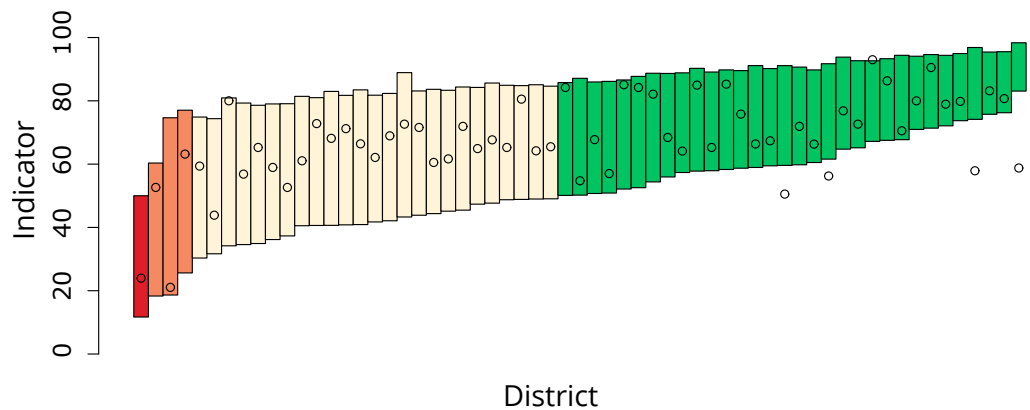


Figure 2.8: Predicted 95% confidence range of district indicators for 2012, compared with the observed value. Priority levels are given by the same colour scheme as in Figure 2.7. Only 4 out of 61 (6.6%) of the observed values were outside the predicted range.

studies reporting wealth and economic access to health care as facilitators behind FBD [27, 28, 31].

This data cannot tell us the mechanism by which the significant covariates we identified in this chapter influence FBD, or indeed if FBD and the significant covariates are caused by some unidentified extra factor. Some covariates are plausible as direct causes: more local health facilities, more roads, and easier terrain, for example. Others, such as education level and wealth index, could be proxies for the ease of accessing facilities (by paying for transport and healthcare fees), but it could also be the case that wealthy and highly educated individuals influence each other and establish a social norm of FBD. The argument could also apply to age: given the improvement in FBD seen over time, the norm among older mothers may have become fixed in previous years when fewer mothers gave birth in health facilities. Although difficult to collect, data on FBD that includes information a mothers' social networks may help to investigate these hypotheses.

The difference between national and subnational results indicated gaps in equitable access to health services. In 2011, for example, the Eastern region district prevalence ranged from 23% in Bukwa district to 82% in Kumi district. Although these two districts contribute to regional and national prevalence, the district prevalence is indicative of the inequitable occurrence of FBD in the districts. A study in Ghana found similar variations, further highlighting the importance of detecting subregional variation when planning health programs and allocating resources to decrease the gaps [82]. To achieve equitable access to services, subregional variation must be detected and addressed [83, 55].

Our study also examined an additional practical use of logistic regression, namely, assigning predicted priorities to districts based on the lowest expected value for FBD. By using the model to construct the predicted range of FBD in each district, we can flag districts, prior to a survey, most likely to need intervention. Such an approach could help lessen the gaps and inequities in maternal health care and help Uganda identify health system changes needed to decrease both maternal and child mortality [55]. This feature of logistic regression suggests it can be used as an early warning system to detect priority districts in need of special attention.

2.7 Summary

Maternal mortality is still very high in certain countries, despite attempts to improve maternal health services. One explanation for this is the lack of health service uptake by mothers in developing countries. In this chapter, we have shown that

LQAS surveys can provide data on the spatial prevalence of health behaviours at regular timepoints. Modelling the association of geographical variables with these health behaviours shows that ease of access is positively correlated with uptake, and allows us to accurately predict poorly performing regions. Through mapping and modelling the areas where the prevalence of positive health behaviours is low, policymakers will have more information available to concentrate resources where they are needed most.

Chapter 3

The role of the health system

3.1 Introduction

In chapter 2 we investigated how individual-level and geographical factors influenced health behaviours. If the goal is to change health behaviours, these types of factors are likely to be hard for policymakers to influence; economic and educational changes, for example, are likely to take many years before influencing behaviour.

Another approach is to study the health systems that individuals interact with, as called for in a paper by Travis et al. [84] which suggests that improving health systems are the only way to meet the Millennium Development Goals. They argue that health systems tend to be neglected by stakeholders, who are drawn more to specific disease interventions such as providing malaria nets and vaccinations. They identify some of the potential problems within poorly performing health systems: the cost of health care, physical inaccessibility, poor quality of care, and the perception of care quality. Removing these barriers should lead to stronger health systems and healthier populations.

3.1.1 Seeking appropriate healthcare

In this chapter, we will investigate health systems as they relate to care seeking for children aged 1-59 months. Roughly 3.7 million of these children die each year [6]: a third are caused by lower respiratory infections and malaria, and diarrhoeal disease accounts for a significant fraction of the rest. The suggested treatments are Oral Rehydration Salts for diarrhoea [1], an antibiotic such as cotrimoxazole, amoxicillin or ampicillin for pneumonia [2], Artemisinin-based Combination Therapies (ACT) for uncomplicated *P. falciparum* malaria [3], and Quinine for complicated or severe malaria [3]. If drugs are available then these diseases are treatable at home [85], but

often parents do not consult with trained professionals to assess the severity of the disease and to ensure the correct medication is prescribed [86]. Parents who do not seek appropriate help often use incorrect dosages or misdiagnose [87]. Underdosing can lead to the spread of drug-resistant *P. falciparum* malaria [88] and antibiotic-resistant respiratory and diarrhoeal diseases [89]. A study by Rutebemberwa et al. in Uganda in 2009 suggests that many parents seek initial treatment from drug shops [87], which often have severe breaches of regulations on the storage and labelling of drugs [90] and provide inappropriate treatment [91, 92].

Community health workers or health facilities are considered to be appropriate places to seek health advice by medical professionals. If parents seek help at a health facility they are not guaranteed to receive correct treatment advice [93], but policymakers have more direct control over the quality of health facilities through staff training, which can lead to better prescription practices [94, 95]. They also have access to Rapid Diagnostic Tests for malaria, which are starting to be recommended over presumptive treatment [96]. There is evidence that seeking treatment at health facilities leads to correct usage of anti-malarials [97], and seeking appropriate help is associated with positive outcome in the case of diarrhoea [98]. Increasing the fraction of individuals who seek initial help at a health facility is therefore likely to improve diagnosis and prescription, and to reduce mortality.

3.1.2 Barriers to access

Previous studies have identified many potential health system barriers to seeking appropriate care. Greater distances to a health facility are strongly associated with treatment at home [99, 100, 101, 102, 103], as are greater costs [104, 105, 106, 101, 102].

When interviewed, parents often identified low health facility quality as a psychological barrier [107, 108, 101], and the lack of availability of drugs [109, 108]. Some studies have shown evidence that higher quality health facilities are more likely to be used [110, 111, 102], where quality is measured either by services offered or perceived quality. Work by Mugisha et al [112] identified a difference between the factors governing patient initiation, or where patients initially seek treatment, and patient retention, or where they seek treatment for subsequent episodes of an illness. Initiation was mainly associated with lower barriers to access, such as cost, whereas retention was associated with higher perceived quality of care. Parents often perceive drug shops and private clinics to be closer, cheaper, and to allow payment by credit, compared with government- or NGO-run facilities [113, 114], despite evidence that objective quality is often worse [115, 116]. It is clear that this difference between

perceived and objective quality, and the link with health facility usage, requires further study.

The studies related to quality mentioned above tend to use descriptive statistics rather than predictive models, and most interview households without linking them to specific health facilities. Akin and Hutchinson [102] did consider both households and the health facilities they visited, and investigated the factors affecting the bypassing health facilities. They found that health facilities with lower quality ratings, higher prices, and more beds were more likely to be bypassed in favour of other facilities, and that closer health facilities were less likely to be bypassed. This study did not address the question of which individuals visit a health facility at all, however, which they identify as being a more complex statistical problem due to the difficulty in linking home-treating individuals with the health facilities they could have visited.

3.1.3 Interventions to improve access

Various interventions have been proposed to lift the major potential barriers identified above, but one of the most comprehensive approaches is Results-based Funding (RBF). RBF is new funding mechanism for health facilities in which facilities are given funds based on meeting certain targets. Traditionally, funding is centrally administered and centrally funded. Meessen et al. [117] claim that this leads to a lack of accountability, does not meet user needs and demands, and may have been partly responsible for the lack of progress made on some Millennium Development Goals. Instead, RBF aims to give local health facilities more direct control of administration and community health by paying a facility for each ‘health unit’, such as a fully immunized child, it produces [118]. Payments are often dependent on quality, as measured using a checklist [117]. RBF is controversial, with some claiming that it is a donor fad that encourages side effects such as false reporting, is inefficient, and is not generalizable to every developing country [119]. There is some evidence that RBF improves healthcare quality and usage, particularly the cluster-randomised trial undertaken by Basinga et al. [111] in Rwanda, but more evidence is required before the debate on RBF can be settled [119, 120].

3.1.4 Evaluating complex interventions

Complex public health interventions are hard to evaluate experimentally, even using Randomised Controlled Trials (RCTs) [121, 122, 123]. Currently, however, this type of trial provides the strongest evidence about interventions. Previous studies of RBF

were observational, with no control group without RBF to compare against. Only one study, by Basinga et al. in 2011 [111], used the RCT method to compare the improvement in quality with and without RBF. This study took into account the clustered nature of the data by calculating robust standard errors in a frequentist framework.

3.1.5 Summary

The prevalence of appropriate treatment-seeking behaviour is still low in many developing countries, often leading to incorrect diagnosis, inappropriate drug dosing, and drug-resistant disease. Policy-makers can have direct impact on this behaviour through improvements to the health system, by reducing barriers to access.

One barrier that has not been fully investigated is the quality of health facilities; both subjective, as perceived by users, and objective, as measured through quality assessments. In particular, a caregiver's seeking of appropriate treatment has not previously been linked to the quality of health facilities in their local area. This can have important implications for choosing number and location of health facilities: will providing many health facilities with low quality improve access, or is it better to provide fewer with high quality?

Policy proposals to improve access must be evaluated, and currently one of the best ways to achieve this is through a Randomised Controlled Trial. These are often analysed in a frequentist framework, which can be difficult to extend to complex study designs. In this chapter we concentrate on one particular policy proposal, Results Based Financing, which is controversial and is lacking in strong evidence. We will use a Bayesian framework to analyse data from an RBF trial in Uganda and compare it to a similar trial in Rwanda that used frequentist methods.

3.2 Methods

3.2.1 Approximate string matching

Throughout this chapter it was necessary to link names given in interview data to names in GIS shapefiles. Due to the lack of consistent local administrative information, the names of geographical regions and health facilities often did not exactly match names provided in interview data. Some examples of the types of discrepancies that occurred are:

- Abbreviations, such as 'TC' instead of 'Town Centre'

- Additional words added to the name, such as ‘GOMA DIVISION’ instead of ‘GOMA’
- Spelling differences, for example ‘KANSHENSHERO’ rather than ‘KASHEN-SHERO’
- Name changes which are not reflected in the GIS shapefile.

This last source of discrepancy could not be corrected without significant data collection on subcounty name changes, but the discrepancy due to the other sources was reduced using approximate string matching. Rather than remove data where the name provided in the interview does not match names in the GIS shapefile exactly, we instead compute a similarity score between each name provided by interview to names in the GIS. Two names are considered ‘matched’ if their similarity score is above a threshold. If a name from interview is similar to more than one GIS name then the GIS name with the highest similarity is chosen. If no GIS names are sufficiently similar then the row is removed from the data.

The similarity between two names was calculated using the token set ratio of the two words. The token set ratio is calculated by first splitting each name into words, removing duplicated words in each name, sorting the words, and finding the intersection of words between the two names. This forms a sorted set of words common to both strings. Two further test strings are created by combining the common words with those that are unique to each of the two names. The similarity score is then the maximum similarity ratio between each of the test strings and set of common words, where similarity ratio is calculated by combining the total length of both strings, T , and the total number of matches, M :

$$\text{ratio}(s_0, s_1) = \frac{2M}{T}.$$

An example helps to illustrate this process:

1. The ‘names’ to be compared are “ABC DEF GHI” and “ABC XYZ DEF”.
2. The sorted common unique set is s_0 = “ABC DEF” (see Figure 3.1)
3. The two further test strings are s_1 = “ABC DEF GHI” and s_2 = “ABC DEF XYZ”
4. The token set ratio is given by $\max(\text{ratio}(s_0, s_1), \text{ratio}(s_0, s_2), \text{ratio}(s_1, s_2))$, which in this case is 78%.

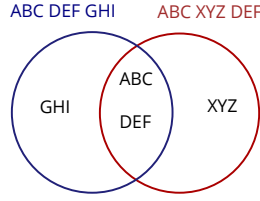


Figure 3.1: Overlap of the token sets for two example ‘names’: “ABC DEF GHI” and “ABC XYZ DEF”. These two names share the tokens “ABC” and “DEF” as a common unique set.

5. Assuming a threshold of 95%, these two ‘names’ would not be considered as being matched.

3.2.2 Hierarchical models

In chapter 2 we have treated data as applying to the same level of the general behavioural model. In this chapter, however, data will be collected at multiple levels: for example, a study which collected data about individuals and the country in which they live. Flat models such as standard logistic regression ignore this extra information about the structure of the data and are thus prone to either underfitting or overfitting. Hierarchical models, by contrast, take this known structure into account, and can therefore lead to better inferences about the data [124]

Hierarchical models are often constructed in a Bayesian framework, in which the parameters θ in the model are random variables and have pre-specified ‘prior’ distributions $p(\theta)$ which represent prior knowledge about their value before data is observed [124]. The model is constructed by specifying the distribution for each parameter conditional on other parameters, and by specifying the likelihood of observed data conditional on the parameters, $p(y|\theta)$ [124]. The priors and likelihood are combined to form a ‘posterior’ distribution for the parameters, which represents the joint probability density function for the value of the parameters [124]:

$$p(\theta|y) = \frac{p(\theta)p(y|\theta)}{\int p(\theta)p(y|\theta)d\theta}. \quad (3.1)$$

Choosing appropriate prior distributions is an important part of Bayesian analysis: priors that are too vague can lead to numerical issues, whereas assuming strong prior data can bias the inference if the prior data is not correct [124]. A general approach to choosing prior distributions when there is little prior knowledge is to use weakly-informative priors, which concentrate parameters around reasonable val-

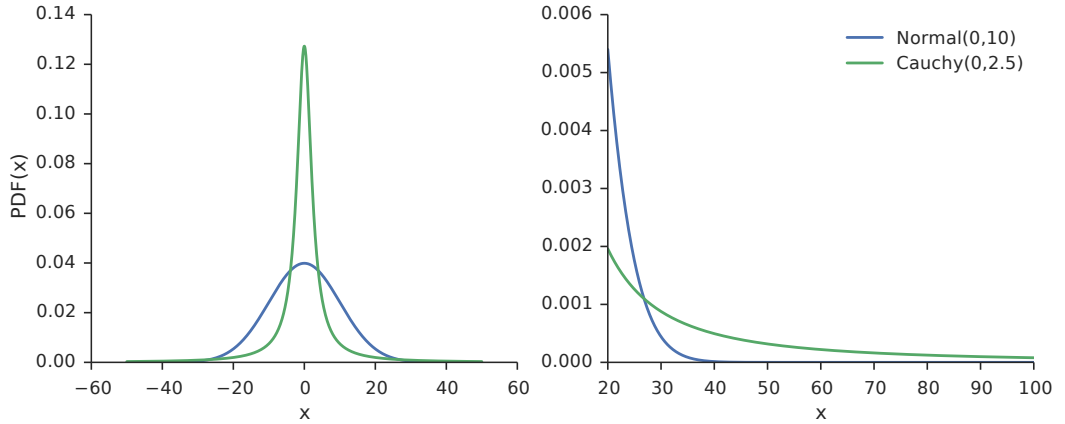


Figure 3.2: Comparison of two weakly-informative priors for logistic regression.

ues without strongly restricting them. For logistic regression parameters, Gelman et al. [125] recommend using a Cauchy distribution with centre 0 and scale 2.5 as a weakly-informative prior distribution, and in this chapter we also use a $\text{Normal}(0, 10)$ (i.e. standard deviation 10) distribution for some parameters to improve numerical stability. These two weakly-informative priors are compared in Figure 3.2, showing that the $\text{Normal}(0, 10)$ prior distribution is less concentrated in the tails of the distribution.

3.2.3 Computational tools

3.2.3.1 MCMC

The posterior density of a hierarchical model is generally a high-dimensional probability distribution without an easily-calculable normalising factor (the denominator in equation (3.1)). This posterior density must therefore be evaluated approximately using computational tools. A popular general computational technique for approximating a posterior density is Markov Chain Monte Carlo (MCMC), which draws values from a distribution approximating the target posterior and gradually corrects those draws after successive iterations to improve the approximation. One particular MCMC algorithm, the Metropolis-Hastings algorithm, is commonly used to sample from posterior distributions; it requires evaluating only the unnormalised posterior density. The algorithm is as follows [124]:

1. Draw initial parameter values θ_0 , preferably in a region of high posterior density
2. Repeat:

- (a) Draw a proposed new value of the parameters, θ^* , from a proposal distribution:

$$\theta^* \sim q(\theta^*|\theta^{t-1})$$

- (b) Calculate the ratio of the posterior density at θ_{t-1} and θ^* as follows:

$$r = \frac{p(\theta^*|y)q(\theta^{t-1}|\theta^*)}{p(\theta^{t-1}|y)q(\theta^*|\theta^{t-1})}$$

- (c) Set

$$\theta_t = \begin{cases} \theta^* & \text{with probability } \alpha = \min(r, 1) \\ \theta_{t-1} & \text{with probability } 1 - \alpha \end{cases}$$

The algorithm is constructed so that the simulated sequence is a Markov chain with a unique stationary distribution, and that this stationary distribution is the same as the target posterior distribution. Therefore, once the sequence has converged to its stationary distribution, the following samples generated from the sequence have the same distribution as the posterior.

3.2.3.2 Hamiltonian Monte Carlo

Although the Metropolis-Hastings algorithm has advantages over other methods of sampling from the posterior, such as importance sampling or Gibbs sampling, it suffers from the problem that highly correlated parameters can lead to slow convergence to the target distribution. This means that the sampler must be run for many iterations and may take a prohibitively long time to perform inference. There has therefore been much recent research on samplers with faster convergence rates for correlated parameters, such as Adaptive Metropolis, Slice Sampling, and Hamiltonian Monte Carlo (HMC). This latter method has been implemented in the convenient software package ‘STAN’ [126], and often provides better convergence than the Metropolis algorithm.

HMC improves on the performance of the Metropolis-Hastings algorithm by modifying its random walk behaviour to move more rapidly through the target distribution, using ideas from dynamics in physics. In HMC, each dimension in the posterior parameter space θ_j has an associated momentum variable χ_j which determines the direction and size of the next proposed parameter. Just as in the Metropolis-Hastings algorithm HMC requires the unnormalised posterior density, but HMC also requires the gradient of the log-posterior density in order to calculate

the momentum variable. One iteration of the HMC algorithm is shown below, and is repeated until the sequence has reached its stationary distribution [124]:

1. Draw χ from its prior distribution, which can be any distribution but is often a multivariate normal with mean 0 and covariance set to a prespecified ‘mass matrix’ M :

$$\chi \sim \text{MVN}(0, M).$$

Often, choosing M to scale with the inverse covariance matrix of the posterior distribution $(\text{var}(\theta|y))^{-1}$, leads to more efficient computation.

2. Repeat L ‘leapfrog’ steps:

- (a) Update χ to make a half-step, based on the log-posterior $\log p(\theta|y)$ defined in equation (3.1):

$$\chi \leftarrow \chi + \frac{1}{2} \frac{d \log p(\theta|y)}{d\theta}$$

- (b) Update θ based on the new value of χ and a ‘mass matrix’ M :

$$\theta \leftarrow \theta + \epsilon M^{-1} \chi$$

3. Calculate the posterior density ratio between the parameters at the start of the leapfrog process, $(\theta_{t-1}, \chi_{t-1})$, and those at the end, (θ^*, χ^*) :

$$r = \frac{p(\theta^*|y)p(\chi^*)}{p(\theta_{t-1}|y)p(\chi_{t-1})}$$

4. Set

$$\theta_t = \begin{cases} \theta^* & \text{probability } \min(r, 1) \\ \theta_{t-1} & \text{otherwise} \end{cases}$$

The parameters L , ϵ , and M , are chosen by tuning during a warmup period. In this warmup period, the parameters are gradually changed until the percentage of steps accepted by the algorithm is near the theoretically-optimal percentage of 65% [127]. The samples generated during the warmup period are then discarded and the tuned parameters used to generate subsequent samples.

We use HMC to sample from the hierarchical models in this chapter because of its superior convergence properties and the convenience of the ‘STAN’ software package.

3.2.3.3 MCMC Diagnostics

An MCMC algorithm provides samples from the posterior distribution once it has converged to its stationary distribution. There are several methods for assessing convergence, but one suggested approach [128] is to use the Gelman-Rubin statistic [129]. This statistic requires multiple independent markov chains to be run from different initial points, and then analyses the variance of samples between chains (B) and within the chains (W). If B and W are sufficiently similar then the chains are likely to have converged. The variances are calculated as follows:

$$B = \frac{n}{m-1} \sum_{j=1}^m (\bar{\theta}_{.j} - \bar{\theta}_{..})^2$$

$$W = \frac{1}{m} \sum_{j=1}^m \left[\frac{1}{n-1} \sum_{i=1}^n (\theta_{ij} - \bar{\theta}_{.j})^2 \right].$$

These variances can be combined to form an estimate of the marginal posterior variance of θ as follows:

$$\hat{\text{Var}}(\theta|y) = \frac{n-1}{n} W + \frac{1}{n} B.$$

The Gelman-Rubin statistic is then defined as:

$$\hat{R} = \sqrt{\frac{\hat{\text{Var}}(\theta|y)}{W}}.$$

As the chain approaches convergence, W and $\hat{\text{Var}}(\theta|y)$ will converge to the true variance of θ . As such, the Gelman-Rubin statistic will approach $\hat{R} \rightarrow 1$ as the chain reaches convergence.

Another important consideration for MCMC sampling is the autocorrelation of the samples, defined as:

$$A(\tau) = \frac{\mathbb{E}[(X_t - \mu)(X_{t+\tau} - \mu)]}{\sigma^2}.$$

If the autocorrelation is significant at lags greater than 1 then the samples can be ‘thinned’ by taking only every k -th sample, thus reducing the autocorrelation.

3.2.3.4 Model comparison

AIC, as described in chapter 2, is appropriate for comparing flat models, such as Logistic Regression, which satisfy asymptotic normality [130]. For these models, the AIC equals the average generalisation error [131]. Hierarchical models, however, generally do not satisfy the asymptotic normality condition and as such the AIC for these models does not equal the average generalisation error. For Bayesian hierarchical models, Gelman et al. [132] suggest the Widely Applicable Information Criterion (WAIC) [133] as a reasonable and convenient estimate of generalisation error. The WAIC for a posterior distribution is calculated by combining the log pointwise predictive density, i.e. the log of the mean likelihood of the data under the posterior distribution for the parameters, with a penalty term specific to WAIC. The log pointwise predictive density is given by:

$$\text{lppd} = \sum_{i=1}^n \log \left(\frac{1}{S} \sum_{s=1}^S p(y_i | \theta^s) \right),$$

where θ^s is the value of θ for the s -th MCMC sample and y_i is the i -th data point. The penalty term, or the effective number of parameters, is calculated as:

$$p_{\text{WAIC}} = \sum_{i=1}^n V_{s=1}^S (\log p(y_i | \theta^s)).$$

WAIC is then given by:

$$\text{WAIC} = -2\text{lppd} + 2p_{\text{WAIC}}.$$

As with AIC, lower values of WAIC indicate lower generalisation error and hence a better model.

3.2.3.5 Summarising posterior distributions

Throughout this chapter we summarise the marginal posterior distributions for parameters of various Bayesian models using the probability that the sign of the parameter is positive, given the data observed \mathbf{y} :

$$\begin{aligned} P(\theta_i > 0 | \mathbf{y}) &= \int \mathbf{I}_{\{\theta_i > 0\}} p(\theta_i | \mathbf{y}) d\theta_i \\ p(\theta_i | \mathbf{y}) &= \int p(\theta | \mathbf{y}) d\theta_{-i}, \end{aligned}$$

where $p(\theta | \mathbf{y})$ is the joint posterior distribution from equation (3.1) and θ_{-i}

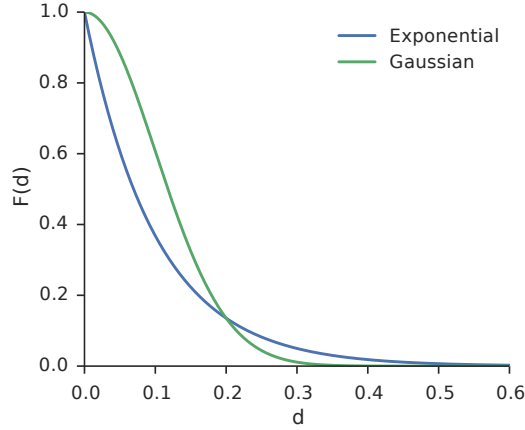


Figure 3.3: Comparison of exponential and Gaussian spatial kernels for the same value of the kernel width parameter $\phi = 0.1$.

are all parameters in the model except θ_i . We use the samples from a Monte Carlo method to estimate this probability as follows:

$$P(\theta_i > 0|\mathbf{y}) = \sum_{s=1}^{N_s} \mathbf{I}_{\{\theta_{i,s} > 0\}}, \quad (3.2)$$

where $\theta_{i,s}$ is the s -th sample of parameter θ_i , and N_s is the total number of samples drawn for this model.

3.2.4 Spatial modelling

Some of the models considered in this chapter include the effect of distance, such as the distance to a health facility, on an outcome. Spatial effects are unlikely to be linear, since after a certain large cutoff distance a spatial effect is likely to be zero. Instead, the effect of distance is assumed to follow a parametric function, which for most spatial applications is large at small distances and decreases with distance. In this chapter we consider a negative exponential distance kernel, which is commonly used to model utilisation rates [100]:

$$F_{exp}(d) = e^{-\frac{d}{\phi}}.$$

The shape of this kernel is given in Figure 3.3, compared to a Gaussian spatial kernel.

Health behaviour outcomes may also be spatially correlated with each other: individuals living near to each other may have similar health behaviours due to the

social influence described in chapter 4, whereas individuals far from each other may have independent outcomes. This is the last aspect of the general behavioural model, and can be included as follows:

$$\mu_{\text{spatial}} \sim \text{MVN}(0, \sigma^2 H(\mathbf{d}, \phi)).$$

H is a spatial correlation function which captures the way in which the correlation between the outcome at two spatial locations varies with the distance between them. One method for investigating the spatial correlation in data is the empirical semivariogram [134], which estimates the difference between the covariance of the data at zero distance and the covariance at other distances. An estimator for the semivariogram at a representative distance lag \mathbf{h}_u is as follows:

$$\hat{\gamma}(\mathbf{h}_u) = \frac{1}{2N(H_u)} \sum_{\mathbf{s}_i - \mathbf{s}_j \in H_u} \{\hat{e}(\mathbf{s}_i) - \hat{e}(\mathbf{s}_j)\}^2 \quad (u = 1, \dots, k),$$

where $\hat{e}(\mathbf{s}_i)$ is the value of the datapoint at location \mathbf{s}_i , H_u is the u 'th bin of distances represented by \mathbf{h}_u , and $N(H_u)$ is the number of datapoints that fall within bin H_u . Figure 3.4 (left) plots an example semivariogram against \mathbf{h}_u for multivariate normal data generated with an exponential distance kernel. This figure shows that closer points have lower semivariance, and that after a certain distance (around $\mathbf{h}_u = 0.5$ for this case) the semivariance saturates and does not increase further. This is a general feature of the semivariograms for spatially-correlated data. Figure 3.4 (right) shows the semivariogram for example data without any spatial correlation. The semivariance does not depend on distance, so the semivariogram is flat.

3.3 The relationship between health centre quality and usage

3.3.1 Data

The data consist of two sets of surveys: Household LQAS surveys and Health Facility Assessment (HFA) surveys. The household surveys, performed in 2012 and 2013 across 78 Ugandan districts, were conducted using the same methodology as described in chapter 2. The HFA surveys assessed 181 health facilities across south-west Uganda in 2012 and 2013, using the methodology described above. In this section we focussed on two HFA modules: Exit Interview and Sick Child Observation. During the Exit Interview six consecutive pharmacy clients at the health

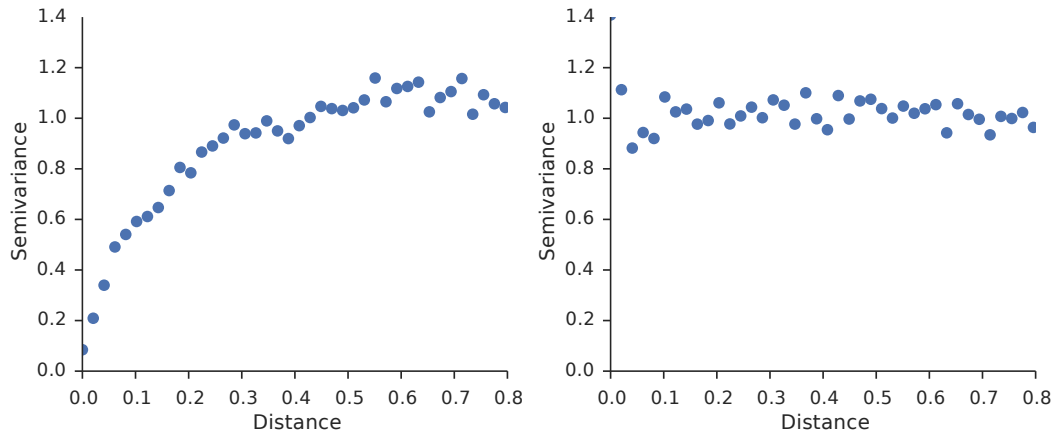


Figure 3.4: Example semivariograms for an exponential distance kernel with $\phi = 0.1$ (left), and no spatial correlation (right).

facility were asked whether they ‘Strongly agree’, ‘Agree’, ‘Not sure’, ‘Disagree’, or ‘Strongly disagree’ with the following statements:

1. All the health workers who reviewed you today carefully listened to you
2. You were given the chance by all the health workers you met during this visit to clearly state your problems and ask questions
3. The health workers who reviewed you today explained well to you about your illness
4. It was easy for you today to get all the drugs that were prescribed for your illness
5. You were handled with respect during this visit
6. You feel you can trust the health workers with your private information
7. You are confident with the privacy you received during review and/or counseling session with the health workers
8. You have no problem today with the sanitation state at this health facility

During the Sick Child Observation the consultations with a health worker were observed for six consecutive sick children from 2-59 months who presented with fever, cough or difficult breathing, or diarrhoea. The health worker’s diagnosis and prescription for each child was recorded and checked using Table 3.1.

