

## Distribution Agreement

In presenting this thesis or dissertation as a partial fulfillment of the requirements for an advanced degree from Emory University, I hereby grant to Emory University and its agents the non-exclusive license to archive, make accessible, and display my thesis or dissertation in whole or in part in all forms of media, now or hereafter known, including display on the world wide web. I understand that I may select some access restrictions as part of the online submission of this thesis or dissertation. I retain all ownership rights to the copyright of the thesis or dissertation. I also retain the right to use in future works (such as articles or books) all or part of this thesis or dissertation.

Signature:

---

Qian An

---

Date

# Models for Statistical Analyses of Infectious Disease Data

By

Qian An

Doctor of Philosophy

Biostatistics

---

Jian Kang, Ph.D.  
Advisor

---

Michael Haber, Ph.D.  
Advisor

---

Howard Chang, Ph.D.  
Committee Member

---

H Irene Hall, Ph.D.  
Committee Member

Accepted:

---

Lisa A. Tedesco, Ph.D.  
Dean of the James T. Laney School of Graduate Studies

---

Date

# Models for Statistical Analyses of Infectious Disease Data

By

Qian An

M.S., University of Minnesota: Twin Cities, 2005

B.E., Beijing Jiaotong University, 1999

Advisors : Jian Kang, Ph.D. and Michael Haber, Ph.D.

An abstract of

A dissertation submitted to the Faculty of the

James T. Laney School of Graduate Studies of Emory University

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

in Biostatistics

2014

## Abstract

### Models for Statistical Analyses of Infectious Disease Data

By Qian An

This dissertation aims at developing new statistical methods to monitor changes in HIV testing behaviors and to evaluate the influenza vaccine effectiveness (VE). The proposed approaches can help evaluate the quality of public health programs and provide guidance for expanding future public health responses.

In the first project, we propose a two-level Bayesian hierarchical model to estimate the HIV testing rate using annual acquired immunodeficiency syndrome (AIDS) and HIV diagnosis data. We introduce a new class of priors for the HIV incidence rate and testing rate taking into account the temporal dependence of these parameters to improve the estimation accuracy. We develop an efficient posterior computation algorithm based on the adaptive rejection metropolis sampling technique (ARMS). The proposed approach is illustrated via simulation studies and the analysis of the national HIV surveillance data in the United States.

In the second project, we propose a novel Bayesian model to estimate the influenza VE using data collected from the test negative design (TND). Given that a person is sampled into TND, the joint probability of this person's vaccination status and influenza infection status is modeled as a function of the influenza VE. To improve the estimation accuracy, subjective priors are elicited from published literatures. We resort to ARMS for efficient posterior computation. To demonstrate the superiority of our approach, we perform simulation studies where model-based estimates of influenza VE are compared with existing odds ratio estimates.

In the third project, we propose an improved nonhomogeneous probability model for evaluating bias and precision of the estimates of influenza VE from traditional case-control design (CCD) and TND. The proposed model describes the data generation process in real life composed of five steps: latent health status, vaccination, acute respiratory illness (ARI) and influenza infection; seeking medical care for ARI and testing for influenza infection. By including a parameter for the latent health status, this model facilitates the evaluation of an important bias resulting from the unobserved variable. We present and compare the numerical results of the bias and the standard error of the VE estimates from the CCD and the TND.

# Models for Statistical Analyses of Infectious Disease Data

By

Qian An

M.S., University of Minnesota: Twin Cities, 2005

B.S., Beijing Jiaotong University, 1999

Advisors: Jian Kang, Ph.D. and Michael Haber, Ph.D.

A dissertation submitted to the Faculty of the  
James T. Laney School of Graduate Studies of Emory University  
in partial fulfillment of the requirements for the degree of  
Doctor of Philosophy  
in Biostatistics

2014

## Acknowledgement

Foremost, I would like to gratefully and sincerely thank my advisors Dr. Jian Kang and Dr. Michael Haber for their patience, motivation, enthusiasm and immense knowledge. Without their continuous guidance and encouragement, my research and dissertation cannot be done. Their knowledge, academic rigor, as well as enthusiasm of statistical research influenced me a lot, enabling me to pursue a continuing career in the statistical field. I feel very fortunate to be their student, and this learning experience under their supervision will be never forgotten.

Besides my advisors, I would also like to thank my committee members, Dr. Howard Chang and Dr. H Irene Hall for their encouragement, insightful comments, detailed review, and active discussion on my dissertation. Their professional background on Bayesian modeling and infectious disease epidemiology provide me constructive advice, making the work integrated and comprehensive. In particular, I sincerely thank Dr. Hall for her generosity in her continuous support and guidance during my Ph.D studies.

My sincere thanks also goes to Dr. Joseph Prejean and Dr. Angela Hernandez from the Centers for Disease Control and Prevention for their consistent support and encouragement during my Ph.D studies.

I thank the Department of Biostatistics and Bioinformatics at Emory University, especially Dr. Limin Peng and Dr. Lance Waller for providing me the opportunity to study and work on exciting projects. I thank all the wonderful faculty, excellent staff, and great classmates for making my graduate study life at Emory delightful

and unforgettable.

Last but not least, I would like to thank my family: my husband Xiangming, my daughter Emma and my parents, Jun and Guangyou, for their unwavering love and support. I cannot image how my life would be without them, and I hope this work could be the best reward to them. Thanks for all your unconditional dedication, faith and encouragement.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>A Bayesian Hierarchical Model with Novel Prior Specifications to estimate HIV Testing Rate</b>	<b>9</b>
2.1	Background . . . . .	9
2.2	The Model . . . . .	14
2.2.1	Data and Notation . . . . .	14
2.2.2	A Hierarchical Model . . . . .	15
2.2.3	Prior Specifications . . . . .	19
2.2.4	Posterior Inference . . . . .	24
2.3	Simulation Study . . . . .	25
2.4	Application . . . . .	30
2.4.1	Analysis of the United States HIV surveillance data . . . . .	30
2.4.2	Model Assessment . . . . .	35
2.5	Discussion . . . . .	36



<b>3</b>	<b>A Bayesian Model to Estimate the Influenza Vaccine Effectiveness from a Test Negative Design</b>	<b>38</b>
3.1	Background . . . . .	38
3.2	The model . . . . .	42
3.2.1	Model representation . . . . .	44
3.2.2	Prior specifications . . . . .	45
3.2.3	Posterior Inference . . . . .	49
3.3	Simulation Study . . . . .	50
3.4	Application . . . . .	53
3.4.1	Analysis of the seasonal 2010-2011 influenza vaccine TND study	54
3.4.2	Model assessment . . . . .	57
3.4.3	Sensitivity Analysis . . . . .	58
3.5	Discussion . . . . .	60
<b>4</b>	<b>A Nonhomogeneous Probability Model for Evaluating Bias and Precision of Estimates of the Influenza Vaccine Effectiveness from Case-Control Studies</b>	<b>62</b>
4.1	Background . . . . .	62
4.1.1	Main sources of bias in case-control studies . . . . .	64
4.1.2	Medically-attended influenza and symptomatic influenza . . . . .	66
4.2	Method . . . . .	68
4.2.1	The study population and designs . . . . .	68

4.2.2	The model . . . . .	69
4.2.3	Outcome of interest and true VE . . . . .	72
4.2.4	VE estimates . . . . .	73
4.2.5	Standard errors of the VE estimates . . . . .	75
4.2.6	Determining the values of the parameters . . . . .	77
4.3	Results . . . . .	78
4.3.1	Analytic results . . . . .	78
4.3.2	Numeric results . . . . .	80
4.3.3	Summary of sources of bias under non-random vaccination . . . . .	95
4.4	Discussion . . . . .	95

# List of Figures

1.1	The natural progression of HIV infection from primary infection through Acute HIV Syndrome to clinical latency. (O'Brien and Hendrickson 2013) . . . . .	3
2.1	Estimated posterior mean and 95% credible intervals for HIV testing rates and time since infection with different choices of $(c^H, c^\lambda)$ in simulation Scenarios 1 and 2.	31
2.2	Estimated posterior mean and 95% credible intervals for the annual HIV testing rates and expected time-since-infection from 1985 to 2010 in United States . . .	34
2.3	Scatterplot of predictive versus realized $\chi^2$ discrepancies under the joint posterior distribution; the $p$ -value is estimated by the proportion of points above the 45 degree line. . . . .	35
4.1	Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation A1 from baseline assumption A while assumptions B and C remain in place. $\rho_\beta$ is the vaccine related ratio in the probability of non-influenza ARI, ranging from 0.25 to 4. The bias is the same for SI and MI. . . . .	83

4.2	Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation A2 from baseline assumption A while assumptions B and C remain in place. $\eta_\beta$ is the health status related ratio in the probability of non-influenza ARI, ranging from 0.25 to 4. The bias is the same for SI and MI. . . . .	84
4.3	Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation B1 from baseline assumption B while assumptions A and C remain in place. $\eta_\gamma$ is the health status related ratio in the probability of influenza ARI, ranging from 0.25 to 4. The bias is the same for SI and MI. . . . .	85
4.4	Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation B2 from baseline assumption B while assumptions A and C remain in place. $\theta_\gamma$ is the inequality in vaccine related ratio in the probability of influenza ARI resulted from health status, ranging from 0.25 to 2. The bias is the same for SI and MI. . . . .	86
4.5	Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation C1 from baseline assumption C while assumptions A and B remain in place. $\eta_\delta$ is the health status related ratio in the probability of seeking medical care for ARI, ranging from 0.25 to 4. The bias is the same for SI and MI in Scenario 1. . . . .	87
4.6	Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation C2 from baseline assumption C while assumptions A and B remain in place. $\theta_\delta$ is the inequality in health status related ratio in the probability of seeking medical care for ARI resulted from the type of infection leading to ARI, ranging from 0.25 to 4. The bias is the same for SI and MI in Scenario 1. . . . .	88

# List of Tables

- 2.1 Illustration of the relationship between observed data and latent quantities in the model. The column totals represent the number of persons diagnosed with AIDS and HIV not AIDS in each year and those undiagnosed by the end of the most recent year. Row totals represent the number of new infections in each year. In each cell,  $A_{it}$  represents the number of persons infected in year  $i$  and diagnosed with AIDS in year  $t$  and  $H_{it}$  represents the number of persons infected in year  $i$  and diagnosed with HIV not AIDS in year  $t$ . Only column totals  $\{A_t\}_{t=1}^T$  and  $\{H_t\}_{t=1}^T$  are observed, and all other quantities are latent. . . . . 15
- 2.2 Values of AIDS diagnosis rate generated from the hazard function of a Gamma distribution with shape parameter of 2 and scale parameter of 4.  $t$  is year of AIDS diagnosis and  $i$  is year of infection. . . . . 26

2.3	Values for the parameters used in the simulation studies. $\lambda$ is the mean annual number of new HIV infections, $p^H$ is the annual HIV testing rate, $A$ is the observed annual number of AIDS diagnoses and $H$ is the observed annual number of HIV diagnoses. Scenario 1 is for a 34-year period with a gradual increasing trend in $p^H$ and Scenario 2 is for a 20-year period with an increasing trend followed by a decreasing trend in $p^H$ . . . . .	27
2.4	Simulation model fitting results for different choices of hyperparameters ( $c^H, c^\lambda$ ). $\log(\text{BF})$ is the estimated log Bayes factor. AMSE is the average mean square error, ACI is the average length of 95% credible intervals. . . . .	30
3.1	Parameters and notation used in the model representation and their value ranges . . . . .	46
3.2	Estimates of $\beta_{1v}$ and $\beta_{2v}$ from the 14 comparisons identified through 5 publications . . . . .	48
3.3	True values for each parameter and each scenario used in the simulation studies . . . . .	50
3.4	Simulation model fitting results for each scenario: model-based estimates versus OR estimates. MSE is the mean square error. . . . .	52
3.5	Simulation model fitting results: absolute relative bias for estimates of all parameters for each scenario . . . . .	53

3.6	Number of patients by influenza status and vaccination status for each age group and each vaccine component in the 2010-2011 TND study	55
3.7	The model-based and unadjusted odds ratio (OR) influenza VE estimates for each age group and each vaccine type in the 2010-2011 TND study . . . . .	56
3.8	The priors for $\beta_{10}$ and $\beta_{20}$ used in the sensitivity analysis . . . . .	58
3.9	The mean, minimum, median, and maximum of VE estimates from sensitivity analyses . . . . .	59
4.1	Random variables used to define the five steps of the latent-class process model and in the calculation . . . . .	72
4.2	List of parameters and other notation used in this article . . . . .	79
4.3	Bias and standard errors of VE estimates for SI and MI from TND and CCD for scenario 1 when vaccination is random ( $\alpha_0 = \alpha_1 = 0.6$ ) . . .	89
4.4	Bias and standard errors of VE estimates for SI and MI from TND and CCD for scenario 2 when healthy persons are twice more likely to get vaccination than frail persons ( $\alpha_0 = 0.4, \alpha_1 = 0.8$ ) . . . . .	90
4.5	Bias and standard errors of VE estimates for SI and MI from TND and CCD for scenario 3 when frail persons are twice more likely to get vaccination than healthy persons ( $\alpha_0 = 0.8, \alpha_1 = 0.4$ ) . . . . .	91
4.6	Situations for severe bias ( $ bias  > 0.1$ ) under non-random vaccination	96

# Chapter 1

## Introduction

Despite the progress in the diagnosis, prevention and treatment of infectious disease over the past three decades, human immunodeficiency virus infection (HIV) and influenza still pose significant public health challenges domestically and globally. The World Health Organization (WHO) reports that HIV has claimed more than 25 million lives over the past three decades and there were approximately 35.3 million people living with HIV in 2012. (WHO Media Centre 2013) In the United States, more than 1.1 million people were living with the HIV infection at the end of 2010. (CDC 2012) Similarly, seasonal and pandemic influenza could also sweep the globe, leading to hospitalization and massive death. In the United States, influenza is responsible for more than 226,000 hospitalizations each year. (Thompson et al. 2004) The estimates of annual influenza-related death during 1976 to 2007 ranged from 3000 to 48000. (Thompson et al. 2010; 2003) In June 2009, WHO declared the 2009 influenza A (H1N1) pandemic and in October 2009, President Obama declared a national emergency. (Larson and Heymann 2010)



Maintaining and expanding public health responses to HIV and influenza require the assessment and evaluation of public health prevention and intervention programs. Evaluating the effectiveness of the public health programs can disseminate information on whether the programs achieve the desired results, pinpoint the areas that need improvement and provide guidance for expanding future public health responses. Although the conceptual framework for public health program evaluation is established and systematic, evaluating specific programs or issues frequently requires sophisticated statistical models. This dissertation aims to propose novel statistical methods for estimating the HIV testing rates to monitor changes in HIV testing behaviors and for evaluating the effectiveness of influenza vaccine. These estimates help public health officials to evaluate the impact of HIV testing initiatives and the effectiveness of the influenza vaccines in the United States, and provide guidance for expanding future HIV testing and influenza vaccine services in the United States.

The first project of this dissertation is motivated by the need to make assessment of the HIV testing behavior changes in the United States. First reported in 1981, HIV has become one of the greatest public health challenges both domestically and globally. Unlike most other viruses, HIV attacks a key part of the human immune system and over time it can cause badly damaged immune system, which puts people at risk for fatal opportunistic illness. Figure 1.1 (O'Brien and Hendrickson 2013) outlines the natural history of the HIV disease: from the primary HIV infection, the acute HIV syndrome, and the HIV-specific immune response through a long period of clinical latency to clinically apparent diseases or AIDS-defining illness, and finally death from AIDS. After HIV infection, the "acute HIV viral syndrome" with influenza-like symptoms develops in 40-70% of patients while the symptoms may be very mild or may

not show up at all in other people with HIV. The initial acute infection stage usually lasts one or two weeks and during this period of time serum testing for HIV antibody may, or may not be positive. HIV antibody test is positive in most individuals within one to three months after primary infection and in about 95% patients within six months. The initial symptoms are followed by a very short period of HIV-specific immune response and a prolonged period called clinical latency. The duration of clinical latency varies widely among individuals. Without treatment, the latency stage can last for about three years to over 20 years with a median of 10 years, during which many people are asymptomatic and others develop different symptoms at different times. In the absence of treatment, around half of people infected with HIV develop significant clinically apparent diseases or AIDS-defining illness. (Moore and Chaisson 1999)

HIV testing is the cornerstone for HIV prevention. It can foster early detection of HIV infection and is the first essential step for entry to clinical care to reduce morbidity and mortality. In addition, persons aware of their infections and on treatment are less likely to transmit the virus to others. Since HIV testing first became available in 1985, the importance of HIV testing was soon recognized, emphasized and promoted (Branson et al. 2006). The importance of HIV testing was recognized as early as in 1987, just two years after the first HIV test became available, when the United States Public Health Service (USPHS) issued guidelines making HIV counseling and testing a priority as a prevention strategy for people with high risk behaviors. In 1993, the Centers for Disease Control and Prevention (CDC) updated the recommendations regarding HIV counseling and testing for patients in acute-care hospital settings. In 1995, the first National HIV Testing Day was observed. Throughout the

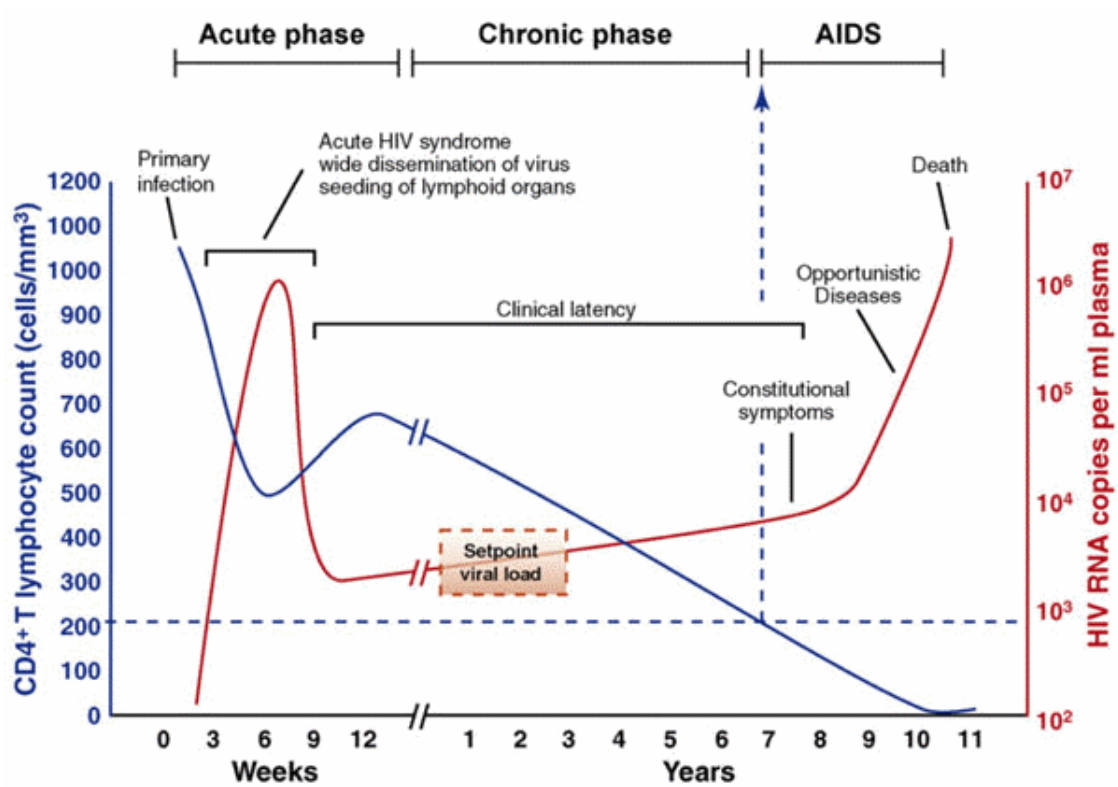


Figure 1.1: The natural progression of HIV infection from primary infection through Acute HIV Syndrome to clinical latency. (O'Brien and Hendrickson 2013)

decade of 2000, a few important recommendations were published in 2001, 2003 and 2006 (Branson et al. 2006) respectively to emphasize routine HIV testing as an important HIV prevention tool for adults, adolescents, and pregnant women in health-care settings. In 2010, CDC implemented the Expanded Testing Initiative to focus on increasing HIV testing among high risk populations, such as African Americans and Latinos as well as gay, bisexual, or other men who have sex with men and injection drug users of all races and ethnicities. In April 2013, the U.S. Preventive Services Task Force (USPSTF) also released HIV testing recommendations that everyone aged 15 to 65 should be screened for HIV infection; teens younger than age 15 and adults older than 65 also should be screened if they are at increased risk for HIV infection; and all pregnant women, including women in labor who do not know if they are infected with HIV, should be screened for HIV infection. Over the years, numerous public health recommendations and initiatives on HIV testing were established. It is important to monitor whether there has been increased HIV testing in the past years and to assess the impact of the public health recommendations and initiatives.

In the first project of this dissertation, we focus on estimating the HIV testing rates in the past thirty years in the United States. We define the HIV testing rate as the probability that an HIV infected person seeks a test and gets diagnosed with HIV, not AIDS, which is an advanced stage of HIV infection with CD4 T-lymphocyte count less than 200/ul or opportunistic illnesses, given no previous positive test has been obtained. We propose a Bayesian hierarchical model with two levels of hierarchy to estimate the HIV testing rate using annual AIDS and HIV diagnoses data. At level one, we model the latent number of HIV infections for each year using a Poisson distribution with the intensity parameter representing the HIV incidence rate. At level

two, the annual number of AIDS and HIV diagnosed cases and all undiagnosed cases stratified by the HIV infections at different years are modeled using a multinomial distribution with parameters including the HIV testing rate. We propose a new class of priors for the HIV incidence rate and HIV testing rate taking into account the temporal dependence of these parameters to improve the estimation accuracy. We develop an efficient posterior computation algorithm based on the adaptive rejection metropolis sampling technique. We demonstrate our model using simulation studies and the analysis of the national HIV surveillance data in United States.

The second and third projects in this dissertation are motivated by the need to assess the effectiveness of influenza vaccine. Influenza is a highly contagious viral infection, one of the most severe illness of the winter season. Seasonal and pandemic influenza can pose a significant threat to human health, leading to hospitalization, disability and death. The most effective way to prevent and control influenza and its complications is to get an annual influenza vaccine. Safe and effective vaccines have been available and used for more than 60 years. (Osterholm et al. 2012) Because the types and strains of influenza virus usually vary from one season to the next, a new vaccine targeting the strains that are expected to circulate during the next season has to be developed and used every year.

The vaccine effectiveness measures the direct protection in a vaccinated individual and the effect of vaccine at the population level (the so-called herd immunity) in the post-licensure phase. (Carrillo-Santistevé et al. 2012) The influenza VE reflects proportionate reduction in the frequency of the influenza illness in those receiving the influenza vaccine compared to individuals who did not receive the influenza vaccine. It is calculated by comparing attack rates in the vaccinated and unvaccinated through

the relative risk (RR) for randomized placebo-controlled clinical trials (RCT) or is estimated by the odds ratios (OR) in case-control studies. (De Serres et al. 2013)

Assessing the influenza vaccine effectiveness (VE) every season is very important in that: (1) it helps to understand the relationship between antigen match or mismatch and VE in order to improve the vaccines developed for future seasons; (2) it enables virologists and public health scientists to evaluate the ongoing impact of vaccination efforts in the setting of antigenic drift and periodic vaccine reformulation; (3) when a new vaccine is introduced it is important to estimate its effectiveness as early as possible; (4) it allows the evaluation of vaccination programs and strategies in terms of individual and population-wide benefits to help inform public health officials of the actual impact of the vaccination program in any given season; (5) it allows the identification of subgroups that should be targeted in the development of future vaccines to increase overall effectiveness; and (6) it allows public health officials to inform healthcare providers and the general public regarding what the benefits of influenza vaccination are in a particular season.

The first population-scale use of the influenza vaccine date back to 1945 in U.S. military personnel. (Meiklejon 1994, Osterholm et al. 2012) In 1960, the public health recommendation of annual influenza vaccination was made for people aged 65 years or older and other high-risk groups, such as individuals with chronic debilitating disease and pregnant women. (Burney 1960) The Advisory Committee on Immunization Practices (ACIP) reaffirmed this recommendation in 1964. (Long 1964) Since 2006 the ACIP has substantially expanded the target population recommended for annual influenza vaccination in the United States, first by including young children aged 6 to 59 months, then all children aged 6 months to 18 years in 2008, and finally all individ-

uals aged 6 months or older. (Ferdinands and Shay 2012, Fiore et al. 2010) For many years, randomized clinical trials (RCTs) that measure laboratory-confirmed influenza virus infection as the outcome were considered to provide the most accurate estimates of influenza VE, although these trials had to exclude individuals for whom vaccination is recommended. The recent universal influenza vaccination recommendation made it unethical and impossible to conduct such trials. For this reason, estimates of influenza VE in the U.S. have increasingly relied on observational studies, particularly case-control studies for reasons of statistical power, logistics and cost. (Ferdinands and Shay 2012) However, case-control studies innately are subject to biases, which can lead to errors in estimates of VE.

In the second project, we propose a novel Bayesian model to estimate the effectiveness of the influenza vaccine using data collected from the newly developed test-negative case-control design (TND). (Orenstein et al. 2007) We model the joint probability of each person’s vaccination status and the infection of influenza status conditional on being sampled into the TND study. We specify subjective priors that can be elicited from published literature. We develop an efficient posterior computation algorithm based on the ARMS technique. We demonstrate the superiority of our model compared with the existing method on the VE estimates via simulation studies and an analysis of real data.

In the third project, we develop an improved probability model for evaluating bias and precision of the VE estimates from the traditional case-control design (CCD) and the TND. The proposed model describes the data generation process in real life composed of five steps: latent health status, vaccination, acute respiratory illness (ARI) and influenza infection, seeking medical care for ARI and testing for influenza infec-

tion. Each step is represented by conditional probability parameters. By including a parameter for the health care seeking behavior, this model facilitates the evaluation of an important bias resulted from the unobserved care-seeking behavior. We present and compare the numerical results of the bias and the standard error of the vaccine effectiveness from the traditional case-control study and the test negative case-control study.



## Chapter 2

# A Bayesian Hierarchical Model with Novel Prior Specifications to estimate HIV Testing Rate

### 2.1 Background

Human immunodeficiency virus (HIV) infection is a severe infectious disease actively spreading globally. The diagnosis of HIV infection or HIV testing is one of the most important tools for HIV prevention and treatment. It can foster early detection of HIV infection and is the first essential step for entry to clinical care to reduce morbidity and mortality. In addition, persons aware of their infections and on treatment are less likely to transmit the virus to others. Since HIV testing first became available in 1985, the importance of HIV testing was soon recognized, emphasized and

promoted (Branson et al. 2006). However, a significant proportion of individuals infected with HIV still remain undiagnosed. As of December 2010, more than 1.1 million people were living with HIV infection in the United States and about 1 in 6 were unaware of their infections (CDC 2012). The occurrence of new HIV diagnoses among a population infected with HIV, i.e., the HIV testing rate, is an important epidemiological parameter for public health. Accurately and timely estimating the HIV testing rates is crucial for public health in that 1) it facilitates monitoring changes in HIV testing behaviors over time; 2) it can be used to assess the effectiveness of the public health prevention and intervention programs and provide guidance for maintaining and expanding future public health responses; and 3) it provides additional information on the average length of time between HIV infection and the first positive HIV test, informing how soon an HIV infection gets diagnosed.