Diagnosis	Correct prescription
Malaria or fever	ACT
Severe Malaria	Quinine
Malaria or fever + other infection (not diarrhoea or ARI)	ACT + antibiotics
ARI / pneumonia	cotrimoxazole or amoxicillin
Diarrhoea without blood	ORS without other antibiotics
Diarrhoea with blood	ORS + antibiotics
No malaria, pneumonia or diarrhoea	No antibiotics or ACT or ORS

Table 3.1: Appropriate prescriptions to treat each diagnosis identified in the Sick Child Observation module [1, 2, 3].

Education level	Classification
Never attended	None
Incomplete primary	None
Complete primary	Primary
Functional adult literacy	Primary
O-level	Secondary
Vocational training	Secondary
A-level	Secondary
Post-secondary	Secondary

Table 3.2: Classification of ‘education level’ based on survey response.

3.3.2 Analysis

In this analysis the study population was limited to mothers with children aged 12-23 months whose child had experienced fever, pneumonia, or diarrhoea in the past two weeks. The outcome variable was a combination of the responses to two questions from the LQAS survey:

1. “Did you seek any advice or treatment for the sickness from any source”
2. “Where did you seek advice or treatment”

Mothers were considered to have used a health facility if they answered ‘Yes’ to question 1 and ‘Hospital/Clinic’ or ‘Pharmacy/drug shop’ to question 2, and assigned an outcome value of 1; otherwise, the outcome value was 0.

The individual-level covariates included in the model were age, education, and whether or not the individual had a regular partner; ‘education level’ is defined in Table 3.2 and ‘regular partner’ is defined in Table 3.3.

In a slight modification from the analysis in chapter 2, we used a mother’s subcounty, rather than supervision area, to identify which cluster she belonged to.

Relationship status	Classified as regular partner?
Single, no partner	No
Single, no regular partner	No
Single, regular partner	Yes
Cohabiting	Yes
Married	Yes
Divorced/Separated	No
Widowed	No

Table 3.3: Classification of ‘regular partner’ based on survey response.

Subcounties, unlike supervision areas, have a direct administrative meaning; we could therefore give a mother an approximate spatial location and hence identify her local health facility. Unfortunately the subcounty names provided in the household survey did not necessarily match the names the GIS data which provided the GPS coordinates of subcounties in Uganda. We therefore used approximate string matching, as described in section 3.2.1, to link the mother’s subcounty name to known subcounties in the GIS: a subcounty in the data was ‘matched’ to a known subcounty if the similarity score between the two strings was greater than 95%. Of the 8232 mothers in the original sample, 6082 were matched to known subcounties. Mothers were linked to health facilities by assuming an exponential spatial kernel based on the distance from the centroid of their subcounty to the location of each health facility. Models were fitted with several values of the kernel width parameter, $\phi = 0.01, 0.025, 0.075, 0.1, 0.2$, and 0.5 , and compared using WAIC.

There were also discrepancies between the names of health facilities in the HFA data and those provided in the GIS data for health facility locations, so we again used approximate string matching, with a similarity score cutoff of 95%, to match health facility names between the two datasets. Of the 181 health facilities in the original HFA data, 164 were matched to a known GIS health facility. There were 303 health facilities in the GIS dataset which were not matched to an HFA facility, either because they were not assessed or because there was too large a discrepancy between the name in the HFA dataset and that in the GIS dataset.

The subjective quality of a health facility was estimated from the pharmacy exit interview component of the HFA. A client was considered to agree with a statement in the exit interview if they answered ‘Agree’ or ‘Strongly agree’; otherwise, they were considered to disagree. The total fraction of agreement for all statements and clients at each health facility was then calculated, and each health facility was given a subjective quality of ‘good’ or ‘bad’ based on whether this total agreement fraction reached a certain threshold. Since a large fraction of health facility clients

agreed with all statements, the threshold in this case was set to 0.9. For this threshold, 49% of health facilities were classified as ‘bad’. The objective quality of a health facility was estimated by calculating the percentage of correct treatments prescribed by the health worker during the Sick Child Observation module, as described above. A health facility was classified as ‘bad’ if this percentage was less than 70% and ‘good’ otherwise, leading to 60% of health facilities being classified as ‘bad’. The quality covariates included in the model were grade (level II, level III, level IV or Hospital), distance to nearest main road, bed capacity, and ownership.

In this section we again used the framework of random-intercept logistic regression, introduced in section 2.3.1. There were the additional complications, however, that the data did not specify which health facility a mother visited, and that not all facilities were assessed for quality. These complications were addressed by building a hierarchical Bayesian model which specified the probability that a mother would have gone to a particular health facility and which imputed the quality of unassessed health facilities based on health facility covariates. A quality estimate was thus obtained for all health facilities, whether or not they were part of the original HFA study. Unlike chapter 2, the additional spatial correlation between mothers within a subcounty, captured by the ξ random intercept, was estimated explicitly. The model is as follows:

$$\begin{aligned}
y_i &\sim \text{Bernoulli}(p_i) \\
Q_j &\sim \text{Bernoulli}(\mu_j) \\
\text{logit}(p_i) &= (\beta^T \Theta)_i + \sum_{j=1}^{N_{hc}} \pi_{k,j} (\gamma^T Q)_j + \xi_k \\
\xi_k &\sim \text{Normal}(0, \sigma_{sc}) \\
\text{logit}(\mu_j) &= (\omega^T \Psi)_j \\
\pi_{k,j} &= \frac{e^{-\frac{d_{k,j}}{\phi}} \mathbf{I}_{d_{k,j} \leq d_c}}{\sum_{j=1}^{N_{hc}} \mathbf{I}_{d_{k,j} \leq d_c} e^{-\frac{d_{k,j}}{\phi}}} \\
\beta, \omega &\sim \text{Cauchy}(0, 2.5) \\
\sigma_{sc} &\sim \text{Cauchy}(0, 25) \\
\gamma &\sim \text{Normal}(0, 10),
\end{aligned} \tag{3.3}$$

where y_i is the i -th outcome datapoint, Q_j is the j -th health facility’s quality (either subjective or objective), Θ is the matrix of household-level covariates, $d_{k,j}$ is the distance between the k -th subcounty and the j -th health facility, $\pi_{k,j}$ is the

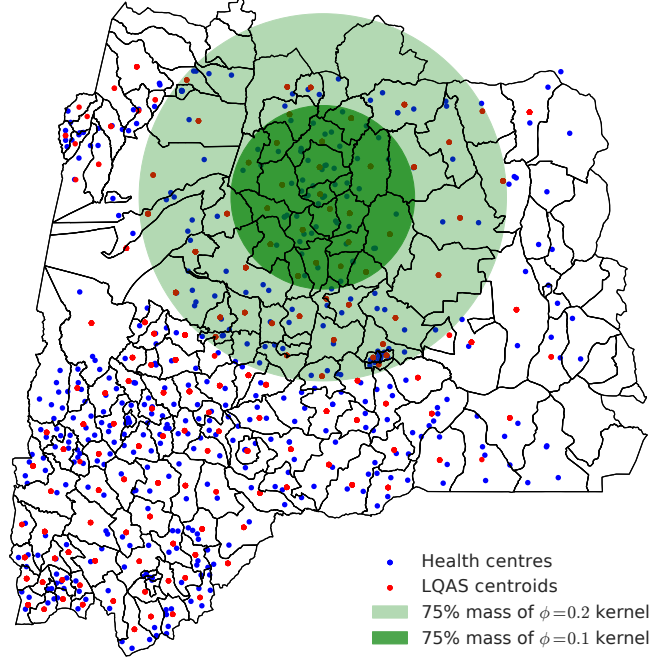


Figure 3.5: Spatial kernel shown in the context of South-West Uganda. The map shows the spatial extent of 75% of the mass of the kernel function evaluated at $\phi = 0.1$ and $\phi = 0.2$, overlaid on a map of the subcounties included in this study. The health facilities and LQAS centroids (where mothers are assumed to be located) are also shown.

probability that health facility j is the local facility of mothers in subcounty k , Ψ is the matrix of covariates for predicting health facility quality, ϕ is the width of the spatial kernel, d_c is the distance cutoff, ξ_k is the unknown random intercept for subcounty k , and β , γ , and ω are the coefficients for household-level covariates, quality covariates, and covariates for predicting health facility quality respectively. σ_{sc} is the standard deviation of the subcounty random intercept.

The spatial extent of the kernel function for two values of the spatial parameter $\phi = 0.1$ and $\phi = 0.2$ is shown in Figure 3.5, overlaid on a map of the study area. This suggests that appropriate values for the spatial parameter should be near these values: much larger values of ϕ will cover a significant portion of the study area and are therefore unlikely to be informative about which health facilities mothers choose.

This general model can be reduced to specific models by setting the value of ϕ , setting $\gamma = 0$ or $\xi_k = 0$. This leads to the following specific models:

- Flat model, No quality: ($\gamma = 0, \xi_k = 0$)
- Flat model, with subjective quality γ : ($\gamma \sim \text{Normal}(0, 10), \xi_k = 0$)
- Flat model, with objective quality γ : ($\gamma \sim \text{Normal}(0, 10), \xi_k = 0$)
- Hierarchical model, No quality: ($\gamma = 0$)
- Hierarchical model, with subjective quality γ : ($\gamma \sim \text{Normal}(0, 10)$)
- Hierarchical model, with objective quality γ : ($\gamma \sim \text{Normal}(0, 10)$)

HMC was used to obtain 4 chains of 3000 samples each (excluding a warmup period of 1000 samples for each chain) from the posterior of these models, for different values of $\phi \in [0.01, 0.5]$. The autocorrelation was plotted and a level of thinning chosen to reduce the autocorrelation. Using the thinned samples the Gelman-Rubin statistic was calculated for each parameter to check convergence diagnostics, and the WAIC calculated to identify the best models for no quality, objective, and subjective quality. The best model for each type of quality was then sampled from again to obtain four chains of 9000 samples each. We then compared the predictive power of the three models by calculating p_i in equation (3.3) for each sample of the parameters and using these to construct a posterior predictive distribution of ROC curves and AUC values, described in section 2.3.4. These distributions were then summarised using the mean and 95% credible interval (given by the $[2.5^{\text{th}}, 97.5^{\text{th}}]$ percentile range).

3.3.3 Results

Neither the subjective nor the objective quality variable showed significant spatial correlation, since the semivariograms shown in Figure 3.6 are flat for both. We therefore fit models without spatial correlation in these variables.

Three types of model were fitted: no quality covariate, a subjective quality covariate, and an objective quality covariate. The latter two model types were compared for several values of the spatial kernel parameter ϕ . The WAIC values for all models are shown in Figure 3.7. All of the flat models performed worse than all of the hierarchical models, even though WAIC takes into account the additional parameters included in the hierarchical models. Most flat models which include quality were worse or as bad as the flat model without quality, and there was no

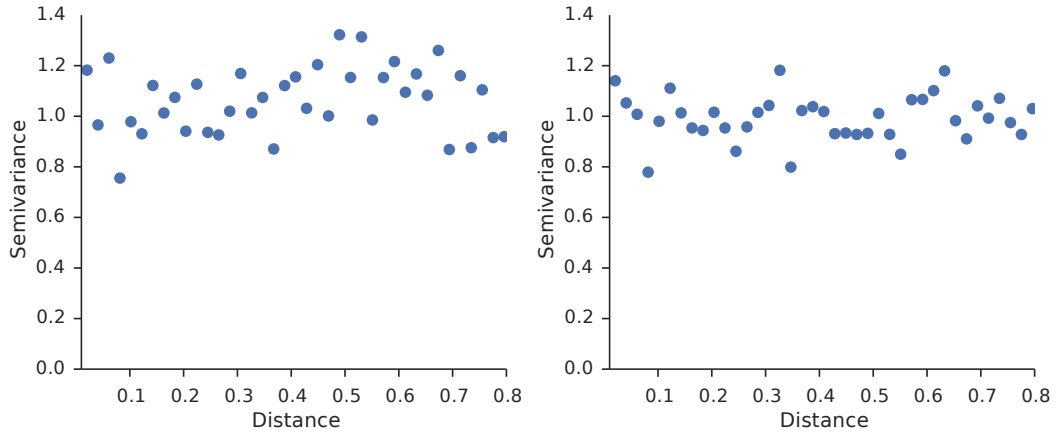


Figure 3.6: Semivariograms for subjective quality (left) and objective quality (right). Neither variable shows a clear semivariance increase with spatial lag, so these variables do not seem to have strong spatial correlation.

clear trend in which values of ϕ provided better models except that models with $\phi < 0.075$ were worse for both subjective and objective quality. For the hierarchical models, almost all of those that included quality were better than the model that did not. For both subjective and objective quality there was a minimum in the WAIC between $\phi = 0.1$ and $\phi = 0.2$, and models below this range performed worse than models above it.

The best performing models with subjective and objective quality, in this case the hierarchical models with $\phi = 0.1$ and $\phi = 0.2$ respectively, were then investigated in more detail. The Gelman-Rubin statistic was calculated for each covariate coefficient parameter in each model and shown in Table 3.4; all were close to 1, consistent with convergence.

An autocorrelation plot for the scalar parameters in the best subjective quality model is given in Figure 3.8, and suggests a thinning of 10 is appropriate. The autocorrelation at this level of thinning is then shown in Figure 3.9.

The marginal posterior distributions for the coefficients of individual-level covariates for the model with subjective quality are shown in Figure 3.10. The strongest association is that older mothers are less likely to seek treatment at a health facility, with $P(\beta_{\text{age}} < 0|\mathbf{y}) = 0.95$. The other covariates are less strongly associated with this outcome, but mothers with regular partners or with secondary education may be more likely to seek health facility treatment ($P(\beta_{\text{regular partner}} > 0|\mathbf{y}) = 0.80$ and $P(\beta_{\text{secondary}} > 0|\mathbf{y}) = 0.76$). Health facility-level coefficient distributions are shown in Figure 3.11: NGO ownership is strongly positively associated with subjective quality ($P(\omega_{\text{NGO}} > 0|\mathbf{y}) = 0.98$), while health facilities of grade III

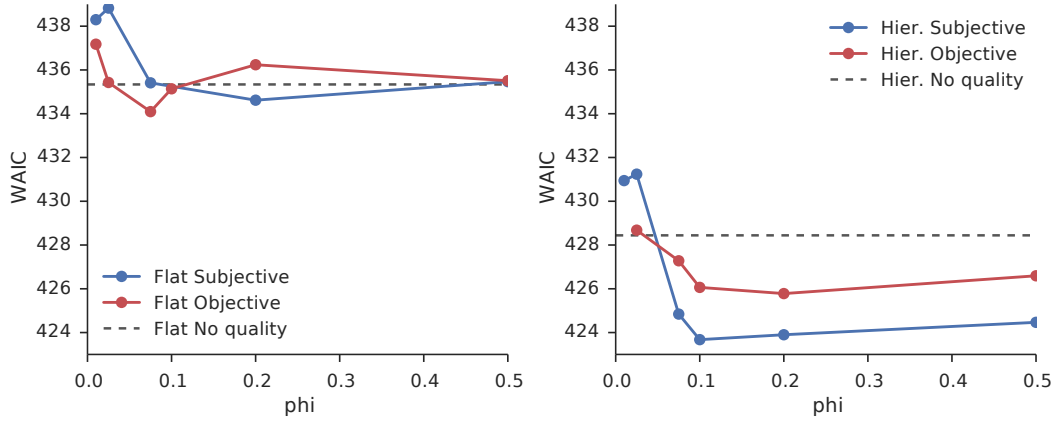


Figure 3.7: WAIC values for models of health facility use including subjective and objective quality, for flat models (left) and for hierarchical models (right). The hierarchical models all perform better than the flat models, and most hierarchical models including quality were better than the model that did not. Hierarchical models below $\phi = 0.1$ all performed worse than those with higher values of ϕ , as did those with values higher than $\phi = 0.2$. There therefore appears to be an optimal value of ϕ between 0.1 and 0.2. Finally, hierarchical models which included subjective quality were almost all better than those including objective quality.

Parameter	\hat{R} for no-quality	\hat{R} for subjective	\hat{R} for objective
grade_HC III	-	1.0	1.0
grade_HC IV	-	1.0	1.0
grade_HOSPITAL	-	1.0	1.0
road_distance	-	1.0	1.0
capacity	-	1.0	1.0
ownership_NGO	-	1.0	1.0
ownership_PRIVATE	-	1.0	1.0
Intercept	-	1.0	1.0
regular_partner	1.0	1.0	1.0
age	1.0	1.0	1.0
education_primary	1.0	1.0	1.0
education_secondary	1.0	1.0	1.0
Intercept	1.0	1.0	1.0
gamma	-	1.0	1.0
sc_sigma	1.0	1.0	1.0

Table 3.4: Gelman-Rubin statistics for the parameters of the best of the three types of model. $\hat{R} \rightarrow 1$ at convergence, so the fact that all of these values are close to 1 is consistent with the models having converged.

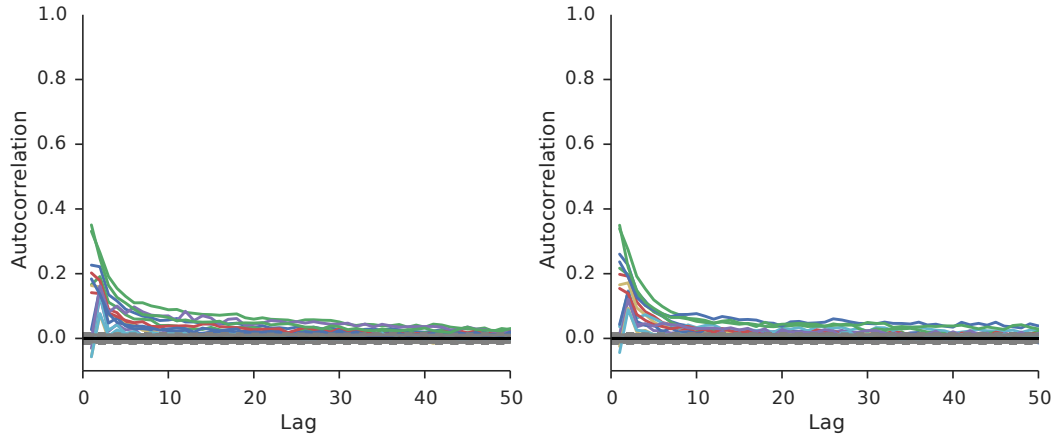


Figure 3.8: Autocorrelations for the subjective (left) and objective (right) model parameters. These plots show that a thinning of 5 to 10 is appropriate.

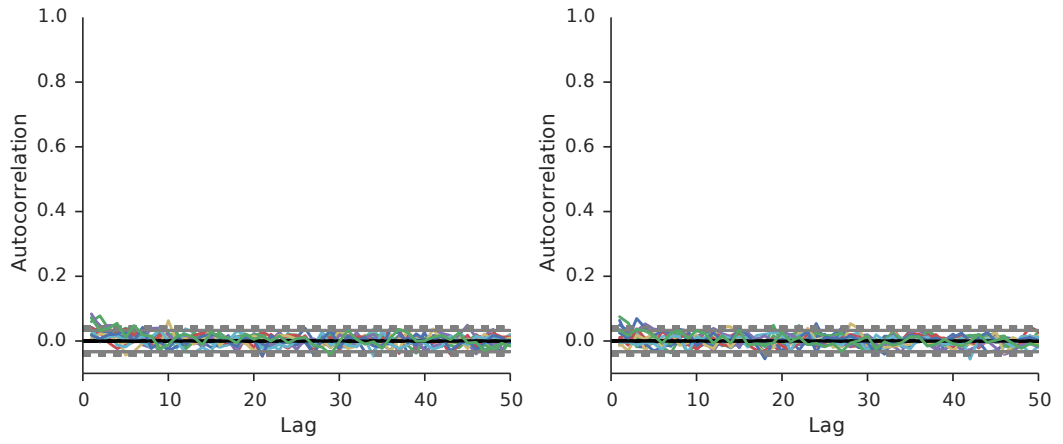


Figure 3.9: Autocorrelations for subjective (left) and objective (right) model parameters with a thinning of 10.

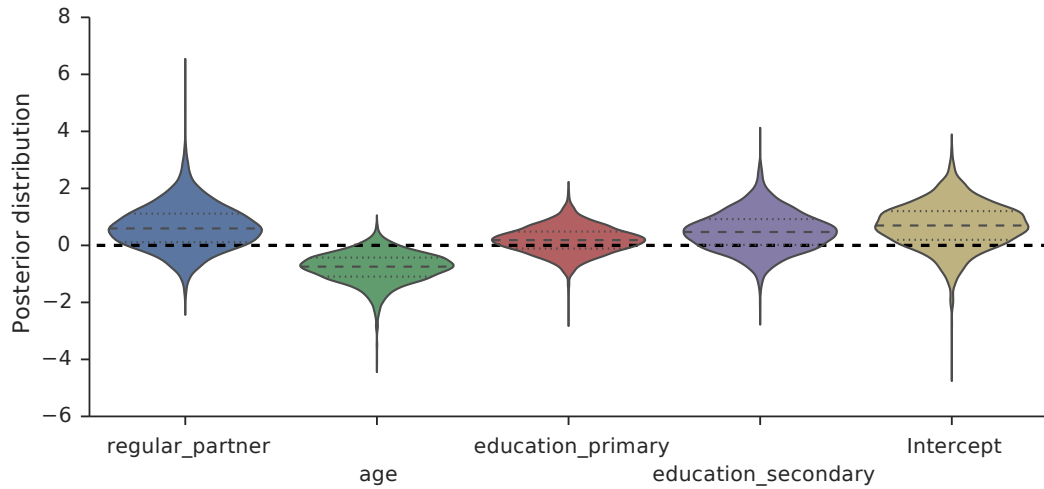


Figure 3.10: Marginal posterior distributions for household-level parameters in the best subjective model. Age is strongly negatively-associated with seeking treatment at a health facility, whereas having a regular partner and secondary education are both weakly positively associated with health facility treatment seeking.

are likely to have a lower subjective quality ($P(\omega_{\text{grade III}} < 0|\mathbf{y}) = 0.95$). None of the other covariates appeared to predict subjective quality. Finally, Figure 3.13 (left) shows the marginal posterior distribution for the quality parameter; higher subjective quality is strongly associated with seeking treatment at a health facility, with $P(\gamma > 0|\mathbf{y}) = 1.0$.

For the best-performing model with objective quality, distributions of coefficients for the individual-level covariates were very similar to those in the subjective model. Distributions for health facility-level coefficients are shown in Figure 3.12: health facilities of grade IV may be more likely to have a higher objective quality ($P(\omega_{\text{grade IV}} < 0|\mathbf{y}) = 0.20$), as may health facilities with lower capacity ($P(\omega_{\text{capacity}} < 0|\mathbf{y}) = 0.87$). NGO ownership does not seem to be associated with objective quality ($P(\omega_{\text{NGO}} < 0|\mathbf{y}) = 0.52$). Figure 3.13 (right) shows that objective quality is also associated with seeking treatment in a health facility, with $P(\gamma > 0|\mathbf{y}) = 0.99$. The standard deviation of the subcounty random intercept for the subjective and objective models are shown in Figure 3.14.

ROC curves were used to compare the specificity and sensitivity of the best-performing models to the model with no quality. Figure 3.15 shows the ROC curves for the best of the three types of model. The area under curve for the three models are $AUC_{\text{noqual}} = 0.73[0.64, 0.80]$, $AUC_{\text{subj}} = 0.76[0.66, 0.82]$, and $AUC_{\text{obj}} = 0.74[0.64, 0.81]$. All credible intervals overlap significantly, with the sub-

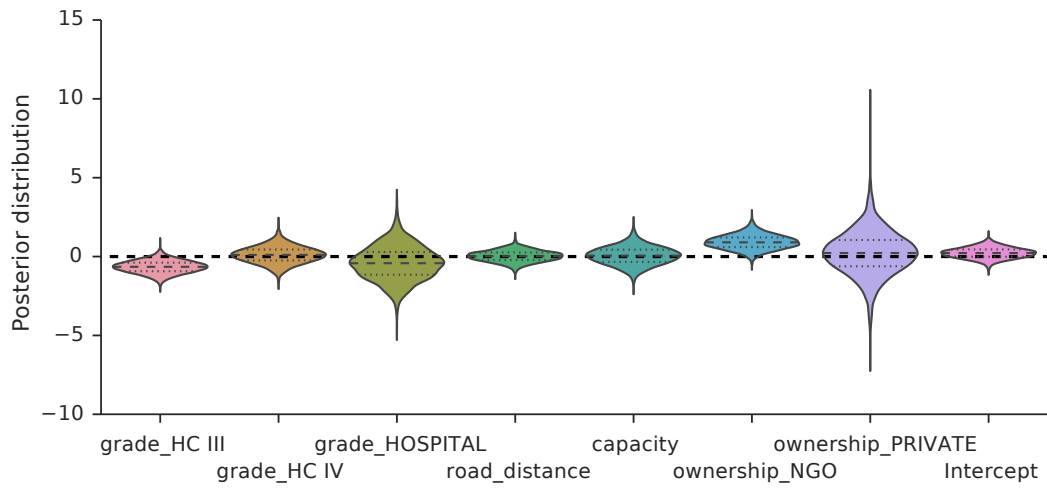


Figure 3.11: Marginal posterior distributions for health facility-level parameters in the subjective quality model. NGO ownership is strongly positively associated with subjective quality, while health facilities of grade III are likely to have a lower subjective quality. None of the other covariates appeared to predict subjective quality.

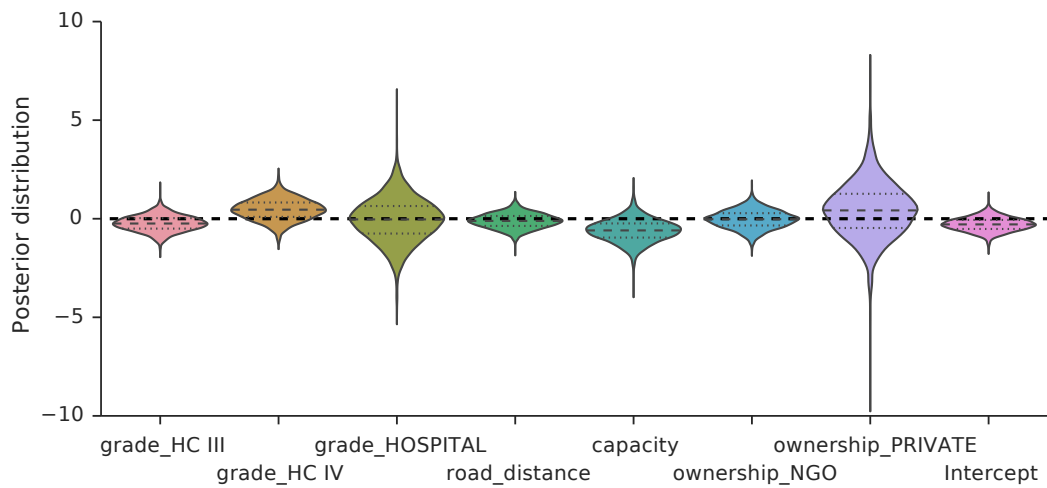


Figure 3.12: Marginal posterior distributions for health facility-level parameters in the objective quality model. NGO ownership is not strongly predictive of objective quality, but health facilities of grade IV are more likely to have higher objective quality as are facilities with lower capacity.

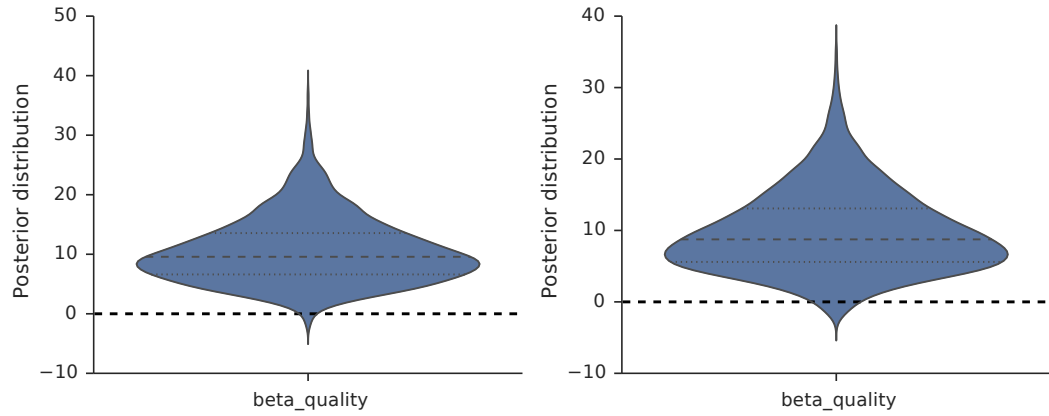


Figure 3.13: Marginal posterior distributions for the subjective (left) and objective (right) quality parameters. Both parameters appear to be strongly associated with seeking treatment in a health facility.

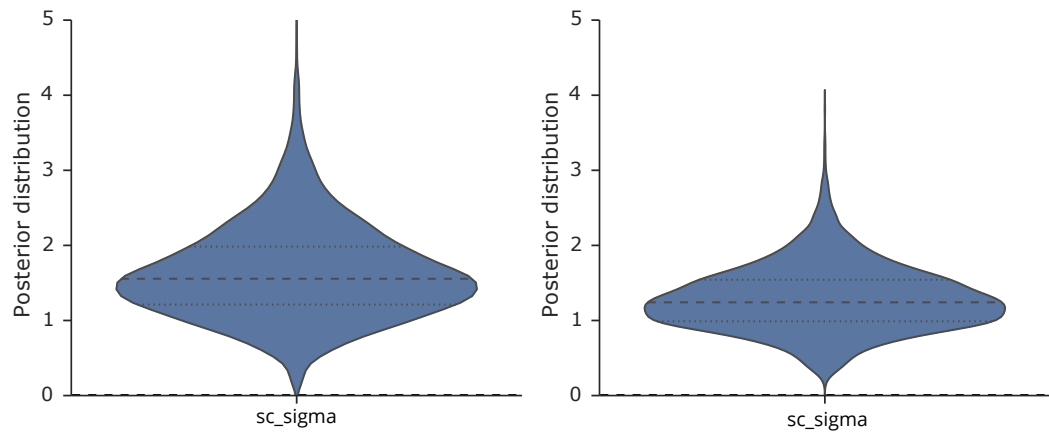


Figure 3.14: Marginal posterior distributions for the standard deviation of the sub-county effect ξ , for the best subjective (left) and objective (right) models.

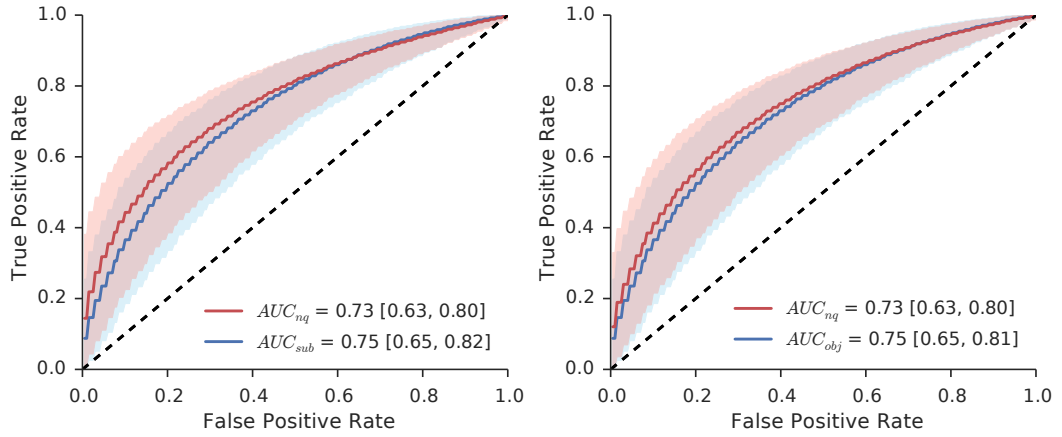


Figure 3.15: ROC curves for the model without quality and the best models with subjective and objective quality. The area under curve for the three models are $AUC_{noqual} = 0.73[0.64, 0.80]$, $AUC_{subj} = 0.76[0.66, 0.82]$, and $AUC_{obj} = 0.74[0.64, 0.81]$. All credible intervals overlap significantly, with the subjective model being slightly better than the others as measured by the area under the curve.

jective model being slightly better than the others using this metric.

3.3.4 Discussion

There are several limitations to this dataset and analysis. Although the match threshold for approximate string matching was set to a high level (95%), and the matches were found on a district-by-district, there is no guarantee that the GIS data has been correctly matched to the HFA and household surveys. In future studies could use a GPS device to record locations for the health facilities and the households, as is the case for the DHS surveys [49]. Privacy could be assured by adding a small amount of noise to household locations. The measures of objective and subjective quality are crude: a cut-off on a percentage score. Further work would investigate the robustness of these results to changes in the cut-off. Distance should also be included in the model, perhaps by including a covariate ‘distance to closest health facility’.

Visiting a pharmacy was included as ‘seeking professional help’ because the subjective quality variable is calculated from a pharmacy exit interview. Although the training of pharmacy staff is likely to be lower than other health facilities, the pharmacy is being recognised as an important healthcare provider [135]. The category in the data indicating a pharmacy, ‘Pharmacy/drug shop’, also includes drug shops; staff at these shops are much less likely to be well-trained and so should not be classified as ‘seeking professional help’. In future studies, separation of these two

distinct types of treatment provider would provide better classification and might lead to stronger results.

The current model is random-intercept: it assumes that the effect of any predictor is the same for each subcounty. There are many reasons why this might not be true in practice, such as a very low-income subcounty containing individuals who cannot afford to visit a health facility even if it is of high quality, but the random-intercept model provides a simple starting point which can be expanded in further work.

The AUC is very high, given the quality of the data. This might indicate some overfitting: the AUC for new data is likely to be lower. Some subcounties contained only one surveyed mother, due to the name discrepancies discussed above, but there is some evidence from other studies that these ‘singletons’ do not significantly affect the results of hierarchical models [136, 137].

A positive outcome for this study is the identification of the spatial range within which mothers in Uganda are likely to visit health facilities: this range (around $\phi = 0.1$) corresponds to visiting neighbouring subcounties but not much further. The average travel speed in Uganda is around 13km/h [138], so converting the spatial kernel parameter into km (where 1 decimal degree corresponds to ~ 110 km at the latitude of Uganda) the probability of selecting a health facility one hour/13km away is 31% lower than selecting one that is very close.

The negligible difference between the ROC curves for the subjective model and the model without quality, implies that although the model with subjective quality is *statistically* better, it is not *practically* better. It is also interesting to note that the WAIC difference between the flat model without quality and the hierarchical model without quality is almost twice that between the hierarchical subjective and no-quality models. This implies that the subcounty effect, which includes social influence, is much more important for predicting usage than quality.

3.4 The effect of funding mechanism on health facility usage

3.4.1 Background

In July 2012, the Liverpool School Tropical Medicine initiated a project in Northern Uganda to evaluate the impact of a novel mechanism for funding health facilities. This new mechanism, Results-Based Financing (RBF), provides funding to health facilities based on measurable improvements in the service they provide. For this

trial, funding is provided by the Department for International Development (DFID) to Private Not-For-Profit (PNFP) health facilities in two seriously deprived regions in Northern Uganda, Acholi and Lango. In 2012, all 31 PNFP facilities in these two regions were funded using traditional input-based financing; then, the 21 facilities in Acholi were funded under RBF while the 10 Lango facilities remained under the traditional mechanism. Of the 31 PNFP facilities in the trial, 29 were evaluated by LSTM before the intervention (in 2012) and two years after it had begun (in 2014). One of the remaining facilities was excluded because it only offered HIV/AIDS services and no maternal or child health services, and the other (Lacor Hospital) provides a greater range of services than normal facilities of its type.

The effect of the funding change was evaluated using two methods. Health facilities were evaluated using a Health Facility Assessment (HFA), which consists of a series of modules assessing quality of care, infrastructure and staff training. The effect on the local community was evaluated using Lot Quality Assurance Sampling (LQAS) for households in catchment areas around each health facility. The LQAS questionnaire provides demographic information such as the age and education of members of the household, health behaviour information such as whether mothers of children with fevers took their child to a health facility, and wealth information such as the average annual income of the household.

The HFA and household data are rich sources, so this analysis was limited to investigating the effect of the intervention on the behaviour of mothers whose child is under 5 years old had a fever in the last two weeks. In particular, we investigated whether changing the funding mechanism of a health facility to RBF leads to a significant increase in the odds of these mothers taking their febrile child to a health facility.

3.4.2 Data

The data consist of Health Facility Assessments (HFA) for 29 Private-Not-For-Profit health facilities and household-level data sampled from the catchment areas of these facilities. The health facilities are located in two northern regions of Uganda: Acholi and Lango. The data was collected before (in 2012) and after (in 2014) a change in the funding mechanism for the facilities in Acholi, with facilities in Lango unchanged and used as a control group.

Catchment areas around each health facility were chosen based on the level of the health facility and the population density around it, such that the majority of users of the health facility would come from this catchment area. Each catchment area was split into four supervision areas (SAs), and within each of those SAs villages

were selected by Probability-Proportional-to-Size (PPS) sampling. This ensured that all households in a SA had the same probability of being sampled.

The HFA data was obtained only for PNFPs, and due to the difficulty of obtaining information in Northern Uganda, data on government health facilities other than their location was very sparse. In addition, some of the PNFP catchment areas overlapped, and mothers did not necessarily go to their local PNFP. These limitations are addressed by the modelling in the analysis stage.

3.4.3 Analysis

The outcome variable was a combination of two questions: “Did you seek advice or treatment for [child name]’s fever from any source?”, and “Where did you first seek advice or treatment?”. If the mother answered ‘Yes’ to the first question and one of ‘Government Health facility’, ‘Government Hospital’, ‘Private Hospital/Clinic’, or ‘PNFP’ to the second then the outcome was set to ‘True’. If she answered ‘No’ to the first question or answered ‘Private Pharmacy’, ‘Private drug shop’, ‘Relative or friend’, or ‘Traditional healer’ the outcome was set to ‘False’. Overall, 69% of the mothers interviewed in 2012 took their child to a health facility compared with 74% in 2014.

The covariates considered can be split into those affecting the household and those affecting the local health facility. The household-level covariates included were the age of the mother, whether she was literate, whether she had a regular partner, and the year she was surveyed. Health facility-level covariates were the average cost paid by six mothers with sick children who were interviewed as part of the HFA, the distance to the nearest level-III or higher government health facility, and the treatment effect (True for health facilities in Acholi in 2014, False otherwise).

The outcome was modelled as a hierarchical Bayesian logistic regression, with covariates included at the household and health facility levels and a random effect introduced at the level of the supervision area. The SA-level random effect was introduced to control for heterogeneity within a catchment area: for example, some parts of the catchment area may be wealthier or further away from the health facility than others. The treatment variable was set to 1 for all health facilities in Acholi district in 2014, and zero otherwise. All coefficient parameters were given weakly-informative priors: Cauchy distributions with centre 0 and scale 2.5. The random effect for each supervision area was drawn from Normal distribution with mean 0 and standard deviation σ_{sa} , with a Cauchy(0,10) prior on σ_{sa} .

The model is shown graphically in Figure 3.16 is specified as follows:

$$\begin{aligned}
y_{i,t} &\sim \text{Binomial}(p_{i,t}) \\
\text{logit}(p_{i,t}) &= \beta_0 + \beta_{\text{time}}t + (\beta^T \boldsymbol{\Theta})_i + (\gamma \mathbf{E})_k + \gamma_{\text{treat}} \mathbf{1}_{\{\text{treated}:k,t\}} + \xi_j \\
\xi_j &\sim \text{Normal}(0, \sigma_{\text{sa}}) \\
\beta_0, \beta_{\text{time}}, \beta &\sim \text{Cauchy}(0, 2.5) \\
\gamma, \gamma_{\text{treat}} &\sim \text{Cauchy}(0, 2.5) \\
\sigma_{\text{sa}} &\sim \text{Cauchy}(0, 25).
\end{aligned} \tag{3.4}$$

Here, $y_{i,t}$ is the outcome value for household i surveyed at time t (0 for 2012, 1 for 2014), which is in the catchment area of health facility k and is in supervision area j . $\boldsymbol{\Theta}$ is the matrix of covariates for households, ξ_j is the supervision area-level effect. \mathbf{E} is the matrix of covariates associated with the health facility, and $\mathbf{1}_{\text{treated}:k,t}$ is an indicator function with value 1 if health facility k was part of the treatment group and $t = 1$.

HMC was used to obtain 4 chains of 9000 samples each (excluding a warmup period of 1000 samples for each chain) from the posterior of this model. The autocorrelation was plotted and a level of thinning chosen to reduce the autocorrelation. Using the thinned samples the Gelman-Rubin statistic was calculated for each parameter to check convergence diagnostics. We estimated the predictive power of the model by calculating $p_{i,t}$ in equation (3.4) for each sample of the parameters and using these to construct a posterior predictive distribution of ROC curves and AUC values, as described in section 2.3.4. The distribution of the AUC was then summarised using the mean and [2.5th, 97.5th] percentiles. Marginal posterior distributions for the parameters were summarised using the probability that they were positive given the data, $P(\beta > 0|\mathbf{y})$, estimated from the HMC samples using equation (3.2).

Finally, we used a hierarchical linear regression model to identify any significant differences in the price or the objective quality of the health facilities between the RBF and non-RBF groups. A significant difference could suggest that part of any observed treatment effect may act through these two variables. The model is as follows:

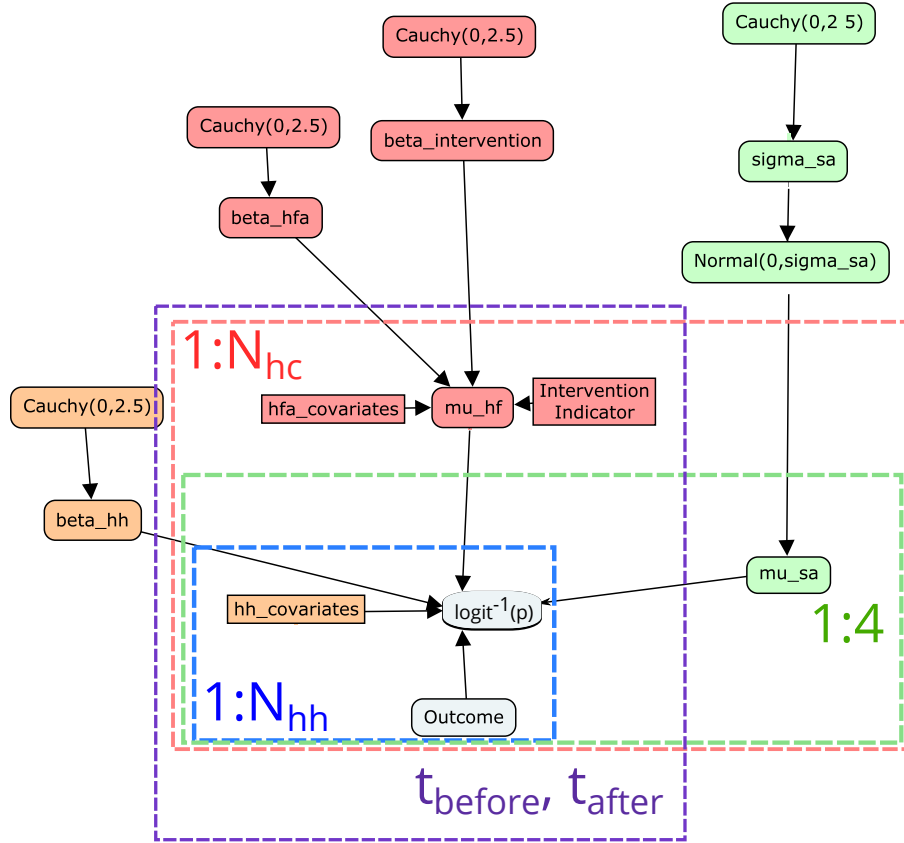


Figure 3.16: Graphical representation of model. N_{hh} is the number of households surveyed in a supervision area (usually 24), N_{sa} is the number of supervision areas per health facility catchment area (always 4 for this data), N_{hc} is the number of health facilities (30 for this data), and $t_{\text{before}}, t_{\text{after}}$ indicate times before and after the intervention (2012 and 2014) respectively.

Parameter	Gelman-Rubin Statistic
regular_partner	1.00
age	1.00
is_literate	1.00
time	1.00
Intercept	1.00
hfa_price	1.00
nearest_govt	1.00
treatment	1.00
sa_sigma	1.00

Table 3.5: Gelman-Rubin convergence diagnostics for each parameter. A value of 1 indicates convergence.

$$\begin{aligned}
x_{k,t} &\sim \text{Normal}(\mu_{k,t}, \sigma) \\
\mu_{k,t} &= \chi_0 + \chi_{\text{treat}} \mathbf{I}_{\{\text{treated:}k,t\}} + \chi_{\text{time}} t + \alpha_k \\
\alpha_k &\sim \text{Normal}(0, \sigma_{hf}) \\
\chi_0, \chi_{\text{treat}}, \chi_{\text{time}} &\sim \text{Cauchy}(0, 25) \\
\sigma &\sim \text{Cauchy}(0, 25) \\
\sigma_{hf} &\sim \text{Cauchy}(0, 25),
\end{aligned} \tag{3.5}$$

where $x_{k,t}$ is either the standardised price or standardised objective quality of health facility k at time t . The treatment term $\mathbf{I}_{\{\text{treated:}k,t\}}$ is an indicator function with value 1 if the health facility k is in Acholi region and $t = 1$, and value 0 otherwise. α_k is a random-intercept term to take into account differences between health facilities that are present at both times. The price and quality models were both sampled from using HMC, with 4 chains of 30000 samples each (excluding 2000 warmup samples for each chain). As before, an autocorrelation plot was used to choose an appropriate level of thinning and Gelman-Rubin convergence diagnostics calculated using the thinned samples.

3.4.4 Results

The autocorrelation shown in Figure 3.17, and implies that a thinning of 8 is appropriate. The Gelman-Rubin statistics for each variable are given in Table 3.5.

The marginal posterior distributions for various parameters are shown in Figs 3.18, 3.19, and 3.20. Figure 3.18 shows the posterior for the coefficients on

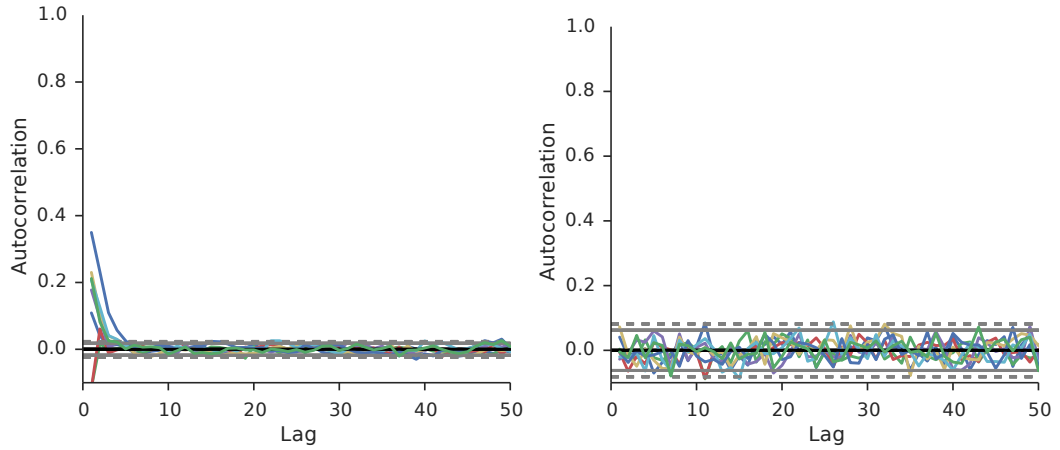


Figure 3.17: Autocorrelation of all household-level covariates (left), and the autocorrelation after a thinning of 8 (right).

health facility-level covariates and the effect of treatment. The mean price paid for taking a child to the local PNFP health facility appears to be weakly negatively associated with the probability of a mother taking her child to a health facility ($P(\gamma_{\text{price}} < 0|\mathbf{y}) = 0.82$). Mothers with local PNFP facilities that are further away from government health facilities are significantly less likely to seek treatment in a health facility ($P(\gamma_{\text{nearest govt.}} < 0|\mathbf{y}) = 1.0$). Finally, mothers in 2014 who live close to a health facility in the treatment group are significantly more likely to seek treatment from a health facility ($P(\gamma_{\text{treat}} > 0|\mathbf{y}) = 1.0$).

Figure 3.19 shows the posterior for household-level covariate coefficients. Older mothers are significantly less likely to seek treatment in a health facility ($P(\beta_{\text{age}} < 0|\mathbf{y}) = 0.99$). Literate and illiterate mothers do not have significantly different probabilities of seeking health facility treatment ($P(\beta_{\text{literate}} < 0|\mathbf{y}) = 0.56$), but mothers with a regular partner have significantly higher probability ($P(\beta_{\text{regular partner}} > 0|\mathbf{y}) = 0.98$). Mothers surveyed before and after the intervention do not have significantly different probabilities after controlling for other factors ($P(\beta_{\text{time}} < 0|\mathbf{y}) = 0.98$).

Figure 3.20 (left) shows posterior distributions for the random effect of selected supervision areas. SA 118, with the lowest mean random effect $\xi_{118} = -0.99$, is in the catchment area of Minakulu Health facility: a small, rural health facility. SA 101, on the other hand, has the highest mean random effect $\xi_{101} = 0.85$; it is located in the catchment area of Boroboro Health facility, on the outskirts of the main town of Lira district, and is close to many government health facilities. Figure 3.20 (right) shows the posterior distribution for the standard deviation of these random

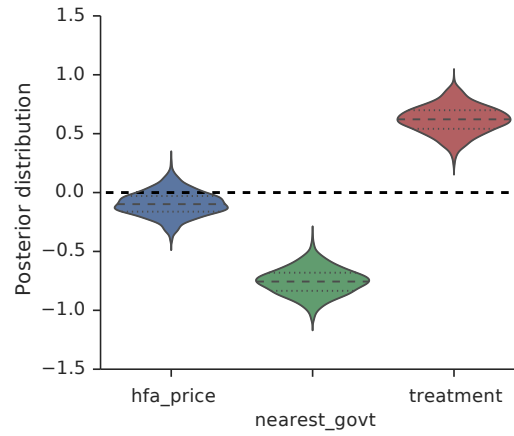


Figure 3.18: Marginal posterior density for coefficients on HFA covariates, showing the parameters for mean price paid by mothers visiting that health facility with sick children, distance to the nearest government health facility of level III or above, and the effect of a health facility being in the treatment group in 2014.

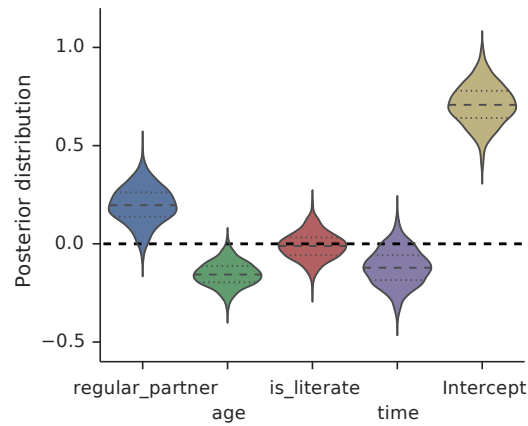


Figure 3.19: Marginal posterior density for coefficients on household-level covariates, showing parameters for regular partner, age, literacy, year surveyed, and the intercept.

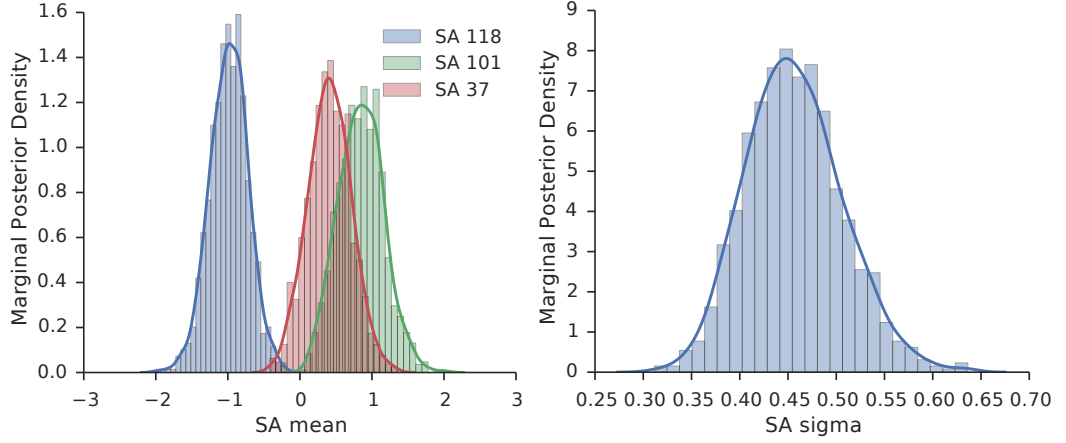


Figure 3.20: (Left) Selected marginal posterior density for Supervision Area-level random effect. Shown are the largest and smallest effects and a middle-range effect. (Right) Marginal posterior density for the standard deviation of supervision area random effects

Parameter	\hat{R} for cost	\hat{R} for objective quality
β_0	1.00	1.00
β_{treat}	1.00	1.00
β_{time}	1.00	1.00
σ	1.00	1.00
σ_{hf}	1.00	1.00

Table 3.6: Gelman-Rubin diagnostics for parameters in the two health facility difference models,, consistent with convergence.

intercepts, and shows that these random effects are similar in size to the treatment effect.

The ROC curve for this model is shown in Figure 3.21, with posterior predictive mean and 95% posterior predictive credible interval for the $AUC = 0.691 [0.685, 0.697]$, which again is similar to good epidemiological models of disease [70, 71].

Autocorrelation plots of the two health facility difference models are shown in Figure 3.22 for a thinning of 25. There is no significant autocorrelation at this level of thinning. The Gelman-Rubin statistics for the parameters of both models are given in Table 3.6, and are consistent with convergence.

The marginal posterior distributions for the regression parameters are shown in Figure 3.23. Treatment (i.e. converting to RBF funding) is associated with a significant decrease in cost, with $P(\chi_{treat} < 0 | \mathbf{x}_{cost}) = 0.98$, but not with a change in quality, with $P(\chi_{treat} < 0 | \mathbf{x}_{quality}) = 0.41$. There was a significant increasing time trend for cost, with $P(\chi_{time} > 0 | \mathbf{x}_{cost}) = 0.96$, but there is slightly less

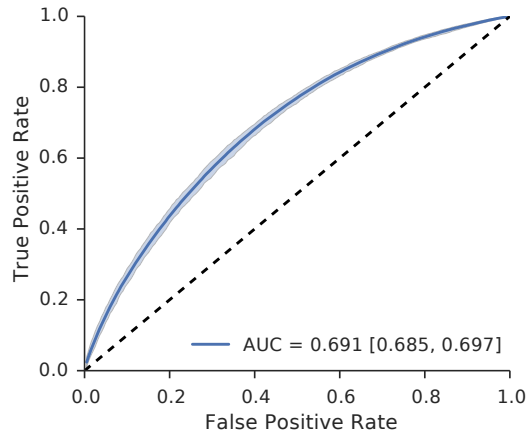


Figure 3.21: ROC curve for this model, with posterior predictive mean and 95% posterior predictive credible interval for the Area Under Curve $AUC = 0.691 [0.685, 0.697]$

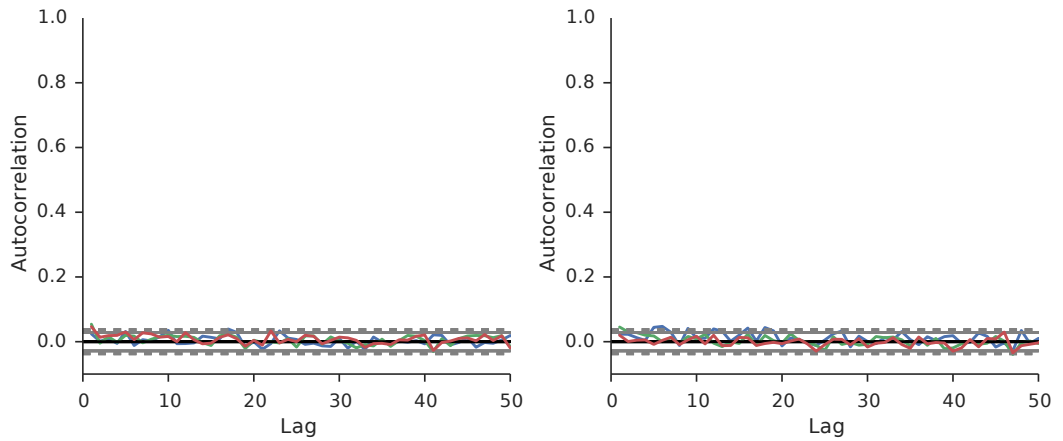


Figure 3.22: Autocorrelation of the beta parameters for the health facility difference models of cost (left) and objective quality (right) at a thinning of 25.

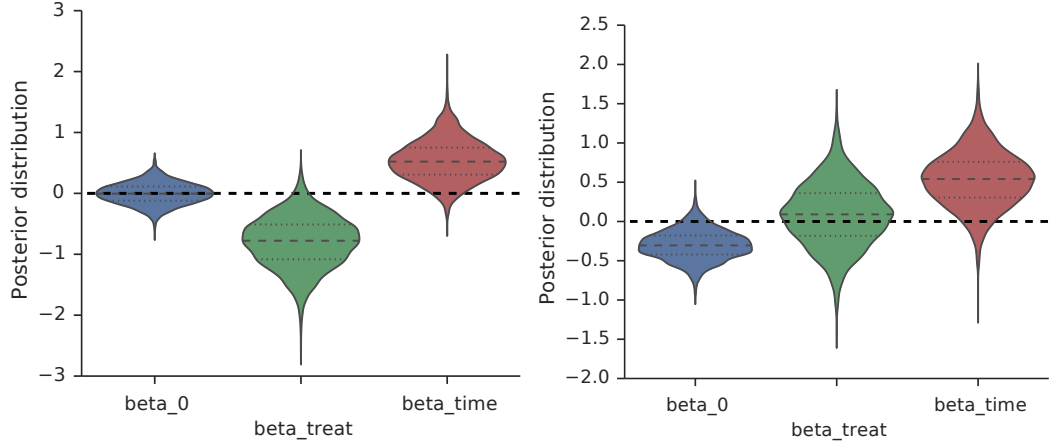


Figure 3.23: Marginal posterior distributions of the regression parameters for the two health facility difference models of cost (left) and objective quality (right).

evidence for an increasing time trend for quality with $P(\chi_{\text{time}} > 0 | \mathbf{x}_{\text{quality}}) = 0.94$.

The marginal posterior distributions for the standard deviation parameters are shown in Figure 3.24. There is more health-facility-level variation compared to time variation for quality, whereas the two sources of variation are similar for cost.