To be more specific, at a particular time point the HIV testing rate is defined as the probability that an HIV infected person seeks a test and gets diagnosed with HIV, not acquired immunodeficiency syndrome (AIDS), which is an advanced stage of HIV infection with CD4 T-lymphocyte count less than 200/ $\mu$ l or opportunistic illnesses, given no previous positive test has been obtained. It is challenging to obtain an accurate estimate of the HIV testing rate in that modeling the probability of HIV testing often involves the complex time lag from HIV infection to HIV diagnosis because the time of HIV infection is rarely observable (Becker et al. 2003, Blaxhult and Svensson 1992). The HIV testing rate was initially introduced by Marschner (1994) and Farewell et al. (1994) independently. In Farewell's work, the HIV testing rate was assumed to have two constants for the pre- and post-1984 periods, respectively. In Marschner's work, the HIV testing rate was characterized by a three-parameter

Weibull distribution. Following their work, various parametric models have been proposed to characterize the HIV testing process using the frequentist approach. Among them, one work models the HIV testing process using a two-parameter exponential distribution incorporating dependence between the time since infection and the calendar time (Bellocco and Marschner 2000). Other methods include the heterogeneous mixed exponential model, leading to a Pareto distribution with a decreasing hazard function for the duration between HIV infection and HIV diagnosis (Wand et al. 2009; 2010, Yan et al. 2011), and the additive hazards model, which partitions the HIV testing process to the constant routine HIV testing and the symptoms-driven testing characterized by an exponential distribution (Becker et al. 2003, Chau et al. 2003, Cui and Becker 2000). Another framework of modeling the HIV disease and diagnosis process is through a multi-state formulation, describing the progression through various disease stages from infection to AIDS. Often various HIV disease stages are characterized by laboratory markers such as the CD4 T-cell count (Aalen et al. 1997). In recent years, Sweeting and Birrell proposed to use a Bayesian formulation of a similar multi-state model for estimation of HIV incidence, in which the natural disease progression probabilities and HIV testing probabilities collectively define the transition probabilities from one state to another. The HIV testing probabilities are estimated from external HIV diagnoses data (Birrell et al. 2013, Sweeting et al. 2005). In addition, observed longitudinal CD4 data are used to model the date of HIV infection using joint linear mixed models (Taffe and May 2008). However, the inclusion of additional information in the model increases the complexity of the model. This potentially introduces bias when the observed laboratory data are not representative of all the new diagnoses.

Although the HIV testing rate has been explored and modeled in various frameworks, it does not serve as a primary parameter of interest. It is an important parameter to the so-called back-calculation or back-projection model, which mainly focuses on reconstructing the past pattern of HIV infections in many countries using the annual numbers of diagnosed HIV and AIDS cases (Chau et al. 2003, Hall et al. 2008, Mallitt et al. 2012, Punyacharoensin and Viwatwongkasem 2009, Wand et al. 2010). In those models, HIV testing rates are assumed to have certain parametric forms or to be a constant for a few years. They are not flexible and do not adequately characterize changes over time. To the best of our knowledge, there is no existing statistical framework to systematically model and estimate the HIV testing rates over time. To fill this gap, in this paper, we particularly focus on the Bayesian modeling of annual HIV testing rates over the calendar years using the annual observed numbers of persons diagnosed with HIV, with or without AIDS at HIV diagnosis. Numerous countries have established national HIV registries or surveillance systems to collect the numbers of HIV and AIDS diagnoses. For example, our motivating dataset comes from the national HIV surveillance system in the United States which started collecting data on the number of persons diagnosed with AIDS since 1981 and the number of HIV diagnoses since 1994.

In this work, we assume that the HIV testing rate is only dependent on the calendar year and these annual probabilities can be considered as the discrete-time analogue of the HIV testing intensity. To introduce the smoothness on the HIV testing rates and the HIV incidence rates over years, we resort to a Bayesian shrinkage approach by developing a new class of priors. In the past decades, the Bayesian shrinkage methods using various priors, such as Laplace priors, have been successful in many

applications (Bae and Mallick 2004, Figueiredo and Nowak 2003, Genkin et al. 2007, Kyung et al. 2010, Park and Casella 2008). These methods are mainly developed for variable selection under a regression framework. In particular, the fused lasso priors (Kyung et al. 2010, Tibshirani et al. 2005), i.e., extended Laplace priors, are used to impose smoothness between the model parameters. In a similar fashion, for our problem, we develop a new class of priors that can smooth the annual HIV testing rates and the HIV incidence rates.

Compared with existing models involving the HIV testing rate, our proposed approach has the following remarkable features: 1) our model is among the very first to propose a structured and systematic model to estimate the HIV testing rate in the United States; 2) unlike other back-calculation models (Hall et al. 2008), our model does not impose the constraints that HIV testing rates remain constant within certain years, neither does it make assumptions on the parametric form of the HIV testing rates; 3) our model has an ability to characterize the temporal dependence and smoothness between the annual HIV testing rates, which substantially improve the model fitting and parameter inference accuracy; and 4) our model is widely applicable in that it only needs the annual numbers of HIV and AIDS diagnoses as data in contrast to other approaches, such as the multi-state formulation which requires good representative laboratory data such as the CD4 data.

This chapter is organized as follows: in Section 2.2, we present the detailed formations of our Bayesian hierarchical model for estimating annual HIV testing rates, where two new Laplace-type prior models are introduced and the corresponding properties are discussed in Section 2.2.3, and the posterior computation strategy are developed in Section 2.2.4. We demonstrate the superiority of our proposed methods

via simulation studies in Section 2.3, and illustrate our methods via analysis of the United States national HIV surveillance data in Section 2.4. We conclude our paper with a discussion of future work in Section 2.5.

## 2.2 The Model

### 2.2.1 Data and Notation

First we outline the problem and the observed data. We use January 1, 1977 as the time origin in this analysis because the earliest diagnosis dates of AIDS cases in U.S. were in 1977. The time unit used in this analysis is calendar year.

We consider HIV infection from year 1 to year  $T$  in this paper. When an individual becomes infected with HIV in year  $i$ , he or she could be (1) diagnosed with HIV but not AIDS (referred to as HIV not AIDS) in year  $t$ , where  $i \leq t \leq T$ ; or (2) diagnosed with AIDS in year  $t$ , where  $i \leq t \leq T$ ; or (3) remain undiagnosed as of the most recent year  $T$ . In other words, one could be diagnosed with HIV, with or without AIDS anytime after being infected with HIV, or remain undiagnosed as of the most recent year. The number of HIV and AIDS diagnoses in a calendar year includes persons infected with HIV any time up to and including that year. Note that there does not exist a diagnosis test to determine when the individual became infected.

Denote by  $A_{it}$  the number of persons infected in year  $i$  and diagnosed with AIDS in year  $t$ , by  $H_{it}$  the number of persons infected in year  $i$  and diagnosed with HIV not AIDS in year  $t$  and by  $U_{iT}$  as the number of persons infected in year  $i$  but remain undiagnosed at the end of year  $T$ . Let  $N_i$  be the total number of new HIV infections

in year  $i$ . We have

$$N_i = \sum_{t=i}^T (A_{it} + H_{it}) + U_{iT}, \quad (2.1)$$

We observe the total number of cases diagnosed with AIDS and the total number of cases diagnosed with HIV not AIDS in year  $t$ , denoted by  $A_t$  and  $H_t$  respectively. We have

$$A_t = \sum_{i=1}^t A_{it} \quad \text{and} \quad H_t = \sum_{i=1}^t H_{it}. \quad (2.2)$$

Table 2.1 illustrates the relationship between  $N_i$ ,  $A_{it}$ ,  $H_{it}$ ,  $U_{iT}$ ,  $A_t$  and  $H_t$ , where columns characterize the number of persons diagnosed with AIDS and HIV not AIDS, and those undiagnosed, and rows represent the number of infections at different years.

Table 2.1: Illustration of the relationship between observed data and latent quantities in the model. The column totals represent the number of persons diagnosed with AIDS and HIV not AIDS in each year and those undiagnosed by the end of the most recent year. Row totals represent the number of new infections in each year. In each cell,  $A_{it}$  represents the number of persons infected in year  $i$  and diagnosed with AIDS in year  $t$  and  $H_{it}$  represents the number of persons infected in year  $i$  and diagnosed with HIV not AIDS in year  $t$ . Only column totals  $\{A_t\}_{t=1}^T$  and  $\{H_t\}_{t=1}^T$  are observed, and all other quantities are latent.

Year of infection (i)	Year of AIDS or HIV diagnosis (t)							Undiagnosed by T	Incidence	
	1		2		...		T			
1	$A_{11}$	$H_{11}$	$A_{12}$	$H_{12}$	...	...	$A_{1T}$	$H_{1T}$	$U_{1T}$	$N_1$
2			$A_{22}$	$H_{22}$	...	...	$A_{2T}$	$H_{2T}$	$U_{2T}$	$N_2$
...					...	...	...	...	...	...
T							$A_{TT}$	$H_{TT}$	$U_{TT}$	$N_T$
Observed	$A_1$	$H_1$	$A_2$	$H_2$	...	...	$A_T$	$H_T$		

### 2.2.2 A Hierarchical Model

The primary interest of this study is to estimate the annual HIV testing rates and re-construct the trend in the HIV testing processes over the years. We propose a Bayesian hierarchical model with two levels of hierarchy. To begin with the top level, we model the latent total number of HIV infections  $N_i$ , for  $i = 1, \dots, T$ , which are assumed to be mutually independent and follow a Poisson distribution with intensity  $\lambda_i$ , i.e.

$$[N_i | \lambda_i] \sim \text{Poisson}(\lambda_i), \quad (2.3)$$

where  $\text{Poisson}(\mu)$  denotes a Poisson distribution with mean  $\mu$ .

At level 2, given  $N_i$ , we model the annual numbers of AIDS and HIV diagnosed cases and all the undiagnosed cases stratified by the HIV infections at different years. Write  $\mathbf{N}_i^{AHU} = (A_{ii}, A_{i,i+1}, \dots, A_{iT}, H_{ii}, H_{i,i+1}, \dots, H_{iT}, U_{iT})$  for  $i = 1, \dots, T$ . It represents a collection of the numbers of persons diagnosed with AIDS and HIV not AIDS and undiagnosed persons, infected with HIV in year  $i$  and diagnosed in different years. Let  $\text{Multinomial}(\mathbf{p}, n)$  represent a multinomial distribution with event probability  $\mathbf{p}$  and number of trials  $n$ , for  $i = 1, \dots, T$ , then we assume

$$[\mathbf{N}_i^{AHU} | \mathbf{q}_i^{AHU}, N_i] \sim \text{Multinomial}(\mathbf{q}_i^{AHU}, N_i), \quad (2.4)$$

where  $\mathbf{q}_i^{AHU} = (q_{ii}^A, q_{i,i+1}^A, \dots, q_{iT}^A, q_{ii}^H, q_{i,i+1}^H, \dots, q_{iT}^H, q_{iT}^U)$ . Given that a person is infected with HIV in year  $i$ ,  $q_{it}^A$  and  $q_{it}^H$  represent the probability of being diagnosed with AIDS and HIV not AIDS in year  $t$  respectively, for  $1 \leq i \leq t \leq T$ , and



$q_{iT}^U$  is the probability of remaining undiagnosed at the end of year  $T$ . Note that  $\sum_{t=i}^T (q_{it}^A + q_{it}^H) + q_{iT}^U = 1$  and  $q_{it}^A, q_{it}^H, q_{iT}^U \geq 0, \forall i$ .

To characterize  $\mathbf{q}_i^{AHU}$  we introduce two types of conditional probabilities: the annual AIDS diagnosis rate denoted by  $p_{t-i}^A$  and the annual HIV testing rate denoted by  $p_t^H$ , which is the primary interest of this study. The annual AIDS diagnosis rate  $p_{t-i}^A$  is the probability that a person infected with HIV in year  $i$  gets diagnosed with AIDS in year  $t$ ,  $t \geq i$  given no previous positive tests have been obtained. We assume that all persons newly diagnosed with AIDS did not receive HIV treatment before diagnosis. The treatment-free AIDS diagnosis rate  $p_{t-i}^A$  can be generated from the known AIDS incubation period, which has been studied and modeled by a Gamma distribution with the shape parameter of 2 and the scale parameter of 4 (Longini et al. 1989; 1991). The AIDS incubation period is only determined by the time interval from HIV infection to AIDS diagnosis, i.e., the value of  $(t - i)$ . This means that without treatment, the AIDS diagnosis rate  $p_{t-i}^A$  only depends on how long a person has been infected with HIV.

The annual HIV testing rate  $p_t^H$  is the probability that an HIV positive person seeks HIV test and gets diagnosed with HIV not AIDS in year  $t$  given no previous positive tests have been obtained. For the HIV testing rate, we assume that persons diagnosed with HIV in the same year have the same HIV testing rate, regardless of when they became infected. That is,  $p_t^H$  is only dependent on the diagnosis year  $t$  and is independent of infection time  $i$ . In the first year of HIV infection, assuming HIV infection happens uniformly in the calendar year, HIV testing can only happen after HIV infection and before the calendar year-end. Therefore, the probability of being diagnosed with HIV not AIDS in the year of HIV infection is proportional to

the time interval from HIV infection and the calendar year-end. On average, the HIV testing rate in the year of HIV infection is half of the annual HIV testing rate. Thus, for  $1 \leq i \leq t \leq T$ , we can represent  $q_{it}^H$  using  $p_{t-i}^A$  and  $p_t^H$  which is given by

$$q_{it}^H = \begin{cases} \frac{1}{2}p_i^H \times (1 - p_0^A) & t = i \\ p_t^H(1 - \frac{1}{2}p_i^H) \prod_{k=i+1}^{t-1}(1 - p_k^H) \times \prod_{k=i}^t(1 - p_{k-i}^A) & t \geq i + 1 \end{cases} \quad (2.5)$$

where we define  $\prod_{k=i+1}^i(1 - p_k^H) = 1$ . The term  $\frac{1}{2}p_i^H$  represents the conditional probability that a person gets HIV infected and tested in year  $i$  given no AIDS diagnosis in the same year. The term  $(1 - p_0^A)$  represents the probability that a person is not diagnosed with AIDS in the same year of HIV infection. The term  $p_t^H(1 - \frac{1}{2}p_i^H) \prod_{k=i+1}^{t-1}(1 - p_k^H)$  represents the conditional probability that a person gets infected with HIV in year  $i$  but is not tested until year  $t$  given no AIDS diagnosis between year  $i$  and year  $t$ . The term  $\prod_{k=i}^t(1 - p_{k-i}^A)$  represents the probability that a person is not diagnosed with AIDS from year  $i$  to year  $t$ . Similarly, we can represent  $q_{it}^A$  as

$$q_{it}^A = \begin{cases} p_0^A & t = i \\ p_{t-i}^A \prod_{k=i}^{t-1}(1 - p_{k-i}^A) \times (1 - \frac{1}{2}p_i^H) \prod_{k=i+1}^{t-1}(1 - p_k^H) & t \geq i + 1. \end{cases} \quad (2.6)$$

This further implies that  $q_{iT}^U$  can be written as a function of  $p_{k-i}^A$  and  $p_i^H$  using  $q_{iT}^U = 1 - \sum_{k=i}^T(q_{ik}^A + q_{ik}^H)$ .

In models (2.1)–(2.6), the only observed data are  $\mathbf{A} = \{A_t\}_{t=1}^T$  and  $\mathbf{H} = \{H_t\}_{t=1}^T$  while all other quantities are latent. The estimated values of parameter  $\mathbf{p}^A = \{p_{t-i}^A\}_{t \geq i}$

can be obtained from published literature (Longini et al. 1989; 1991). Our primary interest is in making inference on the HIV testing rates  $\mathbf{p}^H = \{p_t^H\}_{t=1}^T$ .

Another quantity of great interest in the public health community related to the HIV testing rates  $\mathbf{p}^H$  is the expected time-since-infection. Let  $\xi_t$  denote the time from infection to HIV or AIDS diagnosis for individuals diagnosed at a particular time  $t$ . According to the definition of  $q_{it}^A$  and  $q_{it}^H$ , we can represent the probability mass function for  $\xi_t$  as

$$\Pr(\xi_t = t - i \mid q_{it}^A, q_{it}^H) \propto (q_{it}^A + q_{it}^H)\lambda_i, \text{ for } i = 1, \dots, t.$$

This implies that the expected time-since-infection, denoted  $\eta_t$ , can be represented in terms of  $\{\mathbf{q}_i^{AHU}\}_{i=1}^t$  which are functions of  $\mathbf{p}^H$ , i.e.,

$$\eta_t = \mathbb{E}(\xi_t \mid \mathbf{p}^H) = \frac{\sum_{i=1}^t (t - i)(q_{it}^A + q_{it}^H)\lambda_i}{\sum_{i=1}^t (q_{it}^A + q_{it}^H)\lambda_i}, \text{ for } t = 1, \dots, T.$$

This informs on average how long individuals diagnosed with HIV or AIDS in a certain year have been infected. Since  $\eta_t$  is a deterministic function of  $\mathbf{p}^H$ , the posterior inference on  $\eta_t$  can be obtained straightforwardly through the posterior inference on  $\mathbf{p}^H$ , for which we introduce a class of new priors in Section 2.2.3.

### 2.2.3 Prior Specifications

In this section, we discuss the prior specifications for our proposed hierarchical model. We start from the most important parameter in our analysis, i.e., the HIV testing rates  $\mathbf{p}^H$ . It is generally believed that the HIV testing process usually does not change

dramatically between two succeeding years. Thus, it is meaningful to assume that the HIV testing rate in the current year is associated with the rate in the previous year. To characterize such an association, we propose a new family of probability distributions defined on  $[0, 1]$ , based on which we construct a prior for  $\mathbf{p}^H$  taking into account the temporal dependence between HIV testing rates over the years. Specifically, we introduce the following definition of the new distribution.

**Definition 2.2.1** (Beta-Laplace Distribution). *Let  $\mu \in [0, 1]$ ,  $a, b, c \in \mathbb{R}^+$ . For  $x \in [0, 1]$ , let*

$$\pi(x; \mu, a, b, c) = \frac{1}{L(\mu, a, b, c)} x^{a-1} (1-x)^{b-1} \exp(-c|x-\mu|), \quad (2.7)$$

where  $L(\mu, a, b, c) = \int_0^1 x^{a-1} (1-x)^{b-1} \exp(-c|x-\mu|) dx$ . We refer to  $\pi(x; \mu, a, b, c)$  as the probability density function of a beta-Laplace distribution with the location parameter  $\mu$ , the rate parameter  $c$  and two shape parameters  $a$  and  $b$ . A random variable  $X$  following this distribution is denoted as  $X \sim \text{Beta-Laplace}(\mu, a, b, c)$ .

**Remark:** 1) when  $c = 0$ , the beta-Laplace distribution reduces to a beta distribution with shape parameters  $a$  and  $b$ . 2) When  $a = b = 1$ , the beta-Laplace distribution becomes a truncated Laplace distribution with location parameter  $\mu \in [0, 1]$  and rate parameter  $c$ . Also, the properties of the Beta-Laplace distribution are summarized in the following proposition.

**Proposition 2.2.1.** *Let  $X \sim \text{Beta-Laplace}(\mu, a, b, c)$ . Then*

1.  $1 - X \sim \text{Beta-Laplace}(1 - \mu, b, a, c)$ .

2. The  $n$ th moment of  $X$  is given by

$$\mathbb{E}[X^n] = \frac{L(\mu, a + n, b, c)}{L(\mu, a, b, c)}, \quad \text{for } n = 1, 2, \dots, \quad (2.8)$$

where  $L(\mu, a, b, c)$  is defined in Definition 2.2.1.

3. Given  $a, b \in \mathbb{R}^+$  and  $\mu \in [0, 1]$ , we have

$$\lim_{c \uparrow +\infty} \mathbb{E}[X] = \mu \quad \text{and} \quad \lim_{c \downarrow 0} \mathbb{E}[X] = \frac{a}{a + b}. \quad (2.9)$$

See the proof in Appendix 1.1. Proposition 2.2.1(3) implies that the parameters  $c$  and  $\mu$  control how the mean of Beta-Laplace( $\mu, a, b, c$ ) deviates from the mean of Beta( $a, b$ ). It approximately equals  $\mu$  when  $c$  is sufficiently large, and it gets close to the mean of Beta( $a, b$ ) for a small  $c$ .

Now we assign the priors for the annual HIV testing rates using beta-Laplace distributions. Specifically, we have

$$p_1^H \sim \text{Beta}(a^H, b^H) \quad \text{and} \quad [p_{t+1}^H \mid p_t^H] \sim \text{Beta-Laplace}(p_t^H, a^H, b^H, c_t^H) \quad (2.10)$$

for  $t = 1, \dots, T - 1$ . This implies a priori, the distribution of  $p_{t+1}^H$  borrows information from  $p_t^H$ . According to the property of the beta-Laplace distribution, the parameter  $c_t^H$  controls the difference between  $\mathbb{E}(p_{t+1}^H \mid p_t^H)$  and  $p_t^H$ , which reflects the average change of the HIV testing rate from year  $t$  to  $t + 1$ . The larger  $c_t^H$  is, the closer  $\mathbb{E}(p_{t+1}^H \mid p_t^H)$  gets to  $p_t^H$ . It is generally believed that the variation of the HIV testing rates over the years usually becomes smaller when the HIV testing rate reaches a

certain level. A reasonable choice for  $c_t^H$  is to assume that it is proportional to  $p_t^H$ , i.e.  $c_t^H = c^H p_t^H$ , where  $c^H > 0$  is an important parameter that controls the overall smoothness of the HIV testing rates over the years. For hyperparameters, we choose  $a^H = b^H = 0.5$ . The choice of the  $c^H$  can be determined via Bayes factors (Aitkin 1991, Kass and Raftery 1995). We discuss this in Section 2.2.4.

Similar to the HIV testing rates, the annual HIV incidence which is reflected by parameters  $\boldsymbol{\lambda} = \{\lambda_i\}_{i=1}^T$  in (2.3) is also dependent over the years. To characterize such temporal dependence, we introduce another new family of probability distributions to specify the priors for  $\boldsymbol{\lambda}$ .

**Definition 2.2.2** (Gamma-Laplace Distribution). *Let  $a, b, c, \mu \in \mathbb{R}^+$ . For  $x \in \mathbb{R}^+$ , let*

$$\pi(x; \mu, a, b, c) = \frac{1}{K(\mu, a, b, c)} x^{a-1} \exp(-bx - c|x - \mu|),$$

where  $K(\mu, a, b, c) = \int_0^\infty x^{a-1} \exp(-bx - c|x - \mu|) dx$ . We refer to  $\pi(x; \mu, a, b, c)$  as the probability density function of a gamma-Laplace distribution with the location parameter  $\mu$ , the shape parameter  $a$  and two scale parameters  $b$  and  $c$ . A random variable  $X$  following this distribution is denoted as  $X \sim \text{Gamma-Laplace}(\mu, a, b, c)$ .

**Remark:** 1) when  $c = 0$ , the gamma-Laplace distribution reduces to a gamma distribution with shape  $a$  and rate  $b$ . 2) when  $a = 1$  and  $b = 0$ , the gamma-Laplace distribution becomes a truncated Laplace distribution with location  $\mu$  and rate  $c$ . The properties of the gamma-Laplace distribution are summarized in the following proposition.

**Proposition 2.2.2.** *Let  $X \sim \text{Gamma-Laplace}(\mu, a, b, c)$ . Then*

1. For  $\tau \in \mathbb{R}^+$ ,  $X/\tau \sim \text{Gamma-Laplace}(\mu/\tau, a, \tau b, \tau c)$ .

2. The  $n$ th moment of  $X$  is given by

$$\mathbb{E}[X^n] = \frac{K(\mu, a + n, b, c)}{K(\mu, a, b, c)}, \quad \text{for } n = 1, 2, \dots, \quad (2.11)$$

where  $K(\mu, a, b, c)$  is defined in Definition 2.2.2.

3. Given  $a, b, \mu \in \mathbb{R}^+$ , we have

$$\lim_{c \uparrow +\infty} \mathbb{E}[X] = \mu \quad \text{and} \quad \lim_{c \downarrow 0} \mathbb{E}[X] = \frac{a}{b}. \quad (2.12)$$

The proof is straightforward and similar to proposition 2.1. Proposition (2.2.2)(3) implies that the parameters  $c$  and  $\mu$  reflect how different the mean of  $\text{Gamma-Laplace}(\mu, a, b, c)$  is from the mean of  $\text{Gamma}(a, b)$ . It gets close to  $\mu$  when  $c$  is sufficiently large and it approaches to the mean of  $\text{Gamma}(a, b)$  when  $c$  is very small. Based on this property, we assign the following priors for  $\boldsymbol{\lambda}$ :

$$\lambda_1 \sim \text{Gamma}(a^\lambda, b^\lambda) \quad \text{and} \quad [\lambda_{t+1} \mid \lambda_t] \sim \text{Gamma-Laplace}(\lambda_t, a^\lambda, b^\lambda, c^\lambda), \quad (2.13)$$

for  $t = 1, \dots, T-1$ . As a priori,  $\lambda_{t+1}$  is assumed to follow a distribution characterized by  $\lambda_t$ , where the difference between  $\mathbb{E}[\lambda_{t+1} \mid \lambda_t]$  and  $\lambda_t$  is controlled by  $c^\lambda > 0$ . This prior specification implies that the number of HIV infections in the current year can be dependent on that of the previous year to a certain extent. This further assists better posterior inference on the HIV testing rate. We demonstrate this in simulation studies. For hyperparameters in (2.13), we choose  $a^\lambda = 1, b^\lambda = 0.00002$  so that the

mean of Gamma(a, b) is the average of the annual number of total diagnoses. Similar to the  $c^H$ , we discuss the choice of  $c^\lambda$  in Section 2.2.4.

## 2.2.4 Posterior Inference

To simplify the posterior computation, we consider an equivalent model representation by integrating out the latent quantities  $\mathbf{N}_i^{AHU}$  and  $N_i$  in models (2.3) and (2.4).