3.4.5 Discussion

These data take the form of a partially-controlled trial. By measuring the same health facilities and supervision areas before and after an intervention, we control for heterogeneities between facilities. One concern with this data is that the treatment and control groups were not picked at random - instead, health facilities were assigned based on their geographic region. Using a hierarchical model, however, controls for unknown factors between the two regions by including them in the supervision area effect ξ_k . The other potentially confounding factor is the interaction between time and geographic region: perhaps an unrelated change in Acholi region caused the observed increase in health facility use. This possibility is unlikely, since each PNFP health facility is administered separately; a regional time effect would require PNFPs to act together in each region but separately across regions. To rule it out would require further studies of a similar type but with treatment and control groups randomised across regions. Another limitation with this data is that the difference in total funding between input-based and RBF was not measured. This means that the RBF catchment areas could be performing better because they recieved more total funding than the input-based ones. The Basinga et al study [111] isolated the effect of RBF incentives from the effect of a general increase in resources and still

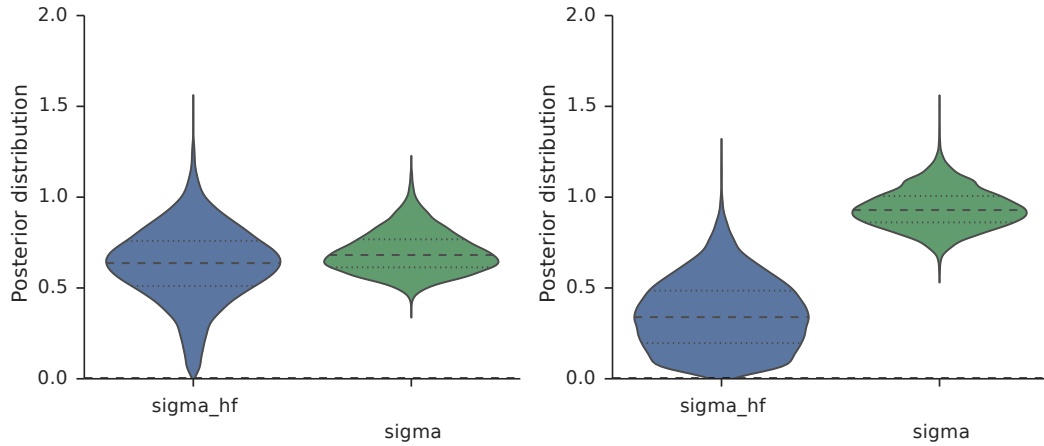


Figure 3.24: Marginal posterior distributions of the standard deviation parameters for the two health facility difference models of cost (left) and objective quality (right).

found a significant effect, so RBF does appear to have a real benefit over input-based funding.

For this section, unlike the previous section, visiting a pharmacy was not considered as visiting a health facility. This is because the funding change applied only to health facilities, not pharmacies. This is unlikely to make a significant difference to the results because only 24 out of 5760 (0.004%) of mothers took their child to a pharmacy in this dataset.

At the health facility level, these results show that distance to the nearest government health facility is strongly predictive of health facility usage. Many of the mothers surveyed took their children to a government health facility rather than their local PNFP, which may explain this result. In addition, this covariate may be a proxy for the size of the local PNFP's catchment area, and hence the distance of each household to their local PNFP.

The random-intercept effect appears to capture the additional variability in access to health facilities not already captured by the covariates explicitly included in the model, and will also include any social influence effect for each supervision area. Interestingly, the largest and smallest of the random effects are of similar magnitude to the treatment effect; although the effects of RBF are significant, it is clear that other factors such as social influence play an equally important role in determining usage.

The results from the health facility difference models show that health facilities had lower costs after a funding change to RBF, whereas health facilities without a funding change did not show a significant change. This may be one of the reasons

for the positive treatment effect of RBF, with lower costs allowing more mothers to take sick children to a health facility. There was no observable difference in the objective quality (percentage of correctly treated children) between RBF and non-RBF funded facilities, which provides more evidence for the conclusions of section §3.3 that objective quality does not seem to strongly influence health facility use. This highlights a possible weakness in the RBF method: although more sick children are being taken to health facilities in RBF catchment areas, if they are receiving poor quality treatment then disease outcomes may not be improving. Further work should investigate the direct effect of RBF on child mortality.

Even if cost does not have a strong influence on health facility usage, there is evidence from other studies that healthcare costs can put low-income families into debt and poverty [139, 140]; the strong evidence that overall costs had increased between 2012 and 2014 is therefore a worrying sign and highlights the need for health system interventions such as RBF which seem to reduce healthcare costs.

3.5 Summary

Understanding the health system barriers to positive health behaviours such as health facility usage, and how to remove these barriers through policy changes, could lead to significant improvements to health in developing countries. In this chapter we investigated a barrier that has previously been difficult to analyse: the quality of local health facilities. By using a Bayesian hierarchical framework we linked households to local health facilities, imputed the objective and subjective qualities of these health facilities, and investigated the association between local health facility quality and usage. We found that NGO-run health facilities were significantly more likely to be rated better by users, but were not significantly better when measured objectively. We also found that both higher objective and higher subjective quality was predictive of higher usage, but that this predictive power was only half that of including other unknown local factors such as social influence in the model. The poor quality of the data means that further studies of health facility quality and usage are required.

We also investigated the effect of a policy change on health facility usage behaviour, using quasi-experimental data collected before and after a change in the funding mechanism of health facilities. We show evidence that changing the funding mechanism of a health facility to results-based funding, in which health facilities receive extra funding for reaching specific targets, causes an increase in health facility usage for mothers of sick children who live in the catchment area of that health

facility. Other factors associated with higher usage were having a regular partner, lower age, and living nearer to a government facility. There was evidence that RBF caused a decrease in the average price paid by users of a health facility, but there was no evidence of a change in the quality of a health facility after RBF was introduced.

The Bayesian hierarchical modelling framework allows us to investigate health behaviours at both the individual and health system level. It is also extendable, and in future could include data from other levels such as local communities or between countries. This multi-level approach makes it possible to evaluate the complex interventions required to understand behaviour-changing policies.

Chapter 4

The spread of health behaviours

4.1 Introduction

Social influence, or the impact of others on our behaviour, plays an important role in understanding health behaviours. So far, in chapters 2 and 3, we assumed that social influence took place on a much faster timescale than the timescales of individual factors, such as education, and health system changes. In this chapter we take a different limit and assume that it is only social influence that affects the dynamics, and that the population is homogeneous. This allows the effect of covariates on behavioural change to be combined into a single rate which applies to all individuals. Using the definitions in chapter 1, we are considering:

$$\begin{aligned}\mathbf{B}(t, f(\mathbf{B}(t)), \mathbf{X}(t)) &\rightarrow \mathbf{B}(t, \mathbf{X}, f(\mathbf{B}(t))) \\ \mathbf{X}_i &= \mathbf{X}_0 \quad \forall i \in 1..N.\end{aligned}$$

The aim of this section is to develop a model for the social influence term $f(\mathbf{B}(t))$ based on the mechanism of ‘complex contagion’, and from there build a general model of the spread of behaviour through a homogeneous population. Various models of social influence have been proposed, but the difficulty of collecting data on behaviours outside the lab has means that these models have not been tested statistically. After investigating the complex contagion model and its deterministic limit, we statistically compare it to the popular ‘simple contagion’ model that is commonly used to model behaviour, using data on the spread of behavioural fads.

4.1.1 Social Influence

There is clear evidence that other people have an influence on our behaviour [141, 142, 143, 144], and health behaviours are no exception [145, 146, 147]. In the Health Belief Model (HBM) introduced in chapter 2, all the dimensions in the model are based on an individual’s perceptions of the susceptibility to disease, the severity of the disease, the benefits of a health behaviour, and the barriers to taking part in the behaviour [23]. Since the beliefs of others influence our own perceptions [148], it is likely that social influence plays a role in changing health behaviours. Social influence may also help to understand the “cue to action” required by the HBM to trigger a decision to change behaviour [23]. By understanding how social influence affects behaviours in a population, we can better predict how health interventions will improve health on a population level. We can also use ideas such as ‘nudges’ to encourage the spread of good health behaviours and reduce the spread of bad ones.

Not all health behaviours seem to be affected by social influence: Keating et al. [149] investigated cancer screening behaviour and found minimal evidence for social influence. There are many studies of other health behaviours, however, which do provide evidence for the importance of social influence on health behaviour. In particular, the work of Christakis and Fowler [147] analysed longitudinal social network and health data from the Framingham Heart Study and showed that if an individual had a friend, sibling, or spouse who had become obese in a given time interval then that individual was significantly more likely to also become obese. Similar results were also found when studying the cessation of smoking [146].

The work of Christakis and Fowler has proved controversial; it has been shown that social influence cannot be distinguished from homophily, or the clustering of individuals who are similar, in observational studies [150]. Aral et al. [151] try to determine an upper bound for the importance of social influence for behaviour spread, and find that for the adoption of a particular social media app at least half of the observed adoption events can be attributed to homophily. This discussion highlights the difficulty of using observational data to distinguish the effect of individual-level factors, in the form of homophily, from social influence. This same difficulty is not present in experimental data, however. Bond et al. performed a randomised controlled trial over Facebook to find evidence for social influence on the decision to vote. By sending direct messages to ‘seed’ nodes in a network, and then tracking the behaviour of their contacts, the experimenters showed that individuals were significantly more likely to vote if one of their close friends had received a message. In a study more closely related to health behaviour, Centola [145] placed individuals in an artificially-structured online community in which users were informed about the

health activities of their assigned contacts. This experiment showed that social signals significantly increased the likelihood of an individual taking part in a behaviour, and that up to three additional social signals significantly increased this likelihood even further. Taken together, these studies show that while individual-level factors are significant, as discussed in chapter 2, social influence is also important in determining health behaviours.

4.1.2 Previous models

Models of social influence have taken three main forms: experimental generalisations, agent-based models, and compartmental models. Experimental generalisations take historical data on the spread of a behaviour and try to find functional forms which match that data. One of the first examples of this approach was by Bass [152], who created a model of product adoption based on the idea of innovators and imitators. More recent attempts include fitting a variety of statistical distributions to the popularity of Internet memes [153]. The main disadvantage to this approach is that it does not provide a mechanistic model for social influence, and hence does not provide much insight into individual-level processes.

Agent-based models take almost the opposite approach to the experimental generalisations mentioned above, in that they simulate all of the individual- (or ‘agent-’) level processes occurring and then try to calibrate the model by matching the aggregate behaviour to data [154, 155]. Agent-based models are useful tools for reproducing the complex phenomena observed in real systems, but it is extremely difficult to fit their parameters to data well.

Compartmental models put each individual in the population into one of a certain number of states, or compartments. Only the number of individuals in each compartment and the transitions between them are tracked, and hence the number of dimensions of the system can be much less than an equivalent agent-based model. This in turn allows a compartmental model to be fitted to data more easily than agent-based models, while remaining a mechanistic description of the underlying system. Treating social influence in this compartmental way has a long history, an example being Dietz in 1967 [156] who developed a model for the spreading of rumours similar to models from epidemiology. In fact, much of the social influence literature using compartmental models has been based on the SIRS model of an epidemic. In the SIRS model there are three compartments: susceptible (S), infectious (I), and recovered (R). Susceptible individuals have not yet been infected with the disease, infected individuals currently have the disease and are spreading it, and recovered individuals have had the disease but are no longer spreading it. In

the canonical SIRS model, individuals move between these compartments with the following rates:

$$\begin{aligned}
(S, I) &\rightarrow (S - 1, I + 1) & \text{rate } \beta I \\
(S, I) &\rightarrow (S, I - 1) & \text{rate } \gamma \\
(S, I) &\rightarrow (S + 1, I) & \text{rate } \delta.
\end{aligned} \tag{4.1}$$

This canonical model can be modified by changing the functions for the rates, and by adding or removing compartments. For models of social influence on behaviour, the ‘infectious’ compartment represents individuals taking part in a behaviour and spreading it, and ‘recovered’ means the individual is no longer influencing others to take part in the behaviour. Many previous studies of social influence modify the canonical model by changing the rates at which individuals move between compartments. Isham et al. [157], for example, developed a model for rumours on a network based on the SIR model modified to include ‘stiflers’ who cause infectious individuals to recover at a faster rate.

Very few compartmental models for social influence modify the form of the infection term in the canonical model. However, as shown in experimental studies [145], there is significant evidence that the form of ‘infection’ in social influence is different to that in a biological epidemic. The important difference is the number of exposures to infection that an individual must receive before becoming infected: in biological infection only one source of infection is required for a non-zero probability of infection, whereas in social influence multiple sources are required. Dodds and Watts [158], for example, generalise the SIS model to allow for infection processes that require multiple exposures.

4.1.3 Population models

Modelling every individual in a population can lead to models with a large number of parameters, which real world datasets on behaviour generally cannot support. Instead, here we consider groups of individuals that are all in a particular state. Individuals change state with a certain rate. This is called ‘compartmental modelling’. A common application of this approach is to epidemic modelling, where individuals are grouped by their infection status. Perhaps the simplest example from epidemiology is the SIR model, which has three compartments: S, or susceptible, containing people who have never caught the disease; I, or infected, containing people who have caught the disease and are infectious; and R, or recovered, who caught the disease in the past but are no longer infectious and cannot catch the disease again. This

model, and extensions to it, have been very successful in capturing the dynamics of real epidemics [15]. One of the advantages of compartmental modelling is that additional compartments can be added to models in order to capture higher levels of detail. The SEIR model, for example, is an extension to SIR which gives an improved description of some diseases by including an ‘exposed’ period during which an individual is infected but is not infectious.

Compartmental models can be treated stochastically or deterministically. For small populations, randomness will have a large effect, and so we have to model the movement of individuals between each compartment using probability distributions. Individual changes of state are assumed to follow a Poisson Process, which means that the waiting time between moves follows an Exponential distribution.

4.1.4 Analysis methods

4.1.4.1 Gillespie Algorithm

One method for simulating these stochastic models is the Gillespie Algorithm [159], which is as follows:

Algorithm 4.1 The Gillespie Algorithm for simulating from a compartmental model

1. Add up the rates of all possible transitions from the current state, $R = \sum_i R_i$
 2. Draw time δt from an Exponential distribution with mean R^{-1}
 3. Calculate the fraction of each transition rate to the total rate, $r_i = \frac{R_i}{R}$
 4. Pick one transition to occur by drawing randomly from the possible transitions weighted by their rates.
 5. Update the states based on the picked transition, and set the time that this transition occurred to be at time δt after the last transition.
-

This is a general technique for simulating the time dynamics of continuous-time Markov chains, which we are considering in this chapter.

4.1.4.2 Bailey’s Method

Bailey’s method calculates the exact final size distribution for the SIR model [160]. Neuts and Li’s implementation [161] of it has been shown to be numerically stable and computationally efficient [160]. This method is convenient as it removes the noise in estimating final size distributions that comes from simulating using the Gillespie Algorithm. It can also be extended to systems with more compartments

and different transition rate functions [162], as long as the system has an absorbing state. The method relies on using the transition matrix to calculate the probability that the system will pass through a particular state, and performing this calculation in an order such that the probability of each state depends only on quantities that have already been calculated. Figure 4.1 explains this process diagrammatically. Figure 4.1 (left) shows one particular realisation of an SIR stochastic process with population size 9, showing the states visited in order and ending with 7 recovered individuals. Since this is an SIR process, infected individuals can never return to the susceptible state. As such, the only allowed transitions in this system are those that increase I by 1 and leave R unchanged, or those that decrease I by 1 and increase R by 1. This allows the calculation order shown in Figure 4.1 (right), which starts from the initial state with probability 1. For the canonical SIR model, the probability of states $(I_0 + 1, R_0)$, $(I_0 + 2, R_0)$, \dots , (N, R_0) are as follows:

$$\begin{aligned}
P(I_0, R_0) &= 1 \\
P(I_0 + 1, R_0) &= \frac{\beta P(I_0, R_0)}{\beta P(I_0, R_0) + \gamma P(I_0, R_0)} \\
P(I_0 + 2, R_0) &= \frac{\beta P(I_0 + 1, R_0)}{\beta P(I_0 + 1, R_0) + \gamma P(I_0 + 1, R_0)} \\
&\vdots
\end{aligned}$$

Once the probabilities of visiting each state in column R_0 have been calculated, there is enough information to determine the probability of state $(I_0 - 1, R_0 + 1)$. From this, values in column $R_0 + 1$ can be calculated:

$$\begin{aligned}
P(I_0 - 1, R_0 + 1) &= \frac{\gamma P(I_0, R_0)}{\beta P(I_0, R_0) + \gamma P(I_0, R_0)} \\
P(I_0, R_0 + 1) &= \frac{\gamma P(I_0 + 1, R_0)}{\beta P(I_0 + 1, R_0) + \gamma P(I_0 + 1, R_0)} + \\
&\quad \frac{\beta P(I_0 - 1, R_0 + 1)}{\beta P(I_0 - 1, R_0 + 1) + \gamma P(I_0 - 1, R_0)} \\
P(I_0 + 1, R_0 + 1) &= \frac{\gamma P(I_0 + 2, R_0)}{\beta P(I_0 + 2, R_0) + \gamma P(I_0 + 2, R_0)} + \\
&\quad \frac{\beta P(I_0, R_0 + 1)}{\beta P(I_0, R_0 + 1) + \gamma P(I_0, R_0)} \\
&\vdots
\end{aligned}$$

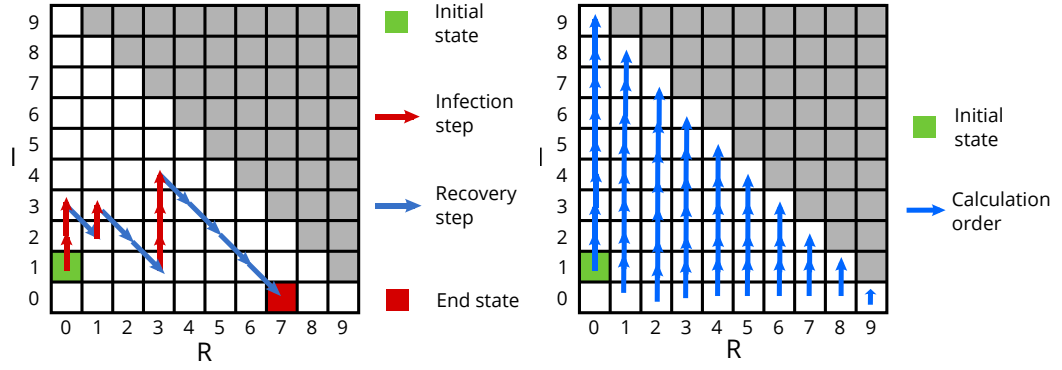


Figure 4.1: One realisation of an SIR stochastic process with population size 9, showing the states visited in order (left). Calculation order such that the probability of visiting each state depends only on previously-calculated states (right).

The calculation then proceeds along columns of R until $R = N$, and the probability of visiting all allowed states has been calculated. The final size distribution is then given by:

$$P(Z) = P(I = 0, R = Z).$$

4.2 General complex contagion model

4.2.1 Introduction

We will treat the change in the population's behaviour over time as an infectious process, with individuals becoming 'infected' with the behaviour of their social contacts. An individual's behavioural state can be one of the following: S (susceptible), not taking part in the behaviour; I (newly infectious), recently converted to the behaviour and is currently trying to spread the it; J (infectious), still spreading the behaviour and has been doing so for some time; and R (recovered), previously taking part in the behaviour but is no longer spreading it. By ignoring individual covariates, we can split the whole population into compartments depending on their behavioural state. Within each compartment individuals are interchangeable and are not separately tracked. The numbers of individuals in each compartment are defined as follows:

$$\begin{aligned}
S &= \sum_i \mathbf{I}_{\{B_i(t)=S\}} \\
I &= \sum_i \mathbf{I}_{\{B_i(t)=I\}} \\
J &= \sum_i \mathbf{I}_{\{B_i(t)=J\}} \\
R &= \sum_i \mathbf{I}_{\{B_i(t)=R\}}.
\end{aligned}$$

We can model behaviour at the population level by specifying the rates at which individuals move between these compartments. In this section we extend the work of House [163], who introduced a population-level form of complex contagion, to include multiple infectious compartments with different thresholds. We then insert this form of complex contagion into a general stochastic compartmental model based on an SIRS epidemic to create a general model of social contagion applicable to many spreading situations.

Using the Gillespie Algorithm to simulate realisations and an implementation of Bailey’s Method to calculate the final size distribution, we study a special case of the general model that represents the spreading of ‘fads’: short-lived behaviours which rapidly take over a large fraction of the population and eventually die out. We find that, for complex contagion, chance is very important in determining the success of a fad. We also find that a small ‘nudge’ (convincing a small number of individuals to take part in a fad who would not otherwise have done so) is sufficient to significantly change the success of the fad.

4.2.2 Social influence term

For this work we use a simple model for complex contagion that captures the nonlinear response to the number infected contacts while remaining analytically tractable. Each individual canvasses C contacts, and if the number of these contacts taking part in a behaviour is greater than some threshold τ then the individual becomes infected with some rate:

$$D(k|\tau) = \begin{cases} 0 & k < \tau \\ \beta & k \geq \tau, \end{cases}$$

where k is the number of infected contacts and τ is the threshold for behavioural spread.

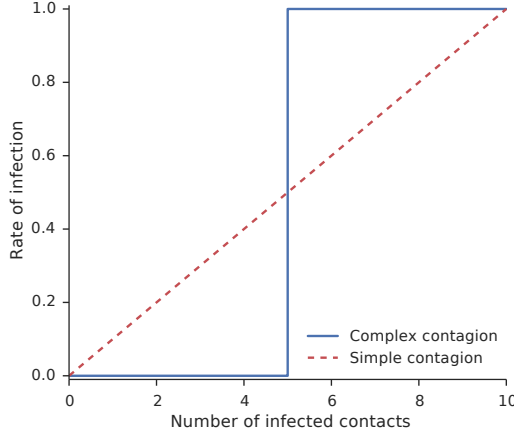


Figure 4.2: Individual infection rate for simple and complex contagion with $\beta = 1$, $\tau = 5$, and $C = 10$. Simple contagion increases linearly with number of infected contacts, whereas this complex contagion model is a step function with threshold τ .

Figure 4.2 compares the transmission rates for simple and complex contagion as a function of the number of infected contacts. For simple contagion, each additional infected contact gives the same increase to the transmission rate, whereas complex contagion does not spread until the threshold is reached.

We follow House [163], and treat this mechanism at the population level. For simplicity we assume an all-to-all contact network, and that all individuals have the same values for C and τ . Repeated exposures to the same infected individual each count towards reaching the threshold for infection, and as such the C contacts are chosen randomly with replacement from the population. Under these conditions, the probability p of an individual having k infectious contacts is given by a Binomial distribution:

$$p(k|I, N, C) = \text{Bin} \left(k \middle| \frac{I}{N}, C \right),$$

where I is the number of infectious individuals in the population and N is the total population size.

The infection rate is therefore given by the total probability that there are at least τ infected individuals in the canvas group:

$$f(I) = \beta \sum_{k=\tau}^C \text{Bin} \left(k \middle| \frac{I}{N}, C \right). \quad (4.2)$$

Figure 4.3 shows the population-level behaviour of the infection rate. Whereas for simple contagion the infection rate increases linearly with the fraction of infected

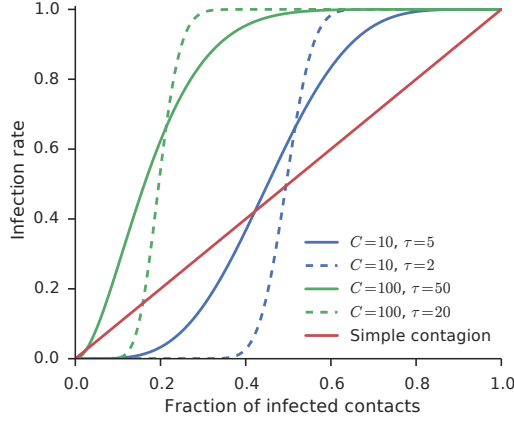


Figure 4.3: Population-level infection rate for simple contagion and four complex contagions with different parameters. The simple contagion term is again linear, but each complex contagion has a sigmoidal infection function which is steeper as C is increased.

contacts, complex contagion has sigmoidal infection rate with the point of inflection at the position of the threshold. Increasing C creates a sharper transition at the threshold point.

The Binomial model given above sufficient for a single complex contagion compartment, but if there are two or more of these compartments then the thresholds must be combined in some way. The simplest model is to treat the contribution of each compartment to its threshold separately, and then calculate the total probability of either of the compartments reaching its threshold. If there are two infectious compartments, the total probability that at least one of the thresholds is reached is given by the following:

$$\begin{aligned}
 f(I, J) = & \beta \sum_{k=\tau_i}^C \sum_{l=0}^C \text{Multi}(k, l | \frac{I}{N}, \frac{J}{N}, C) \\
 & + \beta \sum_{k=0}^{\tau_i} \sum_{l=\tau_j}^C \text{Multi}(k, l | \frac{I}{N}, \frac{J}{N}, C).
 \end{aligned} \tag{4.3}$$

The probability surface for this model is shown in Figure 4.4 for $C = 10$ and $C = 100$, and with $\tau_i = \frac{2}{10}C$ and $\tau_j = \frac{4}{10}C$. As with the single-threshold model, higher values of C lead to sharper transitions between high and low probability. This figure also compares the Multinomial model with an alternative model for multiple compartments with the same threshold, which sums the fraction of the population

in all complex contagion compartments and uses this total infectious fraction in equation (4.2). This model is not easily generalisable to multiple thresholds. Other models for the contributions of the infectious compartments to the threshold could be developed, but the model given in equation (4.3) provides an initial approach.

Much of the analysis for this chapter fixes the number of contacts an individual canvasses to $C = 10$. In part, this number was chosen for computational reasons: higher values of C increase the time taken to calculate the social influence term. This value is also within the range of 10-15 ‘sympathy’ contacts identified by Hill and Dunbar in real human social networks [164]; these contacts are outside immediate family but nonetheless have a strong connection with the individual.

4.2.3 Stochastic Model

The complex contagion term given above was inserted into a compartmental model framework to provide a general model of behaviour spread through a population. We based this general model on the SJIRS model of epidemic spread, and then applied it to specific spreading situations by setting particular parameters to 0 or taking them to infinity.

We consider a fixed population size $N = S + I + J + R$. The transition rates for state (S, I, J) are given by:

$$\begin{aligned}
(S, I, J) &\rightarrow (S - 1, I + 1, J) & \text{rate } f(S, I, J) \\
(S, I, J) &\rightarrow (S, I - 1, J + 1) & \text{rate } g(S, I, J) \\
(S, I, J) &\rightarrow (S, I, J - 1) & \text{rate } h(S, I, J) \\
(S, I, J) &\rightarrow (S + 1, I, J) & \text{rate } r(S, I, J).
\end{aligned} \tag{4.4}$$

This model is ‘SIRS-like’, but if $r \rightarrow \infty$ it becomes ‘SIS-like’, and if $r \rightarrow 0$ it becomes ‘SIR-like’. The general model can therefore be tailored to fit many spreading situations.

The infection rate and decay rates have the following general forms:

$$\begin{aligned}
f(S, I, J) &= \sum_{x=0}^C \sum_{y=0}^C p(x, y, z | S, I, J, N, C) F(x, y, z | \theta) \\
h(S, I, J) &= \sum_{x=0}^C \sum_{y=0}^C p(x, y, z | S, I, J, N, C) G(x, y, z | \theta^*),
\end{aligned} \tag{4.5}$$

where $p(x, y | \dots)$ is the probability that of the C contacts, x are in state S , y are in

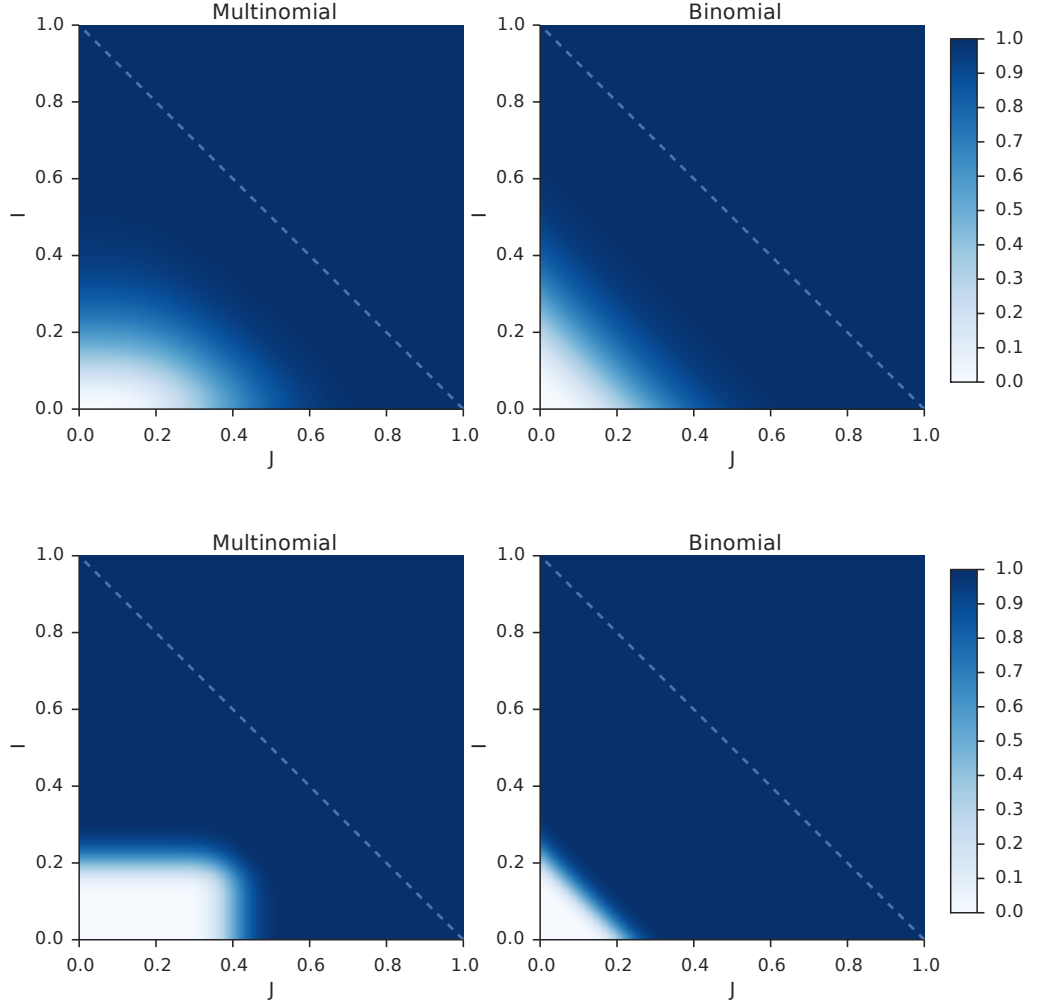


Figure 4.4: Comparison of Multinomial and Binomial social term for two values of C . The Multinomial term can include multiple types with different thresholds. For the Multinomial model, there is an extra region of low probability in the region where the total infected fraction ($I + J$) is above the lowest threshold, but neither threshold individually has been reached. Parameters are $C = 10$, $\tau_I = 2$, $\tau_J = 4$ (left) and $C = 100$, $\tau_I = 20$, $\tau_J = 40$ (right).

state I , and z are in state J . $F(x, y, z|\theta)$ is the social influence term which models the probability of taking part in a behaviour if x contacts are in state S , y are in state I , and z are in state J . $G(x, y|\theta^*)$ models the probability of stopping a behaviour if x contacts are in state S , y are in state I , and z are in state J .