Since  $N_i$  follows a Poisson distribution with mean  $\lambda_i$ , it is straightforward to show that  $A_{it}$  and  $H_{it}$  follow Poisson distributions with means  $q_{it}^A \lambda_i$  and  $q_{it}^H \lambda_i$ , respectively, i.e.

$$[A_{it} \mid \lambda_i, q_{it}^A] \sim \text{Poisson}(q_{it}^A \lambda_i) \quad \text{and} \quad [H_{it} \mid \lambda_i, q_{it}^H] \sim \text{Poisson}(q_{it}^H \lambda_i). \quad (2.14)$$

See the proof in Appendix 1.2. Note that  $A_{it}$  and  $H_{it}$  are mutually independent given  $\lambda$ ,  $q_{it}^A$  and  $q_{it}^H$ . Assuming the annual numbers of HIV infections are independent, in each calendar year  $t$ , the observed total numbers of cases diagnosed with AIDS ( $A_t = \sum_{i=1}^t A_{it}$ ) and HIV ( $H_t = \sum_{i=1}^t H_{it}$ ) follow Poisson distributions with means  $\sum_{i=1}^t q_{it}^A \lambda_i$  and  $\sum_{i=1}^t q_{it}^H \lambda_i$ , respectively, i.e.

$$[A_t \mid \mathbf{p}^H, \lambda] \sim \text{Poisson} \left( \sum_{i=1}^t q_{it}^A \lambda_i \right) \quad \text{and} \quad [H_t \mid \mathbf{p}^H, \lambda] \sim \text{Poisson} \left( \sum_{i=1}^t q_{it}^H \lambda_i \right), \quad (2.15)$$

where both  $q_{it}^A$  and  $q_{it}^H$  are functions of  $\mathbf{p}^H$  according to models (2.6) and (2.5). The joint posterior distribution of  $\mathbf{p}^H$  and  $\lambda$  given data  $\mathbf{A}$  and  $\mathbf{H}$  and hyperparameters



$(c^H, c^\lambda)$  is given by

$$\begin{aligned} & \pi(\mathbf{p}^H, \boldsymbol{\lambda} \mid \mathbf{A}, \mathbf{H}, c^H, c^\lambda) \\ & \propto \prod_{t=1}^T \pi(A_t, H_t \mid \mathbf{p}^H, \boldsymbol{\lambda}) \times \pi(p_1^H) \prod_{t=1}^{T-1} \pi(p_{t+1}^H \mid p_t^H, c^H) \times \pi(\lambda_1) \prod_{t=1}^{T-1} \pi(\lambda_{t+1} \mid \lambda_t, c^\lambda). \end{aligned}$$

The posterior distribution of parameters given the data is complicated and has no closed form solution. Thus, to sample from this posterior distribution for given  $c^\lambda$  and  $c^H$ , we resort to the adaptive rejection metropolis sampling within Gibbs sampling (Gilks et al. 1995). Details of the full conditional distributions of  $p_t^H$  and  $\lambda_t$  are provided in Appendix 1.3. For the choice of hyperparameters  $(c^H, c^\lambda)$ , we maximize the estimated Bayes factors, where the reference model is the case when  $(c^H, c^\lambda) = (0, 0)$ , on a set of pre-specified values  $\{(c^{H(k)}, c^{\lambda(k)})\}_{k=1}^K$ , i.e. we choose

$$(\widehat{c^H}, \widehat{c^\lambda}) = (c^{H(\hat{k})}, c^{\lambda(\hat{k})}) \quad \text{with} \quad \hat{k} = \arg \max_k \left[ \sum_{s=1}^S \pi^{-1}(\mathbf{A}, \mathbf{H} \mid \mathbf{p}^{H(k,s)}, \boldsymbol{\lambda}^{(k,s)}) \right]^{-1},$$

where  $\{\mathbf{p}^{H(k,s)}, \boldsymbol{\lambda}^{(k,s)}\}_{s=1}^S$  are the simulated samples from the posterior distribution  $\pi(\mathbf{p}^H, \boldsymbol{\lambda} \mid \mathbf{A}, \mathbf{H}, c^{H(k)}, c^{\lambda(k)})$ .

## 2.3 Simulation Study

To demonstrate the performance of the proposed model, we conducted simulation studies to estimate the HIV testing rates and the time-since-infection. We specify the AIDS diagnosis rate ( $\mathbf{p}^A$ ) from the hazard function of Gamma(2, 4), shown in Table 2.2. To simulate the observed numbers of HIV and AIDS diagnoses over the years ( $\mathbf{A}$

and  $\mathbf{H}$ ), we specify the true values for the mean numbers of new HIV infections over the years ( $\boldsymbol{\lambda}$ ) and the HIV testing rates ( $\mathbf{p}^H$ ) (see Table 2.3). In particular, we have two scenarios for testing rates. In Scenario 1, we consider a 34-year period with a gradual increasing trend in the HIV testing rate. In Scenario 2, we consider a 20-year period with an increasing trend followed by a decreasing trend in the HIV testing rate.

Table 2.2: Values of AIDS diagnosis rate generated from the hazard function of a Gamma distribution with shape parameter of 2 and scale parameter of 4.  $t$  is year of AIDS diagnosis and  $i$  is year of infection.

$t - i$	$p_{t-i}^A$	$t - i$	$p_{t-i}^A$	$t - i$	$p_{t-i}^A$	$t - i$	$p_{t-i}^A$	$t - i$	$p_{t-i}^A$
0	0.00934	7	0.14701	14	0.17670	21	0.18941	28	0.19648
1	0.04761	8	0.15346	15	0.17910	22	0.19066	29	0.19724
2	0.07934	9	0.15889	16	0.18126	23	0.19181	30	0.19795
3	0.10124	10	0.16351	17	0.18321	24	0.19288	31	0.19863
4	0.11727	11	0.16749	18	0.18498	25	0.19387	32	0.19926
5	0.12952	12	0.17095	19	0.18659	26	0.19480	33	0.19986
6	0.13919	13	0.17400	20	0.18806	27	0.19567	34	0.20043

Given a set of values of  $\boldsymbol{\lambda}$ ,  $\mathbf{p}^H$  and  $\mathbf{p}^A$ , the HIV infections and diagnoses data are simulated from a process that mimics data generation and collection in real-life: a person becomes infected with HIV in year  $i$ , and he or she gets diagnosed in a later year  $t \geq i$  or remains undiagnosed as of the most recent year  $T$ . At the time of diagnosis, he or she can be diagnosed with HIV and AIDS in the same year (AIDS) or in different years (HIV not AIDS). The diagnosis date and disease status are determined and reported to a national surveillance registry. Annual numbers of HIV and AIDS diagnoses are thus summarized. This process is simulated as follows:

**Step 1:** For each year  $i, i = 1, \dots, T$ , the number of new HIV infections,  $N_i$  is generated through a Poisson distribution based on a mean  $\lambda_i$ .

Table 2.3: Values for the parameters used in the simulation studies.  $\lambda$  is the mean annual number of new HIV infections,  $p^H$  is the annual HIV testing rate,  $A$  is the observed annual number of AIDS diagnoses and  $H$  is the observed annual number of HIV diagnoses. Scenario 1 is for a 34-year period with a gradual increasing trend in  $p^H$  and Scenario 2 is for a 20-year period with an increasing trend followed by a decreasing trend in  $p^H$ .

year	Scenario 1				Scenario 2			
	$\lambda$	$p^H$	$A$	$H$	$\lambda$	$p^H$	$A$	$H$
1	24	0.060	0	1	24	0.060	0	1
2	86	0.060	2	4	86	0.060	2	4
3	86	0.060	6	8	86	0.060	6	8
4	244	0.060	14	17	244	0.060	14	17
5	244	0.096	28	46	244	0.096	28	46
6	862	0.096	49	91	862	0.096	49	91
7	862	0.096	92	156	862	0.096	92	156
8	2521	0.120	161	361	2521	0.120	161	361
9	2521	0.120	285	586	2521	0.120	285	586
10	3337	0.120	439	815	3337	0.120	439	815
11	3337	0.156	618	1356	3337	0.156	618	1356
12	4675	0.156	780	1649	4675	0.156	780	1649
13	4675	0.156	965	1970	4675	0.156	965	1970
14	5305	0.156	1151	2262	5305	0.156	1151	2262
15	5305	0.156	1335	2529	5305	0.156	1335	2529
16	5305	0.156	1500	2728	5305	0.144	1500	2518
17	3429	0.156	1619	2730	3429	0.144	1640	2547
18	3429	0.156	1647	2582	3429	0.144	1687	2431
19	3429	0.180	1621	2840	3429	0.120	1678	1944
20	2000	0.180	1522	2542	2000	0.120	1674	1835
21	2000	0.180	1379	2196				
22	1500	0.180	1227	1895				
23	1500	0.192	1074	1739				
24	2346	0.192	934	1596				
25	2346	0.192	852	1576				
26	2346	0.192	808	1569				
27	2346	0.192	784	1568				
28	2642	0.192	775	1597				
29	2642	0.204	783	1750				
30	2642	0.204	786	1772				
31	2400	0.204	789	1763				
32	2400	0.204	785	1733				
33	2400	0.204	777	1711				
34	2400	0.204	769	1694				

**Step 2:** For each case of infection in year  $i$ ,

1. Simulate the time of HIV infection in the year from  $\text{Uniform}(0,1)$ , assuming an HIV infection happens uniformly throughout the year.
2. Simulate the time interval from HIV infection to AIDS diagnosis (i.e., the AIDS incubation period) using a  $\text{Gamma}(2, 4)$  distribution.
3. Determine the year of diagnosis and categorize the case as either “HIV not AIDS” or “AIDS” at the time of diagnosis.
  - (a) If the AIDS incubation period is smaller than one (i.e., AIDS diagnosis happens in the same year of HIV infection), the case is categorized as ”AIDS” and the year of diagnosis is the year of HIV infection.
  - (b) If the AIDS incubation period is greater than one (i.e., AIDS is diagnosed in years after the year of HIV infection), determine whether the case had an HIV test before the year of AIDS diagnosis based on the HIV testing rates in each year before AIDS diagnosis.
    - If a case has an HIV test before AIDS diagnosis, then it is categorized as ”HIV not AIDS” and the year of diagnosis is the year of HIV test;
    - Else if the year of AIDS diagnosis is earlier than the most recent year, the case is an AIDS case and the year of diagnosis is the year of AIDS diagnosis.
    - Otherwise, the case remains undiagnosed as of the most recent year.

**Step 3:** After looping through each infection and each year, summarize the diagnosed

HIV not AIDS cases ( $\mathbf{H}$ ) and AIDS cases ( $\mathbf{A}$ ) over years.

Given the simulated  $\mathbf{A}$  and  $\mathbf{H}$ , we set the initial values for  $\boldsymbol{\lambda}$  as a half of the annual observed numbers of HIV and AIDS diagnoses and set  $\mathbf{p}^H$  as random values between 0 and 1, respectively. We choose the hyperparameters  $(c^H, c^\lambda)$  by maximizing the Bayes factors as discussed in Section 2.2.4 on a set of pre-specified values  $(0,0)$ ,  $(300,0.001)$ ,  $(500,0.001)$ ,  $(800,0.001)$ ,  $(300,0.002)$ ,  $(500,0.002)$ ,  $(800,0.002)$ ,  $(300,0.004)$ ,  $(500,0.004)$ ,  $(800,0.004)$ . We run the proposed posterior simulation algorithm 2,000 iterations with 200 burn-in for each set of  $(c^H, c^\lambda)$  in both scenarios. We check the convergence of the simulated Markov chains using the Gelman and Rubin diagnostic (Gelman and Rubin 1992) by running five additional Markov chains with different initial values. The potential scale reduction factors (PSRF) of the log-likelihood for scenarios 1 and 2 are respectively 0.99 and 1.06, which are both close to 1, indicating the convergence of the posterior simulations.

The selected values of  $(c^H, c^\lambda)$  are  $(300, 0.002)$  for scenarios 1 and 2. The posterior mean and 95% credible interval of  $p_t^H$  and  $\eta_t$  are shown in Figure 2.1. For scenarios 1 and 2, the estimated posterior means of  $p_t^H$  and  $\eta_t$ , for  $t = 1, \dots, T$ , are quite close to the true values of the HIV testing rates and the time-since-infection and the associated 95% credible intervals cover the true values for all years. The results show that our model can provide accurate estimates under different trends and different periods of time. In Table 2.4, for both scenarios, we compare the model fitting results for different choices of hyperparameters  $(c^H, c^\lambda)$ . For both  $p_t^H$  and  $\eta_t$ , Table 2.4 summarizes the average mean square error (AMSE) over years, the average length of 95% credible intervals (ACI) over years and the estimated Bayes factors (BF). A special case (shown in Figure 2.1) is when  $(c^H, c^\lambda) = (0,0)$ , corresponding to a

regular choice of independent beta priors for  $p_t^H$  and independent gamma priors for  $\lambda_t$ , compared to which our method has a better model fitting (larger BF) and provide more accurate estimates and inference on the HIV testing rate and the time-since-infection (smaller AMSE and ACI).

Table 2.4: Simulation model fitting results for different choices of hyperparameters  $(c^H, c^\lambda)$ .  $\log(\text{BF})$  is the estimated log Bayes factor. AMSE is the average mean square error, ACI is the average length of 95% credible intervals.

$(c^H, c^\lambda)$	Scenario 1					Scenario 2				
	$\log(\text{BF})$	$p^H$		$\eta$		$\log(\text{BF})$	$p^H$		$\eta$	
		AMSE	ACI	AMSE	ACI		AMSE	ACI	AMSE	ACI
(0,0)	0	5.7e-05	0.038	0.008	0.36	0	4.4e-05	0.039	0.008	0.37
(300,0.001)	10.2	2.2e-05	0.029	0.005	0.31	6.5	3.3e-05	0.030	0.006	0.32
(500,0.001)	8.7	3.6e-05	0.024	0.005	0.32	7.5	4.4e-05	0.030	0.006	0.35
(800,0.001)	12.2	3.4e-05	0.023	0.004	0.28	7.9	4.8e-05	0.027	0.004	0.33
(300,0.002)	13.8	2.5e-05	0.025	0.003	0.28	8.5	3.3e-05	0.031	0.007	0.33
(500,0.002)	11.9	2.3e-05	0.024	0.003	0.29	5.8	3.2e-05	0.029	0.008	0.32
(800,0.002)	10.8	4.0e-05	0.022	0.004	0.27	4.8	5.3e-05	0.027	0.006	0.34
(300,0.004)	4.1	2.4e-05	0.026	0.004	0.28	6.9	3.6e-05	0.030	0.006	0.31
(500,0.004)	11.1	3.1e-05	0.025	0.004	0.28	5.2	4.0e-05	0.028	0.005	0.31
(800,0.004)	13.6	2.6e-05	0.023	0.004	0.27	-3.9	4.3e-05	0.027	0.006	0.29

## 2.4 Application

In this section, we apply the proposed Bayesian hierarchical model to the data from national HIV surveillance in the United States.

### 2.4.1 Analysis of the United States HIV surveillance data

Since 1982, all 50 states and the District of Columbia have reported AIDS cases to the Centers for Disease Control and Prevention (CDC) using a standardized case report form. In 1994, the CDC implemented data management for national reporting of HIV integrated with AIDS case reporting, at which time 25 states with confiden-

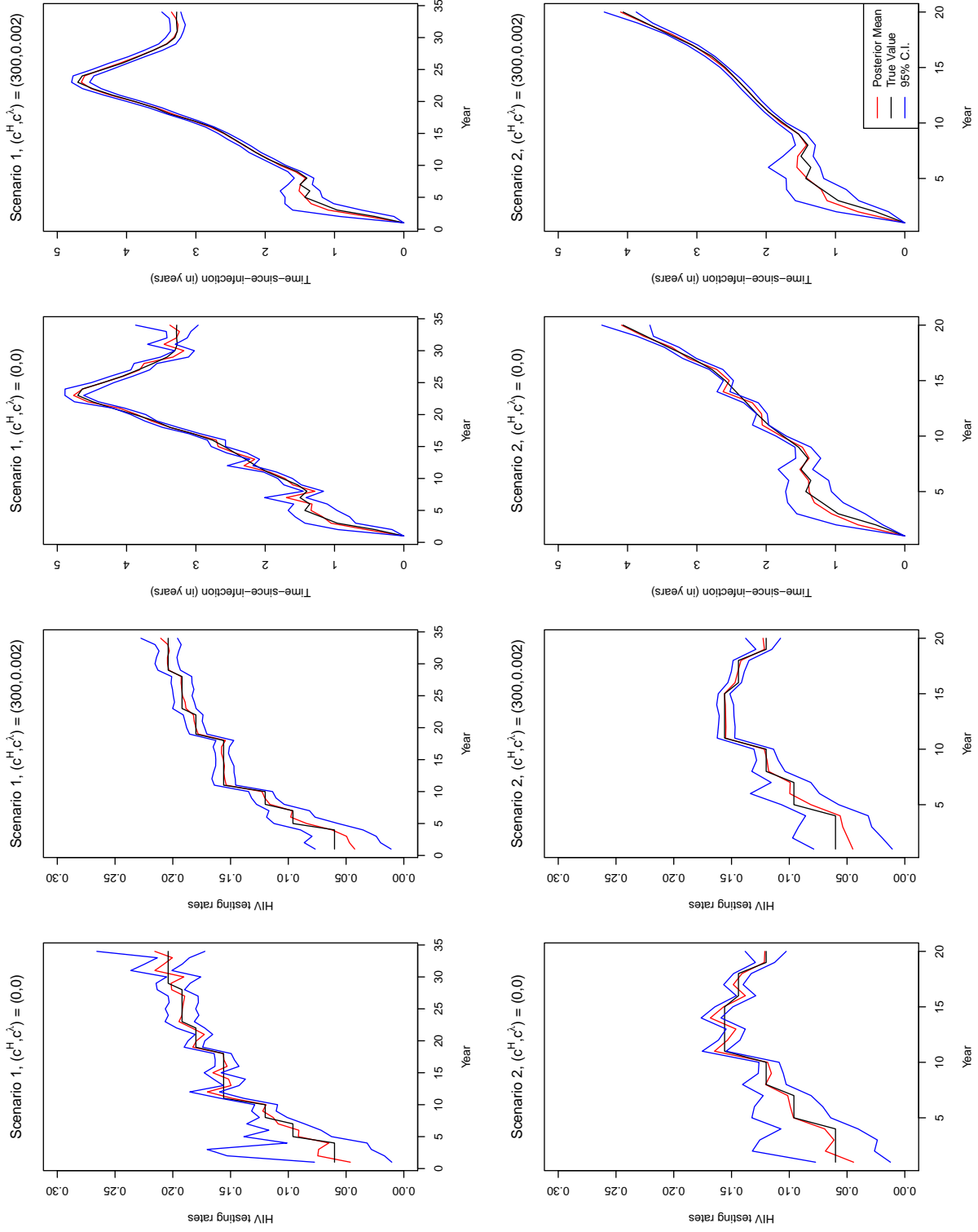


Figure 2.1: Estimated posterior mean and 95% credible intervals for HIV testing rates and time since infection with different choices of  $(c^H, c^\lambda)$  in simulation Scenarios 1 and 2.

tial name-based HIV surveillance started submitting case reports to the CDC. Over time, additional states implemented name-based HIV surveillance, and all states had implemented such surveillance in 2008.

In this study, we use HIV and AIDS data reported to the CDC through June 2012. The data are adjusted for incomplete reporting, reporting delay and misclassification of the diagnosis dates (Green 1998, Song et al. 2005). We estimate the annual numbers of HIV and AIDS diagnosed cases each year from 1977 (the beginning of the HIV epidemic) to 2010.

The initial values for HIV testing rates are randomly assigned to be values between 0 and 1. The annual numbers of new HIV infections are initially assigned a half of the observed number of HIV and AIDS diagnoses. We choose the hyperparameters  $(c^H, c^\lambda)$  as (500, 0.0005), which has the maximal value (3.1) of the Bayes factor on the log scale among a set of pre-specified values (0, 0), (300, 0.0001), (300, 0.0003), (300, 0.0005), (300, 0.0008), (400, 0.0005), (500, 0.0001), (500, 0.0003), (500, 0.0004), (500, 0.0005), (500, 0.0006), (500, 0.0008), (600, 0.0005), (800, 0.0001), (800, 0.0003), (800, 0.0005), (800, 0.0008). We run the posterior simulation algorithm 2,000 iterations with 200 burn-in. The PSRF of the loglikelihood is 1.04 from running five additional chains, indicating the convergence of the posterior simulations.

Figure 2.2 presents the posterior means and 95% credible intervals for the HIV testing rate and the expected time-since-infection (in years) from 1985 when the first HIV test became available in the United States, to 2010. In the first few years after HIV test became available, HIV testing was widely adopted and it continued to increase until 1990. The testing rate went down in early 1990s. This is likely caused by the change in the CDC AIDS definition during that time resulting in



a high proportion of simultaneous AIDS diagnoses, which could indicate low HIV testing rate. Another possible reason could be a sudden increase in the number of new HIV infections in early 1990s resulting in a high number of undiagnosed HIV infections and consequently low HIV testing rate. After that, the HIV testing rate gradually increased and sustained the increasing trend ever since. In the most recent years since 2007, the annual HIV testing rate has been stable around 0.22. As for the expected time interval (in years) since HIV infection to HIV diagnosis among individuals diagnosed in a specific calendar year, there was an increasing trend followed by a decreasing trend. The expected time-since-infection was short among those diagnosed in the early years of the epidemic, which is likely because (1) the early infections were mainly concentrated among men who have sex with men and the targeted HIV testing among this population could result in shorter time interval from infection to diagnoses; and (2) without proper treatment, cases diagnosed in the early years could be fast progressors and the estimated expected time-since-infection might be limited by the short history of the disease. As time went by, with the HIV epidemic spread to a more general population and the treatment improving, the expected time-since-infection gradually increased and reached the peak of 4.2 years among cases diagnosed in 1997. Since 1997, because of the increased HIV testing rate, the estimated expected time-since-infection decreased and was 3.3 years in 2010.

The estimated HIV testing rates and expected time-since-infection reflect the impact of important public health initiatives and recommendations on HIV testing. As shown in the results, the HIV testing rate increased when the first HIV test became available in 1985. In 1987, the United States Public Health Service (USPHS) issued guidelines making HIV counseling and testing a priority as a prevention strategy for

people with high risk behaviors. (CDC et al. 1987) As a result, the HIV testing rate increased in the late 1980s. Though HIV testing went down during 1991 to 1993, the HIV testing rate started increasing since 1993 when CDC updated the recommendations regarding HIV counseling and testing of patients in acute-care hospital settings. (Ward et al. 1993) In 1995, when National HIV Testing Day was observed, the HIV testing rate sustained the upward trend until 2000. Throughout the first decade of 2000, a few important recommendations were published in 2001 (Allen et al. 1999), 2003 (Janssen et al. 2003) and 2006 (Branson et al. 2006) respectively to emphasize routine HIV testing as an important HIV prevention tool. The HIV testing rate increased after each recommendation and it maintained an increasing trend since 2001. As a result of the increased HIV testing, the expected time-since-infection decreased for cases diagnosed since 2001. These results indicate that public health recommendations on HIV testing have a consistently positive impact on people’s HIV testing awareness and testing behavior.

## 2.4.2 Model Assessment

We conduct a posterior predictive model assessment using the  $\chi^2$  discrepancy, which is a summary statistic for the sum of squares of standardized residuals of the data with respect to their expectations under the model (Gelman et al. 1996). For our model, the  $\chi^2$  discrepancy is defined as:

$$\chi^2(\mathbf{H}; \mathbf{p}^H, \boldsymbol{\lambda}) = \sum_{t=1}^{34} \frac{(H_t - E(H_t | \mathbf{p}^H, \boldsymbol{\lambda}))^2}{\text{Var}(H_t | \mathbf{p}^H, \boldsymbol{\lambda})}.$$

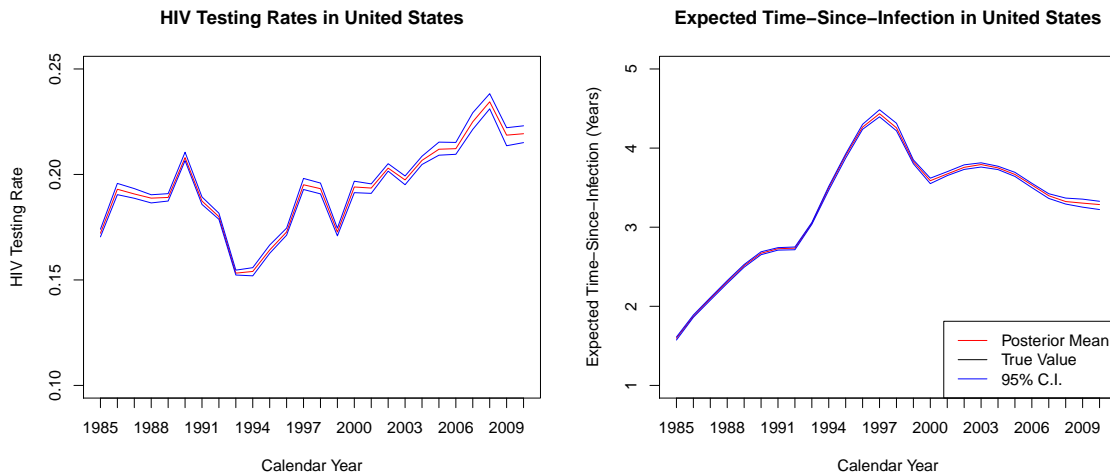


Figure 2.2: Estimated posterior mean and 95% credible intervals for the annual HIV testing rates and expected time-since-infection from 1985 to 2010 in United States

We calculate the posterior predictive  $p$ -value based on  $\chi^2$  as  $p = P(\chi^2(\mathbf{H}^{rep}; \mathbf{p}^H, \boldsymbol{\lambda}) \geq \chi^2(\mathbf{H}; \mathbf{p}^H, \boldsymbol{\lambda}))$ , where  $\mathbf{H}^{rep}$  represents the predictive replication and  $\mathbf{H}$  represents the observed data. The  $p$ -value of the predictive versus realized  $\chi^2$  discrepancies is 0.48. This implies that the model fits the data pretty well.

## 2.5 Discussion

In this paper, we develop a Bayesian hierarchical model to estimate the intensities of HIV testing from 1977 to 2010 using annual numbers of HIV and AIDS diagnosed cases collected through national HIV surveillance in the United States. Our model takes the most general form and makes no parametric assumptions for HIV testing rates. We assume that the HIV testing rate is only dependent on the calendar year and

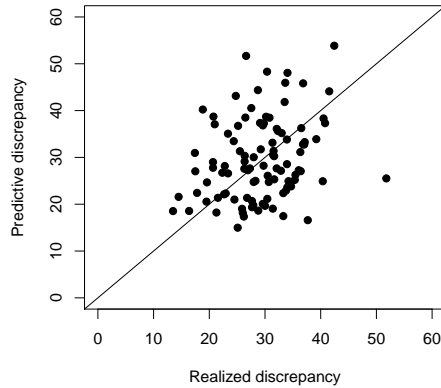


Figure 2.3: Scatterplot of predictive versus realized  $\chi^2$  discrepancies under the joint posterior distribution; the  $p$ -value is estimated by the proportion of points above the 45 degree line.

these annual probabilities can be considered the discrete-time analogue of the HIV testing intensity. We propose beta-Laplace and gamma-Laplace priors to characterize the temporal dependence for annual HIV testing rates and annual HIV incidence rates respectively, which greatly improve the estimation accuracy and the model fitting. The simulation studies show that our model can make much more accurate inference on HIV testing rates and the time-since-infections of different trends for either long or short periods of time compared to a regular choice of priors.

The proposed beta-Laplace and gamma-Laplace priors are general and can have different applications. For example, in spatial statistics, these two prior distributions can impose the smoothness over the spatial parameters, which are alternatives to the use of Gaussian random fields requiring the normality assumption which are not valid in certain cases. In the spatial modeling of disease mapping, the disease infection probabilities over space can be assigned with beta-Laplace priors and the intensity of the disease clusters can be modeled by gamma-Laplace priors. Also, in the analysis

of functional neuro-imaging data, the voxel-wise probabilities of activation can be modeled with beta-Laplace priors and the gamma-Laplace model should be a good choice for the intensity of peak activation locations.

There are several possible future directions that extend our current work. One extension is to include covariates such as sex, race/ethnicity, or transmission category in the model, which can adjust for the confounding factors that might affect the HIV testing rates. Also, in contrast to the current modeling of the whole population in the United States, this extension can provide stratified estimates for sub-populations, e.g., HIV testing rates for demographic groups or different regions in the country, which are of particular interest for public health officials and program evaluation. Another direction is that we can develop alternative strategies to choose the hyper-parameters  $(c^H, c^\lambda)$  which are strongly related to the performance of model fitting and parameter estimations. A fully Bayesian approach can be considered by jointly updating  $(c^H, c^\lambda)$  in the posterior simulation and using the Bayesian model averaging to make the inference. The key steps are to choose appropriate priors for  $(c^H, c^\lambda)$  and to develop an efficient posterior sampling strategy which is worthy of investigation in that this approach would take into account more sources of variation in the model and potentially produce better posterior inference on the model parameters.