In this chapter we concentrate on one specific spreading situation: behavioural ‘fads’ that briefly flourish and then die out. In this situation, fads are guaranteed to eventually cease to spread in any closed population. We therefore set rate $r = 0$ in equation (4.4), which prevents recovered individuals from becoming reinfected. Later in this chapter we consider the behaviour spreading to separate sub-populations.

In this model, both the infection rate and the recovery rate are complex contagions. Individuals may be infected by canvassed individuals in either of the infectious compartments (I or J), but are more strongly influenced by newly infectious contacts; the number of I contacts required for infection is therefore lower than the required number of J contacts. The existence of two different thresholds requires us to use the Multinomial form of the complex contagion mechanism, so the infectious transition rate is given by equation (4.3). The recovery complex contagion rate depends on only one compartment, the number of infectious individuals, and is therefore given by equation (4.2). Individuals are assumed to transition from the newly-infectious state to the infectious state at a constant rate, ϵ .

Substituting these rates into the general spreading model gives the following system:

$$\begin{aligned}
(S, I, J) \rightarrow (S - 1, I + 1, J) \quad \text{rate} \quad & \beta \sum_{k=\tau_i}^C \sum_{l=0}^C \text{Multi} \left(k, l \middle| \frac{I}{N}, \frac{J}{N}, C \right) \\
& + \beta \sum_{k=0}^{\tau_i} \sum_{l=\tau_j}^C \text{Multi} \left(k, l \middle| \frac{I}{N}, \frac{J}{N}, C \right); \\
(S, I, J) \rightarrow (S, I - 1, J + 1) \quad \text{rate} \quad & \epsilon; \\
(S, I, J) \rightarrow (S, I, J - 1) \quad \text{rate} \quad & \gamma \sum_{y=\tau_\gamma}^C \text{Bin} \left(y \middle| \frac{J}{N}, C \right). \tag{4.6}
\end{aligned}$$

This system was simulated using the Gillespie algorithm, and the final size probability distribution (number of recovered individuals when the epidemic dies out) calculated to machine precision using an implementation of Bailey’s method.

4.2.4 Stochastic results

The fad model was simulated for a population size of $N = 500$. Multiple realisations for the same set of parameters are shown in Figure 4.5a with $\tau_j = 5$, and in Figure 4.5b with $\tau_j = 3$. Several of the realisations died out without infecting a significant fraction of the population. Others infected a large fraction, and all of these had very similar timeseries once they had reached a particular threshold. If shifted so that the peak times of each realisation coincided, the timeseries of all the successful fads were extremely similar. The timeseries for each successful fad is characterised by a slow initial increase until a certain threshold is reached, followed by a sharp spike in the number of infected individuals. There is then a fast decay just after the peak which gradually slows until the fad has died out. The effect of decreasing τ_j is to allow more fads which do not spread significantly at the beginning of the fad to eventually reach the threshold for the infectious state and eventually spread to a large fraction of the population.

4.2.5 Analysis using Bailey’s method

Although the Gillespie Algorithm be used to simulate realisations from the social influence compartmental model, it is inefficient for investigating ensemble properties of the model such as the final size distribution. The Gillespie method scales poorly with population size, and many realisations are required to estimate ensemble properties with sufficiently low Monte Carlo error. Instead, we investigated the effect of different parameter values and of ‘nudges’ on the final size distribution of the fad model using an implementation of Bailey’s method.

We extended Neuts and Li’s implementation of Bailey’s method [161] to include the extra newly-infectious compartment and modified the infection and recovery rates to represent complex contagion. Figure 4.6 compares the final size distribution calculated using the new implementation of Bailey’s method with that estimated from 10^5 realisations simulated using the Gillespie Algorithm. Both estimates are very similar, but small differences due to noise can still be observed even with a large number of iterations. Figure 4.7 shows how the time taken to calculate the final size distribution using Bailey’s method scales with population size for this model. Many steps in the calculation involve the same quantities, so this implementation uses caching to store previous function evaluations and avoid recalculating the same value multiple times.

Figure 4.8 shows the relationship between β and ϵ in determining the extent of the fad. It plots the probability that more than half the population take part in

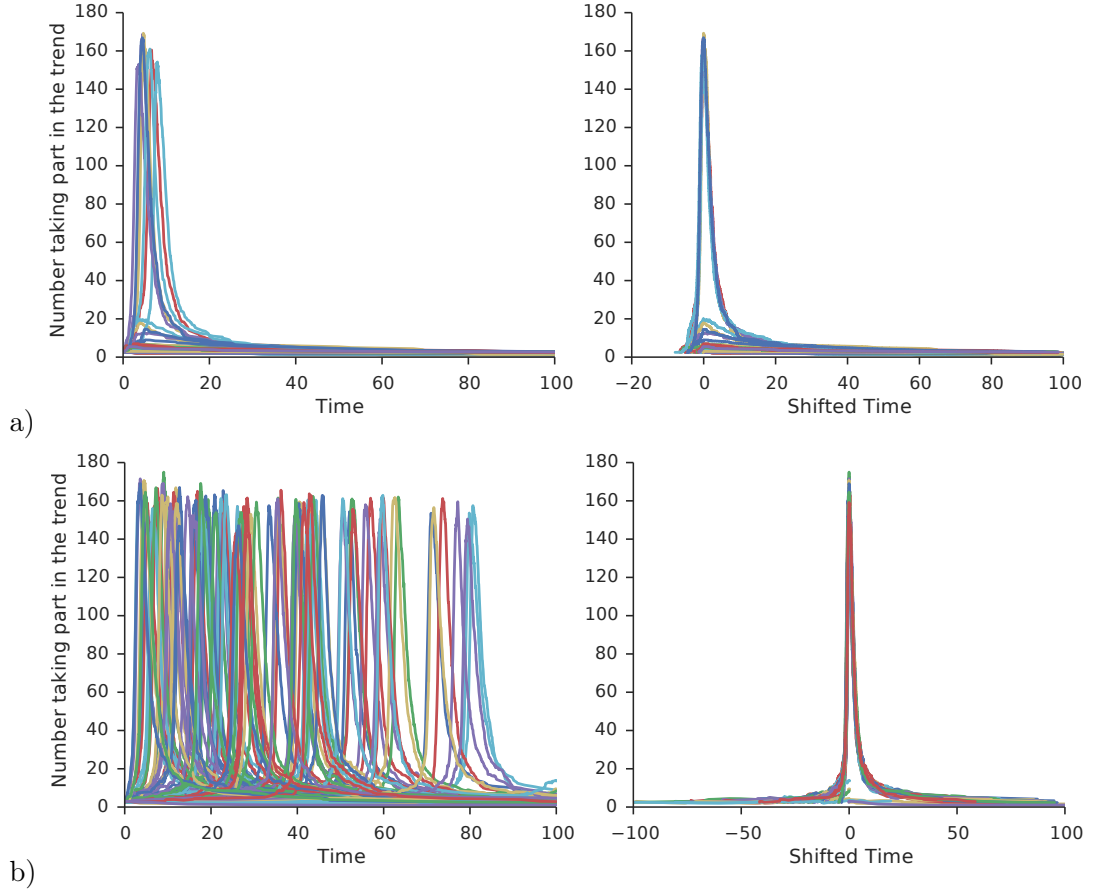


Figure 4.5: Realisations of the stochastic fad system (left) and the same realisations shifted on the time axis to have the same peak time (right), for different values of τ_j : a) $\tau_j = 5$ b) $\tau_j = 3$. Other parameters are: $N = 500$, $\beta = 1.2$, $\epsilon = 1$, $\tau_i = 2$, $\tau_r = 2$, $I(0) = 5$. All successful fads have very similar shape and timescale

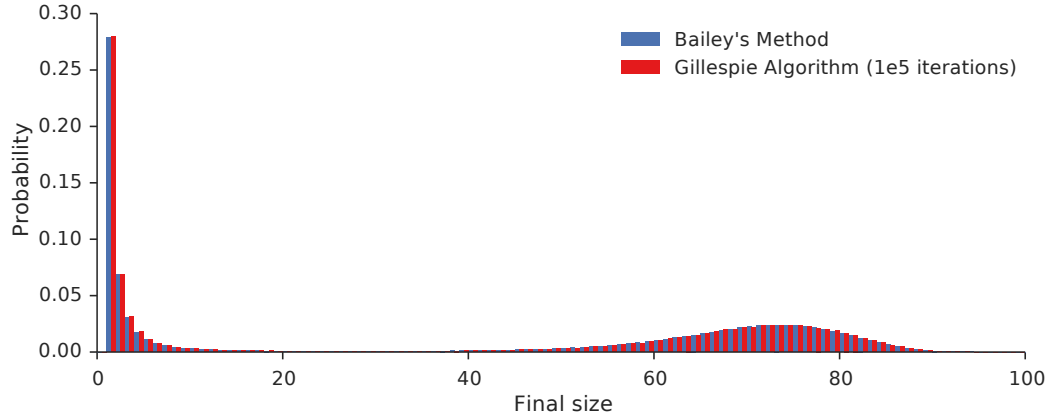


Figure 4.6: Comparison of the final size distribution calculated using Bailey's Method to the distribution estimated by 1×10^5 simulations using the Gillespie Algorithm. Parameter values are: $N = 100$, $\beta = 0.2$, $\epsilon = 1$, $\tau_i = 1$, $\tau_j = 5$, $\tau_r = 4$, $I(0) = 1$. The two distributions match well, although small variations are noticeable due to Monte Carlo error in the simulation estimate.

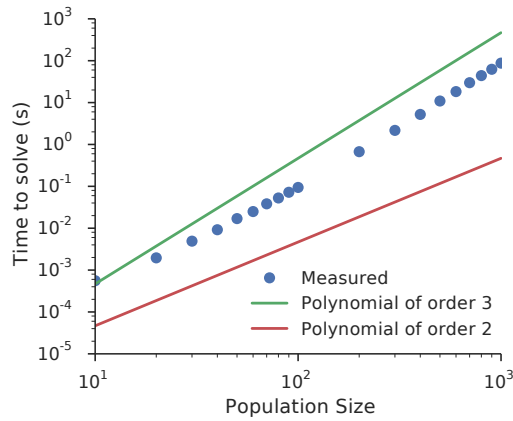


Figure 4.7: Time taken to calculate the final size distribution using this implementation of Bailey's Method against population size N . The time scales as a polynomial between order 2 and order 3.

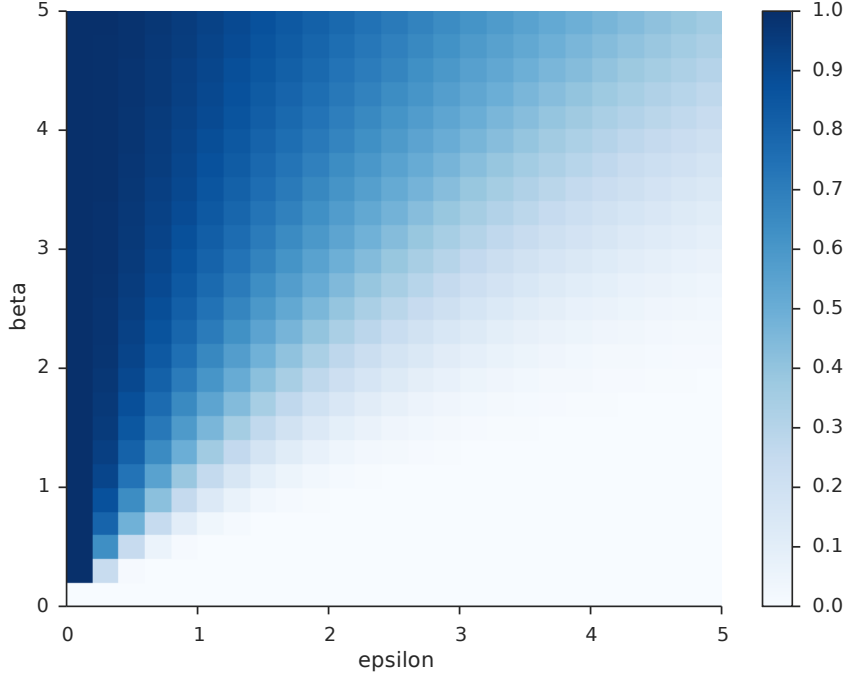


Figure 4.8: Probability of more than 50% of the population taking part in the fad for each value of β and ϵ . Other parameter values are: $N = 100$, $\tau_i = 2$, $\tau_j = 5$, $\tau_r = 3$, $I(0) = 2$.

the fad based on the value of β and ϵ . As ϵ increases, a higher value of β is required for the fad to spread to more than half the population. The boundary (β, ϵ) line below which the fad does not spread appears to be almost linear. Figure 4.9 is a similar plot showing the relationship between β and $I(0)$. Below a certain value of β , in this case around $\beta < 0.5$, the fad will die out whatever the initial value of I . As β increases, a lower initial value of I is required for the fad to spread to most of the population.

Figure 4.10 shows the effect of moving one individual from the S compartment to the I compartment, by comparing the final size probability with and without this ‘nudge’ at the beginning of the fad. Without the nudge, the fad is guaranteed to die out. With the nudge, however, there is a small but non-zero probability that almost the whole population will be affected by the fad. This abrupt change from die-out to widespread uptake is not present in systems with simple contagion.

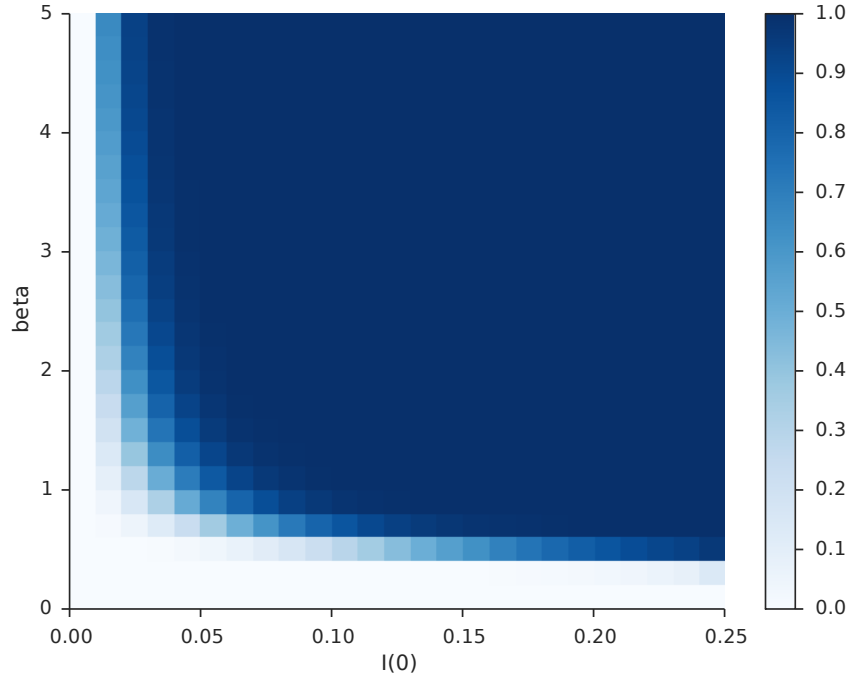


Figure 4.9: Probability of more than 50% of the population taking part in the fad for each value of β and $\frac{I(0)}{N}$. Other parameter values are: $N = 100$, $\epsilon = 1$, $\tau_i = 2$, $\tau_j = 5$, $\tau_r = 3$, $I(0) = 2$.

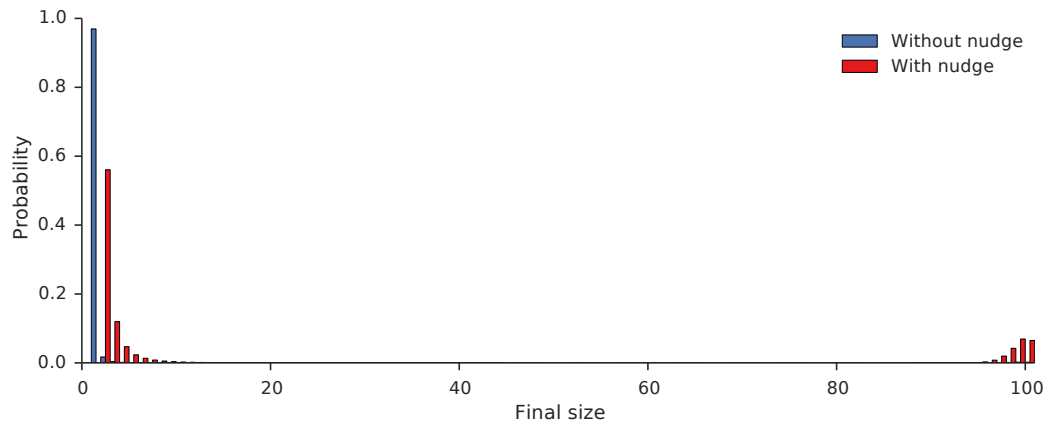


Figure 4.10: Comparison of the final size distributions for fads with identical parameters before and after one individual has been moved from the S compartment to the I compartment. This small ‘nudge’ leads to a significant chance of the fad infecting the entire population.

4.2.6 Discussion

Using the social influence terms given in equation (4.2) and equation (4.3), we can capture the mechanism of complex contagion within a compartmental model. The multinomial form given in equation (4.3) allows for multiple infectious compartments with different thresholds. This form relies on each compartment reaching its threshold separately, with individuals in one compartment not contributing towards reaching the threshold of another. This may be an unrealistic assumption: it is more likely that the total number of infectious individuals also has a threshold. The general features, however, are robust to changes in the form of the complex contagion term if it has a threshold property. Further work is required to investigate the effect of canvas size C on the spread of fads, but the ratio of the thresholds τ_i , τ_j and τ_r to C seems to be more important in determining the general form of the spread.

We have also introduced a general compartmental model which can include a variety of social influence terms, including complex contagion. By taking the limit of the system as some parameters go to either zero or infinity, the model can be used to represent a range of situations in which social influence is important. In particular, we investigate the limit in which ‘recovered’ individuals cannot become re-infected; this corresponds to a ‘fad’, which is briefly popular before eventually dying out. The presence of an absorbing state allows us to implement Bailey’s Method for the fad system, providing additional insight into the social contagion process.

Simulating the fad system using the Gillespie Algorithm showed that a fad’s popularity not only depends on the parameters of the system (such as the ‘infectiousness’ of the fad), but also strongly depends on chance. Realisations of the stochastic system demonstrate high variability and low predictability. Another feature of the model is that those fads which do affect the majority of the population have very similar timeseries shapes, but take become popular at different times. This implies that certain features of fads are predictable once they have been identified as popular; in particular, the time before the fad dies out should be very similar for fads with the same parameters.

Stochastic simulations provide insight into the timeseries of fads, but the stochasticity makes this method inefficient for investigating the effect of each parameter on the fad. Instead, an implementation of Bailey’s method was used to efficiently calculate the final size probability distribution for given parameters. In particular, in Figure 4.8 there appears to be an almost linear boundary between the infectivity (β) and the loss of enthusiasm parameter (ϵ) below which the trend does not spread. There is also a fast transition in the probability of spread for values of β and $I(0)$, shown in Figure 4.9: a small change in either can lead to a large change

in the probability.

Bailey’s method was also used to investigate the effect of a small nudge to the system on the final size distribution. A nudge of one individual moved from the susceptible to the infectious state is sufficient to turn a fad that was guaranteed to die out to one that has a non-zero probability of spreading to the whole population.

In summary, the general social influence model explains the main characteristics of social influence identified by Salganik et al. [141]: unpredictability of which behaviours will succeed, and large variability in the number of people affected by the behaviour. One version of this model, the fad model, also captures the phenomenon of ‘boredom’ with a fad observed by Denrell and Kovacs [165]. Finally, we can use Bailey’s method to understand control measures for behaviour spread; specifically, that only a small ‘nudge’ is required to make a behaviour either die out or take over the population.

4.3 Deterministic limit

4.3.1 Introduction

The stochastic model described in the previous section is important for understanding the effect of chance on the spread of behaviours. A full stochastic model, however, has three main disadvantages: simulation is computationally intensive, many realisations are required to reduce noise, and methods for fitting these models to data are advanced and experimental. In this section we consider the general model, equation (4.4), in the limit of large population size N . This allows us to treat the model as a system of ODEs, and apply methods of analysis that would not be available for the stochastic system. For many real-world situations the population size is large enough to make this a valid approximation.

4.3.2 Model

The stochastic model equation (4.5) can be approximated by the following system of ODEs [166, 167], with error $O\left(N^{-\frac{1}{2}}\right)$ where N is the population size:

$$\begin{aligned}
\frac{dS}{dt} &= -f(S, I, J)S + r(S, I, J)R \\
\frac{dI}{dt} &= f(S, I, J)S - g(S, I, J)I \\
\frac{dJ}{dt} &= g(S, I, J)I - h(S, I, J)J \\
R &= 1 - S - I - J,
\end{aligned} \tag{4.7}$$

where we have redefined $\frac{S}{N} \rightarrow S$, $\frac{I}{N} \rightarrow I$, $\frac{J}{N} \rightarrow J$, and $\frac{R}{N} \rightarrow R$ when converting from the stochastic model to the ODE model.

For the fad model, this system is reduced to the following:

$$\begin{aligned}
\frac{dS}{dt} &= -f(I, J)S \\
\frac{dI}{dt} &= f(I, J)S - \epsilon I \\
\frac{dJ}{dt} &= \epsilon I - h(J)J \\
R &= 1 - S - I - J,
\end{aligned} \tag{4.8}$$

where

$$\begin{aligned}
f(I, J) &= \beta \sum_{k=\tau_i}^C \sum_{l=0}^C \text{Multi}(k, l | I, J, C) \\
&\quad + \beta \sum_{k=0}^{\tau_i} \sum_{l=\tau_j}^C \text{Multi}(k, l | I, J, C)
\end{aligned} \tag{4.9}$$

and

$$h(J) = \sum_{y=\tau_\gamma}^C \text{Bin}\left(y \middle| \frac{J}{N}, C\right). \tag{4.10}$$

4.3.3 Results

We confirm that this ODE system approximates the stochastic system in Figure 4.11. If N is large enough (typically $N \gtrsim 10^4$), then the mean of an ensemble of stochastic realisations lies very close to the timeseries predicted by the ODE model. Interestingly, due to the property of the stochastic fad model that realisations which spread

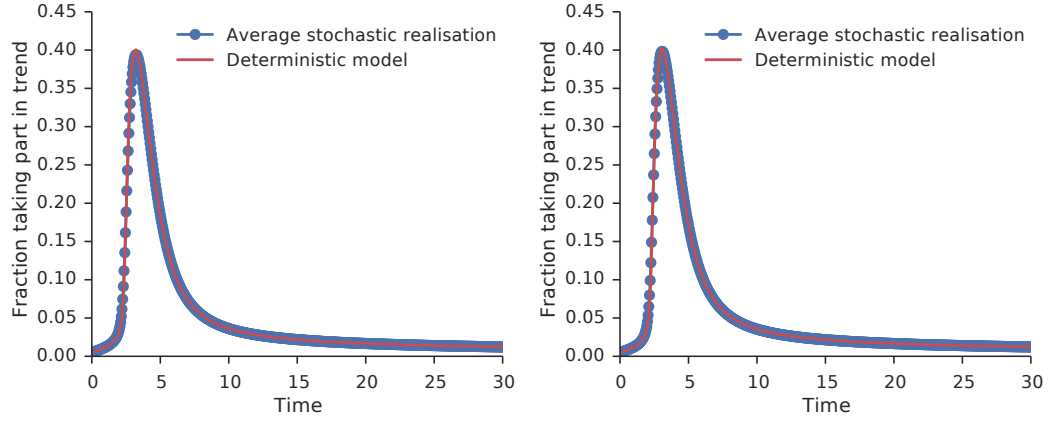


Figure 4.11: Comparison of stochastic and deterministic models for two values of τ_j . Stochastic model is averaged over 100 realisations. Stochastic parameters are $N = 2 \times 10^5$ and $I(0) = 0.01N$. The other parameters are $\beta = 2.7$, $\epsilon = 1$, $\tau_i = 2$, $I(0) = 1$. $\tau_j = 5$ (left) $\tau_j = 3$ (right).

all share a similar shape, individual realisations of the fad model which survive are also very close to the ODE timeseries.

Figure 4.12 explores the effect of parameters β , ϵ , τ_i and τ_r on the time course of the behavioural fad, with all other parameters held constant. There is a sharp threshold for β , below which fads die out and above which they take over a significant part of the population. For fads which do take off, the effect of β is to make the peak in the fad occur faster and die quicker. Increasing the rate of enthusiasm decay parameter, ϵ , has the effect of compressing the width of the peak. Above a certain threshold it also stops the spread of the fad to the majority of the population. The effect of increasing τ_i is to increase the time before the fad spreads to a majority of the population, until reaching a value at which the fad dies out immediately. The effect of increasing τ_r is to increase the time it takes for the fad to die out after the peak has been reached; it does not have an effect on the start of the fad.

The initial number of infectives, $I(0)$, is very important for determining whether the fad will spread. In the example given in Figure 4.13, a difference of 0.001 is sufficient to make one fad take over the population and the other to die out. For simple contagion, however, significant spreading is determined by the parameters rather than by the fraction of initial infectives, and the final size of the epidemic increases smoothly with this initial fraction.

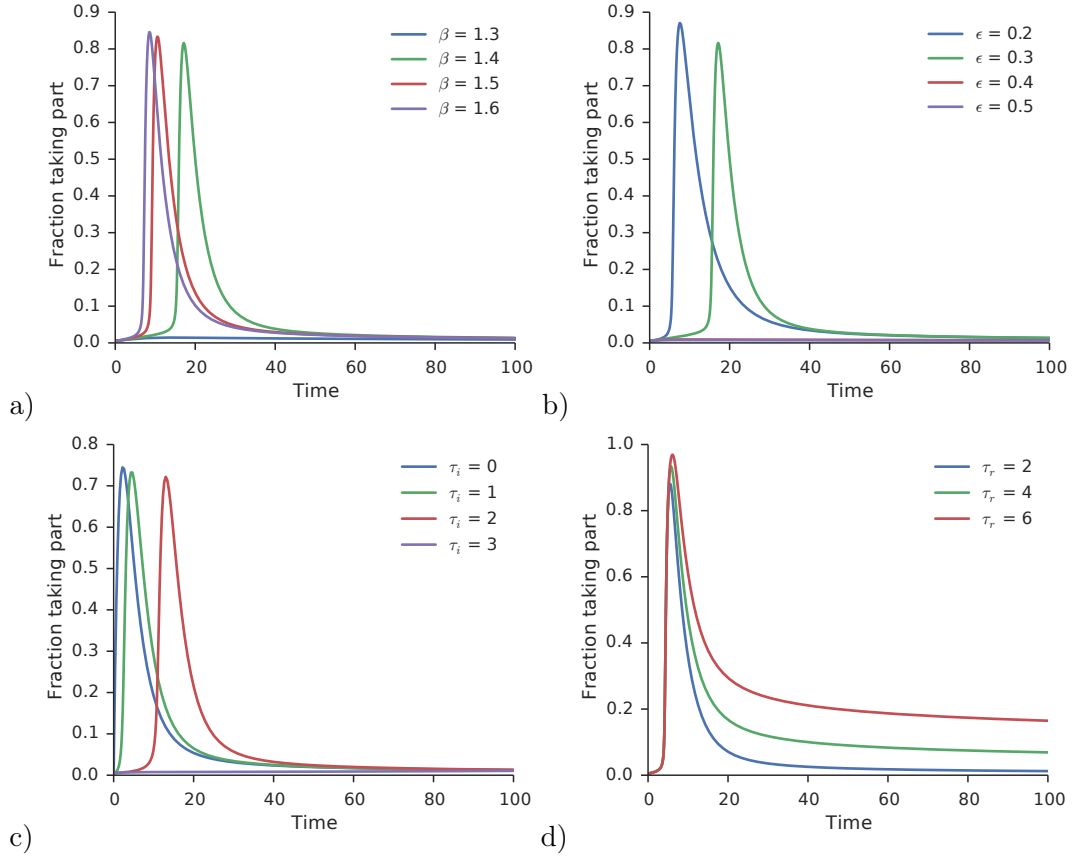


Figure 4.12: Effect of a) β , b) ϵ , c) τ_i , and d) τ_r on the solution to the ODE system. The two continuous parameters, β and ϵ , both exhibit a sharp threshold between fads that are successful and those that are not. The thresholds τ_i and τ_r control the beginning and end of the fad respectively.

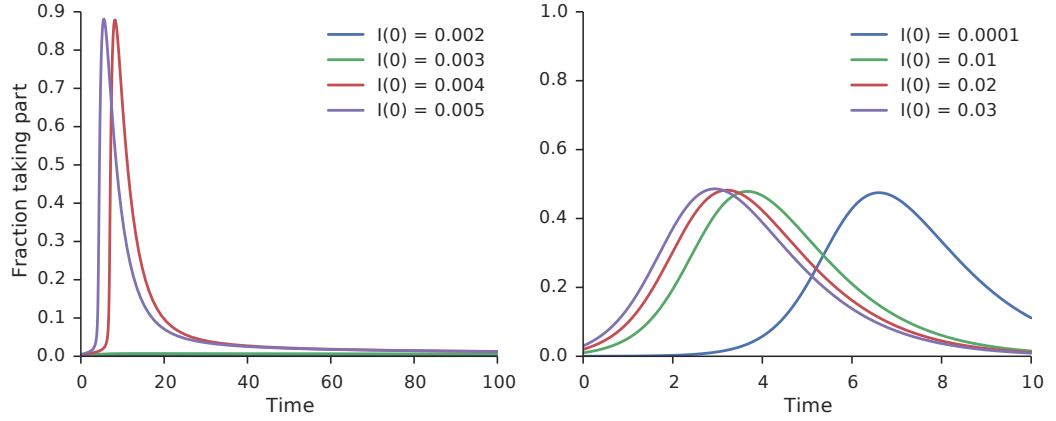


Figure 4.13: The effect of initial infected fraction on complex contagion (left) and simple contagion (right). For complex contagion, $I(0)$ is very important for determining the success of the fad; for simple contagion $I(0)$ does not affect the success of the fad.

4.3.4 Discussion

Approximating the general stochastic social influence model as an ODE system is appropriate when considering population sizes $N \gtrsim 10^4$, and it provides additional tools to analyse the spread of behaviours. In particular, it shows the effect of the parameters on the timeseries of the spread more efficiently than in the stochastic system. The continuous parameters β , ϵ and $I(0)$ have critical values which separate fads which spread and those which immediately die out, and the threshold parameters control either the start or the tail of the fad.

While the general ODE model captures the average behaviour of the general stochastic model, in the case of the fad model even the behaviour of individual realisations of the process are captured by the ODE. This is encouraging, as it implies that if it is known that a fad is starting to become popular in a large population then the ODE model will predict its behaviour.

4.4 Evidence for complex contagion

4.4.1 Introduction

The ODE model introduced in the previous section has the advantage that it allows us to fit to data. We used the fad version of the model to compare simple and complex contagion to real-world behaviours, and statistically decide which was better. We then tested the model's predictive power by using it to predict the spread of a fad

before it had reached its peak.

4.4.2 Data

The data chosen was Google search volumes for a particular category of Internet meme: photo fads. These fads consist of users uploading photos of themselves in a particular pose; for example, the photo fad ‘planking’ involved participants uploading photos of themselves lying face-down in public places. Photo fads were chosen because they tended to have distinctive names, allowing them to be clearly identified in search data; they involved real-world behaviours that were spread by and reported on the Internet; and they were undertaken for no ostensive reason beyond their online popularity. These photo fads tended to be global phenomena, and hence took place in a population large enough to satisfy the assumptions of the ODE model.

We avoided selection bias by taking all 37 Photo Fads listed on the website KnowYourMeme.com (a comprehensive source of information on internet memes). The search data was obtained from Google Trends, and consisted of search volumes quoted in terms of a percentage of the peak value, and aggregated weekly. We fitted models to the 26 fads with sufficient (greater than 15) non-zero datapoints.

4.4.3 Statistical model

The data take the form of a set of real-valued Google Trends at discrete time points $\mathbf{y} := (y_t)_{t=1}^T$. Search data was assumed to be a proxy for the number of people taking part in the trend: infected individuals search for information about these fads at a constant rate. The noise in the data was therefore modelled as arising from overdispersed sampling with mean $\mu(t) := I(t) + J(t)$, where $I(t)$ and $J(t)$ are solutions to the ODE fad model in equation (4.8). For known count data the Negative Binomial distribution would be appropriate to model this overdispersed sampling, but the data provided by Google Trends is instead given as a percentage of the peak and is therefore real-valued. As such we use the Gamma distribution, which approximates the Negative Binomial in the limit of large population size and is defined on the positive real numbers, to model the noise around the mean. This gives the following likelihood function:

$$L(\mathbf{y}|\theta) = \prod_{t=1}^T \text{Gamma}(y_t | A\mu(t + \Delta t), r),$$

AIC difference	Evidence
0 - 2	Weak
2 - 6	Positive
6 - 10	Strong
> 10	Very strong

Table 4.1: Grades of evidence provided by the AIC difference between two models, based on Stylianou et al. [4].

where we use the ‘mean-shape’ parameterisation of the Gamma distribution. This likelihood contains three additional ‘nuisance’ parameters: A is the relative amplitude term to adjust for the fact that Google Trends data is quoted in terms of the fraction of the peak; Δt is an additive time shift to match model time with real time; and r is the Gamma shape parameter to capture overdispersal.

To fit the model, L was maximized with respect to all parameters. For each set of fad data we calculated the Akaike Information Criterion (AIC), as described in chapter 2, for the simple and complex contagion models:

$$AIC = 2k - 2\ln L^*,$$

where k is the number of parameters for each model (8 for simple contagion and 9 for complex contagion) and L^* is the maximum value of the likelihood. We then classified the difference in AIC between the two models into different grades of evidence, based on the suggestions of Stylianou et al. [4], as shown in Table 4.1.

Some fads showed two clear peaks in the data. We therefore fitted a model in which two separate sub-populations become infected, with the total infected fraction being the sum of infected in the sub-populations. The parameters for each population were fitted independently, except for the thresholds in the complex contagion model that were assumed constant. The AIC was again used to select between one-population and two-population versions of both contagion mechanisms.

4.4.4 Results

Of these fads, 22 of 26 showed significant evidence that complex contagion was a better model for the data than simple contagion. The fitted timeseries for all fads are provided in Figure 4.14, ordered by log-likelihood difference. Most fads showed similar characteristics: a fast uptake, a drop in interest after the peak that was almost as fast, and then a long tail of activity taking a long time to die out.

The complex contagion model’s threshold for social influence allows it to capture the fast increase in popularity seen in most of the trends. The linear force

of influence in the simple contagion model, however, means that it is slower to build to peak popularity. After the peak, the simple contagion model has a constant rate for individuals leaving the fad, leading to exponential decay in popularity. The complex contagion initially shows a fast drop in popularity as individuals see that their contacts are already taking part in the fad, but once most of the population has stopped taking part the few individuals remaining take longer to give it up. This correctly captures the ‘long tail’ of popularity seen in the data.

For a minority of fads, the simple contagion model was also adequate, but this was typically linked to few datapoints and / or poor signal quality.

Table 4.2 shows the log-likelihood difference, $\Delta\mathcal{L} = \ln L_c^* - \ln L_s^*$, between the complex contagion and the simple contagion models (the difference in number of parameters is constant for the single population models and for the double population models) and the AIC evidence grade for each fad. For 22 out of 26 fads the complex contagion model is significantly better than simple contagion. The three fads with no positive evidence for either model were noisier and had higher background search volumes than the other fads. The names of these fads (‘caught me sleeping’, ‘people eating money’, ‘playing dead’) are phrases that could appear in searches unrelated to photo fads, leading to higher noise. It is interesting that the one case where simple contagion was a significantly better model, ‘horsemanning’, was the only one started by the Internet news site ‘BuzzFeed’ in an attempt to create a fad artificially. This suggests that a strong external driver not included in the model, such as mass media influence, can have a significant effect on the spread of a fad.

4.4.5 Prediction

The complex contagion model was used to predict the future spread of another fad, ‘ALS Icebucket Challenge’. This was a charity campaign that spread in a viral manner, with friends nominating each other to take part. A previous fad, ‘Neknomination’, had spread in a similar way, and so we used the parameters fitted from that fad to predict the future spread of ‘ALS Icebucket Challenge’. We made a verifiable prediction at the start of the campaign, shown in Figure 4.16, and overlaid the final data when the campaign had finished. The data are generally within the 95% prediction interval of the model, and the time and duration of interest in the campaign were predicted well: the peak occurred in the week predicted by the model, and the campaign was popular for the same length of time as the model.

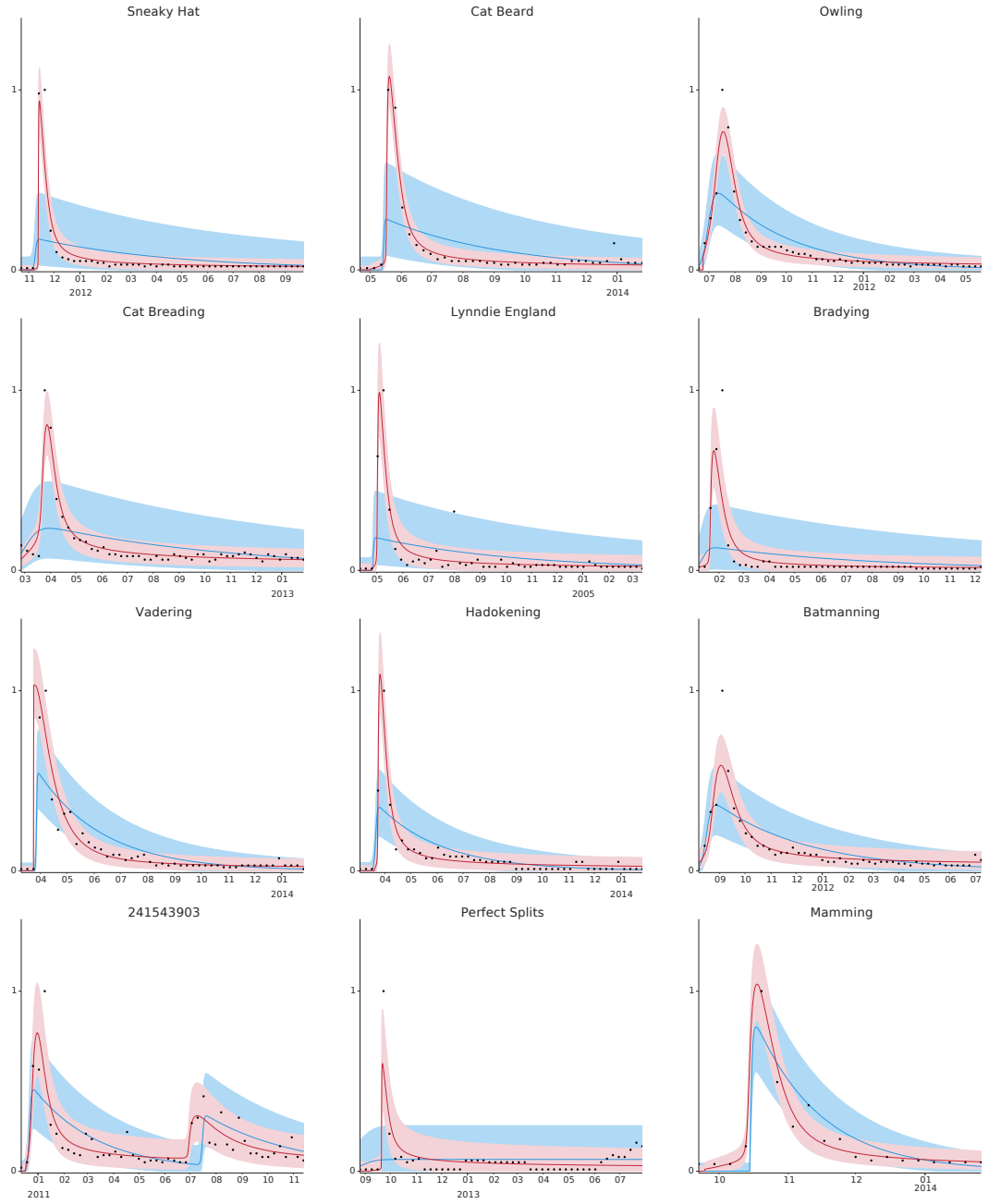


Figure 4.14: Simple and complex contagion model fits for all fads with sufficient data, ordered by log-likelihood difference. (continued on next page)

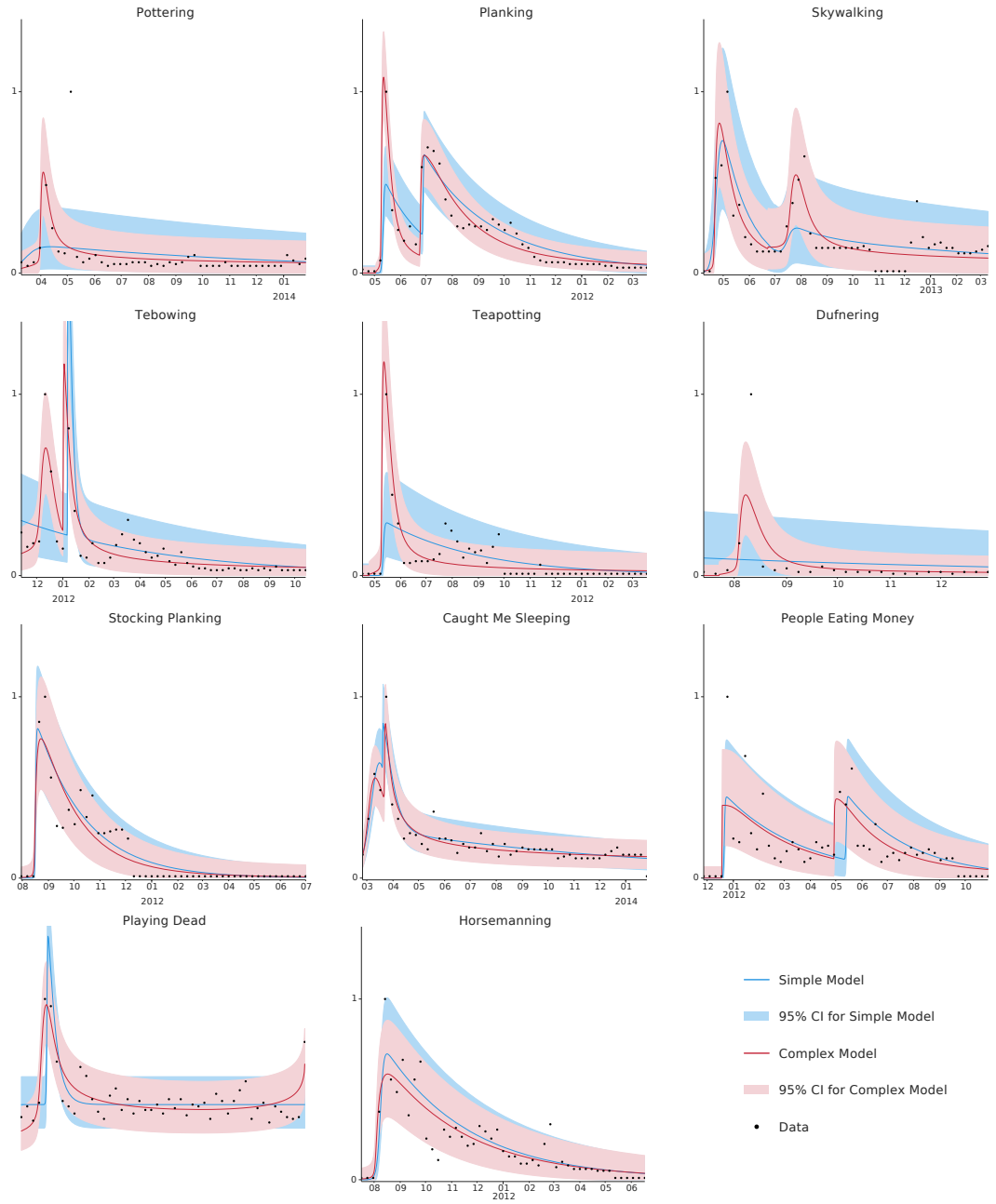


Figure 4.15: (Continued) Simple and complex contagion model fits for all fads with sufficient data, ordered by log-likelihood difference.

Photo Fad	Log-likelihood difference	AIC Evidence	
Sneaky Hat	47.4	***	
Cat Beard	44.0	***	
Owling	39.3	***	
Cat Breeding	38.9	***	
Lynndie England	27.3	***	
Bradying	26.6	***	
Vadering	24.6	***	
Hadokening	24.2	***	
Batmanning	24.2	***	
Lying Down Game	17.8	***	
Leisure Diving	16.8	***	
Sleeveface	12.1	***	
241543903	11.7	***	†
Perfect Splits	11.7	***	
Mamming	10.2	***	
Pottering	9.5	***	
Planking	8.5	***	†
Skywalking	7.1	***	†
Tebowing	6.9	***	†
Teapotting	6.0	***	
Dufnering	4.7	**	
Stocking Planking	2.0	*	
Caught Me Sleeping	-0.0	.	†
People Eating Money	-1.7	.	†
Playing Dead	-1.9	.	†
Horsemanning	-3.2	-	

Table 4.2: The log-likelihood difference between the simple and complex contagion models. (***) is very strong evidence, (**) is strong evidence, (*) is positive evidence, (.) is no significant evidence for either model, (-) is strong evidence against.

† AIC selected models with two peaks.

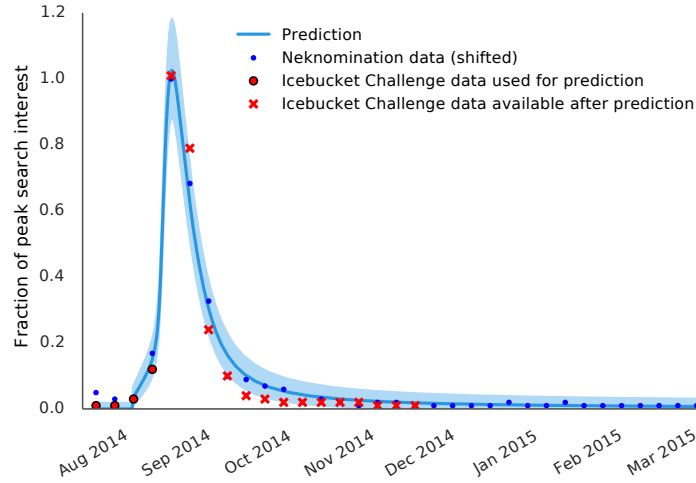


Figure 4.16: Prediction for the ‘Icebucket Challenge’ fad, using a model fitted to a similar fad. The plot shows the data available when the prediction was made, and the subsequent course of the fad. The model successfully predicts the timing of the peak and the duration of the fad.

4.4.6 Discussion

The general social influence model provides a framework in which we can test competing models of social influence. In this case we used data on real-world behaviours driven by online spreading to compare simple contagion and complex contagion statistically. For all but one fad the complex contagion model before at least as well as simple contagion, and in most cases was significantly better. The complex contagion was also capable of predicting the peak time and duration of a fad that had not yet reached its peak. This demonstrates that complex contagion provides a predictive modelling framework for real-world behaviours which spread online.

4.5 Summary

In this chapter we have developed a general model of behaviour spread based on the psychologically motivated and experimentally confirmed mechanism of complex contagion. This general model can be simplified to suit a variety of spreading situations, including behavioural fads which are briefly popular and then die out. A convenient feature of the fad model is that it can be analysed using Bailey’s method, a computationally efficient technique for calculating the final size distribution of a fad to machine precision. This in turn highlighted an interesting feature of complex

contagion models: that a small ‘nudge’ of a few individuals can mean the difference between a behaviour that spreads or one that dies out, even for behaviours with the same parameter values.

We then investigated the deterministic limit of the general spread model, and used it to investigate the effect of the parameters on the mean behaviour of the model. The deterministic limit also allowed a statistical comparison of complex contagion with the simple contagion model that is currently used in most other studies, using observational data for the online spread of real-world behaviours. The predictive power of the complex contagion model was then demonstrated by predicting the spread of a fad before it had reached its peak.

Chapter 5

Towards a combined social influence and covariate model

5.1 Introduction

So far in this thesis we have separately tackled each health behaviour component identified in chapter 1. An ideal model for health behaviour, however, would combine all of these components into a predictive model that could be statistically fitted if given appropriate data. It might also attempt to model a disease spreading at the same time as the behaviour, each influencing the other. In this overall model defined in chapter 1, the behaviour of the population $\mathbf{B}(t)$ is a function of time, covariates $\mathbf{X}(t)$, and a separate dynamic process, such as a disease, $\mathbf{Z}(t, \mathbf{B}(t))$. This is an ambitious goal, which will not be reached in this thesis. Instead, in this Chapter we discuss some possible approaches which could lead to this ideal model.

5.2 Coupling behaviour and disease spread

5.2.1 Background

We will use the complex contagion model described above to investigate the dynamics of health behaviours that spread through a homogeneous population while an epidemic is taking place:

$$\mathbf{B}(t) \rightarrow \mathbf{B}(t, \mathbf{Z}(t, \mathbf{B}(t))).$$

Using a condom during sex, for example, can greatly reduce the chance of contracting HIV [7], and individuals are more likely to use a condom if they know that an HIV epidemic is occurring [18]. Other preventative health behaviours have also been

shown to become more prevalent during an epidemic, such as vaccinations during a measles outbreak [19] and avoiding infectious contact [20]. Some studies of disease outbreak data use behavioural changes to explain sudden decreases in transmission rate [168, 169], but there are few studies which model this mechanism directly. Funk et al. [170] review and classify recent theoretical attempts to investigate human behaviour and disease. They identify the properties as ways to group models of behaviour-disease dynamics: if information about the disease is transmitted globally to the whole population, or locally through word of mouth; if individuals change their behaviour due to the prevalence of the disease, or the prevalence of the behaviour; and if the effect of the behaviour is to change the disease state, the model parameters, or the contact structure of the population. Developing a general model of health behaviour linked to disease is beyond the scope of this thesis, but an aspect lacking in the studies reviewed by Funk et al. is the principled social influence term we introduced in chapter 4. In this section we will investigate a model linking disease to a complex behavioural contagion.

5.2.2 Model

We chose to base our approach on the work of Funk et al. [5], which considers both the spread of a disease and an awareness of the disease, with aware individuals taking part in a behaviour which changes their infection rate and recovery rate. Awareness is transmitted locally, by word of mouth, and individuals change their behaviour based on the prevalence of behaviour around them rather than the prevalence of the disease. The effect of behaviour is to modify the model parameters. The model consists of six compartments: S_u, I_u, R_u , representing unaware individuals who are Susceptible, Infected or Recovered respectively, and S_a, I_a, R_a , representing aware individuals in the corresponding states. Both disease and awareness are treated as simple SIRS contagions with the rates between each compartment given in Figure 5.1.

Four mean-field equilibria are identified: disease and awareness both spread, only awareness spreads, only disease spreads, and nothing spreads. Each of these equilibria are determined by the parameters of the model, for all non-trivial initial conditions.

We simplified the structure of the Funk model by removing the Recovered

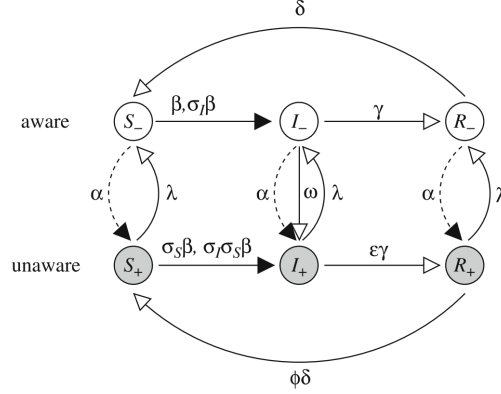


Figure 5.1: Diagram showing the rates between the compartments in the Funk model [5]. Arrows represent possible transitions between the compartments

states and setting $\omega = 0$. The states are defined in terms of $\mathbf{B}(t)$ and $\mathbf{Z}(t)$ as follows:

$$\begin{aligned}
 S_u &= \sum_i \mathbf{I}_{\{B_i(t)=u\}} \mathbf{I}_{\{Z_i(t)=S\}} \\
 I_u &= \sum_i \mathbf{I}_{\{B_i(t)=u\}} \mathbf{I}_{\{Z_i(t)=I\}} \\
 S_a &= \sum_i \mathbf{I}_{\{B_i(t)=a\}} \mathbf{I}_{\{Z_i(t)=S\}} \\
 I_a &= \sum_i \mathbf{I}_{\{B_i(t)=a\}} \mathbf{I}_{\{Z_i(t)=I\}},
 \end{aligned}$$

where $B_i(t) \in \{u, a\}$ and $Z_i(t) \in \{S, I\}$.

We separated the rate α into a rate α_I for infected individuals and a rate α_S for susceptible individuals, and similarly for λ . We concentrated on a situation where information about the risks of certain behaviours are less important than peer pressure/social norms. In our model, equation (5.1), there is no direct interaction between the disease and the rate of becoming aware of the behaviour; instead, awareness spreads through the population as a complex contagion (‘peer pressure’) with rates dependent on whether an individual’s contacts are infected or susceptible.

$$\begin{aligned}
\frac{dS_u}{dt} &= -\beta_u(I_u + I_a)S_u + \gamma_u I_u - \alpha_S D(S_u, I_a, I_u)S_u + \lambda_S S_a \\
\frac{dI_u}{dt} &= \beta_u(I_u + I_a)S_u - \gamma_u I_u - \alpha_I D(S_u, I_a, I_u)I_u \\
\frac{dI_a}{dt} &= \beta_a(I_u + I_a)S_a - \gamma_a I_a + \alpha_I D(S_u, I_a, I_u)I_u \\
S_a &= 1 - S_u - I_u - I_a.
\end{aligned} \tag{5.1}$$

In this model, susceptible individuals take steps to change their probability of becoming infected if they are ‘aware’ of the behaviour. Unaware individuals become infected with rate β_u , and aware individuals with rate β_a . Individuals stop being infectious at rates γ_u and γ_a if they are unaware or aware respectively, at which point unaware individuals return to the S_u state and aware return to the S_a state. Awareness spreads through the mechanism of complex contagion, with individuals canvassing C contacts and becoming aware with a particular rate if the number of aware infectious individuals around them is greater than a threshold τ_i , or if the number of aware susceptible individuals around them is greater than a threshold τ_s . The rate of awareness spreading if one of these thresholds is reached depends on whether the individual itself is susceptible or infected: susceptible individuals become aware at rate α_S and infectious individuals at rate α_I . The loss of awareness happens at constant rate λ_S for susceptible individuals, but aware infectious individuals are assumed not to lose their awareness ($\lambda_I = 0$).

5.2.3 Methods

We analyse the system of ODEs given in equation (5.1) using Linear Stability Analysis [171]. This is a method for investigating the local behaviour of an ODE system around its fixed points. First the fixed points are found by solving:

$$\begin{aligned}
f(S_u^*, I_u^*, I_a^*) &= \frac{dS_u}{dt}(S_u^*, I_u^*, I_a^*) = 0 \\
g(S_u^*, I_u^*, I_a^*) &= \frac{dI_u}{dt}(S_u^*, I_u^*, I_a^*) = 0 \\
h(S_u^*, I_u^*, I_a^*) &= \frac{dI_a}{dt}(S_u^*, I_u^*, I_a^*) = 0.
\end{aligned}$$

To find the local behaviour around the fixed point, we consider a small perturbation from the fixed point:

$$(S_u^* + \Delta S_u, I_u^* + \Delta I_u, I_a^* + \Delta I_a).$$

Using a Taylor series expansion and linearizing (removing all $O(\Delta S_u^2)$ etc terms), we arrive at the following linear system of equations for the dynamics of a small perturbation:

$$\begin{pmatrix} \frac{\Delta S_u}{dt} \\ \frac{\Delta I_u}{dt} \\ \frac{\Delta I_a}{dt} \end{pmatrix} = \begin{pmatrix} \frac{df}{dS_u} & \frac{df}{dI_u} & \frac{df}{dI_a} \\ \frac{dg}{dS_u} & \frac{dg}{dI_u} & \frac{dg}{dI_a} \\ \frac{dh}{dS_u} & \frac{dh}{dI_u} & \frac{dh}{dI_a} \end{pmatrix} \begin{pmatrix} \Delta S_u \\ \Delta I_u \\ \Delta I_a \end{pmatrix}. \quad (5.2)$$

The matrix in equation (5.2) is the Jacobian matrix for the model, and its eigenvalues determine the local behaviour of the perturbations: if any of the eigenvalues have real part greater than 1 then perturbations grow and the system is unstable, otherwise perturbations shrink and any small perturbation from the fixed point will return to it [171]. If some of the eigenvalues are greater than one and others less than one then the fixed point is a ‘saddle point’, with some directions unstable and others stable. The stability of fixed points can change as the values of the model parameters change.

5.2.4 Results

For this model, the parameter values can be grouped into two categories: parameter values for which all fixed points have endemic disease, and values for which there is a stable disease-free fixed point. Both of these categories have one endemic fixed point with low awareness, and a second fixed point with higher awareness. Figure 5.2 shows the system reaching both fixed points for the case where there is a disease-free fixed point. These parameter values have $R_0 < 1$ among aware individuals. Figure 5.3 shows the system reaching both fixed points for the case with no non-trivial disease-free fixed point. For parameter values in this first category, $R_0 > 1$ even among aware individuals.

We used linear stability analysis to investigate the fixed points of this system. Fixed points were found numerically, and the Jacobian matrix calculated using numerical differentiation. We explored the behaviour of the system away from these fixed points by plotting $(\frac{dS_u}{dt}, \frac{dI_u}{dt}, \frac{dI_a}{dt})$ as a function of (S_u, I_u, I_a) in the form of a flow diagram. Figure 5.4 shows a 2d cross-section of this diagram for $(S_u, I_u, 0)$ in the case where there is a stable disease-free fixed point. In this system there are two fixed points, corresponding to a state with high awareness and one with low awareness. The high awareness state has no infected individuals, and is thus disease-free,

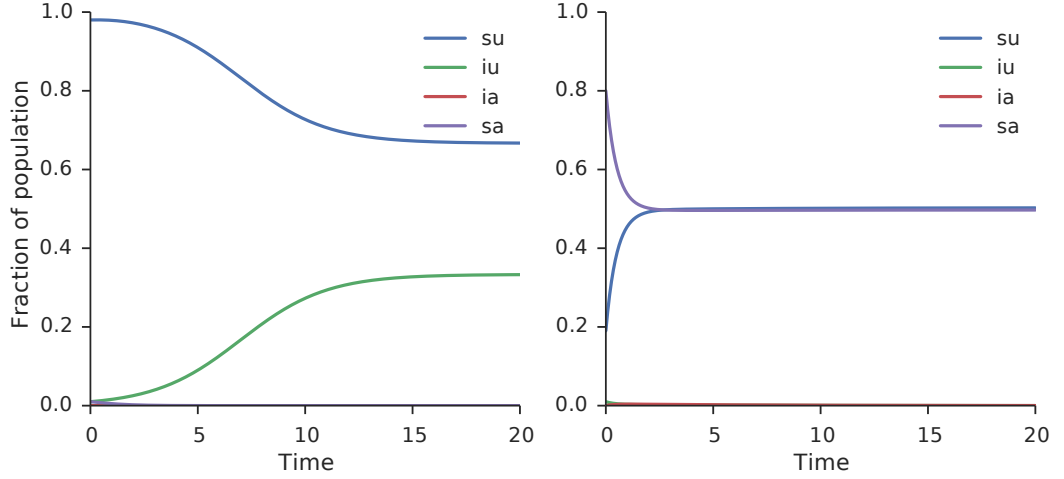


Figure 5.2: Timeseries of disease and awareness for the parameter regime with a stable disease-free fixed point, approaching the low-awareness (left) and high-awareness (right) states. In this case the parameter values were: $\beta_u = 1.5$, $\gamma_u = 1.0$, $\beta_a = 0.2$, $\gamma_a = 1.0$, $\lambda_S = 1.0$, $\alpha_S = 1.0$, $\alpha_I = 1.0$, $\tau_S = 3$, $\tau_I = 3$. The initial condition for low awareness was $(S_u(0) = 0.98, I_u(0) = 0.01, I_a(0) = 0.0)$, and for high awareness was $(S_u(0) = 0.19, I_u(0) = 0.01, I_a(0) = 0.0)$.

and the low awareness state has no aware individuals: awareness has totally died out. The other two fixed points are the saddle points at $(1, 0, 0)$, corresponding to the case where there are no infected or aware individuals in the system, and at $(0.67, 0.3, 0.01)$. This second saddle point separates trajectories that lead to the high awareness fixed point and those that lead to the low-awareness fixed point.

Figure 5.5 gives the same flow diagram for the case without any stable disease-free fixed points. At the stable low-awareness fixed point awareness has again completely died out, but the disease-free fixed point has become a saddle node and a new stable fixed point with endemic disease has appeared. There is again a saddle point separating low-awareness and high-awareness trajectories.