## Chapter 3

# A Bayesian Model to Estimate the Influenza Vaccine Effectiveness from a Test Negative Design

### 3.1 Background

A highly contagious viral infection, influenza can pose a significant threat to human health and a great cost to the economy. In the United States, influenza is responsible for more than 226,000 hospitalizations each year. (Thompson et al. 2004) During 1976 to 2007, the estimates of annual influenza-related death ranged from 3000 to 48000. (Thompson et al. 2010; 2003) It was estimated that the annual influenza epidemics could lead to a \$87.1 billion economic burden in the U.S. (Molinari et al. 2007). Annual influenza vaccination is the most effective way to prevent influenza

and its complications. Because the types and strains of influenza virus usually vary from one season to the next, a new vaccine targeting the strains that are expected to circulate during the next season has to be developed and used every year.

Accurate and precise assessment of the influenza vaccine effectiveness (VE) is important for (1) understanding the relationship between antigen match or mismatch and VE to improve the vaccines developed for future seasons; (2) assessing the ongoing impact of vaccination efforts in the setting of antigenic drift and periodic vaccine reformulation; (3) evaluating of vaccination programs and strategies in terms of individual and population-wide benefits; (4) identifying risk factors for vaccine failure to assist in determining strategies to improve effectiveness; and (5) identifying subgroups that should be targeted in future vaccination campaigns to increase overall effectiveness.

The test negative design (TND) is a type of observational study that was recently developed by scientists in the British Columbia Centre for Disease Control (BCCDC) as a convenient approach to assessing the influenza vaccine effectiveness (VE) using the available sentinel physician network. (Skowronski et al. 2005; 2007) Since 2007, the TND has been popularly used to estimate the influenza VE by investigators in Europe (Hardelid et al. 2011, Kissling et al. 2009, Valenciano et al. 2011), the United States (Belongia et al. 2009, Treanor et al. 2012) and Australia (Fielding et al. 2011, Kelly et al. 2009). This design is popular because of its ease to implement and its advantage to reduce the risk of a particular type of confounding due to underlying differences in health and health care-seeking behavior between persons who choose to receive influenza vaccine and persons who do not (Jackson et al. 2006a;b). In the TND, patients seeking health care for an acute respiratory illness (ARI) are recruited

into the study and tested for influenza using a highly sensitive and specific laboratory test, usually a polymerase chain reaction (PCR) assay. Those tested positive are cases and those tested negative are controls. Patients' vaccination status can be ascertained through patient's report or medical record review. Therefore, the observed data from a TND study can be summarized by four numbers:  $N_{11}$ , the number of persons who are influenza-negative and vaccinated;  $N_{21}$ , the number of persons influenza-positive and vaccinated;  $N_{10}$ , the number of persons who are influenza-negative and not vaccinated; and  $N_{20}$ , the number of persons influenza-positive and not vaccinated.

The TND is different from the traditional case-control design (CCD) in three ways: (1) persons are sampled into the CCD based on their case status (i.e. cases vs. controls), however, persons are sampled into the TND before their case status is known; (2) the marginal ratio of cases to controls is often known and specified in the CCD but not in the TND; and (3) the TND only includes persons with ARI who seek medical care, while the CCD includes both persons with and without ARI. Influenza VE is estimated the same way in the TND and the CCD by 1 minus the ratio of the odds of vaccination in cases to that in controls, i.e.  $1 - \frac{N_{10}N_{21}}{N_{11}N_{20}}$ .

The first application of the TND as a pilot study for assessing the 2004/2005 season influenza VE in Canada was reported in 2005 (Skowronski et al. 2005) and subsequently in 2007 for the 2005/2006 season influenza VE (Skowronski et al. 2007). Since then, the TND has been used in Canada to evaluate the influenza VE within the existing surveillance structures annually. (Janjua et al. 2012, Skowronski et al. 2013; 2010; 2009; 2011; 2012) Since 2007, researchers began theoretical work to examine the validity, accuracy, and assumptions of the TND for estimating influenza VE. Among them, one work examined the impact of the sensitivity and specificity of diagnostic



tests on VE and found that test specificity is the most critical factor influencing the VE estimate in the TND. (Orenstein et al. 2007). Another study verified a core assumption of the TND that the vaccine has no effect on other respiratory viruses that may cause influenza-like-illness and compared the VE estimates from TND to randomized placebo-controlled clinical trials. (De Serres et al. 2013) Findings showed that as long as the core assumption is met, VE estimates from TND have comparable accuracy and precision to those from randomized clinical trial analysis of the same data sets, however, assessment for bias and confounding is still needed when used in observational studies. Following this result, two studies assessed the bias and confounding for TND used in observational studies. (Foppa et al. 2013, Jackson and Nelson 2013) Jackson found that compared to traditional CCD and cohort study, TND is less susceptible to bias due to misclassification of infection and to confounding by health care-seeking behavior. Foppa found that VE estimates from TND are valid and unbiased under a wide range of assumptions. However, if vaccinated cases are less severely ill and seek care less frequently than unvaccinated cases, then adjustment for illness severity or general health status is required to avoid bias in VE estimates.

Although researches have shown that estimates of the influenza VE from the TND are generally valid and less susceptible to bias due to misclassification of infection status and health care seeking behavior, all the studies treat the TND the same as the CCD and estimate the influenza VE from TND using the odds ratio (OR) estimate. However, a TND study is not a strict case-control study and to our knowledge, there does not exist a study that assesses whether the OR estimate is appropriate and the amount of bias of the OR estimate under the TND framework. To fill this gap, in this work we propose to estimate the influenza VE from the TND using a novel

Bayesian model. We model the probability of each person's status determined jointly by vaccination status and the influenza infection status conditional on being sampled into the TND. We take the most basic approach to estimate the influenza VE from the TND. We compare the model-based estimate to OR estimate and assess whether both estimates are accurate and precise.

## 3.2 The model

Consider a study population from which a TND will be conducted during an influenza season. Some persons in the study population develop ARI symptoms and seek medical care, while others either develop ARI symptoms but do not seek medical care or simply do not develop ARI symptoms. Note that ARI symptoms may be due to infection with the influenza virus or with another pathogen. For persons with ARI symptoms seeking medical care, assume they are tested for influenza using PCR test which is assumed to be 100% sensitive and specific.

Denote by  $V$  a person's vaccination status, where  $V = 1$  for vaccination and  $V = 0$  for un-vaccination. We define

$$\alpha = P(V = 1).$$

Denote by  $E$  a variable that contains the information about a person's ARI symptoms and influenza infection status, where  $E = 0$  for persons not developing ARI symptoms,  $E = 1$  for symptomatic persons who are influenza-negative and  $E = 2$  for symptomatic persons who are influenza-positive. We assume that whether a person develops ARI symptoms and his infection status depend on his vaccination status.

We define

$$\beta_{ev} = P(E = e|V = v), \text{ where } e = 0, 1, 2 \text{ and } v = 1, 0$$

$$\text{and } \sum_{e=0}^2 \beta_{ev} = 1.$$

Denote by  $M$  a variable indicating whether a person seek medical care for ARI, where  $M = 1$  representing a person seeking medical care and  $M = 0$  for persons not seeking medical care. We assume that a person's medical care seeking behavior depends on his ARI/influenza infection status and his vaccination status. We define

$$\delta_{ev} = P(M = 1|E = e, V = v), \text{ where } e = 1, 2 \text{ and } v = 1, 0.$$

Normally, a person with no ARI symptoms won't seek medical care for ARI, therefore we define  $\delta_{0v} = P(M = 1|E = 0) = 0$ .

Consider a study population of size  $N$ , from which a TND study will be conducted. Only persons who develop ARI symptoms and seek medical care (i.e.  $M = 1$ ) will be enrolled in a TND study. The probability that a person with ARI symptoms and seeking medical care is:

$$P(M = 1) = \sum_{v,e} P(M = 1|V = v, E = e)P(E = e|V = v)P(V = v)$$

where  $e = 1, 2$  and  $v = 0, 1$ .

Each subject enrolled in the TND study will fall into one of the four categories determined by his/her infection status and vaccination status, i.e. the values of  $E$  and  $V$ . Let  $E_i$  be subject  $i$ 's ARI and influenza infection status,  $V_i$  be his vaccination

status and  $M_i$  be his medical care seeking status. For subject  $i$  in the TND, we have

$$\begin{aligned} P(E_i = e, V_i = v | M_i = 1) &= P(E_i = e | V_i = v, M_i = 1) P(V_i = v | M_i = 1) \\ &= \frac{\alpha \delta_{ev} \beta_{ev}}{\alpha(\delta_{11}\beta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})} \end{aligned}$$

where  $i = 1, \dots, N$ ,  $e = 1, 2$  and  $v = 0, 1$ . The above probabilities can be calculated using the law of total probability and conditional probability. All the details can be found in the Appendix 2.1.

The likelihood function for the observed data  $\mathbf{N} = (N_{10}, N_{11}, N_{20}, N_{21})$  is:

$$\begin{aligned} \pi(N_{11}, N_{21}, N_{20}, N_{10} | \beta_{ev}, \alpha, \delta_{ev}) &= \prod_{i=1}^N P(E_i = e, V_i = v | M_i = 1) \\ &= \frac{\alpha^{N_{11}+N_{21}} (1 - \alpha)^{N_{10}+N_{20}} (\delta_{11}\beta_{11})^{N_{11}} (\delta_{21}\beta_{21})^{N_{21}} (\delta_{10}\beta_{10})^{N_{10}} (\delta_{20}\beta_{20})^{N_{20}}}{[\alpha(\delta_{11}\beta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})]^N} \end{aligned}$$

### 3.2.1 Model representation

The primary interest of this study is the influenza VE against symptomatic influenza. Based on a model we developed in a previous paper (Haber et al. 2014), the true VE against the symptomatic influenza from a TND study is  $1 - \frac{\beta_{21}}{\beta_{20}}$ . We define a few probability ratios from the base parameters:  $\rho_{\beta 1} = \frac{\beta_{11}}{\beta_{10}}$  is the vaccine-related ratio in the probability of non-influenza ARI;  $\rho_{\beta 2} = \frac{\beta_{21}}{\beta_{20}}$  is the vaccine-related ratio in the probability of influenza ARI;  $\rho_{\delta 1} = \frac{\delta_{11}}{\delta_{10}}$  is the vaccine-related ratio in the probability

of seeking medical care for non-influenza ARI;  $\rho_{\delta 2} = \frac{\delta_{21}}{\delta_{20}}$  is the vaccine-related ratio in the probability of seeking medical care for influenza ARI;  $\theta_{\delta} = \frac{\rho_{\delta 2}}{\rho_{\delta 1}}$  measures the inequality in  $\rho_{\delta}$  resulted from the type of infection leading to ARI. With the definition of the probability ratios, the influenza VE =  $1 - \rho_{\beta 2}$ , making  $\rho_{\beta 2}$  the primary interest of this study. A list of all the model parameters and the range of their values are provided in the Table 3.1.

With the newly defined ratios, the likelihood of observed data can be represented as:

$$\frac{\alpha^{N_{11}+N_{21}}(1-\alpha)^{N_{10}+N_{20}}(\rho_{\delta 1}\delta_{10}\rho_{\beta 1}\beta_{10})^{N_{11}}(\delta_{20}\theta_{\delta}\rho_{\delta 1}\rho_{\beta 2}\beta_{20})^{N_{21}}(\delta_{10}\beta_{10})^{N_{10}}(\delta_{20}\beta_{20})^{N_{20}}}{[\alpha(\rho_{\delta 1}\delta_{10}\rho_{\beta 1}\beta_{10} + \delta_{20}\theta_{\delta}\rho_{\delta 1}\rho_{\beta 2}\beta_{20}) + (1-\alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})]^N}$$

### 3.2.2 Prior specifications

In this section, we discuss the prior specifications for the proposed model.

According to the most recent publication (CDC 2013), influenza vaccine coverage in the US in the 2011-12 season ranged between 30% to 70%. We assigned the proportion of influenza vaccine coverage  $\alpha$  as 0.6 for this study. The influenza vaccine coverage will be treated as a constant in this study as it does not affect the influenza VE estimates.

Priors for other parameters were determined based on knowledge from published literature. To determine the potential ranges of these parameters, we used data from randomized clinical trials (RCT) from a recent review paper (Osterholm et al. 2012) and other sources. We found five publications (Beran et al. 2009, Edwards et al. 1994, Frey et al. 2010, Jackson et al. 2010, Madhi et al. 2011) where the numbers of

Table 3.1: Parameters and notation used in the model representation and their value ranges

Parameter	Definition	Values
$\alpha$	Probability of being vaccinated (vaccine coverage)	0.4-0.8
$\beta_{1v}$	Probability of non-influenza ARI for a person of vaccination status $v$	From randomized clinical trials
$\beta_{2v}$	Probability of influenza ARI for a person of vaccination status $v$	From randomized clinical trials
$\rho_{\beta 1} = \frac{\beta_{11}}{\beta_{10}}$	Ratio comparing vaccinees and non-vaccinees w.r.t. probability of non-influenza ARI	0.25-4
$\rho_{\beta 2} = \frac{\beta_{21}}{\beta_{20}}$	Ratio comparing vaccinees and non-vaccinees w.r.t. probability of influenza ARI	0.2-0.7
$\delta_{ev}$	Probability of seeking medical care for ARI for a person of illness/infection status $e$ and vaccination status $v$	0-0.5
$\rho_{\delta 1} = \frac{\delta_{11}}{\delta_{10}}$	Ratio comparing vaccinees and non-vaccinees w.r.t. probability of seeking care for non-influenza ARI	0.25-4
$\rho_{\delta 2} = \frac{\delta_{21}}{\delta_{20}}$	Ratio comparing vaccinees and non-vaccinees w.r.t. probability of seeking care for influenza ARI	0.25-4
$\theta_{\delta} = \frac{\rho_{\delta 2}}{\rho_{\delta 1}}$	Ratio of the 2 ratios defined above	0.25-4

vaccinated and unvaccinated RCT participants who developed ARI with and without influenza infection could be determined. In all these RCTs, culture or RT-PCR was used to confirm influenza infection. Because some of the publications included RCT data from more than one season or RCTs with more than one active vaccine, we identified 14 comparisons of an active influenza vaccine and a placebo in a specific influenza season from the five publications. (Haber et al. 2014) For each of the comparisons, we obtained estimates of  $\beta_{1v}$  and  $\beta_{2v}$  ( $V = 0, 1$ ) from the numbers of influenza and non-influenza cases of ARI stratified by vaccination status. A list of these comparisons and the corresponding estimates for  $\beta_{1v}$  and  $\beta_{2v}$  is given in Table 3.2. Based on the results, we assigned  $\beta_{10}$  and  $\beta_{20}$  both uniform priors with the ranges slightly larger than the minimum and maximum values regardless of vaccination status from the 14 comparisons.

The influenza VE varies by demographic groups. But it's generally believed that the mean influenza VE is around 0.6, with a range between 0.3 and 0.8 and skewed to the left. We therefore assigned a lognormal prior distribution for  $\rho_{\beta_2}$ . We defined  $y = \log \rho_{\beta_2}$  and selected -0.96 as the mean and 0.32 as the standard deviation for  $y$  respectively so that the mean of  $\rho_{\beta_2}$  is 0.4 and the 95% CI is (0.2, 0.7). This makes the VE has a mean of 0.6 and the 95% CI of (0.3, 0.8).

The probability of seeking medical care for ARI in the US has been estimated to be between 0.2 and 0.5 (Ferdinands and Shay 2012). We therefore assigned a uniform prior with ranges (0, 0.5) for both  $\delta_{10}$  and  $\delta_{20}$ .

For the other three parameters,  $\rho_{\beta_1}$ ,  $\rho_{\delta_1}$  and  $\theta_\delta$ , they could range from 0.25 to 4, with a mean around 1. We therefore assigned a lognormal prior distribution for each of them. We defined  $x = \log \rho_{\beta_1}$ ,  $z_{\delta_1} = \log \rho_{\delta_1}$  and  $z_\delta = \log \theta_\delta$ . We selected the

Table 3.2: Estimates of  $\beta_{1v}$  and  $\beta_{2v}$  from the 14 comparisons identified through 5 publications

Comparison	Group	Total	ARI	influenza infection	$\hat{\beta}_{1v}$	$\hat{\beta}_{2v}$
		N	N	N		
1	Vaccine	872	89	6	0.0952	0.0069
	Placebo	878	92	28	0.0729	0.0319
2	Vaccine	878	75	6	0.0786	0.0068
	Placebo	878	92	28	0.0729	0.0319
3	Vaccine	1029	103	12	0.0884	0.0117
	Placebo	1064	125	29	0.0902	0.0273
4	Vaccine	1060	122	9	0.1066	0.0085
	Placebo	1064	125	29	0.0902	0.0273
5	Vaccine	1114	95	3	0.0826	0.0027
	Placebo	1125	119	32	0.0773	0.0284
6	Vaccine	1126	89	8	0.0719	0.0071
	Placebo	1125	119	32	0.0773	0.0284
7	Vaccine	999	78	8	0.0701	0.0080
	Placebo	1016	93	18	0.0738	0.0177
8	Vaccine	1016	75	4	0.0699	0.0039
	Placebo	1016	93	18	0.0738	0.0177
9	Vaccine	4011	254	28	0.0563	0.0070
	Placebo	2003	120	18	0.0509	0.0090
10	Vaccine	3776	189	42	0.0389	0.0111
	Placebo	3843	353	140	0.0554	0.0364
11	Vaccine	3638	243	49	0.0533	0.0135
	Placebo	3843	353	140	0.0554	0.0364
12	Vaccine	1703	181	19	0.0951	0.0112
	Placebo	1725	233	38	0.1130	0.0220
13	Vaccine	2011	181	11	0.0845	0.0055
	Placebo	2043	194	22	0.0842	0.0108
14	Vaccine	255	51	3	0.1882	0.0118
	Placebo	251	57	12	0.1793	0.0478



mean and standard deviation to be 0 and 0.7 respectively so that the mean of the probability ratio is about 1.28 and the 95% CI is (0.25, 3.94).

The above priors are summarized as below:

$$\begin{aligned}
 \beta_{10} &\sim \text{Uniform}\{0, 0.2\} \\
 \beta_{20} &\sim \text{Uniform}\{0, 0.1\} \\
 x &\sim \text{Normal}\{0, 0.7^2\} \\
 y &\sim \text{Normal}\{-0.96, 0.32^2\} \\
 \delta_{10} &\sim \text{Uniform}\{0, 0.5\} \\
 \delta_{20} &\sim \text{Uniform}\{0, 0.5\} \\
 z_{\delta_1} &\sim \text{Normal}\{0, 0.7^2\} \\
 z_{\delta} &\sim \text{Normal}\{0, 0.7^2\}
 \end{aligned}$$

### 3.2.3 Posterior Inference

The posterior distribution of parameters given the data is complicated and has no closed form solution. Thus, to sample from the posterior distribution for given parameters, we resort to the adaptive rejection metropolis sampling within Gibbs sampling (Gilks et al. 1995). Details of the full conditional posterior distributions of  $\beta_{10}$ ,  $\beta_{20}$ ,  $x$ ,  $y$ ,  $\delta_{10}$ ,  $\delta_{20}$ ,  $z_{\delta_1}$  and  $z_{\delta}$  are provided in Appendix 2.2. The model-based estimate is the posterior mean.

### 3.3 Simulation Study

To demonstrate the performance of the proposed model, we conducted simulation studies to estimate  $\rho_{\beta_2}$  for a few different scenarios. Based on previous studies (De Serres et al. 2013, Haber et al. 2014), three factors can affect the bias of the influenza VE from TND study: sample size, vaccine-related ratios in the probability of non-influenza ARI,  $\rho_{\beta_1}$  and the ratio in vaccine-related probability ratio of seeking medical care for ARI,  $\theta_\delta$ . We considered two different sample sizes (large versus small), two different values for  $\rho_{\beta_1}$  (equal to 1 versus unequal to 1, we chose 1.012 as the value close to 1 and 0.709 as the value not close to 1), and two different values for  $\theta_\delta$  (1 versus 1.5). The values of  $\rho_{\beta_1}$  and  $\theta_\delta$  equal to 1 make the VE estimates unbiased. Therefore, there were 8 scenarios in total. The true values of the parameters for each scenario are shown in Table 3.3.

Table 3.3: True values for each parameter and each scenario used in the simulation studies

Scenario	Sample	$\rho_{\beta_1}$	$\theta_\delta$	$\beta_{10}$	$\beta_{20}$	$\beta_{11}$	$\beta_{21}$	$\delta_{10}$	$\delta_{11}$	$\delta_{20}$	$\delta_{21}$
1	large	1.012	1	0.084	0.011	0.085	0.005	0.2	0.3	0.4	0.6
2	large	1.012	1.5	0.084	0.011	0.085	0.005	0.2	0.3	0.4	0.9
3	large	0.709	1	0.055	0.036	0.039	0.011	0.2	0.3	0.4	0.6
4	large	0.709	1.5	0.055	0.036	0.039	0.011	0.2	0.3	0.4	0.9
5	small	1.012	1	0.084	0.011	0.085	0.005	0.2	0.3	0.4	0.6
6	small	1.012	1.5	0.084	0.011	0.085	0.005	0.2	0.3	0.4	0.9
7	small	0.709	1	0.055	0.036	0.039	0.011	0.2	0.3	0.4	0.6
8	small	0.709	1.5	0.055	0.036	0.039	0.011	0.2	0.3	0.4	0.9

Assume a study population from which a TND will be conducted have a population size of  $N$ . The probability that a person will be enrolled into a TND study is the probability that a person seeking medical care, i.e.  $P(M = 1)$ . Given the values of the parameters, we first simulated the number of persons enrolled in a TND study

based on a binomial distribution with sample size  $N$  and the probability  $P(M = 1) = \alpha(\delta_{11}\beta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})$ . The selected parameters making  $P(M = 1)$  around 0.03, we considered the population size of  $N = 50000$  and  $N=10000$  respectively, corresponding to the large sample size of about 1500 persons and the small sample size of about 300 persons in the TND study.

For each person enrolled in the TND study, he or she is assigned to one of the four categories determined by the vaccination status ( $V$ , Yes/1 or No/0) and influenza infection status ( $E$ , uninfected/1 or infected/2) based on a multinomial distribution with sample size one and probabilities  $(p_{11}, p_{10}, p_{21}, p_{20})$ . The probabilities of being in each category are calculated as following:

$$\begin{aligned}
 p_{11} &= \frac{\alpha\beta_{11}\delta_{11}}{\alpha(\beta_{11}\delta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})} \\
 p_{10} &= \frac{(1 - \alpha)\beta_{10}\delta_{10}}{\alpha(\beta_{11}\delta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})} \\
 p_{21} &= \frac{\alpha\beta_{21}\delta_{21}}{\alpha(\beta_{11}\delta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})} \\
 p_{20} &= \frac{(1 - \alpha)\beta_{20}\delta_{20}}{\alpha(\beta_{11}\delta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})}
 \end{aligned}$$

We summarized the simulated numbers of persons in each of the four categories  $\mathbf{N} = (N_{11}, N_{10}, N_{21}, N_{20})$ .

For each scenario, we simulated 500  $\mathbf{N}$ . For each simulated  $\mathbf{N}$ , we set the initial values for  $\beta_{10}$ ,  $\beta_{20}$ ,  $\delta_{10}$  and  $\delta_{20}$  random values in their corresponding value ranges and the initial values for  $x$ ,  $y$ ,  $z_{\delta_1}$  and  $z_{\delta}$  random values between -1 and 1. We run the proposed posterior simulation algorithm 5000 iterations with 1000 burn-in.

We calculated the odds ratio (OR) estimate for  $\rho_{\beta_2} = \frac{N_{10}N_{21}}{N_{11}N_{20}}$  and the model-based estimate for  $\rho_{\beta_2}$  as the posterior mean. We calculated the average bias and the mean square error (MSE) across the 500 OR and model-based estimates. The average bias was calculated as  $\frac{1}{500} \sum_{i=1}^{500} (\hat{\rho}_{\beta_2}^{(i)} - \rho_{\beta_2}^{(i)})$ . The MSE was calculated as  $\frac{1}{500} \sum_{i=1}^{500} (\hat{\rho}_{\beta_2}^{(i)} - \rho_{\beta_2}^{(i)})^2$ .

Table 3.4: Simulation model fitting results for each scenario: model-based estimates versus OR estimates. MSE is the mean square error.

Scenario	Sample	$\rho_{\beta_1}$	$\theta_\delta$	True	OR		Model-based	
				$\rho_{\beta_2}$	Bias	MSE	Bias	MSE
1	large	1.012	1	0.455	0.002	0.006	-0.022	0.001
2	large	1.012	1.5	0.455	0.230	0.068	0.006	0.001
3	large	0.709	1	0.305	0.124	0.018	0.113	0.013
4	large	0.709	1.5	0.305	0.346	0.126	0.131	0.018
5	small	1.012	1	0.455	0.050	0.075	-0.025	0.001
6	small	1.012	1.5	0.455	0.294	0.211	0.002	0.001
7	small	0.709	1	0.305	0.133	0.034	0.114	0.013
8	small	0.709	1.5	0.305	0.360	0.165	0.139	0.020

Table 3.4 presents the average bias and the MSE for the OR and model-based estimates for  $\rho_{\beta_2}$  (true VE =  $1 - \rho_{\beta_2}$ ) for all the scenarios. The results show that in all 8 scenarios, the model-based estimates have smaller average bias and MSE than the OR estimates, implying that the model-based estimates are more accurate and reliable than the OR estimates. Comparison of the estimates for  $\rho_{\beta_1}$  equal to 1 versus unequal to 1 show that both estimators have larger average bias and MSE when  $\rho_{\beta_1}$  is unequal to 1, confirming that  $\rho_{\beta_1} = 1$ , that is equal probability of non-influenza ARI among vaccinees and nonvaccinees is the core assumption for TND. Sample size and  $\theta_\delta$  seemed to have different impact on the two estimators. For OR estimates, both the average bias and the MSE are bigger for small sample size as compared to big sample size and are bigger when  $\theta_\delta = 1.5$  as compared to 1. However, neither of

them seem to affect the model-based estimates.

In addition to the estimates for the parameter of primary interest, the proposed model also provides estimates for all other parameters. Table 3.5 presents the absolute relative bias for model-based estimates for all the parameters of each scenario. The results show that the relative bias are pretty large for some parameters and this is likely due to the small scales of the parameters' true values.

Table 3.5: Simulation model fitting results: absolute relative bias for estimates of all parameters for each scenario

Parameter	Scenario							
	1	2	3	4	5	6	7	8
$\beta_{10}$	0.616	0.614	0.245	0.226	0.638	0.635	0.396	0.183
$\beta_{20}$	1.466	1.513	0.738	0.667	1.791	1.899	0.791	0.553
$\beta_{11}$	1.380	1.235	1.429	1.136	1.449	1.316	1.595	1.006
$\beta_{21}$	1.333	1.527	1.401	1.412	1.579	1.843	1.459	1.247
$\delta_{10}$	0.472	0.481	0.237	0.185	0.561	0.527	0.895	0.806
$\delta_{20}$	0.015	0.009	0.475	0.430	0.051	0.029	0.928	0.867
$\delta_{11}$	0.138	0.492	0.522	0.355	0.177	0.442	0.918	0.816
$\delta_{21}$	0.927	0.858	0.212	0.225	0.957	0.839	0.110	0.070

The simulation studies show that the model-based estimates are more accurate and reliable than the OR estimates. The model-based estimates also have the advantage over the OR estimates in that the former is not affected by sample size or  $\theta_\delta$ . In addition, the proposed model can provide estimates for other parameters that cannot be estimated through any other existing methods.