The relationship between the infectivity of aware and unaware individuals in determining whether there is a stable disease-free fixed point is shown in Figure 5.6. In this figure, the total infected fraction at the lowest infected stable fixed point is shown for each combination of β_a and β_u . It is clear that for each value of β_u there is a critical value of β_a below which there is a stable disease-free fixed point. Plotting this critical value β_a^* against β_u , as shown in Figure 5.7, displays a linear relationship. The critical value of β_a^* for these parameters is given by

$$\beta_a^* = -1.1\beta_u + 2.1.$$

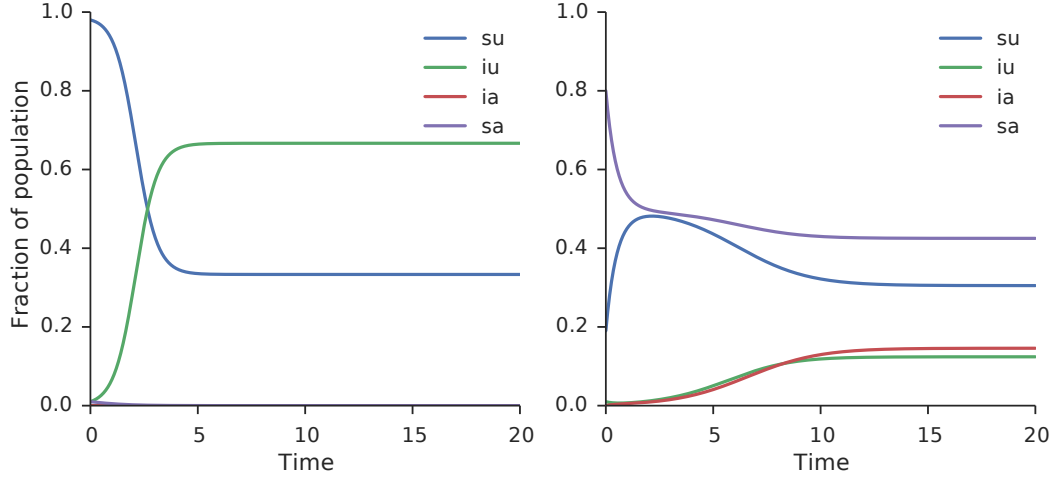


Figure 5.3: Timeseries of disease and awareness for the parameter regime with no stable disease-free fixed point, approaching the low-awareness (left) and high-awareness (right) states. In this case the parameter values were: $\beta_u = 3.0, \gamma_u = 1.0, \beta_a = 0.2, \gamma_a = 1.0, \lambda_S = 1.0, \alpha_S = 1.0, \alpha_I = 1.0, \tau_S = 3, \tau_I = 3$. The initial condition for low awareness was $(S_u(0) = 0.98, I_u(0) = 0.01, I_a(0) = 0.0)$, and for high awareness was $(S_u(0) = 0.19, I_u(0) = 0.01, I_a(0) = 0.0)$.

For unaware infectivity greater than 1.91, even behaviours with $\beta_a = 0$ will have no disease-free stable fixed point. This linear relationship defines a basic reproduction number for the low-disease state:

$$R_0 = \frac{I_a^* \frac{\beta_a}{\gamma_a} + I_u^* \frac{\beta_u}{\gamma_u}}{I_u^* + I_a^*},$$

where I_a^* is the fraction of infected aware individuals at the fixed point, and I_u^* is the number of infected unaware individuals at the fixed point. If $R_0 > 1$ then there is no stable disease-free fixed point.

5.2.5 Discussion

By coupling disease dynamics with behaviour we can include an aspect of disease spreading which is lacking in many epidemic models. In particular, we studied a simple SIS disease coupled with a peer pressure complex social contagion. This model could be applicable to the study of gonorrhea prevention in adolescents, for example, since gonorrhea does not confer significant immunity to future gonorrhea infections [172] and adolescents may be hesitant to take part in a behaviour unless a significant number of their friends are also taking part.

The coupled disease-behaviour model has important implications for the con-

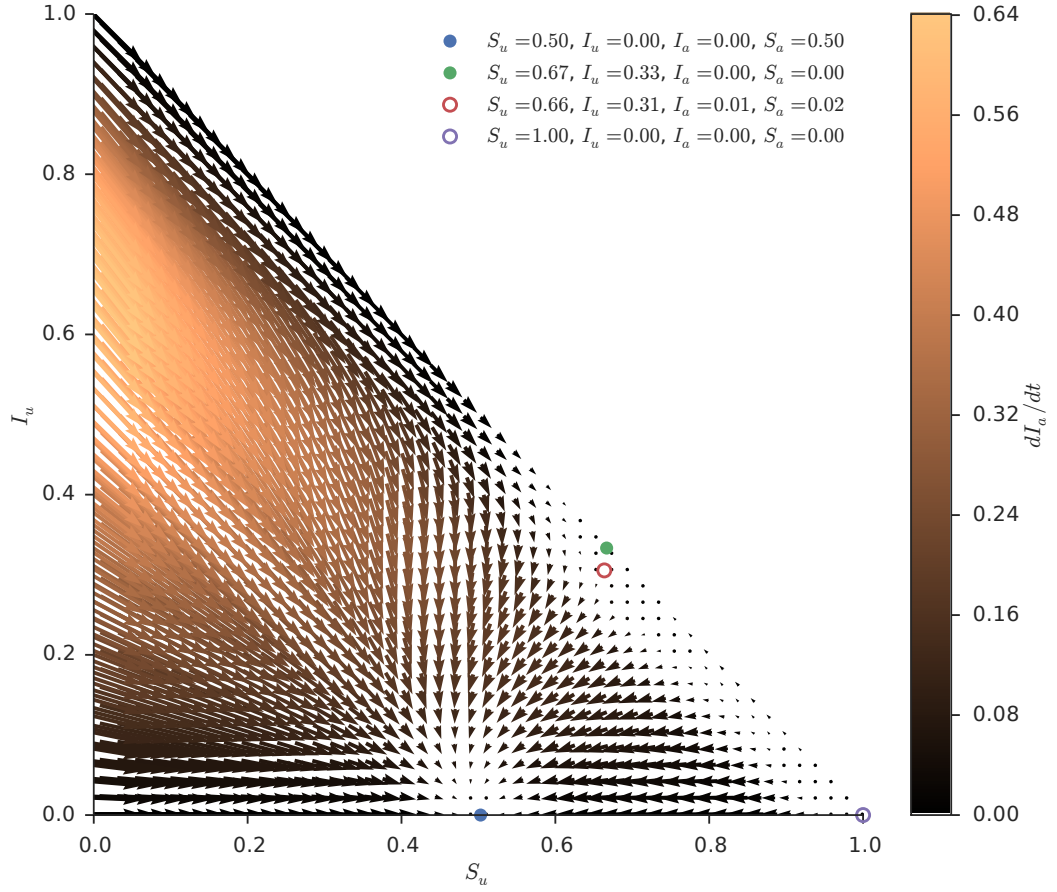


Figure 5.4: Vector flow diagram cross-section of $(\frac{dS_u}{dt}, \frac{dI_u}{dt}, \frac{dI_a}{dt})$ at $(S_u, I_u, 0)$, for the parameter regime with a stable disease-free fixed point. Stable fixed points are shown as filled circles, and unstable points with open circles.

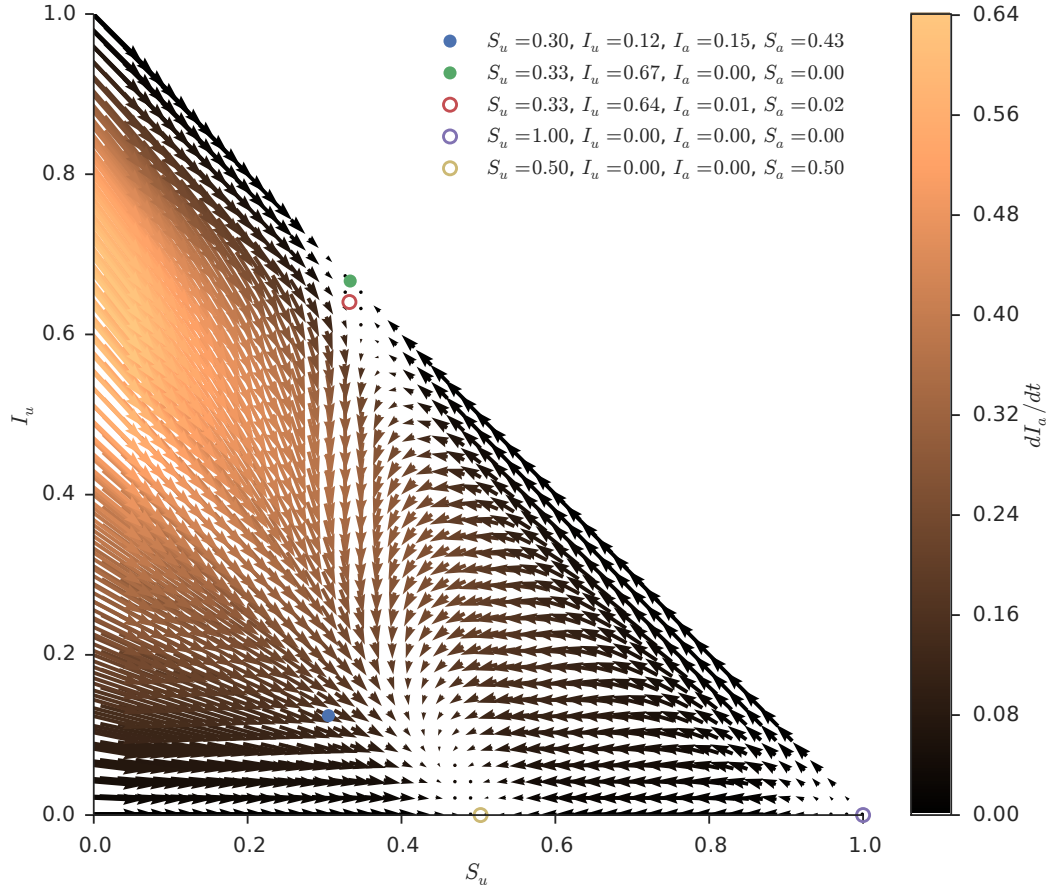


Figure 5.5: Vector flow diagram cross-section of $(\frac{dS_u}{dt}, \frac{dI_u}{dt}, \frac{dI_a}{dt})$ at $(S_u, I_u, 0)$, for the parameter regime with no stable disease-free fixed point. Stable fixed points are shown as filled circles, and unstable points with open circles.

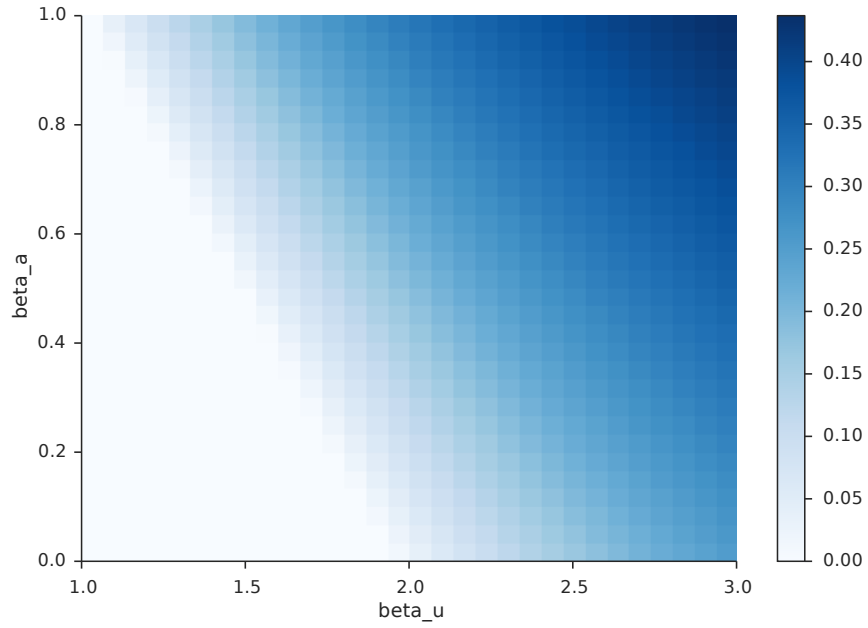


Figure 5.6: Lowest stable endemic disease for parameters β_a and β_u . There is a linear boundary between combinations of β_a and β_u that lead to disease-free states and those that do not.

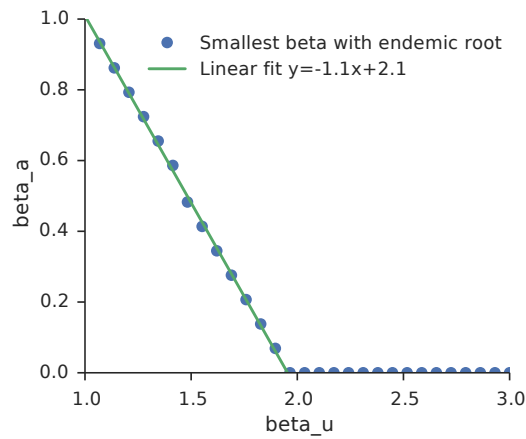


Figure 5.7: Critical value of β_a below which there is a disease-free state. This is a linear relation, with slope and intercept dependent on the other parameters. In this case, for $\beta_u > 1.91$ there is no value of β_a that gives a disease-free state.

trol of disease. Even a behaviour which reduces infectivity to zero will still lead to endemic disease if the number of individuals taking part in the behaviour is not sufficient for the behaviour to spread to a significant fraction of the population. As with the fad model in the previous section, there is a sharp threshold in the number of aware individuals required for the behaviour to spread to enough of the population to ensure a disease-free state. This is due to a saddle-point which is not present in an SIS model with simple contagion. The saddle-point dynamics of this model may help to explain the sharp decline in the rates of some sexually transmitted diseases in certain Western countries in recent decades. Hiltunen-Back et al. [173], for example, could not find an explanation for the rapid decrease of endemic gonorrhea in Finland from 1990 to 1995, which occurred despite a constant rate of imported gonorrhea. A possible explanation based on this coupled model would suggest that the number of aware individuals taking part in preventative behaviour became sufficiently large that a threshold was reached and the population quickly reached a high-awareness state with an associated drop in endemic disease.

The assumption that infected individuals do not lose awareness is unlikely to be correct, but is motivated by research suggesting that individuals who are diagnosed with HIV often do take steps to reduce their risk of infecting others [174, 175, 176], particularly if counseled to do so [177], and often stop having intercourse with uninfected partners [178]. These studies continued from three months to a year after diagnosis, so if the disease and spreading of awareness acts on much faster timescales than this then the rate of awareness loss for infected individuals can be ignored.

The high-awareness and low-awareness states are associated with different values for the basic reproduction number of the disease. The R_0 value is given by the average of $R_{0,a}$ and $R_{0,u}$ (the basic reproduction numbers for the aware and unaware categories respectively) weighted by the fraction of the population in each category. This feature, where the basic reproduction number depends on the initial state of the system, did not occur in the model of Funk et al. [5].

5.3 Towards a combined social and spatial model

5.3.1 Background

In chapter 4 we investigated various versions of a general model for social influence. This general model assumed all-to-all contact network connectivity, with each individual having the same potential social influence as any other. This is an unrealistic assumption when modelling health behaviours across a large spatial area, since in

real human contact networks individuals are more likely to be connected if they are geographically closer [179, 180]. We now suggest a model which includes a very simple heterogeneity in the population.

5.3.2 Social model with a spatial component

We consider each individual to have two spatial regimes: ‘near’ and ‘far’. Individuals canvas C_n ‘near’ contacts and C_f ‘far’ contacts, with social influence thresholds τ_n and τ_f respectively. A simple model for ‘SIS’-like behavioural spread could therefore be:

$$\begin{aligned}\frac{dI}{dt} &= \beta D(I, J, \tau_n, \tau_f, C_n, C_f) S_I - \gamma I \\ \frac{dJ}{dt} &= \beta D(J, I, \tau_n, \tau_f, C_n, C_f) S_J - \gamma J \\ \frac{dS_I}{dt} &= -\beta D(I, J, \tau_n, \tau_f, C_n, C_f) S_I + \gamma I \\ \frac{dS_J}{dt} &= -\beta D(J, I, \tau_n, \tau_f, C_n, C_f) S_J + \gamma J,\end{aligned}$$

where I and J are the fraction of individuals taking part in the behaviour in two distinct spatial areas, $D(I, J, \tau_n, \tau_f, C_n, C_f)$ is the multinomial form of the social influence term given in equation (4.3), and other parameters have the same meanings as in section §4.2.

5.3.3 Statistical model with spatial social influence

Approaching the problem from the statistical side, one way to include a more accurate social influence term than that used in chapters 2 and 3 is to take into account social influence between spatial regions. This can be achieved by adding a spatially-correlated random effect to a Bayesian hierarchical model with individual- and health system-level covariates; a technique suggested by Gelfand [181]. An example of such a model, where data is collected within supervision areas as in chapters 2 and 3, is as follows:

$$\begin{aligned}\mathbf{Y}(\mathbf{X}(t)) &\sim \text{Binomial}(\mathbf{p}) \\ \text{logit}(\mathbf{p}) &= \beta^T \mathbf{\Theta} + \gamma^T \mathbf{E} + \mu_{\text{spatial}} \\ \mu_{\text{spatial}} &\sim \text{MVN}(0, \sigma_{SA}^2 H(\mathbf{d}, \phi)),\end{aligned}$$

where σ_{SA} is the within-supervision area standard deviation, $H(\mathbf{d}, \phi)$ is a covari-

ance matrix which specifies the covariance between supervision areas given distance matrix \mathbf{d} between them and with spatial scale parameters ϕ . The spatial scale parameters can be estimated, and this will provide information about the spatial extent of factors such as social influence that might act between supervision areas.

5.3.4 Data

The weak link in connecting these two regimes is the absence of high spatial- and temporal- resolution data. The lack of spatial resolution is easier to remedy: better spatial data collection, such as using GPS devices during household surveys, would allow more exact measurement of the distances between sites, and hence allow more precise fitting of the spatial kernel parameters ϕ .

Temporal data is more difficult to obtain. Even LQAS surveys are too expensive to run more than once per year, and social influence is likely to act on faster timescales than this. There have been studies on contact networks in small communities [182] in which contacts are tracked between individuals over the course of each day; perhaps a health behaviour intervention could be introduced to some individuals and the resulting spread analysed. The internet provides high temporal-resolution data, as described in chapter 4, but it often lacks spatial data. As increasing numbers of people use smartphones with inbuilt GPS, however, it is likely that high-quality data with both dimensions will become easier to obtain.

5.3.5 Discussion

The two models above represent attempts to approach the problem of social influence in heterogeneous populations from the dynamical and the statistical side. The first contains very simple heterogeneity, in which the population is split into two spatial components with strong social interactions within them and weak interactions between them. This could then be generalised to include further spatial components, some of which do not interact at all, by adding additional compartments to the model. The second captures social influence between spatial regions, an effect that was previously ignored. In principle, including many additional spatial components in the first model and reducing the effect of the β parameters in the second model should lead to the two models corresponding.

High quality data is required to fit these combined models. Including GPS locations for households, and combining these with detailed geographic covariates, could allow us to control for the spatial variation due to geography and be more confident in isolating the residual correlation due to social influence. In principle this

could be used to construct a continuous ‘social influence surface’ across a country, using techniques such as latent spatial Gaussian Processes [183, 184, 134]. Such a surface could highlight areas of a country that have unexpectedly poor health behaviours due to social influence; these areas may be more susceptible to nudges.

5.4 Summary

In this chapter we have suggested some ways of relaxing the limits of the model discussed in previous chapters, by linking behaviour with disease and by moving towards a combined social influence and covariate model. The linked behaviour-disease model shows a phenomenon not present in traditional models of disease: a basic reproduction number that depends non-trivially on the fraction of the population that is infected. We approached the problem of combining social influence dynamics and a heterogeneous population from both the dynamic and statistical sides, in the hope that they can be fitted to high quality data as it becomes available.

Chapter 6

Conclusions and Further Work

In chapter 2, we explained and predicted local variations in health behaviour by taking into account geographical and demographic factors. Our results are consistent with the claims of the Health Belief Model that ease of access plays a role in understanding which individuals take part in certain health behaviours. The predictive power of this model means that it could become a useful policy tool for prioritising poorly-performing areas of a country. Possible extensions for this model include treating it in the Bayesian framework introduced in chapter 3: this would allow estimation of regional effects such as social influence that are not currently modelled. The surveys used in chapter 2 provide rich data on many health behaviours, and further work is needed to apply the model to these.

In chapter 3, we showed that Bayesian hierarchical modelling can provide new ways of studying the health systems that shape individual behaviour. We showed that one aspect of the health system that had previously been difficult to study, health facility quality, is associated with health facility usage, but that including quality in a model of health facility usage behaviour did not provide as much predictive benefit as taking into account unknown regional effects such as social influence. The Bayesian framework also aided the evaluation of a complex policy intervention to improve health behaviours. The evaluation provided evidence that funding a health facility based on its results, rather than its inputs, led to lower costs for users and an increase in health facility usage in its catchment area. This was despite there being no evidence for an increase in the objective quality of facilities provided with results-based funding. This was again a very rich dataset, and it would be interesting to study the wealth and transport data available from the household surveys.

In chapter 4, we developed a psychologically motivated model for behavioural spread. The model captured an important feature of systems with social influence:

success was highly unpredictable for individual realisations, and small nudges could mean the difference between a successful behaviour and one that died out. The deterministic limit of the model allowed us to compare complex contagion with simple contagion using data for the online spread of real-world behavioural fads, and showed that the complex contagion model performed at least as well as simple contagion for all but one of the datasets studied. The deterministic limit model was also used to predict the timing and duration of another fad before it had peaked, with implications for predicting health behaviours such as fad diets and other behaviours such as charity campaign donations. Applying this model to data other than Google Trends, and using it to predict more fads, would give more confidence in the model's general applicability.

Finally, in chapter 5, we introduced some ways of combining these limits into a general model, and to link behaviour to the spread of communicable disease. In particular, we investigated a model in which peer pressure is the driving force behind a positive health behaviour change which reduces the infectivity of a disease. The complex contagion mechanism of social spread leads to two equilibrium states depending on the initial fraction of infected: a high disease state and a low disease state. For some parameter values the disease can be eradicated, whereas for others even the low disease state is above zero; the basic reproduction number of the system is dependent on the initial fraction of infected individuals. Even this basic model shows phenomena not present in disease models that do not take into account complex contagion.

As more health behaviour data becomes available from online interactions, and from public health Monitoring and Evaluation projects in developing countries, more complex limiting cases of the general health behaviour model can be fitted and used predictively. We believe, however, that the limiting cases we consider in this thesis apply well to the data currently available and provide insight into important real-world situations.

By considering individual barriers to health behaviour, the role of the health system in allowing and promoting those behaviours, and the spread of those behaviours in the population, we believe that this model has the potential to inform policies that encourage good health behaviours, and through this fight disease.

Bibliography

- [1] O’Ryan, M, Prado, V, & Pickering, L. K. (2005) A millennium update on pediatric diarrheal illness in the developing world. *Seminars in Pediatric Infectious Diseases* **16**, 125–136.
- [2] World Health Organization (WHO). (1990) Acute respiratory infections in children: case management in developing countries., (World Health Organisation, Geneva), Technical report.
- [3] World Health Organization. (2010) *Guidelines for the treatment of malaria*. (Geneva), Second edition, p. 197p.
- [4] Stylianou, C, Pickles, A, & Roberts, S. a. (2013) Using Bonferroni, BIC and AIC to assess evidence for alternative biological pathways: covariate selection for the multilevel Embryo-Uterus model. *BMC medical research methodology* **13**, 73.
- [5] Funk, S, Gilad, E, & Jansen, V. a. a. (2010) Endemic disease, awareness, and local behavioural response. *Journal of theoretical biology* **264**, 501–509.
- [6] GBD 2013 Collaborators. (2015) Global , regional and national levels of age-specific mortality and 240 causes of death , 1990-2013 : A systematic analysis for the Global Burden of Disease Study 2013. *The Lancet* **385**, 1990–2013.
- [7] Holmes, K. K, Levine, R, & Weaver, M. (2004) Effectiveness of condoms in preventing sexually transmitted infections. *Bulletin of the World Health Organization* **82**, 454–461.
- [8] Alter, M. (2007) Epidemiology of hepatitis C virus infection. *World journal of gastroenterology : WJG* **13**, 2436–2441.
- [9] Danaei, G, Vander Hoorn, S, Lopez, A. D, Murray, C. J. L, & Ezzati, M. (2005) Causes of cancer in the world: Comparative risk assessment of nine behavioural and environmental risk factors. *The Lancet* **366**, 1784–1793.

- [10] Eriksson, K & Lindgärde, F. (1991) Prevention of type 2 (non-insulin-dependent) diabetes mellitus by diet and physical exercise: The 6-year Malmö feasibility study. *Diabetologia* **2**, 891–898.
- [11] Fouque, D & Laville, M. (2009) Low protein diets for chronic kidney disease in non diabetic adults. *Cochrane Database of Systematic Reviews*.
- [12] Frost, L, Hune, L. J, & Vestergaard, P. (2005) Overweight and obesity as risk factors for atrial fibrillation or flutter: the Danish Diet, Cancer, and Health Study. *The American journal of medicine* **118**, 489–495.
- [13] Sladek, R, Rocheleau, G, Rung, J, Dina, C, Shen, L, Serre, D, Boutin, P, Vincent, D, Belisle, A, Hadjadj, S, Balkau, B, Heude, B, Charpentier, G, Hudson, T. J, Montpetit, A, Pshezhetsky, A. V, Prentki, M, Posner, B. I, Balding, D. J, Meyre, D, Polychronakos, C, & Froguel, P. (2007) A genome-wide association study identifies novel risk loci for type 2 diabetes. *Nature* **445**, 881–885.
- [14] Dabelea, D, Hanson, R. L, Lindsay, R. S, Pettitt, D. J, Imperatore, G, Gabir, M. M, Roumain, J, Bennett, P. H, & Knowler, W. C. (2000) Intrauterine exposure to diabetes conveys risks for type 2 diabetes and obesity: A Study of Discordant Sibships. *Diabetes* **49**, 2208–2211.
- [15] Keeling, M. J & Rohani, P. (2008) *Modeling Infectious Diseases in Humans and Animals*. (Princeton University Press).
- [16] Jerrett, M, Burnett, R. T, Ma, R, Pope, C. A, Krewski, D, Newbold, K. B, Thurston, G, Shi, Y, Finkelstein, N, Calle, E. E, & Thun, M. J. (2005) Spatial analysis of air pollution and mortality in Los Angeles. *Epidemiology (Cambridge, Mass.)* **16**, 727–736.
- [17] Giles-Corti, B & Donovan, R. J. (2002) The relative influence of individual, social and physical environment determinants of physical activity. *Social science & medicine* **54**, 1793–1812.
- [18] Ahituv, A, Hotz, V. J, & Philipson, T. (1996) The Responsiveness of the Demand for Condoms to the Local Prevalence of AIDS. *The Journal of Human Resources* **31**, 869–897.
- [19] Philipson, T. (1996) Private Vaccination and Public Health : An Empirical Examination for U . S . Measles. *The Journal of Human Resources* **31**, 611–630.

- [20] Campbell, G. L & Hughes, J. M. (1995) Plague in India: A new warning from an old nemesis. *Annals of Internal Medicine* **122**, 151–153.
- [21] Beaulac, J, Kristjansson, E, & Cummins, S. (2009) A systematic review of food deserts, 1966-2007. *Preventing chronic disease* **6**, A105.
- [22] Glanz, K, Rimer, B. K, & Viswanath, K. (2008) *Health Behaviour and Health Education: theory, research and practice*.
- [23] Janz, N. K & Becker, M. H. (1984) The Health Belief Model: a decade later. *Health education quarterly* **11**, 1–47.
- [24] Ahmed, S. M, Adams, A. M, Chowdhury, M, & Bhuiya, A. (2003) Changing health-seeking behaviour in Matlab, Bangladesh: Do development interventions matter? *Health Policy and Planning* **18**, 306–315.
- [25] Lagarde, E, Caraël, M, Glynn, J. R, Kanhonou, L, Abega, S. C, Kahindo, M, Musonda, R, Auvert, B, & Buvé, A. (2001) Educational level is associated with condom use within non-spousal partnerships in four cities of sub-Saharan Africa. *AIDS* **15**, 1399–1408.
- [26] Gabrysch, S, Cousens, S, Cox, J, & Campbell, O. M. R. (2011) The influence of distance and level of care on delivery place in rural Zambia: A study of linked national data in a geographic information system. *PLoS Medicine* **8**.
- [27] Afsana, K & Rashid, S. F. (2001) The challenges of meeting rural Bangladeshi women’s needs in delivery care. *Reproductive Health Matters* **9**, 79–89.
- [28] Amooti-Kaguna, B & Nuwaha, F. (2000) Factors influencing choice of delivery sites in Rakai district of Uganda. *Social Science and Medicine* **50**, 203–213.
- [29] Anyait, A, Mukanga, D, Oundo, G. B, & Nuwaha, F. (2012) Predictors for health facility delivery in Busia district of Uganda: a cross sectional study. *BioMed Central Pregnancy and Childbirth* **12**.
- [30] Metcalfe, R & Adegoke, A. a. (2013) Strategies to increase facility-based skilled birth attendance in South Asia: a literature review. *International health* **5**, 96–105.
- [31] Essendi, H, Mills, S, & Fotso, J.-C. (2011) Barriers to formal emergency obstetric care services’ utilization. *Journal of Urban Health : Bulletin of the New York Academy of Medicine* **88**, S356—S369.