## 3.4 Application

In this section, we apply the bayesian model to estimate the influenza VE to data from the a large TND study conducted during 2010-2011 in the United States reported by Treanor et al. (2012). This study was conducted to assess the 3 types of seasonal influenza vaccines that were antigenically similar to circulating influenza viruses during the 2010-2011 US influenza season.

### 3.4.1 Analysis of the seasonal 2010-2011 influenza vaccine TND study

The study population included children and adults who resided in the counties and zip codes surrounding the study centers: Marshfield Clinic and St. Joseph’s Hospital, Marshfield, Wisconsin; the University of Michigan Health System, Ann Arbor, and Henry Ford Health System, Detroit, Michigan; the Strong Memorial and Rochester General Hospitals from the University of Rochester, Rochester, New York; and Vanderbilt University, Summit, St. Thomas, and Baptist Hospitals, Nashville, Tennessee. More than 5000 subjects who were aged  $\geq 6$  months and had an ARI with a duration of  $\leq 7$  days with documented fever or history of feverishness or cough were enrolled in the study and about 4700 were included in the analysis.

Participants’ receipt of seasonal 2010-2011 influenza vaccine was ascertained through patient or parental report and confirmed by medical record review, state registries or a real-time internet-based vaccine registry for Wisconsin participants. Those who received at least 1 dose of seasonal influenza vaccine at least 14 days before illness

onset were defined as being vaccinated. Cases of medically-attended lab-confirmed influenza were defined as individuals meeting the medically attended ARI definition with rRT-PCR-confirmed influenza. Controls were medically-attending ARI individuals whose rRT-PCR was negative for influenza.

Since influenza vaccine effectiveness may vary by age group, especially it may be lower among individuals aged  $\geq 65$  years, we assessed the influenza VE for each age group. Data from the US 2010-2011 case-control study for each type of the seasonal influenza vaccine for each age group are presented in Table 3.6. There are 10 groups in total.

Table 3.6: Number of patients by influenza status and vaccination status for each age group and each vaccine component in the 2010-2011 TND study

Age Group	Influenza positive		Influenza negative	
	Vaccinated	unvaccinated	Vaccinated	unvaccinated
Any seasonal vaccine				
6 months-2 years	47	46	465	210
3-8 years	79	190	438	318
9-49 years	104	375	534	891
50-64 years	47	77	263	207
$\geq 65$ years	40	23	258	100
Inactivated seasonal vaccine				
2-8 years	66	217	443	390
9-49 years	91	375	477	891
$\geq 50$ years	79	100	491	307
Live attenuated seasonal vaccine				
2-8 years	22	217	128	390
9-49 years	9	375	34	891

We applied the proposed model to estimate the influenza VE for each age group and vaccine type. The initial values for all the parameters were assigned random values in their corresponding value ranges. We calculated the posterior mean and

95% credible intervals (CI) for the influenza VE. We checked the convergence of the simulated Markov chains using the Gelman and Rubin diagnostic (Gelman and Rubin 1992) by running five additional Markov chains with different initial values. The potential scale reduction factors (PSRF) of the log likelihood for the simulated scenarios ranged between 1 and 1.04, which are close to 1, indicating the convergence of posterior simulations.

Table 3.7: The model-based and unadjusted odds ratio (OR) influenza VE estimates for each age group and each vaccine type in the 2010-2011 TND study

Age Group	Unadjusted OR		Model-based		
	VE	95% CI <sup>1</sup>	VE	95 % CI <sup>2</sup>	Model checking P-value <sup>3</sup>
Any seasonal vaccine					
6 months-2 years	0.54	0.28-0.70	0.58	0.26-0.78	0.07
3-8 years	0.70	0.59-0.78	0.62	0.34-0.81	0.06
9-49 years	0.54	0.41-0.64	0.62	0.39-0.80	0.07
50-64 years	0.52	0.28-0.68	0.58	0.28-0.78	0.07
>= 65 years	0.33	-0.18-0.62	0.56	0.25-0.77	0.09
Inactivated seasonal vaccine					
2-8 years	0.73	0.64-0.80	0.62	0.35-0.80	0.06
9-49 years	0.55	0.42-0.65	0.62	0.38-0.79	0.06
>= 50 years	0.51	0.31-0.64	0.59	0.28-0.77	0.00
Live attenuated seasonal vaccine					
2-8 years	0.69	0.50-0.81	0.65	0.41-0.81	0.08
9-49 years	0.37	-0.32-0.70	0.68	0.45-0.83	0.05

1. CI stands for confidence interval.
2. CI stands for credible interval.
3. Model checking P-value refers to the posterior predictive p-values based on the chi-square discrepancy test described in section 3.4.2.

Table 3.7 presents the OR influenza VE estimate and 95% confidence interval and the model-based influenza VE estimate and 95% credible interval for each age group and each vaccine type. The results show that the model-based and the OR



estimates are comparable for most groups except for the any seasonal vaccine  $\geq 65$  years group and the live attenuated seasonal vaccine 9-49 years group. The model-based 95% credible intervals overlap with the OR 95% confidence intervals for all groups and the former is slightly larger than the latter for most groups except for those two groups. For the two groups, the model-based estimates are higher than the OR estimates and the model-based 95% credible intervals are smaller than the OR 95% confidence intervals. The number of influenza positive nonvaccinees in the any seasonal vaccine  $\geq 65$  years group is 23 and the number of influenza positive vaccinees in the live attenuated seasonal vaccine 9-49 years group is only 9. The small numbers of participants in these groups may affect the accuracy and the uncertainty for both the model-based and OR the estimates. However, the model-based estimates may be less affected by the small sample size as demonstrated in the simulation studies.

### 3.4.2 Model assessment

We conducted a posterior predictive model assessment using the  $\chi^2$  discrepancy for each age group. The  $\chi^2$  discrepancy is a summary statistic for the sum of squares of standardized residuals of the data with respect to their expectations under the model (Gelman et al. 1996). For our model, the  $\chi^2$  discrepancy is defined as:

$$\chi^2(\mathbf{N}; \beta_{10}, \beta_{20}, x, y, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}) = \sum_{e=1,2;v=1,0} \frac{(\mathbf{N} - \mathbf{E}(\mathbf{N}))^2}{\text{Var}(\mathbf{N}|\beta_{10}, \beta_{20}, x, y, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta})}$$

where  $\mathbf{N} = (N_{11}, N_{10}, N_{21})$ . We calculated the posterior predictive  $p$ -value based on  $\chi^2$  as  $P = P(\chi^2(\mathbf{N}^{rep}; \beta_{10}, \beta_{20}, x, y, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}) \geq \chi^2(\mathbf{N}; \beta_{10}, \beta_{20}, x, y, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}))$ , where  $\mathbf{N}^{rep}$  represents the predictive replication and  $\mathbf{N}$  represents the observed data.

The model checking p-values of the predictive versus realized  $\chi^2$  discrepancies are presented in Table 3.7. The p-values are greater than or equal to 0.05 among 9 out of the 10 age groups, indicating the difference between observed and predictive  $\mathbf{N}$  are not significant. This implies that the model fits the data pretty well except for the inactivated seasonal vaccine  $\geq 50$  years group. Potential reasons for the big discrepancy in this group may be: (1) the influenza VE may differ between persons aged 50-64 and persons aged 65 years and older, therefore, a single VE estimate for this group may not be sufficient and (2) the specified priors may not be specific enough for this age group, which may result in biased estimate for this group.

### 3.4.3 Sensitivity Analysis

We also performed a sensitivity analyses to assess the impact of priors on the model-based VE estimates. Because the influenza  $VE = 1 - \frac{\beta_{21}}{\beta_{20}}$ ,  $\beta_{10}$  and  $\beta_{20}$  are important factors that will affect VE estimates. We therefore specified different priors for  $\beta_{10}$  and  $\beta_{20}$  while keeping the priors for the other parameters unchanged. We estimated the influenza VE with different prior specifications for  $\beta_{10}$  and  $\beta_{20}$  and assessed how different priors would affect the model-based estimates for influenza VE.

In the proposed model, we assigned non-informative priors for  $\beta_{10}$  and  $\beta_{20}$  through Uniform(0, 0.2) and Uniform(0, 0.1), respectively. In the sensitivity analyses, we considered two additional prior specifications for  $\beta_{10}$  and  $\beta_{20}$  through Beta distributions with wider or smaller ranges. For  $\beta_{10}$ , we assigned Beta (2, 20), which has a mean of 0.17 and 95% CI of (0.01, 0.24) and Beta (2, 8), which has a mean of 0.2 and 95% CI of (0.03, 0.48). For  $\beta_{20}$ , we assigned Beta (2, 40), which has a mean of 0.09

Table 3.8: The priors for  $\beta_{10}$  and  $\beta_{20}$  used in the sensitivity analysis

Analysis	Prior for $\beta_{10}$	Prior for $\beta_{20}$
1	Unif(0, 0.2)	Unif(0, 0.1)
2	Unif(0, 0.2)	Beta(2, 40)
3	Unif(0, 0.2)	Beta(2, 16)
4	Beta(2, 20)	Unif(0, 0.1)
5	Beta(2, 20)	Beta(2, 40)
6	Beta(2, 20)	Beta(2, 16)
7	Beta(2, 8)	Unif(0, 0.1)
8	Beta(2, 8)	Beta(2, 40)
9	Beta(2, 8)	Beta(2, 16)

and 95% CI of (0.006, 0.13) and Beta (2, 16), which has a mean of 0.1 and 95% CI of (0.015, 0.29). Together with the non-informative priors specified in the proposed model, there are 9 prior specifications for  $\beta_{10}$  and  $\beta_{20}$ . The nine prior specifications are summarized in the Table 3.8.

For each prior specification, we estimated the influenza VE estimates for each age group. Table 3.9 presents the mean, minimum, median and maximum of the VE estimates from the sensitivity analyses for each age group. The results show that the VE estimates from using nine different prior specifications are very close and they only vary in a very tight range, implying that the influenza VE estimates are not affected by the prior specifications of  $\beta_{10}$  and  $\beta_{20}$ .

### 3.5 Discussion

In this paper, we develop a novel Bayesian model to estimate the influenza VE using data collected from the test negative study. Our model is based on the joint probability of a person's vaccination status and influenza infection status conditional on the

Table 3.9: The mean, minimum, median, and maximum of VE estimates from sensitivity analyses

Name	Mean	Minimum	Median	Maximum
Any seasonal vaccine				
6 months-2 years	0.577	0.564	0.578	0.590
3-8 years	0.613	0.602	0.611	0.622
9-49 years	0.611	0.594	0.612	0.629
50-64 years	0.595	0.582	0.599	0.607
$\geq 65$ years	0.557	0.550	0.556	0.565
Inactivated seasonal vaccine				
2-8 years	0.621	0.609	0.621	0.636
9-49 years	0.615	0.594	0.619	0.630
$\geq 50$ years	0.584	0.563	0.582	0.601
Live attenuated seasonal vaccine				
2-8 years	0.645	0.641	0.646	0.650
9-49 years	0.676	0.670	0.676	0.680

person seeking medical care for ARI symptoms and is tested for influenza infection. We elicit the subjective priors from published literature. Through the simulation studies, we demonstrate that the estimates from the proposed model are superior to the OR estimates in the following aspects: (1) the model-based estimates are more accurate and reliable than the OR estimates; (2) the model-based estimates are not affected by the sample size,  $\theta_\delta$ , the inequality of the vaccine related probability ratios of seeking medical care for ARI, while the OR estimates can be affected by both factors; and (3) in addition to the influenza VE estimates, the proposed model can provide estimates for other parameters that are of interest to researchers in the area of influenza research but have not been estimated through any existing methods from observational studies previously.

Our model is the very first that estimates the influenza VE through a modeling approach instead of the odds ratio estimate. The priors for our model are chosen

subjectively based on current literature. As shown in the application, the influenza VE estimates from the proposed model are mostly comparable to the OR estimates, indicating that the simple OR approach can usually provide good estimates. However, when sample sizes are small, the OR estimates can be more biased and less precise than the model-based estimates. The 95% credible intervals of the model-based estimates are slightly larger than the 95% confidence intervals of the OR estimates, this might be due to the prior setup. In addition, the limited input data (only four data points) can affect the precision and power of the model-based estimates.

Our current work can be extended in several possible directions. One extension is to develop a hierarchical Bayesian model that includes covariates such as age effect and health status in the model. Assuming parameters among different age groups can be characterized through a certain relationship, inclusion of the age effect in the model has the advantage of making use of all the data points in one analysis instead of separate stratified analyses. Another issue is the prior setup. Although the priors in the current model are selected subjectively based on literature, they are general priors. Our sensitivity analyses demonstrate that the VE estimates are not affected by different priors for  $\beta_{10}$  and  $\beta_{20}$ , however, impact of different priors for other parameters on the VE estimates are unclear and worthy of investigation. For example, influenza VE differs by age group, therefore, a specific prior for influenza VE of each age group could help improve the estimates.

## Chapter 4

# A Nonhomogeneous Probability Model for Evaluating Bias and Precision of Estimates of the Influenza Vaccine Effectiveness from Case-Control Studies

### 4.1 Background

Many factors make accurate estimation of influenza VE very challenging. Firstly, the predominant influenza virus types, subtypes and phenotypes change from one season to the next, necessitating a new vaccine targeting different strains and a new

VE estimate in every season. Secondly, for many years, randomized clinical trials (RCTs) that measure laboratory-confirmed influenza virus infection as the outcome have been considered to provide the most accurate estimates of influenza VE. However, such trials have become unethical to conduct because influenza vaccination is now recommended for all persons 6 months of age and older in the U.S. and in other countries. (CDC 2013) Lately, estimates of influenza VE have increasingly relied on observational studies, especially case-control studies for reasons of statistical power, logistics and cost. However, observational studies innately are subject to more biases. Thirdly, it is not easy to find all or most influenza patients in a given community, because influenza symptoms can be mild and many patients do not seek medical care to alleviate them. Fourth, symptoms of influenza are not specific, hence many patients who develop an acute respiratory illness (ARI) are not infected with the influenza virus. Fifth, special laboratory tests are required to confirm influenza infection, and these tests are not 100% sensitive and specific, causing misclassification bias. For these reasons, there is increasing need to carefully design observational studies to avoid, or at least to minimize the various sources of bias in the estimates of influenza VE.

Because of the advantages in cost and logistics, case-control studies are commonly used observation studies to estimate influenza VE. Two commonly used observational study designs for estimating the effectiveness of influenza vaccine against seasonal and pandemic influenza illness are: traditional case-control design (CCD) and the test negative design (TND). In both study designs, individuals of the study population who seek care for an ARI and test positive for influenza infection are considered cases. In the CCD, controls are asymptomatic persons randomly chosen from all members

of the study population who have not reported to the clinic because of ARI either at the same time when the case was reported or at the end of the study. In the TND, ARI patients who test negative for influenza infection serve as controls.

#### **4.1.1 Main sources of bias in case-control studies**

Both CCD and TND are based on subjects who have ARI that could result from an influenza infection and seek medical care. This restriction, along with the fact that vaccination is not randomized, underlies many sources of bias in influenza VE estimates from case-control studies. Below we list a summary of the main sources of bias:

1. Ascertainment of cases (selection bias): A person who develops an ARI may or may not decide to seek medical care. In both CCD and TND, only a person who seeks medical care for ARI can be tested and be considered as a case. In other words, the cases are drawn from a subset of the population (persons who seek care for ARI) that may not be a representative sample of all cases.
2. Confounding by propensity of seeking medical care: The likelihood of seeking medical care may be related to a person's vaccination status, as a vaccinated individual may be more health-conscious so that her/his probability of seeking medical care for ARI may be different from that of an unvaccinated person. In CCD, only persons who seek medical care for ARI can be considered cases, while controls are selected from the entire population. This may confound the association between vaccination status and being considered a case and lead to falsely low estimates of VE. This source of confounding is eliminated in TND, as



both cases and controls are persons seeking care for ARI. However, the selection bias discussed above may become even more pronounced in TND.

3. Probabilities of non-influenza ARI may depend on vaccination status: In TND, persons with non-influenza ARI serve as controls. Therefore one of the assumptions underlying this study design is that vaccinees and non-vaccinees have the same probability of developing ARI as a result of a non-influenza infection. This assumption has not yet been verified. In fact, a recent randomized controlled influenza vaccine trial (Cowling et al. 2012) found that vaccinees had a significantly increased risk of virologically-confirmed non-influenza infection that may lead to ARI.
4. Other confounders: such as health status (Hak et al. 2002, Jackson et al. 2006a), age, exposure, education, SES, may be associated with both the likelihood of being vaccinated and the likelihoods of becoming infected, developing ARI and seeking medical care.
5. Misclassification of influenza infection status and/or vaccination status: As already mentioned, even the best diagnostic tests for influenza infection are not 100% sensitive and specific. Vaccination status may also be misclassified.

Because the two case-control studies are subject to various sources of bias as described above, it is essential to carefully evaluate the properties (bias and precision) of the VE estimates from these study designs. The evaluation of the TND (Orenstein et al. 2007) is especially important, because this design is relatively new and is becoming very popular due to its simplicity. In fact, over 90% of all observational influenza VE studies conducted since 2008 have followed this design. However, the properties

of the TND (bias and precision) have never been evaluated while accounting for all potential sources of bias. The comparison of the TND and the traditional CCD is also of special importance. On one hand, using only ARI patients seeking medical care for their symptoms as controls in the TND should eliminate the impact of an important confounder, namely the propensity to seek medical care, since both cases and controls are selected from the same sub-population of persons seeking care for ARI. On the other hand, this sub-population on which the TND is based may not represent the entire population, resulting in selection bias. In addition, VE may be different for persons who are more likely and less likely to seek care for an ARI. For example, persons with underlying medical conditions may be more likely to have an inferior immune response to vaccination than healthy persons, and they may also be more or less likely to seek care. Very few field studies (Belongia et al. 2009) used both study designs to obtain concurrent estimates of influenza VE in the same population in the same season. The results of this evaluation can help guide the development of new study design or improve existing study designs.

#### **4.1.2 Medically-attended influenza and symptomatic influenza**

Another important issue related to the influenza VE estimates and evaluation of the estimation biases is the outcome measures, i.e. the outcome against which the vaccine is supposed to protect.

In both CCD and TND, a study participant becomes a case if and only if she/he seeks medical care for an ARI and tests positive for influenza infection. Therefore these studies are considered appropriate for estimating VE against medically-attended

influenza (MI), where we define a true case of MI as a truly influenza-infected person who reports to a clinic because of her/his ARI. (Note that testing positive for influenza infection is not required to be a true case of MI). It can be argued that MI is indeed the outcome of interest, as individuals who do not seek medical care when they develop an ARI do not increase health-care costs. On the other hand, persons with ARI who do not seek medical care are still capable of infecting others, missing work or school and developing severe or even deadly complications. Therefore we will also consider in this work a more general definition of influenza illness, namely symptomatic influenza (SI). This definition, which includes everyone who is infected with the influenza virus and develops an ARI, is more appropriate from the public health perspective. Since both CCD and TND are based on persons seeking medical care for an ARI as cases, they miss cases of SI who do not seek medical care. This may further contribute to the bias of VE estimates when SI is the outcome of interest (Orenstein et al. 2007). One should note that when evaluating the biases of VE estimates resulting from CCD and TND, the distinction between MI and SI only affects the true value of VE. The study designs and VE estimates are considered the same, regardless of the outcome of interest.

Recently, several models have been developed to evaluate and compare influenza VE study designs (De Serres et al. 2013, Ferdinands and Shay 2012, Foppa et al. 2013, Haber et al. 2014, Jackson and Nelson 2013, Orenstein et al. 2007). However, these models suffer from the following limitations: (1) they assume random vaccination; (2) they assume that the population is homogeneous with respect to the parameters determining probabilities of infection, illness, seeking medical care, etc.; and (3) they do not consider the precision of VE estimates. Our proposed probability model accounts

for these and other factors that may impact the bias and precision of VE estimates from observational studies.

The goal of this article is to evaluate and compare the bias and precision of estimates of influenza VE resulting from TND and CCD. Specifically, we will (a) evaluate the bias of each of the VE estimates by comparing the expected value of the estimate with the true VE for each of the outcomes of interest, MI and SI, (b) evaluate and compare the standard errors of the VE estimates under equal sample sizes. To conduct these evaluations and comparisons we will develop a detailed stepwise probability model of the process involved in collecting data in these studies and obtaining VE estimates. The model will allow us derive both general and numerical results under different scenarios.

## **4.2 Method**

We first describe the real-life process involved in conducting the two kinds of studies and obtaining the estimates of VE. We then describe the model that will be used to mimic the process.

### **4.2.1 The study population and designs**

The source population for both CCD and TND consists of all individuals who receive most of the health care at a single clinic or at a given network of clinics. Since influenza VE varies by age, we assume that the model pertains to a subpopulation corresponding to a specific age group.

When a member of the study population develops an ARI, she/he may decide to report to a clinic for treatment. At the clinic, the health care provider will ask the person to be tested for influenza infection. If the person agrees, a swab is taken and sent to a laboratory for testing. The test may not be 100% sensitive or specific. In both study designs, a person who tests positive for influenza is eligible to be considered a case. In the traditional CCD, controls are usually selected when cases are identified by randomly contacting members of the population and asking each of them if he has had an ARI since the beginning of the season. Those who have not developed an ARI are eligible to be included as controls. In the TND, an individual who reports to the clinic, is tested and the test result is negative is eligible to be considered as a control. In both study designs, the vaccination status of every case or control is determined from manual or electronic records.

#### **4.2.2 The model**

The model we develop and use for comparing the estimates from the two study designs follows the scheme described above with a few simplifications. We include a latent variable  $X$  for health status. The probabilities of vaccination, ARI, influenza infection and seeking medical care may depend on  $X$ . We assume that (1) when a person seeks medical care for ARI, the probability of being tested for influenza infection does not depend on health status, vaccination status or influenza infection status; (2) given a person's health status and influenza infection status, the probability of seeking medical care no longer depends on the vaccination status; (3) given a person's symptoms and influenza infection status, the sensitivity and specificity of the test do not depend on the tested person's health status or vaccination status; (4) a person's

vaccination status is determined without error; and (5) controls in the CCD are a random sample from the subpopulation of all asymptomatic individuals.

Our model consists of five steps, where the value of a single variable is determined at each step. The distribution of this variable may depend on the values of the variables from the previous steps. Below we define the five steps, the associated variables and the probabilities determining each variable's distribution.

**Step 1: Health status.** We assume the members of the population may be classified as "healthy" or "frail". Define a binary variable  $X$ , where  $X = 1$  for a "healthy" person, and denote  $\pi = P(X = 1)$ . Note that  $X$  (health status) is considered latent (unobservable); we are unable to estimate separate VE in healthy and frail persons or adjust for health status when estimating overall VE.

**Step 2: Vaccination.** A person may be vaccinated against influenza. The probability of vaccination may depend on  $X$ . Define a binary variable  $V$ , where  $V = 1$  for a vaccinated person. Denote  $\alpha_x = P(V = 1 | X = x)$ ,  $x = 0, 1$ .

**Step 3: Influenza infection and ARI.** During the influenza season, a person may become infected with an influenza virus. Both influenza infected and uninfected individuals may develop an ARI. Since our outcome measures only involve symptomatic individuals, we do not distinguish between asymptomatic individuals who are or are not influenza-infected. Therefore, we define a variable  $E$  for the illness/infection status with 3 levels as follows:  $E = 0$  indicating no ARI,  $E = 1$  for ARI without influenza infection (i.e. an ARI resulting from a different pathogen) and  $E = 2$  for ARI and influenza in-

fection (symptomatic influenza). The distribution of  $E$  may depend on  $V$  and  $X$ . Denote  $\beta_{vx} = P(E = 1 \mid V = v, X = x)$ ,  $v = 0, 1$   $x = 0, 1$  and  $\gamma_{vx} = P(E = 2 \mid V = v, X = x)$ ,  $v = 0, 1$   $x = 0, 1$ .  $\beta_{vx} + \gamma_{vx} \leq 1$  for all  $v, x$ .

**Step 4: Seeking medical care for ARI.** A person with ARI may seek medical care and get tested for influenza infection. Define a binary variable  $M$  with  $M = 1$  for a person seeking medical care for his/her ARI. The probability of this event depends on  $E$  (only individuals with an ARI seek medical care) and  $X$ . Denote  $\delta_{ex} = P(M = 1 \mid E = e, X = x)$ ,  $e = 1, 2$   $x = 0, 1$ . Note that  $P(M = 1 \mid E = 0) = 0$ . Because of assumption (2),  $P(M = 1)$  does not directly depend on  $V$  when  $X$  and  $E$  are given.

**Step 5: Testing for influenza infection.** Although only individuals who seeks medical care for ARI are tested for influenza infection, it will be convenient to define a binary variable  $T$  as the (possibly unobserved) test result for any person with an ARI, regardless of whether or not he is actually tested. Define  $T = 1$  ( $T = 0$ ) if a person would test positive (negative) for influenza if tested. Because of assumption (3) above, the probability of testing positive given the person's influenza infection status does not depend on  $X$ ,  $V$  or  $M$ . Denote  $\tau_e = P(T = 1 \mid E = e)$  for  $e = 1, 2$ . Note that  $\tau_1$  is one minus the test's specificity and  $\tau_2$  is the test's sensitivity. In this study, we assume the test has 100% sensitivity and 100% specificity, i.e.  $P(T = 1 \mid E = 1) = \tau_1 = 0$  and  $P(T = 1 \mid E = 2) = \tau_2 = 1$ .

Our model has 17 parameters, which specify the conditional distribution of each variable in terms of the values of the variables determined in the previous steps. The

list of the random variables are provided in Table 4.1.

Table 4.1: Random variables used to define the five steps of the latent-class process model and in the calculation

Variable	Definition	Values
$X$	Health status	0 - frail persons 1 - healthy persons
$V$	Vaccination status	0 - unvaccinated 1 - vaccinated
$E$	Influenza infection and ARI status	0 - no ARI 1 - ARI, not influenza infected 2 - ARI, influenza infected
$M$	Seeking medical care for ARI	0 - no 1 - yes
$T$	Result of test for influenza infection	0 - negative 1 - positive
$C_A$	Case/control status in TND study	0 - control 1 - case
$C_B$	Case/control status in CCD study	0 - control 1 - case

### 4.2.3 Outcome of interest and true VE

As mentioned previously, we consider both medically-attended influenza (MI) and symptomatic influenza (SI) as outcomes of interest. A person is considered a true case of MI if he/she is influenza-infected, develops an ARI and seeks medical care. A person is considered a true case of SI if he/she is influenza-infected and developed an ARI because of the infection. The true VE is defined as one minus the ratio of the probability of the outcome in vaccinees and non-vaccinees.

The true VE against MI is:

$$VE_{TMI} = 1 - RR_{TMI} \quad \text{where} \quad RR_{TMI} = \frac{P(E = 2, M = 1 \mid V = 1)}{P(E = 2, M = 1 \mid V = 0)}$$



The true VE against SI is:

$$VE_{TSI} = 1 - RR_{TSI} \quad \text{where} \quad RR_{TSI} = \frac{P(E = 2 | V = 1)}{P(E = 2 | V = 0)}$$

Using the parameters defined above, and  $VE_{TMI}$  and  $VE_{TSI}$  can be written as:

$$VE_{TMI} = 1 - RR_{TMI} = 1 - \frac{[\delta_{20}\gamma_{10}\alpha_0(1-\pi) + \delta_{21}\gamma_{11}\alpha_1\pi][1 - \alpha_0(1-\pi) - \alpha_1\pi]}{[\alpha_0(1-\pi) + \alpha_1\pi][\delta_{20}\gamma_{00}(1-\alpha_0)(1-\pi) + \delta_{21}\gamma_{01}(1-\alpha_1)\pi]}$$

$$VE_{TSI} = 1 - RR_{TSI} = 1 - \frac{[\gamma_{10}\alpha_0(1-\pi) + \gamma_{11}\alpha_1\pi][1 - \alpha_0(1-\pi) - \alpha_1\pi]}{[\alpha_0(1-\pi) + \alpha_1\pi][\gamma_{00}(1-\alpha_0)(1-\pi) + \gamma_{01}(1-\alpha_1)\pi]}$$

In Appendix 3.1, we show how we derive these expressions.

#### 4.2.4 VE estimates

We only consider estimates of VE that are not adjusted for possible confounders. In both study designs, VE is estimated as one minus odds ratio (OR) in the  $C \times V$  table corresponding to the individuals included in the study, where  $C$  is a binary indicator of case/control status ( $C = 1$  for a case). The bias is defined as the difference between the expectation of the estimated VE and the true VE. For convenience, the TND and CCD will be represented by the letters  $A$  and  $B$ , respectively.

In the TND, the case/control variable is denoted  $C_A$ , where  $\{C_A = 1\} = \{M = 1, T = 1\}$  and  $\{C_A = 0\} = \{M = 1, T = 0\}$ . Then the estimate of VE is:

$$VE_A = 1 - OR_A \quad \text{where} \quad OR_A = \frac{P(C_A = 1, V = 1 | M = 1)P(C_A = 0, V = 0 | M = 1)}{P(C_A = 1, V = 0 | M = 1)P(C_A = 0, V = 1 | M = 1)}$$

Note that all the probabilities conditional on  $M = 1$  as only individuals who seek medical care for ARI can be included in the TNC study.  $OR_A$  can be further written as:

$$OR_A = \frac{P(M = 1, T = 1, V = 1)P(M = 1, T = 0, V = 0)}{P(M = 1, T = 1, V = 0)P(M = 1, T = 0, V = 1)}$$

Using the parameters defined above,  $VE_A$  can be written as:

$$VE_A = 1 - \frac{[\delta_{20}\gamma_{10}\alpha_0(1 - \pi) + \delta_{21}\gamma_{11}\alpha_1\pi][\delta_{10}\beta_{00}(1 - \alpha_0)(1 - \pi) + \delta_{11}\beta_{01}(1 - \alpha_1)\pi]}{[\delta_{20}\gamma_{00}(1 - \alpha_0)(1 - \pi) + \delta_{21}\gamma_{01}(1 - \alpha_1)\pi][\delta_{10}\beta_{10}\alpha_0(1 - \pi) + \delta_{11}\beta_{11}\alpha_1\pi]}$$

The proof can be found in Appendix 3.2.

In the CCD, denote  $C_B$  as the case/control variable. Cases are defined the same way as in the TNC study, i.e.  $\{C_B = 1\} = \{C_A = 1\} = \{M = 1, T = 1\}$ . Controls are those included in a random sample drawn from all the asymptomatic individuals. In other words,  $\{C_B = 0\}$  is a random subset of  $\{E = 0\}$ . In addition, we define a binary variable  $B$  indicating whether or not a person is included in the CCD study, i.e.  $\{B = 1\} = \{C_B = 1 \text{ or } C_B = 0\}$ . The VE estimates is based on the OR in the  $C_B \times V$  table when all the probabilities conditional on  $B = 1$ . The VE estimate from CC study is:

$$VE_B = 1 - OR_B \text{ where } OR_B = \frac{P(C_B = 1, V = 1 | B = 1)P(C_B = 0, V = 0 | B = 1)}{P(C_B = 1, V = 0 | B = 1)P(C_B = 0, V = 1 | B = 1)}$$

$OR_B$  can be further written as:

$$OR_B = \frac{P(M = 1, T = 1, V = 1)P(E = 0, V = 0)}{P(M = 1, T = 1, V = 0)P(E = 0, V = 1)}$$

Using the parameters defined above,  $VE_B$  can be written as:

$$VE_B = 1 - \frac{[\delta_{20}\gamma_{10}\alpha_0(1-\pi) + \delta_{21}\gamma_{11}\alpha_1\pi][(1-\beta_{00}-\gamma_{00})(1-\alpha_0)(1-\pi) + (1-\beta_{01}-\gamma_{01})(1-\alpha_1)\pi]}{[\delta_{20}\gamma_{00}(1-\alpha_0)(1-\pi) + \delta_{21}\gamma_{01}(1-\alpha_1)\pi][(1-\beta_{10}-\gamma_{10})\alpha_0(1-\pi) + (1-\beta_{11}-\gamma_{11})\alpha_1\pi]}$$

The proof can be found in Appendix 3.2.

#### 4.2.5 Standard errors of the VE estimates

We use approximations to the standard errors (SE) of odds ratios (Agresti, 2013) based on the delta method to derive expressions for the SE of both VE estimates in terms of the parameters and the corresponding sample size(s). For evaluating the standard errors we consider the observed odds ratios, where the probabilities are replaced by the observed relative frequencies.

For the TND, the approximate standard error of  $VE_A$  is:

$$\begin{aligned} SE(VE_A) &= SE(OR_A) \approx (OR_A) \times SE(\log(OR_A)) \\ &\approx \frac{p_{11}^A(1-p_{01}^A)}{p_{01}^A(1-p_{11}^A)} \sqrt{\frac{1}{N^A} \left[ \frac{1}{p_{V1}^A p_{11}^A} + \frac{1}{(1-p_{V1}^A) p_{01}^A} + \frac{1}{p_{V1}^A (1-p_{11}^A)} + \frac{1}{(1-p_{V1}^A) (1-p_{01}^A)} \right]} \end{aligned}$$

where  $N^A$  is the number of persons who were tested for influenza ( $M = 1$ ), i.e., the total sample size for the TND. The probabilities ( $p_{V1}^A, p_{11}^A, p_{01}^A$ ) can be written in terms of the parameters defined earlier:

$$\begin{aligned} p_{V1}^A &= P(V = 1 | M = 1) = \frac{P(M = 1, V = 1)}{P(M = 1)} \\ &= \frac{\alpha_0(1-\pi)(\delta_{20}\gamma_{10} + \delta_{10}\beta_{10}) + \alpha_1\pi(\delta_{21}\gamma_{11} + \delta_{11}\beta_{11})}{(1-\pi)[\alpha_0(\delta_{20}\gamma_{10} + \delta_{10}\beta_{10}) + (1-\alpha_0)(\delta_{20}\gamma_{00} + \delta_{10}\beta_{00})] + \pi[\alpha_1(\delta_{21}\gamma_{11} + \delta_{11}\beta_{11}) + (1-\alpha_1)(\delta_{21}\gamma_{01} + \delta_{11}\beta_{01})]} \end{aligned}$$

$$\begin{aligned}
p_{11}^A &= \frac{P(M = 1, T = 1 | V = 1)}{P(M = 1 | V = 1)} = \frac{P(M = 1, T = 1, V = 1)}{P(M = 1, V = 1)} \\
&= \frac{\delta_{20}\gamma_{10}\alpha_0(1 - \pi) + \delta_{21}\gamma_{11}\alpha_1\pi}{\alpha_0(1 - \pi)(\delta_{20}\gamma_{10} + \delta_{10}\beta_{10}) + \alpha_1\pi(\delta_{21}\gamma_{11} + \delta_{11}\beta_{11})} \\
p_{01}^A &= \frac{P(M = 1, T = 1 | V = 0)}{P(M = 1 | V = 0)} = \frac{P(M = 1, T = 1, V = 0)}{P(M = 1, V = 0)} \\
&= \frac{\delta_{20}\gamma_{00}(1 - \alpha_0)(1 - \pi) + \delta_{21}\gamma_{01}(1 - \alpha_1)\pi}{(1 - \alpha_0)(1 - \pi)(\delta_{20}\gamma_{00} + \delta_{10}\beta_{00}) + (1 - \alpha_1)\pi(\delta_{21}\gamma_{01} + \delta_{11}\beta_{01})}
\end{aligned}$$

In the CCD, the approximate standard error of  $VE_B$  is:

$$\begin{aligned}
SE(VE_B) &= SE(OR_B) \approx (OR_B) \times SE(\log(OR_A)) \\
&\approx \frac{p_{11}^B(1 - p_{10}^B)}{p_{10}^B(1 - p_{11}^B)} \sqrt{\frac{1}{N_{C1}^B p_{11}^B} + \frac{1}{N_{C1}^B(1 - p_{11}^B)} + \frac{1}{N_{C0}^B p_{10}^B} + \frac{1}{N_{C0}^B(1 - p_{10}^B)}}
\end{aligned}$$

where  $N_{C1}^B$  is the number of cases and  $N_{C0}^B$  is the number of controls. The probabilities ( $p_{11}^B, p_{10}^B$ ) can be written in terms of the parameters defined earlier:

$$\begin{aligned}
p_{11}^B &= P(V = 1 | C_B = 1, B = 1) = \frac{P(M = 1, T = 1, V = 1)}{P(M = 1, T = 1)} \\
&= \frac{\delta_{20}\gamma_{10}\alpha_0(1 - \pi) + \delta_{21}\gamma_{11}\alpha_1\pi}{\delta_{20}(1 - \pi)[\gamma_{00} + \alpha_0(\gamma_{10} - \gamma_{00})] + \delta_{21}\pi[\gamma_{01} + \alpha_1(\gamma_{11} - \gamma_{01})]}
\end{aligned}$$

$$\begin{aligned}
p_{10}^B &= P(V = 1 | C_B = 0, B = 1) = \frac{P(E = 0, V = 1)}{P(E = 0)} \\
&= \frac{(1 - \beta_{10} - \gamma_{10})\alpha_0(1 - \pi) + (1 - \beta_{11} - \gamma_{11})\alpha_1\pi}{(1 - \pi)[(1 - \beta_{00} - \gamma_{00}) + \alpha_0(\beta_{00} - \beta_{10} + \gamma_{00} - \gamma_{10})] + \pi[(1 - \beta_{01} - \gamma_{01}) + \alpha_1(\beta_{01} - \beta_{11} + \gamma_{01} - \gamma_{11})]}
\end{aligned}$$

## 4.2.6 Determining the values of the parameters

We distinguish between biological and non-biological parameters. The biological parameters are the probabilities of non-influenza and influenza ARIs in non-vaccinees and vaccinees, namely  $\beta_{0x}$ ,  $\beta_{1x}$ ,  $\gamma_{0x}$  and  $\gamma_{1x}$ . Baseline values for these parameters among healthy persons (i.e.  $\beta_{01}$ ,  $\beta_{11}$ ,  $\gamma_{01}$  and  $\gamma_{11}$ ) were estimated from randomized clinical trials (RCT) data using the numbers of influenza and non-influenza cases of ARI in vaccinees and non-vaccinees. More details about the determination of these values can be found in Chapter 3 and Haber et al. (2014).

The non-biological parameters include proportion of healthy persons, proportion of vaccinees and probability of seeking medical care for ARI. Baseline values for these parameters were determined from published literature. In this study, we assume the proportion of healthy persons is 0.7. According to the most recent publication (CDC 2013), influenza vaccine coverage in the US in the 2011-12 season ranged between 30% to 70%. The probability of seeking medical care for ARI has been estimated to be between 0.2 and 0.5 (Ferdinands and Shay 2012). We used 0.3 as the baseline value of this probability for influenza-infected ARI cases and 0.2 as the baseline value of this probability for non-influenza ARI cases. The sensitivity and specificity of the test for influenza infection were assumed to be 100%.

We also define a few probability ratios from the base parameters, to simplify the expressions for the bias of the VE estimates and to examine the impact of the latent health status on the VE estimates:  $\rho_{\beta_x} = \frac{\beta_{1x}}{\beta_{0x}}$  is the vaccine-related relative increase or decrease in the probability of non-influenza ARI;  $\eta_{\beta_v} = \frac{\beta_{v1}}{\beta_{v0}}$  is the health status-related relative increase or decrease in the probability of non-influenza ARI;  $\rho_{\gamma_x} = \frac{\gamma_{1x}}{\gamma_{0x}}$

is the vaccine-related relative increase or decrease in the probability of influenza ARI;  $\theta_\gamma = \frac{\rho_{\gamma 1}}{\rho_{\gamma 0}}$  measures the inequality in  $\rho_\gamma$  resulted from health status;  $\eta_{\delta 1} = \frac{\delta_{11}}{\delta_{10}}$  is the health status-related relative increase or decrease in the probability of seeking medical care for non-influenza ARI;  $\eta_{\delta 2} = \frac{\delta_{21}}{\delta_{20}}$  is the health status-related relative increase or decrease in the probability of seeking medical care for influenza ARI;  $\theta_\delta = \frac{\eta_{\delta 2}}{\eta_{\delta 1}}$  measures the inequality in  $\eta_\delta$  resulted from the type of infection leading to ARI.

A list of all the model parameters and their values is provided in the Table 4.2.

## 4.3 Results

The results presented below are based on the assumption that the influenza test has perfect sensitivity and specificity, i.e.  $\tau_1 = 0, \tau_2 = 1$ .

### 4.3.1 Analytic results

The analytic results are studied under an additional assumption: the probability of non-influenza ARI is the same for all persons regardless of their vaccine status and health status, i.e.  $\beta_{00} = \beta_{01} = \beta_{10} = \beta_{11}$ . We found that the conditions for TND estimates to be unbiased are different for SI and MI.

For SI, the TND estimate is unbiased if  $\theta_\delta = 1$  (i.e.  $\eta_{\delta 1} = \eta_{\delta 2}$ ). For MI, the TND estimate is unbiased if  $\eta_{\delta 1} = 1$ .

Table 4.2: List of parameters and other notation used in this article

Parameter	Definition	Values
$\pi$	Probability of having better health status (i.e. healthy persons)	0.7
$\alpha_x$	Probability of being vaccinated for a person of health status $x$	0.4-0.8
$\beta_{vx}$	Probability of non-influenza ARI for a person of vaccination status $v$ and health status $x$	From randomized clinical trials
$\rho_{\beta_x} = \frac{\beta_{1x}}{\beta_{0x}}$	Ratio comparing vaccinees and non-vaccinees w.r.t. probability of non-influenza ARI for health status $x$	0.25-4
$\eta_{\beta_v} = \frac{\beta_{v1}}{\beta_{v0}}$	Ratio comparing healthy and frail persons w.r.t. probability of non-influenza ARI for vaccination status $v$	0.25-4
$\gamma_{vx}$	Probability of influenza ARI for a person of vaccination status $v$ and health status $x$	From randomized clinical trials
$\rho_{\gamma_x} = \frac{\gamma_{1x}}{\gamma_{0x}}$	Ratio comparing vaccinees and non-vaccinees w.r.t. probability of influenza ARI for health status $x$	0-1
$\theta_\gamma = \frac{\rho_{\gamma 1}}{\rho_{\gamma 0}}$	Ratio of the 2 ratios defined above	0.25-2
$\eta_{\gamma v} = \frac{\gamma_{v1}}{\gamma_{v0}}$	Ratio comparing healthy and frail persons w.r.t. probability of influenza ARI for vaccination status $v$	0.25-4
$\delta_{ex}$	Probability of seeking medical care for ARI for a person of illness/infection status $e$ and health status $x$	0.2-0.5
$\eta_{\delta 1} = \frac{\delta_{11}}{\delta_{10}}$	Ratio comparing healthy and frail persons w.r.t. probability of seeking care for non-influenza ARI	0.25-4
$\eta_{\delta 2} = \frac{\delta_{21}}{\delta_{20}}$	Ratio comparing healthy and frail persons w.r.t. probability of seeking care for influenza ARI	0.25-4
$\theta_\delta = \frac{\eta_{\delta 2}}{\eta_{\delta 1}}$	Ratio of the 2 ratios defined above	0.25-4
$\tau_e$	Probability that a person of illness/infection status $e$ tests positive for influenza infection	$\tau_1 = 0, \tau_2 = 1$

### 4.3.2 Numeric results

Numerical results of the bias and standard errors of VE estimates are presented in Table 4.3 - Table 4.5 and Figure 4.1 - Figure 4.6. The bias was defined as the difference between the estimated and true VE. For example, if the true VE is 0.6 (60%) and the estimated VE is 0.56 (56%), bias is -0.04(-4%). The standard error was calculated for the VE estimates from TND and CCD respectively.

We considered three scenarios with different probabilities of vaccination for healthy and frail persons. In scenario 1, we consider random vaccination with  $\alpha_0 = \alpha_1 = 0.6$ . In scenario 2, we assume that healthy persons are twice more likely to get vaccinated than frail persons with  $\alpha_0 = 0.4, \alpha_1 = 0.8$ . In scenario 3, we assume that frail persons are twice more likely to get vaccinated than healthy persons with  $\alpha_0 = 0.8, \alpha_1 = 0.4$ .

For each scenario, we examine the effects of various deviations from the baseline assumptions on the bias and standard errors of the VE estimates. The baseline assumptions are listed below:

**A:** Probabilities of non-influenza ARI do not depend on vaccination status or health status, i.e.  $\beta_{00} = \beta_{10} = \beta_{01} = \beta_{11}$ . Here the baseline values are  $\beta_{00} = \beta_{10} = \beta_{01} = \beta_{11} = 0.1$ .

**B:** Probabilities of influenza ARI do not depend on health status, i.e.  $\eta_{\gamma_0} = 1, \eta_{\gamma_1} = 1$ . This implies that  $\rho_{\gamma_1} = \rho_{\gamma_0}$ , i.e. VE is the same for healthy and frail persons. The baseline values are  $\gamma_{00} = \gamma_{01} = 0.5, \gamma_{10} = \gamma_{11} = 0.2$ .

**C:** Probabilities of seeking medical care do not depend on health status, i.e.  $\eta_{\delta_1} = \eta_{\delta_2} = 1$ . The baseline values are  $\delta_{10} = \delta_{11} = 0.2, \delta_{20} = \delta_{21} = 0.3$ .



We first examine deviations from baseline assumption A, while assumptions B and C remain in place.

**A1:** Healthy and frail persons have equal probabilities of non-influenza ARI (i.e.  $\eta_\beta = 1$ ), but vaccinees and non-vaccinees have unequal probabilities of non-influenza ARI (i.e.  $\rho_\beta \neq 1$ ). We examine how vaccine-related ratio in the probability of non-influenza ARI (i.e.  $\rho_\beta$  ranging 0.25-4) affects the bias and standard error of VE estimates.

**A2:** Vaccinees and non-vaccinees have equal probabilities of non-influenza ARI (i.e.  $\rho_\beta = 1$ ), but healthy and frail persons have unequal probabilities of non-influenza ARI (i.e.  $\eta_\beta \neq 1$ ). We examine how health status-related ratio in the probability of non-influenza ARI (i.e.  $\eta_\beta$  ranging 0.25-4) affects the bias and standard error of VE estimates.

Next, we examine deviations from baseline assumption B, while assumptions A and C remain in place.

**B1:** Probabilities of influenza ARI may depend on health status but the ratio of these probabilities is the same in vaccinees and nonvaccinees, i.e.  $\eta_{\gamma_1} = \eta_{\gamma_0} = \eta_\gamma \neq 1$ . This still implies  $\rho_{\gamma_1} = \rho_{\gamma_0}$ , i.e. VE does not depend on health status. We examine how health status related ratio in the probability of influenza ARI (i.e.  $\eta_\gamma$  ranging 0.25-4) affects the bias and standard error of VE estimates.

**B2:** We now drop the assumption  $\eta_{\gamma_1} = \eta_{\gamma_0}$ . In this case,  $\rho_{\gamma_1} \neq \rho_{\gamma_0}$ , i.e. VE may depend on health status. We set the ratio of probabilities of influenza ARI comparing vaccinees and nonvaccinees for frail persons at 0.4 and examine how

the ratio  $\theta_\gamma = \frac{\rho_{\gamma 1}}{\rho_{\gamma 0}}$  ranging 0.25-2 affects the bias and standard error of VE estimates.

Finally, we examine deviations from baseline assumption C, while assumptions A and B remain in place.

**C1:** Probabilities of seeking medical care for ARI may depend on health status, but the ratio of these probabilities are the same for ARI patients with and without influenza infection (i.e.  $\theta_\delta = 1$  but  $\eta_\delta \neq 1$  or  $\eta_{\delta 1} = \eta_{\delta 2} = \eta_\delta \neq 1$ ). We examine how the common value of  $\eta_\delta$  ( $\eta_\delta$  ranging 0.25-4) affects the bias and standard error of VE estimates.

**C2:** Assuming non-homogeneity of the probability ratios of seeking medical care for ARI, (i.e.  $\theta_\delta \neq 1$  or  $\eta_{\delta 1} \neq \eta_{\delta 2}$ ). We examine how the inequality in health status-related ratio in the probability of seeking medical care for ARI (i.e.  $\eta_\delta$ ) by the type of infection leading to ARI (i.e.  $\theta_\delta = \frac{\eta_{\delta 2}}{\eta_{\delta 1}}$  ranging 0.25-4) affects the bias and standard error of VE estimates.

Figure 4.1 corresponds to deviations defined in A1. It presents the bias of VE estimates from TND and CCD for SI and MI with  $\rho_\beta$ , vaccine-related ratio in the probability of non-influenza ARI, ranging from 0.25 to 4 in Scenarios 1, 2 and 3. The bias plot is the same for all three scenarios. For both TND and CCD, the bias of VE estimates for SI and MI are the same. For TND, the bias is negative when  $\rho_\beta < 1$  and becomes positive when  $\rho_\beta > 1$ . The absolute bias of TND due to unequal probabilities of non-influenza ARI may be quite substantial, especially when  $\rho_\beta < 1$ , i.e. when vaccinated persons have a lower probability of non-influenza ARI compared

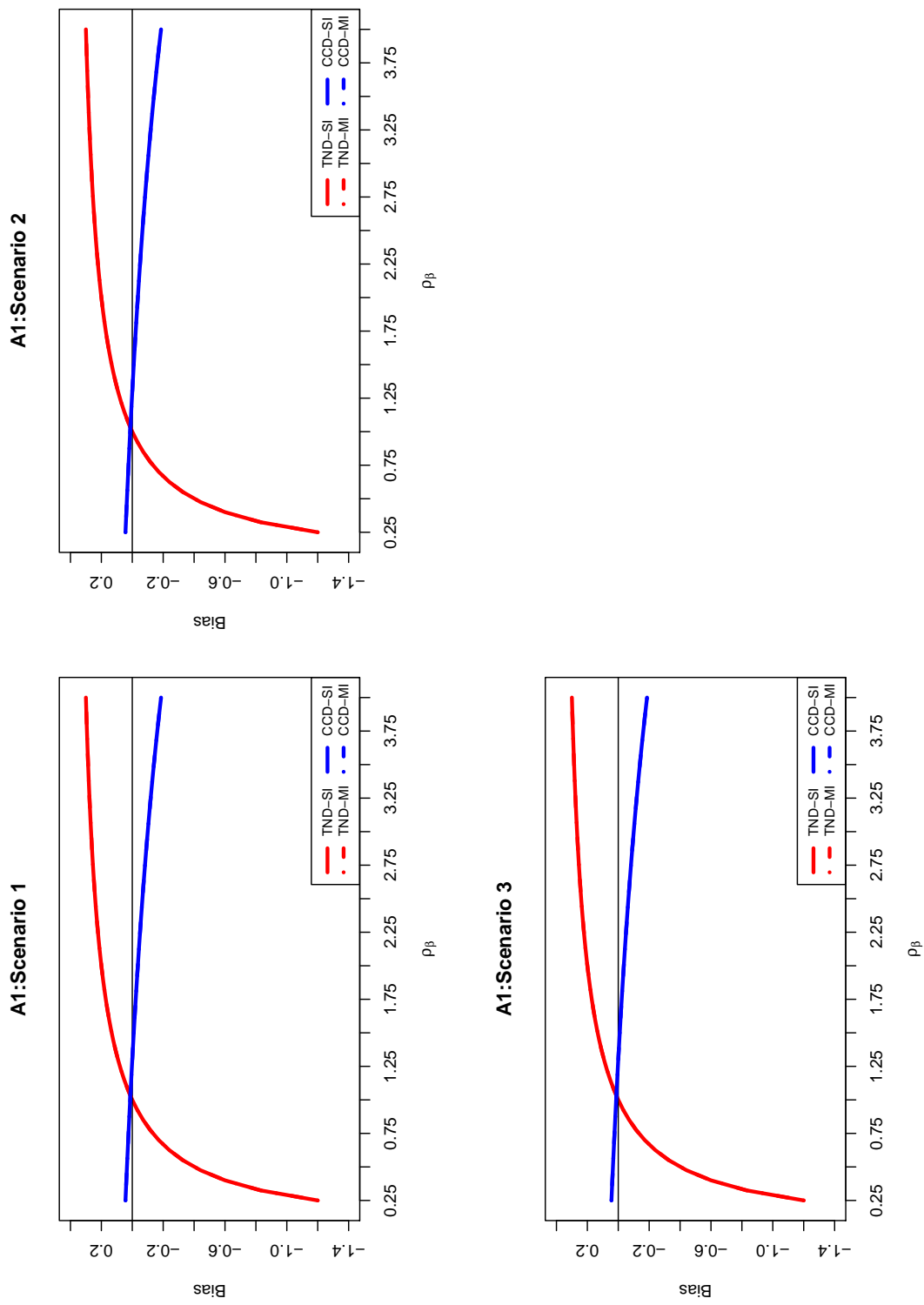


Figure 4.1: Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation A1 from baseline assumption A while assumptions B and C remain in place.  $\rho_\beta$  is the vaccine related ratio in the probability of non-influenza ARI, ranging from 0.25 to 4. The bias is the same for SI and MI.