- [32] Hodge, A, Byrne, A, Morgan, A, & Jimenez-Soto, E. (2014) Utilisation of Health Services and Geography: Deconstructing Regional Differences in Barriers to Facility-Based Delivery in Nepal. *Maternal and Child Health Journal* pp. 566–577.
- [33] Robertson, S. E & Valadez, J. J. (2006) Global review of health care surveys using lot quality assurance sampling (LQAS), 1984-2004. *Social science & medicine (1982)* **63**, 1648–60.
- [34] Rani, M, Bonu, S, Jha, P, Nguyen, S. N, & Jamjoum, L. (2003) Tobacco use in India: prevalence and predictors of smoking and chewing in a national cross sectional household survey. *Tobacco control* **12**, e4.
- [35] Ezzati, M, Lopez, A. D, Rodgers, A, Vander Hoorn, S, & Murray, C. J. (2002) Selected major risk factors and global and regional burden of disease. *The Lancet* **360**, 1347–1360.
- [36] Black, R. E, Morris, S. S, & Bryce, J. (2003) Where and why are 10 million children dying every year? *The Lancet* **361**, 2226–2234.
- [37] Jones, G, Steketee, R. W, Black, R. E, Bhutta, Z. A, & Morris, S. S. (2003) How many child deaths can we prevent this year? *The Lancet* **362**, 65–71.
- [38] United Nations. (2014) The Millennium Development Goals Report, Technical report.
- [39] Ronsmans, C & Graham, W. J. (2006) Maternal mortality: who, when, where, and why. *Lancet* **368**, 1189–1200.
- [40] Darmstadt, G, Bhutta, Z, Cousens, S, Adam, T, Walker, N, & de Bernis, L. (2005) Evidence-based, cost-effective interventions: how many newborn babies can we save? . *The Lancet* **365**, 977–988.
- [41] Ministry of Health. (2008) Situation analysis of newborn health in Uganda: current status and opportunities to improve care and survival, (Kampala), Technical report.
- [42] United Nations Population Fund. (2008) Skilled attendance at birth, (New York), Technical Report Web Page.
- [43] United Nations Children’s Fund, World Health Organization, & The World Bank. (2012) Levels & Trends in Child Mortality: Report 2012 Estimates

Developed by the UN Inter-agency Group for Child Mortality Estimation, (New York, New York), Technical Report Report.

- [44] Campbell, O. M & Graham, W. J. (2006) Strategies for reducing maternal mortality: getting on with what works. *The Lancet* **368**, 1284–1299.
- [45] Uganda Bureau of Statistics & ICF International Inc. (2012) Uganda Demographic and Health Survey 2011, (Kampala, Uganda and Calverton, Maryland), Technical Report Report.
- [46] Tann, C, Kizza, M, Morison, L, Mabey, D, Muwanga, M, Grosskurth, H, & Elliot, A. (2007) Use of antenatal services and delivery care in Entebbe, Uganda: a community survey. *BioMed Central Pregnancy and Childbirth* **7**.
- [47] Joharifard, S, Rulisa, S, Niyonkuru, F, Weinhold, A, Sayinzoga, F, Wilkinson, J, Ostermann, J, & Thielman, N. (2012) Prevalence and predictors of giving birth in health facilities in Bugesera district, Rwanda. *BioMed Central Public Health* **12**.
- [48] Parkhurst, J. O, Rahman, S. A, & Ssengooba, F. (2006) Overcoming access barriers for facility-based delivery in low-income settings: insights from Bangladesh and Uganda. *Journal of Health, Population and Nutrition* **24**, 438–445.
- [49] Demographic and Health Survey. (2013) Guide to DHS Statistics. Demographic and Health Surveys Methodology. pp. 1–168.
- [50] Bhuiya, A, Hanifi, S. M. a, Roy, N, & Streatfield, P. K. (2007) Performance of the lot quality assurance sampling method compared to surveillance for identifying inadequately-performing areas in Matlab, Bangladesh. *Journal of Health, Population and Nutrition* **25**, 37–46.
- [51] Henderson, R. H & Sundaresan, T. (1982) Cluster sampling to assess immunization coverage: A review of experience with a simplified sampling method. *Bulletin of the World Health Organization* **60**, 253–260.
- [52] Valadez, J. J, Hage, J, & Vargas, W. (2005) Understanding the relationship of maternal health behavior change and intervention strategies in a Nicaraguan NGO network. *Social science & medicine (1982)* **61**, 1356–1368.
- [53] Lehtonen Risto and Pahkinen Erkki. (2011) *Practical Methods for Design and Analysis of Complex Surveys*.

- [54] Robertson, S. E, Anker, M, Roisin, A. J, Macklai, N, Engstrom, K, & Laforce, F. M. (1997) The lot quality technique: a global review of applications in in the assessment of health services and disease surveillance. *Wld hlth statist. quart.* **50**.
- [55] Requejo, J, Bryce, J, & Victora, C. (2012) *Building a Future for Women and Children The 2012 Report*.
- [56] Valadez, J, Weiss, W, Leburg, C, & Davis, R. (2007) *Assessing Community Health Programs: A Trainer's Guide. Using LQAS for Baseline Surveys and Regular Monitoring*. (Teaching-aids At Low Cost, St Albans) Vol. second edi.
- [57] UNOCHA. (2011) Humanitarian Response: Common and Fundamental Operational Datasets Registry Uganda.
- [58] US Department of the Interior. (2010) USGS Science for a Changing World: Earth & Resources Observation and Science (EROS) Center, (Washington, DC), Technical Report Web Page.
- [59] Davison, A. C & Hinkley, D. V. (1997) *Bootstrap Methods and Their Applications*. (Cambridge University Press).
- [60] Field, C. A & Welsh, A. H. (2007) Bootstrapping clustered data. *Journal of the Royal Statistical Society B* **69**, 369–390.
- [61] Fox, J. (2008) *Applied regression analysis and generalized linear models*. Third edition.
- [62] Agresti, A. (2002) *Categorical data analysis*. Second edition.
- [63] Akaike, H. (1974) A new look at the statistical model identification. *Automatic Control, IEEE Transactions on*.
- [64] Wu, S, Flach, P, & Ferri, C. (2007) An Improved Model Selection Heuristic for AUC. pp. 478–489.
- [65] Jones, C. M & Athanasiou, T. (2005) Summary receiver operating characteristic curve analysis techniques in the evaluation of diagnostic tests. *The Annals of thoracic surgery* **79**, 16–20.
- [66] Pedregosa, F, Varoquaux, G, Gramfort, A, Michel, V, Thirion, B, Grisel, O, Blondel, M, Prettenhofer, P, Weiss, R, Dubourg, V, Vanderplas, J, Passos,

- A, Cournapeau, D, Brucher, M, Perrot, M, & Duchesnay, E. (2011) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **12**, 2825—2830.
- [67] Arlot, S & Celisse, A. (2010) A survey of cross-validation procedures for model selection. *Statistics Surveys* **4**, 40–79.
- [68] R. Development Core Team. (2013) The R Project for Statistical Computing.
- [69] Bivand, R & Ribeiro, P. J. (2015) Software for Spatial Statistics. **63**, 1–8.
- [70] McDonagh, T. a, Holmer, S, Raymond, I, Luchner, a, Hildebrandt, P, & Dargie, H. J. (2004) NT-proBNP and the diagnosis of heart failure: A pooled analysis of three European epidemiological studies. *European Journal of Heart Failure* **6**, 269–273.
- [71] Turcios, R. M, Widdowson, M.-A, Sulka, A. C, Mead, P. S, & Glass, R. I. (2006) Reevaluation of epidemiological criteria for identifying outbreaks of acute gastroenteritis due to norovirus: United States, 1998-2000. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America* **42**, 964–969.
- [72] Ministry of Health. (2012) Annual Health Sector Performance Report 2011-2012.
- [73] Bohren, M. a, Hunter, E. C, Munthe-Kaas, H. M, Souza, J. a, Vogel, J. P, & Gülmezoglu, a. (2014) Facilitators and barriers to facility-based delivery in low- and middle-income countries: a qualitative evidence synthesis. *Reproductive Health* **11**, 71.
- [74] Bolam, A, Manandhar, D. S, Shrestha, P, Ellis, M, Malla, K, & Costello, A. M. (1998) Factors affecting home delivery in the Kathmandu Valley, Nepal. *Health Policy and Planning* **13**, 152–158.
- [75] Gabrysch, S & Campbell, O. (2009) Still too far to walk: Literature review of the determinants of delivery service use. *BioMed Central Pregnancy and Childbirth* **9**.
- [76] Gage, A. J & Guirlène Calixte, M. (2006) Effects of the physical accessibility of maternal health services on their use in rural Haiti. *Population studies* **60**, 271–88.

- [77] Uganda Ministry of Finance, P & Development, E. (2010) Millennium Development Goals Report for Uganda 2010, Special theme: accelerating progress towards improving maternal health, (Kampala, Uganda), Technical report.
- [78] Interagency Health Team. (2011) Global Health Initiative Uganda: a strategy for accelerating reductions in maternal and neonatal mortality, (Kampala), Technical report.
- [79] Woldemicael, G & Tenkorang, E. Y. (2010) Women's autonomy and maternal health-seeking behavior in Ethiopia. *Maternal and child health journal* **14**, 988–998.
- [80] Celik, Y & Hotchkiss, D. R. (2000) The socio-economic determinants of maternal health care utilization in Turkey. *Social science & medicine (1982)* **50**, 1797–1806.
- [81] Mrisho, M, Schellenberg, J. a, Mushi, A. K, Obrist, B, Mshinda, H, Tanner, M, & Schellenberg, D. (2007) Factors affecting home delivery in rural Tanzania. *Tropical medicine & international health : TM & IH* **12**, 862–872.
- [82] Asante, A. D, Zwi, A. B, & Ho, M. T. (2006) Equity in resource allocation for health: A comparative study of the Ashanti and Northern Regions of Ghana. *Health Policy* **78**, 135–148.
- [83] Valadez, J. J. (2009) Learning to be creative with HIV/AIDS studies: Looking for the variation - not only the average. *International Journal of Epidemiology* **38**, 214–216.
- [84] Travis, P, Bennett, S, Haines, A, Pang, T, Bhutta, Z, & Hyder, a. a. (2004) Overcoming health-system constraints to achieve the millennium development goals. *The Lancet* **364**, 900–906.
- [85] Deming, M. S, Gayibor, a, Murphy, K, Jones, T. S, & Karsa, T. (1989) Home treatment of febrile children with antimalarial drugs in Togo. *Bulletin of the World Health Organization* **67**, 695–700.
- [86] Ruebush, T. K, Kern, M. K, Campbell, C. C, & Oloo, a. J. (1995) Self-treatment of malaria in a rural area of Western Kenya. *Bulletin of the World Health Organization* **73**, 229–236.
- [87] Rutebemberwa, E, Pariyo, G, Peterson, S, Tomson, G, & Kallander, K. (2009) Utilization of public or private health care providers by febrile children after user fee removal in Uganda. *Malaria journal* **8**, 45.

- [88] Deressa, W, Ali, a, & Enqusellassie, F. (2003) Self-treatment of malaria in rural communities, Butajira, southern Ethiopia. *Bulletin of the World Health Organization* **81**, 261–268.
- [89] Okeke, I. N, Lamikanra, a, & Edelman, R. (1999) Socioeconomic and behavioral factors leading to acquired bacterial resistance to antibiotics in developing countries. *Emerging infectious diseases* **5**, 18–27.
- [90] Goodman, C, Kachur, S. P, Abdulla, S, Bloland, P, & Mills, A. (2007) Drug shop regulation and malaria treatment in Tanzania - Why do shops break the rules, and does it matter? *Health Policy and Planning* **22**, 393–403.
- [91] Awor, P, Wamani, H, Bwire, G, Jagoe, G, & Peterson, S. (2012) Private sector drug shops in integrated community case management of malaria, pneumonia, and diarrhea in children in Uganda. *American Journal of Tropical Medicine and Hygiene* **87**, 92–96.
- [92] Wachter, D. a, Joshi, M. P, & Rimal, B. (1999) Antibiotic dispensing by drug retailers in Kathmandu, Nepal. *Tropical Medicine and International Health* **4**, 782–788.
- [93] McCombie, S. C. (2002) Self-treatment for malaria: the evidence and methodological issues. *Health policy and planning* **17**, 333–344.
- [94] Laing, R, Hogerzeil, H, & Ross-Degnan, D. (2001) Ten recommendations to improve use of medicines in developing countries. *Health policy and planning* **16**, 13–20.
- [95] Gouws, E, Bryce, J, Habicht, J. P, Amaral, J, Pariyo, G, Schellenberg, J. a, & Fontaine, O. (2004) Improving antimicrobial use among health workers in first level facilities: results from the multi-country evaluation of the integrated management of childhood illness strategy. *Bull World Health Organ* **82**, 509–515.
- [96] D’Acremont, V, Lengeler, C, Mshinda, H, Mtasiwa, D, Tanner, M, & Genton, B. (2009) Time to Move from Presumptive Malaria Treatment to Laboratory-Confirmed Diagnosis and Treatment in African Children with Fever. *PLoS Medicine* **6**, 7–9.
- [97] Buabeng, K. O, Duwiejua, M, Doodoo, A. N. O, Matowe, L. K, & Enlund, H. (2007) Self-reported use of anti-malarial drugs and health facility management of malaria in Ghana. *Malaria journal* **6**, 85.

- [98] Olango, P & Aboud, F. (1990) Determinants of mothers' treatment of diarrhea in rural Ethiopia. *Social Science and Medicine* **31**, 1245–1249.
- [99] Tanser, F, Hosegood, V, Benzler, J, & Solarsh, G. (2001) New approaches to spatially analyse primary health care usage patterns in rural South Africa. *Tropical medicine & international health : TM & IH* **6**, 826–838.
- [100] Stock, R. (1983) Distance and the utilization of health facilities in rural Nigeria. *Social science & medicine* **17**, 563–570.
- [101] Mwabu, G, Ainsworth, M, & Nyamete, A. (1993) Quality of Medical Care and Choice of Medical Treatment in Kenya : An Empirical Analysis. *The Journal of Human Resources* **28**, 838–862.
- [102] Akin, J. S & Hutchinson, P. (1999) Health-care facility choice and the phenomenon of bypassing. *Health Policy and Planning* **14**, 135–151.
- [103] Onwujekwe, O, Uzochukwu, B, Eze, S, Obikeze, E, Okoli, C, & Ochonma, O. (2008) Improving equity in malaria treatment: relationship of socio-economic status with health seeking as well as with perceptions of ease of using the services of different providers for the treatment of malaria in Nigeria. *Malaria journal* **7**, 5.
- [104] Shaikh, B. T & Hatcher, J. (2005) Health seeking behaviour and health service utilization in Pakistan: Challenging the policy makers. *Journal of Public Health* **27**, 49–54.
- [105] Ching, P. (1995) User fees, demand for children's health care and access across income groups: the Philippine case. *Social Science & Medicine* **41**, 37–46.
- [106] Schuler, S. R, Bates, L. M, & Islam, M. K. (2002) Paying for reproductive health services in Bangladesh: intersections between cost, quality and culture. *Health policy and planning* **17**, 273–280.
- [107] Stekelenburg, J, Jager, B. E, Kolk, P. R, Westen, E. H. M. N, Kwaak, a. V. D, & Wolffers, I. N. (2005) Health care seeking behaviour and utilisation of traditional healers in Kalabo, Zambia. *Health Policy* **71**, 67–81.
- [108] Sahn, D. E & Younger, S. D. (2002) The Demand for Health Care Services in Rural Tanzania. *Oxford Bulletin of Economics and statistics* **65**, 241—59.
- [109] Konde-lule, J, Okuonzi, S, Matsiko, C, Mukanga, D, Onama, V, & Gitta, S. N. (2006) The Potential of the Private Sector to Improve Health Outcomes in

Uganda T, (Makerere University Institute of Public Health, Kampala), Technical Report July.

- [110] El Arifeen, S, Blum, L. S, Hoque, D. M. E, Chowdhury, E. K, Khan, R, Black, R. E, Victora, C. G, & Bryce, J. (2004) Integrated Management of Childhood Illness (IMCI) in Bangladesh: early findings from a cluster-randomised study. *The Lancet* **364**, 11–13.
- [111] Basinga, P, Gertler, P. J, Binagwaho, A, Soucat, A. L, Sturdy, J, & Vermeersch, C. M. (2011) Effect on maternal and child health services in Rwanda of payment to primary health-care providers for performance: An impact evaluation. *The Lancet* **377**, 1421–1428.
- [112] Mugisha, F, Bocar, K, Dong, H, Chepng'eno, G, & Sauerborn, R. (2004) The two faces of enhancing utilization of health-care services: Determinants of patient initiation and retention in rural Burkina Faso. *Bulletin of the World Health Organization* **82**, 572–579.
- [113] Auer, C, Sarol, J, Tanner, M, & Weiss, M. (2000) Health seeking and perceived causes of tuberculosis among patients in Manila, Philippines. *Tropical Medicine and International Health* **5**, 648–656.
- [114] Pathania, V, Almeida, J, & Kochi, A. (1997) TB Patients and Private For-Profit Health Care Providers in India, (World Health Organisation), Technical report.
- [115] Uplekar, M. W, Juvekar, S. K, Parande, S. D, & Dalai, D. B. (1996) Tuberculosis Management in Private Practice and Its Implications. *Indian Journal of Tuberculosis* pp. 6–9.
- [116] Bhatia, J. C & Cleland, J. (2001) Health-care seeking and expenditure by young Indian mothers in the public and private sectors. *Health policy and planning* **16**, 55–61.
- [117] Meessen, B, Soucat, A, & Sekabaraga, C. (2011) Performance-based financing: just a donor fad or a catalyst towards comprehensive health-care reform? *Bulletin of the World Health Organization* **89**, 153–156.
- [118] Meessen, B, Musango, L, Kashala, J. P. I, & Lemlin, J. (2006) Reviewing institutions of rural health centres: The Performance Initiative in Butare, Rwanda. *Tropical Medicine and International Health* **11**, 1303–1317.

- [119] Ireland, M, Paul, E, & Dujardin, B. (2011) Can performance-based financing be used to reform health systems in developing countries? *Bulletin of the World Health Organization* **89**, 695–698.
- [120] Oxman, A. D & Fretheim, A. (2009) Can paying for results help to achieve the Millennium Development Goals? Overview of the effectiveness of results-based financing. *Journal of Evidence-Based Medicine* **2**, 70–83.
- [121] Victora, C. G, Habicht, J. P, & Bryce, J. (2004) Evidence-Based Public Health: Moving Beyond Randomized Trials. *American Journal of Public Health* **94**, 400–405.
- [122] Campbell, M. K & Grimshaw, J. M. (1998) Cluster randomised trials: time for improvement. *BMJ : British Medical Journal* **317**, 1171–1172.
- [123] Campbell, M, Fitzpatrick, R, Haines, A, Kin outh, A, Sandercock, P, Spiegelhalter, D, & Tyrer, P. (2000) Framework for design and evaluation of complex interventions to improve health. *BMJ* **321**, 694–696.
- [124] Gelman, A, Carlin, J. B, Stern, H. S, Dunson, D. B, Vehtari, A, & Rubin, D. B. (2014) *Bayesian Data Analysis*. (Taylor & Francis, Boca Raton), Third edition.
- [125] Gelman, A, Jakulin, A, Pittau, M. G, & Su, Y. S. (2008) A weakly informative default prior distribution for logistic and other regression models. *Annals of Applied Statistics* **2**, 1360–1383.
- [126] Stan Development Team. (2014) PyStan: the Python interface to Stan, Version 2.5.0.
- [127] Beskos, a, Pillai, N. S, Roberts, G. O, Sanz-Serna, J. M, & Stuart, a. M. (2010) The acceptance probability of the hybrid Monte Carlo method in high-dimensional problems. *AIP Conference Proceedings* **1281**, 23–26.
- [128] Cowles, M. K & Carlin, B. P. (1996) Markov Chain Monte Carlo Convergence Diagnostics : A Comparative Review. *J. Am. Stat. Assoc.* **91**, 883–904.
- [129] Gelman, A & Rubin, D. B. (1992) Inference from Iterative Simulation Using Multiple Sequences. *Statistical Science* **7**, 457–472.
- [130] Fahrmeir, L & Kaufmann, H. (1985) Consistency and Asymptotic Normality of the Maximum Likelihood Estimator in Generalized Linear Models. *The Annals of Statistics* **13**, 342–368.

- [131] Hagiwara, K. (2002) On the problem in model selection of neural network regression in overrealizable scenario. *Neural computation* **14**, 1979–2002.
- [132] Gelman, A, Hwang, J, & Vehtari, A. (2013) Understanding predictive information criteria for Bayesian models. *Statistics and Computing* pp. 1–20.
- [133] Watanabe, S. (2010) Asymptotic Equivalence of Bayes Cross Validation and Widely Applicable Information Criterion in Singular Learning Theory. *Journal of Machine Learning Research* **11**, 3571–3594.
- [134] Gelfand, A. E, Diggle, P. J, Fuentes, M, & Guttorp, P. (2010) *Handbook of spatial statistics*. (CRC Press, Boca Raton).
- [135] Azhar, S, Hassali, M. A, Izham, M, Ibrahim, M, Ahmad, M, Masood, I, & Shafie, A. A. (2009) The role of pharmacists in developing countries: the current scenario in Pakistan. *Human Resources for Health* **7**, 1–6.
- [136] Bell, B, Morgan, G, Kromrey, J, & Ferron, J. (2010) The impact of small cluster size on multilevel models: a Monte Carlo examination of two-level models with binary and continuous predictors. *JSM Proceedings, Section on Survey Research Methods* pp. 4057–4067.
- [137] Clarke, P. (2008) When can group level clustering be ignored ? Multilevel models versus single-level models with sparse data. *Journal of epidemiology and community health* pp. 752–758.
- [138] Bryceson, D. F, Mbara, T. C, & Maunder, D. (2003) Livelihoods, daily mobility and poverty in sub-saharan Africa. *Transport Reviews* **23**, 177–196.
- [139] Van Damme, W, Van Leemput, L, Por, I, Hardeman, W, & Meessen, B. (2004) Out-of-pocket health expenditure and debt in poor households: Evidence from Cambodia. *Tropical Medicine and International Health* **9**, 273–280.
- [140] Russells, S. (2004) Economic burden of illness for household in developing countries: a review of studies focusing on malaria, tuberculosis and human immunodeficiency virus /acquired immunodeficiency syndrome. *Am. J.Trop. Med. Hyg* **71**, 147–155.
- [141] Salganik, M. J, Dodds, P. S, & Watts, D. J. (2006) Experimental Study of Inequality and Unpredictability in an Artificial Cultural Market. *Science* **311**.

- [142] Grinblatt, M, Keloharju, M, & Ikäheimo, S. (2008) Social influence and consumption: evidence from the automobile purchases of neighbors. *The Review of Economics and Statistics* **90**, 735–753.
- [143] Bond, R. M, Fariss, C. J, Jones, J. J, Kramer, A. D. I, Marlow, C, Settle, J. E, & Fowler, J. H. (2012) A 61-million-person experiment in social influence and political mobilization. *Nature* **489**, 295–298.
- [144] Kahan, D. M. (1997) Social influence, social meaning, and deterrence. *Virginia Law Review* **83**, 349–395.
- [145] Centola, D. (2010) The spread of behavior in an online social network experiment. *Science (New York, N.Y.)* **329**, 1194–1197.
- [146] Christakis, N. A & Fowler, J. H. (2008) The collective dynamics of smoking in a large social network. *New England journal of medicine* **358**, 2249–2258.
- [147] Christakis, N. A & Fowler, J. H. (2007) The spread of obesity in a large social network over 32 years. *New England journal of medicine*.
- [148] Cialdini, R. B & Goldstein, N. J. (2004) Social influence: compliance and conformity. *Annual review of psychology* **55**, 591–621.
- [149] Keating, N. L, O'Malley, a. J, Murabito, J. M, Smith, K. P, & Christakis, N. a. (2011) Minimal social network effects evident in cancer screening behavior. *Cancer* **117**, 3045–3052.
- [150] Shalizi, C. R & Thomas, A. C. (2011) Homophily and Contagion Are Generically Confounded in Observational Social Network Studies. *Sociological Methods & Research* **40**, 211–239.
- [151] Aral, S, Muchnik, L, & Sundararajan, A. (2009) Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks. *Proceedings of the National Academy of Sciences of the United States of America* **106**, 21544–21549.
- [152] Bass, F. M. (1969) A New Product Growth for Model Consumer Durables. *Management Science* **15**, 215–227.
- [153] Bauckhage, C, Kersting, K, & Hadiji, F. (2013) Mathematical Models of Fads Explain the Temporal Dynamics of Internet Memes. *Proc. ICWSM, AAAI* pp. 22–30.

- [154] Gleeson, J. P, Cellai, D, Onnela, J.-P, Porter, M. a, & Reed-Tsochas, F. (2014) A simple generative model of collective online behavior. *Proceedings of the National Academy of Sciences of the United States of America* **111**, 10411–10415.
- [155] Bentley, R. A, Ormerod, P, & Batty, M. (2010) Evolving social influence in large populations. *Behavioral Ecology and Sociobiology* **65**, 537–546.
- [156] Dietz, K. (1967) Epidemics and rumours: A survey. *Journal of the Royal Statistical Society. Series A* **130**, 505–528.
- [157] Isham, V, Harden, S, & Nekovee, M. (2010) Stochastic epidemics and rumours on finite random networks. *Physica A: Statistical Mechanics and its Applications* **389**, 561–576.
- [158] Dodds, P & Watts, D. (2004) Universal Behavior in a Generalized Model of Contagion. *Physical Review Letters* **92**, 218701.
- [159] Gillespie, D. T. (1977) Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry* **93555**, 2340–2361.
- [160] House, T, Ross, J. V, & Sirl, D. (2013) How big is an outbreak likely to be? Methods for epidemic final-size calculation. *Proceedings of the Royal Society A*.
- [161] Neuts, M. F & Li, J. M. (1996) in *Athens Conference on Applied Probability and Time Series Analysis*, eds. Heyde, C. C, Prohorov, Y. V, Pyke, R, & Rachev, S. T. (Springer, Berlin), pp. pp 295–306.
- [162] Black, A. J & Ross, J. V. (2015) Computation of epidemic final size distributions.
- [163] House, T. (2011) Modelling behavioural contagion. *Journal of the Royal Society, Interface* **8**, 909–912.
- [164] Hill, R. A & Dunbar, R. I. M. (2003) Social network size in humans. *Human Nature* **14**, 53–72.
- [165] Denrell, J. C & Kovacs, B. (2011) The Dynamics of Popularity. *Academy of Management Journal* pp. 1–32.
- [166] Kurtz, T. G. (1970) Solutions of Ordinary Differential Equations as Limits of Pure Jump Markov Processes. *Journal of applied probability* **7**, 49–58.

- [167] Kurtz, T. G. (1971) Limit Theorems for Sequences of Jump Markov Processes Approximating Ordinary Differential Processes. *Journal of Applied Probability* **8**, 344–356.
- [168] Riley, S, Fraser, C, Donnelly, C. a, Ghani, A. C, Abu-Raddad, L. J, Hedley, A. J, Leung, G. M, Ho, L.-M, Lam, T.-H, Thach, T. Q, Chau, P, Chan, K.-P, Lo, S.-V, Leung, P.-Y, Tsang, T, Ho, W, Lee, K.-H, Lau, E. M. C, Ferguson, N. M, & Anderson, R. M. (2003) Transmission dynamics of the etiological agent of SARS in Hong Kong: impact of public health interventions. *Science (New York, N.Y.)* **300**, 1961–1966.
- [169] Nishiura, H. (2007) Time variations in the transmissibility of pandemic influenza in Prussia, Germany, from 1918-19. *Theoretical biology & medical modelling* **4**, 20.
- [170] Funk, S, Salathé, M, & Jansen, V. a. a. (2010) Modelling the influence of human behaviour on the spread of infectious diseases: a review. *Journal of the Royal Society, Interface / the Royal Society* **7**, 1247–1256.
- [171] Strogatz, S. H. (2014) *Nonlinear dynamics and chaos: with applications to physics, biology, chemistry, and engineering* ed. Press, W. 2 edition.
- [172] Hethcote, H. W, Yorke, J. a, & Nold, A. (1982) Gonorrhea modeling: a comparison of control methods. *Mathematical Biosciences* **58**, 93–109.
- [173] Hiltunen-Back, E, Rostila, T, Kautiainen, H, Paavonen, J, & Reunala, T. (1998) Rapid decrease of endemic gonorrhea in Finland. *Sexually transmitted diseases* **25**, 181–186.
- [174] Steward, W. T, Remien, R. H, Higgins, J. a, Dubrow, R, Pinkerton, S. D, Sikkema, K. J, Truong, H. H. M, Johnson, M. O, Hirsch, J, Brooks, R. a, & Morin, S. F. (2009) Behavior change following diagnosis with acute/early HIV infection-a move to serosorting with other HIV-infected individuals. the NIMH multisite acute HIV infection study: III. *AIDS and Behavior* **13**, 1054–1060.
- [175] Fox, J, White, P. J, Macdonald, N, Weber, J, McClure, M, Fidler, S, & Ward, H. (2009) Reductions in HIV transmission risk behaviour following diagnosis of primary HIV infection: A cohort of high-risk men who have sex with men. *HIV Medicine* **10**, 432–438.
- [176] Marks, G, Crepaz, N, Senterfitt, J. W, & Janssen, R. S. (2005) Meta-analysis of high-risk sexual behavior in persons aware and unaware they are infected

- with HIV in the United States: implications for HIV prevention programs. *Journal of acquired immune deficiency syndromes (1999)* **39**, 446–453.
- [177] Allen, S, Meinzen-Derr, J, Kautzman, M, Zulu, I, Trask, S, Fideli, U, Musonda, R, Kasolo, F, Gao, F, & Haworth, A. (2003) Sexual behavior of HIV discordant couples after HIV counseling and testing. *AIDS (London, England)* **17**, 733–740.
 - [178] de Vincenzi, I. (1994) A longitudinal study of human immunodeficiency virus transmission by heterosexual partners. *The New England journal of medicine* **331**, 341–346.
 - [179] Apicella, C. L, Marlowe, F. W, Fowler, J. H, & Christakis, N. a. (2012) Social networks and cooperation in hunter-gatherers. *Nature* **481**, 497–501.
 - [180] Liben-Nowell, D, Novak, J, Kumar, R, Raghavan, P, & Tomkins, A. (2005) Geographic routing in social networks. *Proceedings of the National Academy of Sciences of the United States of America* **102**.
 - [181] Gelfand, A. E. (2012) Hierarchical Modeling for Spatial Data Problems. *Spatial statistics* **1**, 30–39.
 - [182] Kiti, M. C, Kinyanjui, T. M, Koech, D, Munywoki, P. K, Medley, G. F, & Nokes, D. J. (2014) Quantifying age-related rates of social contact using diaries in a rural coastal population of Kenya. *PLoS ONE* **9**.
 - [183] Schliep, E & Hoeting, J. (2013) Multilevel Latent Gaussian Process Model for Mixed Discrete and Continuous Multivariate Response Data. *Journal of Agricultural, Biological, and Environmental Statistics* pp. 1–25.
 - [184] Weir, I. S & Pettitt, A. N. (2000) Binary probability maps using a hidden conditional autoregressive Gaussian process with an application to Finnish common toad data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **49**, 473–484.