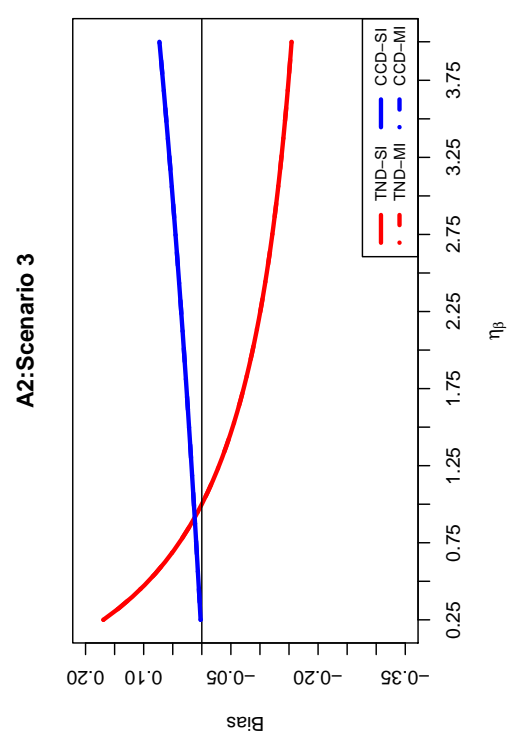
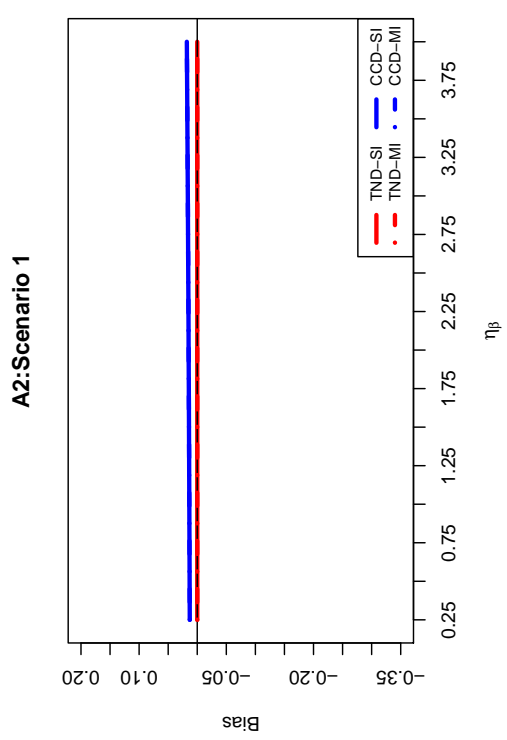
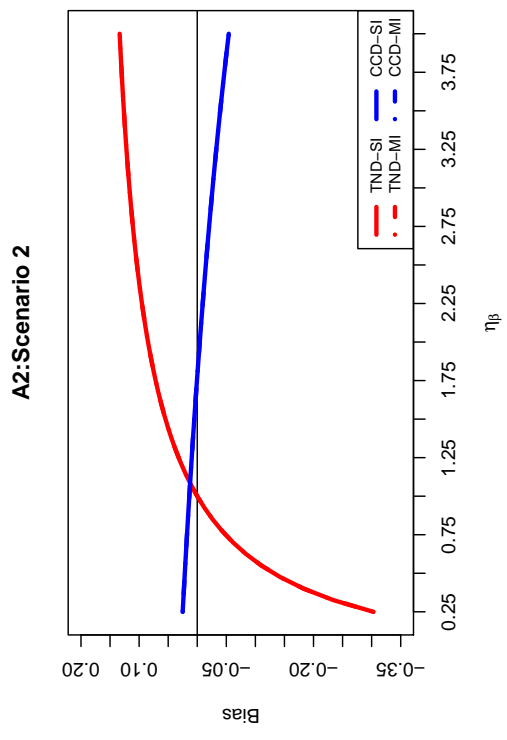


Figure 4.2: Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation A2 from baseline assumption A while assumptions B and C remain in place.  $\eta\beta$  is the health status related ratio in the probability of non-influenza ARI, ranging from 0.25 to 4. The bias is the same for SI and MI.

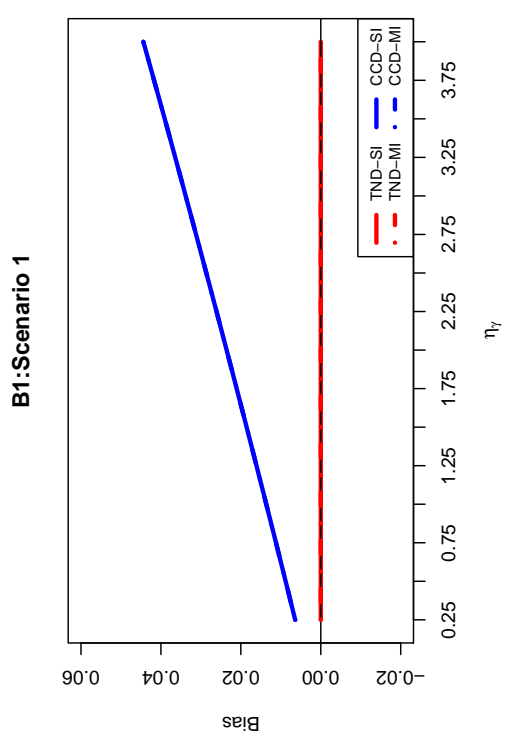
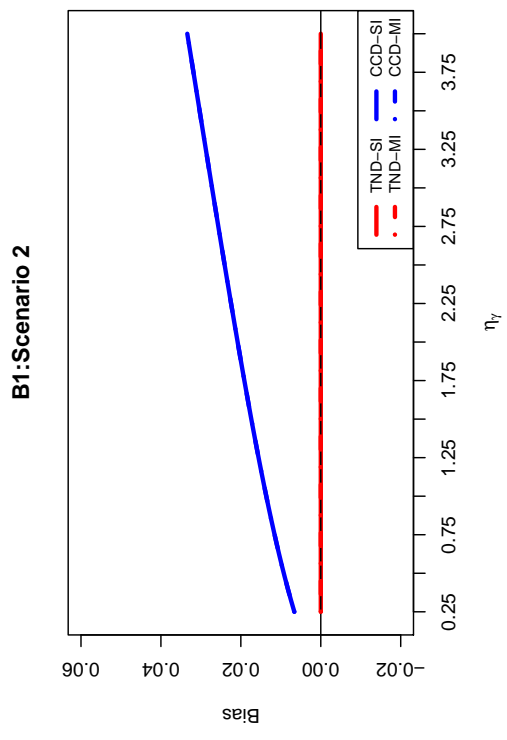


Figure 4.3: Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation B1 from baseline assumption B while assumptions A and C remain in place.  $\eta_T$  is the health status related ratio in the probability of influenza ARI, ranging from 0.25 to 4. The bias is the same for SI and MI.

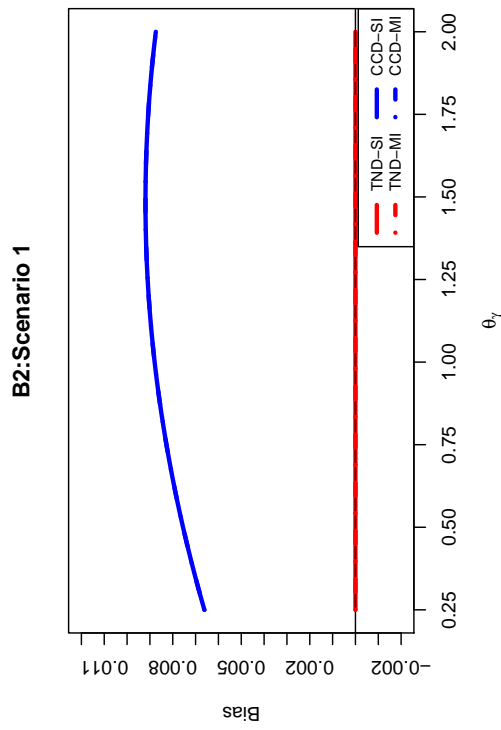
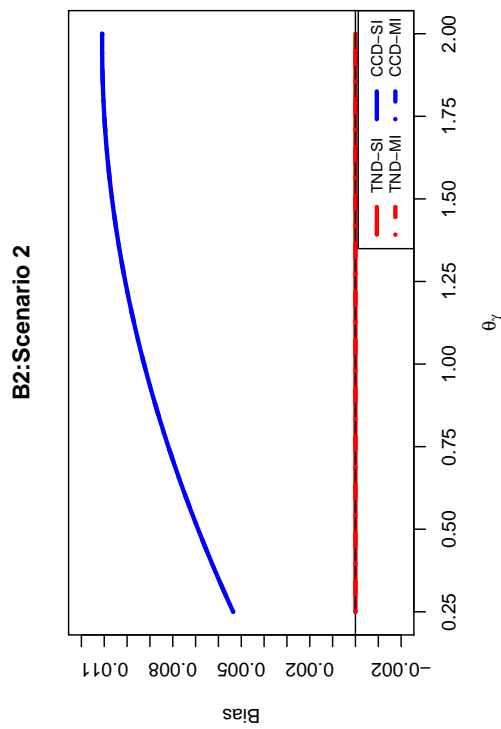


Figure 4.4: Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation B2 from baseline assumption B while assumptions A and C remain in place.  $\theta_\gamma$  is the inequality in vaccine related ratio in the probability of influenza ARI resulted from health status, ranging from 0.25 to 2. The bias is the same for SI and MI.

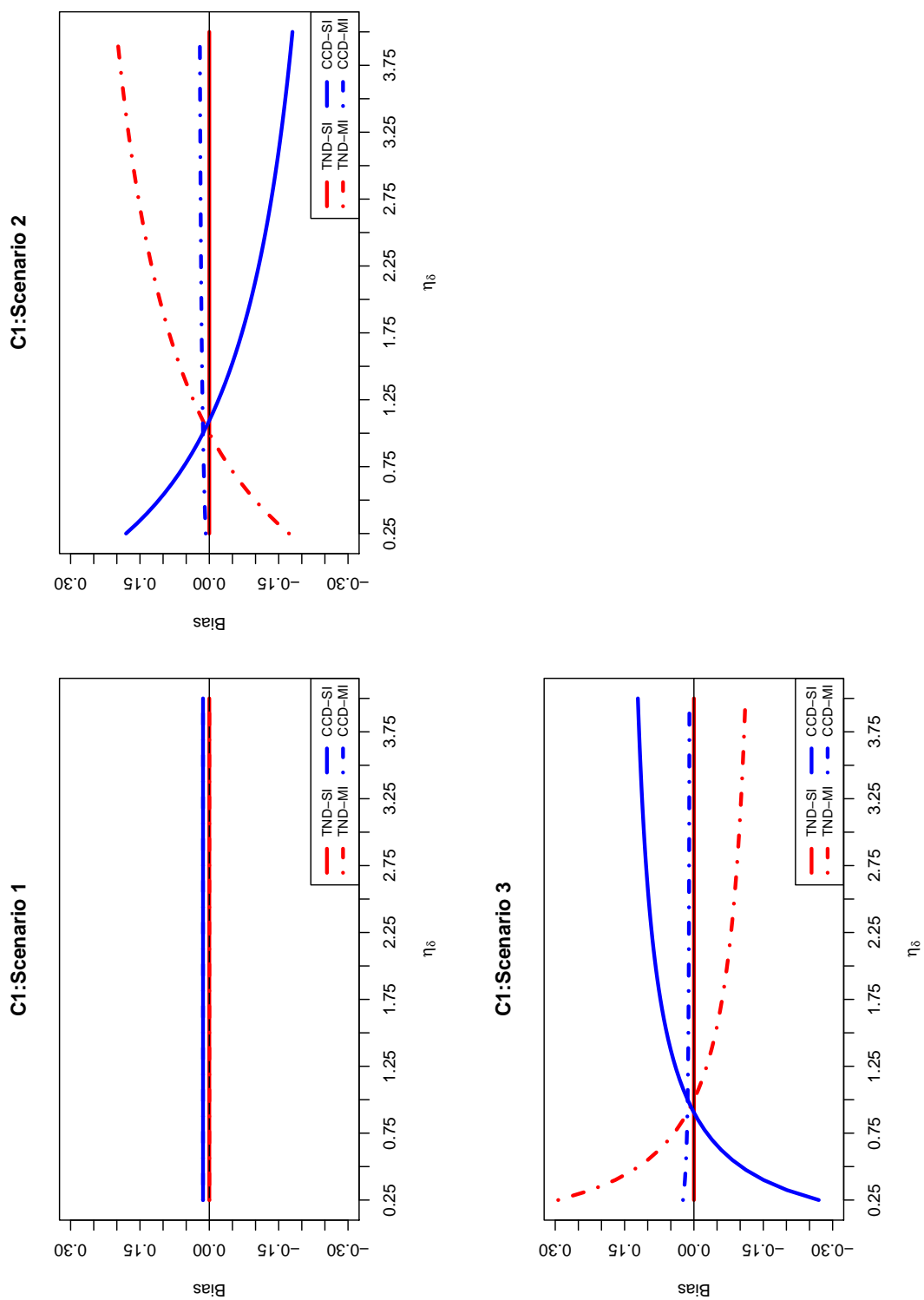


Figure 4.5: Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation C1 from baseline assumption C while assumptions A and B remain in place.  $\eta_\delta$  is the health status related ratio in the probability of seeking medical care for ARI, ranging from 0.25 to 4. The bias is the same for SI and MI in Scenario 1.

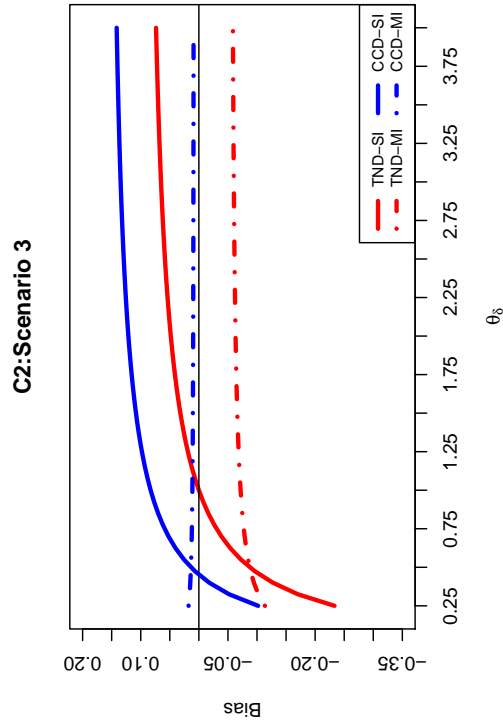
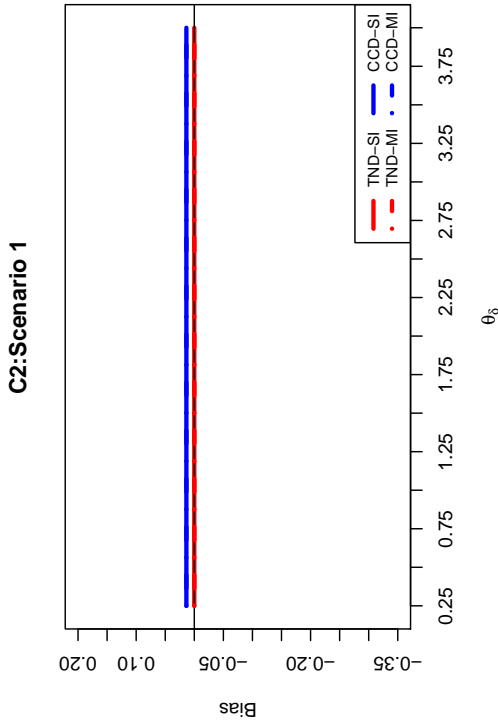
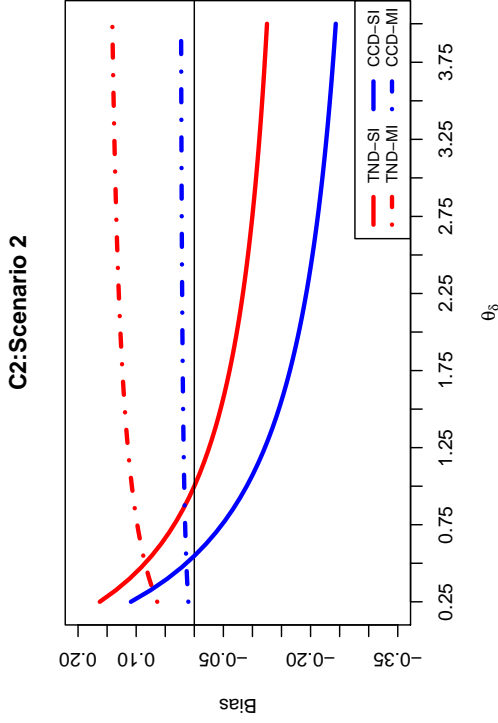


Figure 4.6: Bias of VE estimates for SI and MI from TND and CCD corresponding to effect of deviation C2 from baseline assumption C while assumptions A and B remain in place.  $\theta_6$  is the inequality in health status related ratio in the probability of seeking medical care for ARI resulted from the type of infection leading to ARI, ranging from 0.25 to 4. The bias is the same for SI and MI in Scenario 1.



Table 4.3: Bias and standard errors of VE estimates for SI and MI from TND and CCD for scenario 1 when vaccination is random ( $\alpha_0 = \alpha_1 = 0.6$ )

Effect	Parameter	Bias				Standard Error	
		TND-SI	CCD-SI	TND-MI	CCD-MI	TND	CCD
A1	$\rho_\beta = 0.25$	-1.200	0.044	-1.200	0.044	0.218	0.047
	$\rho_\beta = 0.5$	-0.400	0.034	-0.400	0.034	0.105	0.048
	$\rho_\beta = 1$	0.000	0.014	-0.000	0.014	0.056	0.050
	$\rho_\beta = 2$	0.200	-0.036	0.200	-0.036	0.032	0.056
	$\rho_\beta = 4$	0.300	-0.186	0.300	-0.186	0.020	0.075
A2	$\eta_\beta = 0.25$	-0.000	0.013	-0.000	0.013	0.052	0.050
	$\eta_\beta = 0.5$	0.000	0.013	-0.000	0.013	0.053	0.050
	$\eta_\beta = 1$	0.000	0.014	-0.000	0.014	0.056	0.050
	$\eta_\beta = 2$	0.000	0.015	-0.000	0.015	0.063	0.050
	$\eta_\beta = 4$	0.000	0.018	-0.000	0.018	0.077	0.050
B1	$\eta_\gamma = 0.25$	-0.000	0.006	0.000	0.006	0.067	0.051
	$\eta_\gamma = 0.5$	-0.000	0.009	0.000	0.009	0.061	0.051
	$\eta_\gamma = 1$	-0.000	0.014	-0.000	0.014	0.056	0.050
	$\eta_\gamma = 2$	-0.000	0.024	-0.000	0.024	0.052	0.049
	$\eta_\gamma = 4$	-0.000	0.044	0.000	0.044	0.053	0.046
B2	$\theta_\gamma = 0.25$	0.000	0.007	0.000	0.007	0.041	0.032
	$\theta_\gamma = 0.5$	0.000	0.008	0.000	0.008	0.048	0.038
	$\theta_\gamma = 1$	0.000	0.009	0.000	0.009	0.061	0.051
	$\theta_\gamma = 2$	0.000	0.009	0.000	0.009	0.088	0.078
C1	$\eta_\delta = 0.25$	0.000	0.014	-0.000	0.014	0.030	0.050
	$\eta_\delta = 0.5$	-0.000	0.014	-0.000	0.014	0.042	0.050
	$\eta_\delta = 1$	0.000	0.014	-0.000	0.014	0.056	0.050
	$\eta_\delta = 2$	0.000	0.014	-0.000	0.014	0.067	0.050
	$\eta_\delta = 4$	-0.000	0.014	-0.000	0.014	0.075	0.050
C2	$\theta_\delta = 0.25$	0.000	0.014	-0.000	0.014	0.058	0.050
	$\theta_\delta = 0.5$	0.000	0.014	-0.000	0.014	0.063	0.050
	$\theta_\delta = 1$	0.000	0.014	-0.000	0.014	0.067	0.050
	$\theta_\delta = 2$	-0.000	0.014	-0.000	0.014	0.078	0.050
	$\theta_\delta = 4$	0.000	0.014	0.000	0.014	0.120	0.050

Note: See text for definition of effects of deviations from baseline assumptions. Shaded rows correspond to baseline.

Table 4.4: Bias and standard errors of VE estimates for SI and MI from TND and CCD for scenario 2 when healthy persons are twice more likely to get vaccination than frail persons ( $\alpha_0 = 0.4, \alpha_1 = 0.8$ )

Effect	Parameter	Bias				Standard Error	
		TND-SI	CCD-SI	TND-MI	CCD-MI	TND	CCD
A1	$\rho_\beta = 0.25$	-1.200	0.044	-1.200	0.044	0.208	0.047
	$\rho_\beta = 0.5$	-0.400	0.034	-0.400	0.034	0.103	0.048
	$\rho_\beta = 1$	0.000	0.014	-0.000	0.014	0.056	0.051
	$\rho_\beta = 2$	0.200	-0.036	0.200	-0.036	0.033	0.057
	$\rho_\beta = 4$	0.300	-0.186	0.300	-0.186	0.021	0.075
A2	$\eta_\beta = 0.25$	-0.303	0.025	-0.303	0.025	0.089	0.050
	$\eta_\beta = 0.5$	-0.131	0.021	-0.131	0.021	0.069	0.050
	$\eta_\beta = 1$	0.000	0.014	-0.000	0.014	0.056	0.051
	$\eta_\beta = 2$	0.085	-0.004	0.085	-0.004	0.051	0.053
	$\eta_\beta = 4$	0.133	-0.054	0.133	-0.054	0.052	0.059
B1	$\eta_\gamma = 0.25$	-0.000	0.007	-0.000	0.007	0.039	0.030
	$\eta_\gamma = 0.5$	0.000	0.009	0.000	0.009	0.047	0.039
	$\eta_\gamma = 1$	0.000	0.014	0.000	0.014	0.056	0.051
	$\eta_\gamma = 2$	-0.000	0.021	0.000	0.021	0.067	0.064
	$\eta_\gamma = 4$	-0.000	0.033	-0.000	0.033	0.080	0.075
B2	$\theta_\gamma = 0.25$	0.000	0.005	0.000	0.005	0.026	0.020
	$\theta_\gamma = 0.5$	0.000	0.007	0.000	0.007	0.033	0.026
	$\theta_\gamma = 1$	0.000	0.009	0.000	0.009	0.047	0.039
	$\theta_\gamma = 2$	0.000	0.011	0.000	0.011	0.074	0.066
C1	$\eta_\delta = 0.25$	0.000	0.180	-0.172	0.008	0.037	0.030
	$\eta_\delta = 0.5$	0.000	0.109	-0.099	0.010	0.047	0.039
	$\eta_\delta = 1$	0.000	0.014	-0.000	0.014	0.056	0.051
	$\eta_\delta = 2$	0.000	-0.090	0.107	0.017	0.064	0.065
	$\eta_\delta = 4$	0.000	-0.180	0.200	0.020	0.070	0.077
C2	$\theta_\delta = 0.25$	0.163	0.109	0.064	0.010	0.038	0.039
	$\theta_\delta = 0.5$	0.085	0.014	0.085	0.014	0.051	0.051
	$\theta_\delta = 1$	0.000	-0.090	0.107	0.017	0.064	0.065
	$\theta_\delta = 2$	-0.073	-0.180	0.127	0.020	0.080	0.077
	$\theta_\delta = 4$	-0.125	-0.243	0.141	0.023	0.109	0.085

Note: See text for definition of effects of deviations from baseline assumptions. Shaded rows correspond to baseline.

Table 4.5: Bias and standard errors of VE estimates for SI and MI from TND and CCD for scenario 3 when frail persons are twice more likely to get vaccination than healthy persons ( $\alpha_0 = 0.8, \alpha_1 = 0.4$ )

Effect	Parameter	Bias				Standard Error	
		TND-SI	CCD-SI	TND-MI	CCD-MI	TND	CCD
A1	$\rho_\beta = 0.25$	-1.200	0.044	-1.200	0.044	0.234	0.047
	$\rho_\beta = 0.5$	-0.400	0.034	-0.400	0.034	0.110	0.048
	$\rho_\beta = 1$	-0.000	0.014	0.000	0.014	0.057	0.051
	$\rho_\beta = 2$	0.200	-0.036	0.200	-0.036	0.032	0.058
	$\rho_\beta = 4$	0.300	-0.186	0.300	-0.186	0.019	0.078
A2	$\eta_\beta = 0.25$	0.169	0.002	0.169	0.002	0.031	0.053
	$\eta_\beta = 0.5$	0.092	0.006	0.092	0.006	0.041	0.052
	$\eta_\beta = 1$	-0.000	0.014	0.000	0.014	0.057	0.051
	$\eta_\beta = 2$	-0.087	0.031	-0.087	0.031	0.078	0.049
	$\eta_\beta = 4$	-0.154	0.073	-0.154	0.073	0.108	0.043
B1	$\eta_\gamma = 0.25$	-0.000	0.004	0.000	0.004	0.116	0.088
	$\eta_\gamma = 0.5$	0.000	0.008	0.000	0.008	0.079	0.066
	$\eta_\gamma = 1$	0.000	0.014	0.000	0.014	0.057	0.051
	$\eta_\gamma = 2$	-0.000	0.024	0.000	0.024	0.045	0.041
	$\eta_\gamma = 4$	-0.000	0.044	0.000	0.044	0.040	0.034
B2	$\theta_\gamma = 0.25$	0.000	0.007	0.000	0.007	0.062	0.049
	$\theta_\gamma = 0.5$	0.000	0.008	0.000	0.008	0.068	0.055
	$\theta_\gamma = 1$	0.000	0.008	0.000	0.008	0.079	0.066
	$\theta_\gamma = 2$	0.000	0.007	0.000	0.007	0.103	0.090
C1	$\eta_\delta = 0.25$	-0.000	-0.270	0.294	0.024	0.025	0.085
	$\eta_\delta = 0.5$	-0.000	-0.102	0.120	0.018	0.038	0.065
	$\eta_\delta = 1$	-0.000	0.014	0.000	0.014	0.057	0.051
	$\eta_\delta = 2$	-0.000	0.083	-0.072	0.011	0.078	0.043
	$\eta_\delta = 4$	-0.000	0.121	-0.111	0.010	0.099	0.038
C2	$\theta_\delta = 0.25$	-0.233	-0.102	-0.114	0.018	0.083	0.065
	$\theta_\delta = 0.5$	-0.087	0.014	-0.087	0.014	0.078	0.051
	$\theta_\delta = 1$	-0.000	0.083	-0.072	0.011	0.078	0.043
	$\theta_\delta = 2$	0.048	0.121	-0.063	0.010	0.098	0.038
	$\theta_\delta = 4$	0.074	0.141	-0.059	0.009	0.250	0.036

Note: See text for definition of effects of deviations from baseline assumptions. Shaded rows correspond to baseline.

to unvaccinated persons. For CCD, the bias is positive when  $\rho_\beta < 1$  and becomes negative when  $\rho_\beta > 1$ . The effect of departure of  $\rho_\beta$  from one on the bias of CCD estimates is much smaller than the effect on the bias from TND estimates, and the impact of  $\rho_\beta$  does not depend on the proportion of vaccinated persons.

Figure 4.2 corresponds to the deviations defined in A2. It presents the bias of VE estimates for SI and MI from TND and CCD with  $\eta_\beta$ , health status-related ratio in the probability of non-influenza ARI, ranging from 0.25 to 4 in Scenarios 1, 2 and 3. For all three scenarios, the VE estimates for SI and MI are the same for both TND and CCD. In Scenario 1 (random vaccination), the TND estimate remains unbiased, while the bias of CCD is positive, small and slightly increases with the increase in  $\eta_\beta$ . In Scenario 2, the bias of TND is negative when  $\eta_\beta < 1$  and becomes positive when  $\eta_\beta > 1$ , while the bias of CCD is positive when  $\eta_\beta < 1$  and becomes negative when  $\eta_\beta > 1$ . The absolute bias of TND is more substantial than that of CCD. In Scenario 3, the bias of TND is positive when  $\eta_\beta < 1$  and becomes negative when  $\eta_\beta > 1$ . The bias of CCD remains positive and increases as  $\eta_\beta$  increases. The absolute bias of TND is again more substantial than that of CCD. In summary, when vaccination is random, TND is always unbiased and the bias of CCD is also close to 0. However, when vaccination is not random, the effect of changing  $\eta_\beta$  on the bias from CCD is much smaller than that from TND.

Figure 4.3 corresponds to the deviations defined in B1. It presents the bias of VE estimates for SI and MI from TND and CCD with  $\eta_\gamma$ , health status-related ratio in the probability of influenza ARI, ranging from 0.25 to 4 in Scenarios 1, 2 and 3. For all three scenarios, the VE estimates for both TND and CCD studies are the same for SI and MI. The bias of TND remains 0 while the bias of CCD is positive

and increases as  $\eta_\gamma$  increases. The bias of CCD are the same for Scenario 1 and 3, i.e., when healthy persons have equal or smaller probabilities of vaccination than frail persons. Compared to Scenarios 1 and 3, the bias of CCD are somewhat smaller in Scenario 2 when healthy persons have higher probabilities of vaccination than frail persons. The health status-related ratio in the probability of influenza ARI has no impact on the VE estimates from TND, which is always unbiased, but a positive impact on the bias of VE estimates from CCD. Although the VE estimates from TND is always unbiased, the standard errors of the TND estimates are slightly bigger than those of the CCD estimates (Table 4.4).

Figure 4.4 corresponds to the deviations defined in B2. It presents the bias of VE estimates for SI and MI from TND and CCD with  $\theta_\gamma$ , inequality in vaccine related ratio in the probability of influenza ARI resulted from health status, ranging from 0.25 to 2 in Scenarios 1, 2 and 3. For all three scenarios, the VE estimates for SI and MI are the same for both TND and CCD, and the bias of TND remains 0 while the bias of CCD is positive, close to 0 and changes slightly as  $\theta_\gamma$  increases. The bias of CCD differ slightly across the three scenarios. The inequality in vaccine-related ratio in the probability of influenza ARI resulting from health status has no impact on VE estimates from TND, which is always unbiased, and a small impact on CCD VE estimates. Although the VE estimates from TND is always unbiased, the standard errors of the TND estimates are slightly bigger than those of the CCD estimates (Table 4.3-4.5).

Figure 4.5 corresponds to the deviations defined in C1. It presents the bias of VE estimates for SI and MI from TND and CCD as functions of  $\eta_\delta$ , health status-related ratio in the probability of seeking medical care for ARI, ranging from 0.25 to 4 in

Scenarios 1, 2 and 3. The TND-SI estimate remains unbiased in all three scenarios. The TND-MI estimate is only unbiased in Scenario 1 or when  $\eta_\delta = 1$  in Scenarios 2 and 3. Otherwise, the bias of TND-MI is more substantial. The bias of CCD-MI estimates is small and positive in all three scenarios. The bias of CCD-SI estimates is the same as that of CCD-MI in Scenario 1, but more substantial than that of CCD-MI in Scenarios 2 and 3. With the increase of  $\eta_\delta$ , the bias of TND-MI change from negative to positive in Scenario 2 and vice versa in Scenario 3; and the bias of CCD-SI change from positive to negative in Scenario 2 and vice versa in Scenario 3. The absolute bias of TND-MI and CCD-SI is about the same for the same value of  $\eta_\delta$ . The health status-related ratio in the probability of seeking medical care for ARI has no impact on the bias of TND-SI estimate and little impact on the bias of CCD-MI estimate. When the vaccination is random, there is no effect of  $\eta_\delta$  on the bias of VE estimates. However, when the vaccination is not random, departure of  $\eta_\delta$  from one increases the biases of TND-MI estimate and the CCD-SI estimate.

Figure 4.6 corresponds to the deviations defined in C2. It presents the bias of VE estimates for SI and MI from TND and CCD with  $\theta_\delta$ , inequality in health status-related ratio in the probability of seeking medical care for ARI resulted from the type of infection leading to ARI, ranging from 0.25 to 4 in Scenarios 1, 2 and 3. In Scenario 1, the TND-SI and TND-MI are both unbiased and the bias of CCD is a positive and close to 0. In both Scenarios 2 and 3, TND-SI is only unbiased when  $\theta_\delta = 1$  and the bias of CCD-MI are small and positive. The biases of both TND-SI and CCD-SI estimates change from positive to negative in Scenario 2 and from negative to positive in Scenario 3 with the increase of  $\theta_\delta$  and the absolute bias of TND-SI is bigger than that of CCD-SI when  $\theta_\delta < 1$  but less when  $\theta_\delta > 1$ . The inequality in health status-

related ratio in the probability of seeking medical care for ARI resulted from the type of infection leading to ARI has no impact on the bias of VE estimates when there is random vaccination. When the vaccination probability depends on health status, the effect of changing  $\theta_\delta$  on the bias of CCD-MI is minimal. The departure of  $\theta_\delta$  from one increases the bias of TND-SI and CCD-SI estimates. The impact of changing  $\theta_\delta$  on the bias of CCD-MI depends on the vaccination status. The TND estimates are less precise than CCD estimates. For SI, the absolute bias of CCD is smaller than that of TND when  $\theta_\delta < 1$ , and the reverse is true when  $\theta_\delta > 1$ . For MI, CCD has smaller bias than TND for all values of  $\theta$ .

### 4.3.3 Summary of sources of bias under non-random vaccination

Table 4.6 lists all the situations where severe bias occurs under non-random vaccination scenarios (i.e. scenarios 2 and 3) for each study type. We define severe bias as the absolute value of bias greater than 0.1 (i.e.  $|bias| > 0.1$ ).

## 4.4 Discussion

In this study, we developed a step-wise latent-class model that generates data for an observational study aimed at estimating VE against medically attended influenza and symptomatic influenza. The process is composed of five steps: health status, vaccination, developing of infection and illness, seeking medical care for ARI and testing for influenza infection. The bias and standard error of VE estimates for both

Table 4.6: Situations for severe bias ( $|bias| > 0.1$ ) under non-random vaccination

Source of bias/Scenario	Study type	Situations for severe bias
<b><math>A_1</math> (<math>\rho_\beta \neq 1</math>)</b>		
2	TND	$\rho_\beta < 0.8$ or $\rho_\beta > 1.4$
	CCD	$\rho_\beta > 3.0$
3	TND	$\rho_\beta < 0.8$ or $\rho_\beta > 1.4$
	CCD	$\rho_\beta > 3.0$
<b><math>A_2</math> (<math>\eta_\beta \neq 1</math>)</b>		
2	TND	$\eta_\beta < 0.6$ or $\eta_\beta > 2.4$
	CCD	No severe bias for $0.25 \leq \eta_\beta \leq 4.0$
3	TND	$\eta_\beta < 0.5$ or $\eta_\beta > 2.3$
	CCD	No severe bias for $0.25 \leq \eta_\beta \leq 4.0$
<b><math>B_1</math> (<math>\eta_\gamma \neq 1</math>)</b>		
2	TND	Unbiased regardless of $\eta_\gamma$
	CCD	No severe bias regardless $0.25 \leq \eta_\gamma \leq 4$
3	TND	Unbiased regardless of $\eta_\gamma$
	CCD	No severe bias regardless $0.25 \leq \eta_\gamma \leq 4$
<b><math>B_2</math> (<math>\theta_\gamma \neq 1</math>)</b>		
2	TND	Unbiased regardless of $\theta_\gamma$
	CCD	No severe bias regardless $0.25 \leq \theta_\gamma \leq 2$
3	TND	Unbiased regardless of $\eta_\gamma$
	CCD	No severe bias regardless $0.25 \leq \theta_\gamma \leq 2$
<b><math>C_1</math> (<math>\eta_\delta \neq 1</math>)</b>		
2	TND-SI	Unbiased regardless of $\eta_\delta$
	TND-MI	$\eta_\delta < 0.5$ or $\eta_\delta > 1.9$
	CCD-SI	$\eta_\delta < 0.5$ or $\eta_\delta > 2.2$
	CCD-MI	No severe bias regardless of $\eta_\delta$
3	TND-SI	Unbiased regardless of $\eta_\delta$
	TND-MI	$\eta_\delta < 0.5$ or $\eta_\delta > 3.2$
	CCD-SI	$\eta_\delta < 0.5$ or $\eta_\delta > 2.6$
	CCD-MI	No severe bias regardless of $\eta_\delta$
<b><math>C_2</math> (<math>\theta_\delta \neq 1</math>)</b>		
2	TND-SI	$\theta_\delta < 0.4$ or $\theta_\delta > 2.8$
	TND-MI	$\theta_\delta > 0.8$
	CCD-SI	$\theta_\delta > 1.1$
	CCD-MI	No severe bias regardless of $\theta_\delta$
3	TND-SI	$\theta_\delta < 0.4$
	TND-MI	$\theta_\delta < 0.3$
	CCD-SI	$\theta_\delta > 1.3$
	CCD-MI	No severe bias regardless of $\theta_\delta$



MI and SI based on ordinary case-control study and on test-negative study can be written in terms of the model parameters. This model facilitates the evaluation and comparison of the accuracy and precision of two types of estimates from the two study designs.

Several models and methods for evaluating the bias of influenza VE estimates from TND studies have been proposed in the past (De Serres et al. 2013, Ferdinands and Shay 2012, Foppa et al. 2013, Haber et al. 2014, Jackson and Nelson 2013, Orenstein et al. 2007). Compared to the previous studies, the current model has the following advantages: (1) we studied the impact of non-random vaccination; (2) we assumed that the population may be nonhomogeneous with respect to the parameters determining probabilities of vaccination, infection, illness, and seeking medical care; (3) we considered two outcomes of interest: SI and MI; and (4) we evaluated the precision of the VE estimates.

Under the basic assumption that the test has perfect sensitivity and specificity, we tried to find conditions under which the TND estimates are unbiased. Consider the following 5 conditions: (1) probability of vaccination does not depend on health status; (2) probability of non-influenza ARI does not depend on health status ( $\eta_\beta = 1$ ); (3) probability of non-influenza ARI does not depend on vaccination status ( $\rho_\beta = 1$ ); (4) probability of seeking medical care against non-influenza ARI does not depend on health status ( $\eta_{\delta 1} = 1$ ); and (5) the ratio comparing healthy and frail persons of proportion of seeking medical care is the same for influenza and non-influenza ARI ( $\theta_\delta = 1$ ).

We found that:

- if condition (1) is satisfied, then only condition (3) is needed for unbiasedness of the TND estimates for both SI and MI;
- if condition (1) is not satisfied, then TND-MI will be unbiased when (2), (3) and (4) are met and TND-SI will be unbiased when (2), (3) and (5) are met.

When condition (2), which is the rationale for the TND (De Serres et al. 2013), is violated, the TND-based estimate is always biased (even when vaccination is random) and it can be severely biased. The CCD-based estimates are also biased. Under this situation, the absolute bias of CCD-based estimate is usually much smaller than that of the TND-based estimates. In addition to the non-homogeneity of the probability of non-influenza ARI, factors related to the probability of influenza ARI, probability of seeking medical care for ARI and vaccination status affect the bias of CCD-based estimates. The probability of seeking medical care for ARI affects the bias of CCD-SI more than that of CCD-MI. The CCD-based estimates may be unbiased under an unrealistic situation that vaccine does not affect the probability of having ARI.

This study has a few limitations: (1) we assume the diagnostic test has perfect sensitivity and specificity; (2) we assume the probability of seeking medical care does not depend on vaccination status in persons of the same health status; (3) we assume that the vaccination status is determined without an error; and (4) all the parameters in our model remain unchanged throughout the influenza season. We could eliminate the first three limitations by including additional parameters in the model, but it would be very difficult to determine the values or reasonable ranges of these parameters and make the interpretation of results more difficult. Addressing

limitation (4) would involve assumptions on the infection contact and transmission processes and on the temporal trends in the values of the parameters. Future work that involve modeling the influenza infection transmission process and additional real-life factors will be developed to address this limitation.

# Appendix

## Appendix 1.1: Proofs for Proposition 2.2.1

Note that  $X \sim \text{Beta-Laplace}(\mu, a, b, c)$ .

1. The density function of  $Y = 1 - X$  is

$$f(y) = \pi(1 - x; \mu, a, b, c) \propto (1 - x)^{a-1} x^{b-1} \exp(-c|1 - x - \mu|),$$

which implies that  $1 - X \sim \text{Beta-Laplace}(1 - \mu, b, a, c)$ .

- 2.

$$\mathbb{E}(X^n) = \int_0^1 \frac{1}{L(\mu, a, b, c)} x^{n+a-1} (1-x)^{b-1} \exp(-c|x-\mu|) dx = \frac{L(\mu, a+n, b, c)}{L(\mu, a, b, c)}$$

3. Given  $a, b > 0$ , when  $c \rightarrow \infty$ ,  $\pi(x; a, b, c) \rightarrow I[x = \mu]$ , thus,  $\mathbb{E}(X) \rightarrow \mu$ .

## Appendix 1.2: Derivations for marginal distributions for $A_{it}$ and $H_{it}$

Noting that if  $(X_1, X_2, X_3) \sim \text{Multinomial}(p_1, p_2, p_3, N)$  and  $N \sim \text{Poisson}(\lambda)$ , then the joint distribution of  $(X_1, X_2, X_3)$  and  $N$  is given by

$$P(X_1, X_2, X_3, N) = \frac{N!}{X_1!X_2!X_3!} p_1^{X_1} p_2^{X_2} p_3^{X_3} \frac{\lambda^N}{N!} e^{-\lambda}.$$

Integrating out  $N$  from the above joint probability density, we have

$$\begin{aligned} P(X_1, X_2) &= \sum_{N=X_1+X_2}^{\infty} \frac{p_1^{X_1} p_2^{X_2} (1-p_1-p_2)^{N-X_1-X_2}}{X_1!X_2!(N-X_1-X_2)!} \lambda^{N-X_1-X_2} e^{-\lambda} \lambda^{X_1+X_2} \\ &= \sum_{N=X_1+X_2}^{\infty} \frac{[(1-p_1-p_2)\lambda]^{N-X_1-X_2}}{(N-X_1-X_2)!} e^{-\lambda(1-p_1-p_2)} e^{(p_1+p_2)\lambda} \frac{p_1^{X_1} p_2^{X_2}}{X_1!X_2!} \lambda^{X_1+X_2} \\ &= \frac{(p_1\lambda)^{X_1}}{X_1!} e^{-p_1\lambda} \frac{(p_2\lambda)^{X_2}}{X_2!} e^{-p_2\lambda} \end{aligned}$$

This implies that  $X_1$  and  $X_2$  follow independent Poisson distributions. In a similar fashion of the derivation, it is straightforward to show that if

$$(A_{ii}, \dots, A_{iT}, H_{ii}, \dots, H_{it}, U_{iT}) \sim \text{Multinomial}(q_{ii}^A, \dots, q_{iT}^A, q_{ii}^H, \dots, q_{iT}^H, q_{iT}^U, N_i)$$

and

$$N_i \sim \text{Poisson}(\lambda_i),$$

Then

$$A_{it} \sim \text{Poisson}(q_{it}^A \lambda_i), \text{ and } H_{it} \sim \text{Poisson}(q_{it}^H \lambda_i), t = i, \dots, T,$$

where  $A_{it}$  and  $H_{it}$  are mutually independent.

## Appendix 1.3: Full conditional distributions

- Full conditional distribution of  $\mathbf{p}^H$

$$\begin{aligned} & \pi(\mathbf{p}^H \mid \boldsymbol{\lambda}, \mathbf{A}, \mathbf{H}) \\ & \propto \prod_{t=1}^T \frac{e^{-\sum_{i=1}^t q_{it}^A \lambda_i}}{A_t!} \left( \sum_{i=1}^t q_{it}^A \lambda_i \right)^{A_t} \frac{e^{-\sum_{i=1}^t q_{it}^H \lambda_i}}{H_t!} \left( \sum_{i=1}^t q_{it}^H \lambda_i \right)^{H_t} \\ & \quad \times \prod_{t=1}^T \frac{(p_t^H)^{a_2-1} (1-p_t^H)^{b_2-1}}{B(a_2, b_2)} \times \prod_{t=2}^T \exp\{-c p_t^H | p_t^H - p_{t-1}^H|\} \end{aligned}$$

- Full conditional distribution of  $\boldsymbol{\lambda}$

$$\begin{aligned} & \pi(\boldsymbol{\lambda} \mid \mathbf{p}^H, \mathbf{A}, \mathbf{H}) \\ & \propto \prod_{t=1}^T \frac{e^{-\sum_{i=1}^t q_{it}^A \lambda_i}}{A_t!} \left( \sum_{i=1}^t q_{it}^A \lambda_i \right)^{A_t} \frac{e^{-\sum_{i=1}^t q_{it}^H \lambda_i}}{H_t!} \left( \sum_{i=1}^t q_{it}^H \lambda_i \right)^{H_t} \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda_t^{a_0-1} \exp(-b_0 \lambda_t) \\ & \quad \times \prod_{t=2}^T \exp\{-a |\lambda_t - \lambda_{t-1}|\} \end{aligned}$$

## Appendix 2.1: Details for the calculation of the probabilities in Chapter 3

$$\begin{aligned}
P(M = 1|V = v) &= \sum_e P(M = 1, E = e|V = v) \\
&= \sum_e P(M = 1|E = e, V = v)P(E = e|V = v) \\
&= 0 + \sum_{e=1,2} P(M = 1|E = e, V = v)P(E = e|V = v) \\
&= \delta_{1v}\beta_{1v} + \delta_{2v}\beta_{2v}
\end{aligned}$$

$$\begin{aligned}
P(V = v|M = 1) &= \frac{P(M = 1|V = v)P(V = v)}{P(M = 1)} \\
&= \frac{P(M = 1|V = v)P(V = v)}{\sum_{v=0,1} P(M = 1|V = v)P(V = v)} \\
&= \frac{\alpha_v(\delta_{1v}\beta_{1v} + \delta_{2v}\beta_{2v})}{\alpha(\delta_{11}\beta_{11} + \delta_{21}\beta_{21}) + (1 - \alpha)(\delta_{10}\beta_{10} + \delta_{20}\beta_{20})}
\end{aligned}$$

$$\begin{aligned}
P(E = e|M = 1, V = v) &= \frac{P(E = e, M = 1|V = v)}{P(M = 1|V = v)} \\
&= \frac{P(M = 1|E = e, V = v)P(E = e|V = v)}{P(M = 1|V = v)} \\
&= \frac{\delta_{ev}\beta_{ev}}{P(M = 1|V = v)} \\
&= \frac{\delta_{ev}\beta_{ev}}{\delta_{1v}\beta_{1v} + \delta_{2v}\beta_{2v}}
\end{aligned}$$

## Appendix 2.2: Full conditional distributions

$$\log(\pi(\beta_{10}|\beta_{20}, x, y, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}))$$

$$\begin{aligned} \propto & (N_{11} + N_{10}) \log \beta_{10} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\ & + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \end{aligned}$$

$$\log(\pi(\beta_{20}|\beta_{10}, x, y, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}))$$

$$\begin{aligned} \propto & (N_{21} + N_{20}) \log \beta_{20} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\ & + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \end{aligned}$$

$$\log(\pi(x|\beta_{10}, \beta_{20}, \rho_{\beta_2}, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}))$$

$$\begin{aligned} \propto & N_{11}x - \frac{x^2}{2 \times 0.7^2} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\ & + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \end{aligned}$$

$$\log(\pi(y|\beta_{10}, \beta_{20}, x, \delta_{10}, \delta_{20}, z_{\delta_1}, z_{\delta}))$$

$$\begin{aligned} \propto & N_{21}y - \frac{(y + 0.96)^2}{2 \times 0.32^2} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\ & + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \end{aligned}$$

$$\log(\pi(\delta_{10}|\beta_{10}, \beta_{20}, x, y, \delta_{20}, z_{\delta_1}, z_{\delta}))$$

$$\begin{aligned} \propto & (N_{11} + N_{10}) \log \delta_{10} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\ & + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \end{aligned}$$

$$\log(\pi(\delta_{20}|\beta_{10}, \beta_{20}, x, y, \delta_{10}, z_{\delta_1}, z_{\delta}))$$

$$\begin{aligned} \propto & (N_{21} + N_{20}) \log \delta_{20} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\ & + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \end{aligned}$$



$$\begin{aligned}
& \log(\pi(z_{\delta_1} | \beta_{10}, \beta_{20}, x, y, \delta_{10}, z, z_{\delta})) \\
& \propto (N_{11} + N_{21})z_{\delta_1} - \frac{z_{\delta_1}^2}{2 \times 0.7^2} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \\
& \quad \exp(y)\beta_{20}] + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}]) \\
& \log(\pi(z_{\delta} | \beta_{10}, \beta_{20}, x, y, \delta_{10}, z, z_{\delta_1})) \\
& \propto N_{21}z_{\delta} - \frac{z_{\delta}^2}{2 \times 0.7^2} - N \log(\alpha[\exp(z_{\delta_1})\delta_{10} \exp(x)\beta_{10} + \delta_{20} \exp(z_{\delta}) \exp(z_{\delta_1}) \exp(y)\beta_{20}] \\
& \quad + (1 - \alpha)[\delta_{10}\beta_{10} + \delta_{20}\beta_{20}])
\end{aligned}$$

### Appendix 3.1: True $VE_{TMI}$ and $VE_{TSI}$

The true VE against SI is:

$$VE_{TSI} = 1 - RR_{TSI} \quad \text{where} \quad RR_{TSI} = \frac{P(E = 2 | V = 1)}{P(E = 2 | V = 0)}$$

$$\begin{aligned}
& P(E = 2 | V = 1) \\
& = \sum_{x=0,1} P(E = 2 | V = 1, X = x)P(X = x | V = 1) \\
& = \gamma_{10} \frac{\alpha_0(1 - \pi)}{\alpha_0(1 - \pi) + \alpha_1\pi} + \gamma_{11} \frac{\alpha_1\pi}{\alpha_0(1 - \pi) + \alpha_1\pi} \\
& = \frac{\gamma_{10}\alpha_0(1 - \pi) + \gamma_{11}\alpha_1\pi}{\alpha_0(1 - \pi) + \alpha_1\pi}
\end{aligned}$$

where  $P(X = x | V = v) = \frac{P(V=v|X=x)P(X=x)}{P(V=v)}$  and  $P(V = 1) = \sum_{x=0,1} P(V = 1 | X = x)P(X = x)$ .

$$\begin{aligned}
P(E = 2 | V = 0) &= \sum_{x=0,1} P(E = 2 | V = 0, X = x)P(X = x | V = 0) \\
&= \gamma_{00} \frac{(1 - \alpha_0)(1 - \pi)}{1 - \alpha_0(1 - \pi) - \alpha_1\pi} + \gamma_{01} \frac{(1 - \alpha_1)\pi}{1 - \alpha_0(1 - \pi) - \alpha_1\pi} \\
&= \frac{\gamma_{00}(1 - \alpha_0)(1 - \pi) + \gamma_{01}(1 - \alpha_1)\pi}{1 - \alpha_0(1 - \pi) - \alpha_1\pi}
\end{aligned}$$

Therefore,

$$VE_{TSI} = 1 - RR_{TSI} = 1 - \frac{[\gamma_{10}\alpha_0(1 - \pi) + \gamma_{11}\alpha_1\pi][1 - \alpha_0(1 - \pi) - \alpha_1\pi]}{[\alpha_0(1 - \pi) + \alpha_1\pi][\gamma_{00}(1 - \alpha_0)(1 - \pi) + \gamma_{01}(1 - \alpha_1)\pi]}$$

The true VE against medically-attended influenza is:

$$VE_{TMI} = 1 - RR_{TSI} \quad \text{where} \quad RR_{TMI} = \frac{P(E = 2, M = 1 | V = 1)}{P(E = 2, M = 1 | V = 0)}$$

$$\begin{aligned} P(E = 2, M = 1 | V = 1) &= \sum_{x=0,1} P(E = 2, M = 1 | V = 1, X = x)P(X = x | V = 1) \\ &= \sum_{x=0,1} P(M = 1 | E = 2, V = 1, X = x)P(E = 2 | V = 1, X = x)P(X = x | V = 1) \\ &= \delta_{20}\gamma_{10} \frac{\alpha_0(1 - \pi)}{\alpha_0(1 - \pi) + \alpha_1\pi} + \delta_{21}\gamma_{11} \frac{\alpha_1\pi}{\alpha_0(1 - \pi) + \alpha_1\pi} \\ &= \frac{\delta_{20}\gamma_{10}\alpha_0(1 - \pi) + \delta_{21}\gamma_{11}\alpha_1\pi}{\alpha_0(1 - \pi) + \alpha_1\pi} \end{aligned}$$

$$\begin{aligned} P(E = 2, M = 1 | V = 0) &= \sum_{x=0,1} P(E = 2, M = 1 | V = 0, X = x)P(X = x | V = 0) \\ &= \sum_{x=0,1} P(M = 1 | E = 2, V = 0, X = x)P(E = 2 | V = 0, X = x)P(X = x | V = 0) \\ &= \delta_{20}\gamma_{00} \frac{(1 - \alpha_0)(1 - \pi)}{1 - \alpha_0(1 - \pi) - \alpha_1\pi} + \delta_{21}\gamma_{01} \frac{(1 - \alpha_1)\pi}{1 - \alpha_0(1 - \pi) - \alpha_1\pi} \\ &= \frac{\delta_{20}\gamma_{00}(1 - \alpha_0)(1 - \pi) + \delta_{21}\gamma_{01}(1 - \alpha_1)\pi}{1 - \alpha_0(1 - \pi) - \alpha_1\pi} \end{aligned}$$

Therefore,

$$VE_{TMI} = 1 - RR_{TMI} = 1 - \frac{[\delta_{20}\gamma_{10}\alpha_0(1 - \pi) + \delta_{21}\gamma_{11}\alpha_1\pi][1 - \alpha_0(1 - \pi) - \alpha_1\pi]}{[\alpha_0(1 - \pi) + \alpha_1\pi][\delta_{20}\gamma_{00}(1 - \alpha_0)(1 - \pi) + \delta_{21}\gamma_{01}(1 - \alpha_1)\pi]}$$

## Appendix 3.2: Model-based estimates of VE

The model based estimates from TND study is:

$$VE_A = 1 - OR_A \quad \text{where} \quad OR_A = \frac{P(C_A=1, V=1|M=1)P(C_A=0, V=0|M=1)}{P(C_A=1, V=0|M=1)P(C_A=0, V=1|M=1)}$$

Here  $\{C_A = 1\} = \{M = 1, T = 1\}$  and  $\{C_A = 0\} = \{M = 1, T = 0\}$  and  $OR_A$  can be written as:

$$OR_A = \frac{P(M = 1, T = 1, V = 1)P(M = 1, T = 0, V = 0)}{P(M = 1, T = 1, V = 0)P(M = 1, T = 0, V = 1)}$$

$$\begin{aligned} & P(M = 1, T = 1, V = 1) \\ &= \sum_{x=0,1} P(M = 1, T = 1, V = 1 | X = x)P(X = x) \\ &= \sum_{x=0,1} \left[ \sum_{e=1,2} P(M = 1, T = 1, V = 1 | E = e, X = x)P(E = e | X = x) \right] P(X = x) \\ &= \sum_{x=0,1} \left[ \sum_{e=1,2} P(M = 1, T = 1 | V = 1, E = e, X = x)P(V = 1 | E = e, X = x) \right. \\ &\quad \left. P(E = e | X = x) \right] P(X = x) \\ &= \sum_{x=0,1} \left[ \sum_{e=1,2} P(M = 1, T = 1 | V = 1, E = e, X = x) \frac{P(E = e | V = 1, X = x)}{P(E = e | X = x)} \right. \\ &\quad \left. P(V = 1 | X = x)P(E = e | X = x) \right] P(X = x) \end{aligned}$$

In the above equation, M and T are independent given V and E. In addition, with the assumption of perfect sensitivity and specificity,  $P(T = 1 | E = 1, V = v, X = x) = 0$  and  $P(T = 1 | E = 2, V = v, X = x) = 1$ . So the above equation can be

written as:

$$\begin{aligned}
& P(M = 1, T = 1, V = v) \\
&= \sum_{x=0,1} P(T = 1 \mid E = 2)P(M = 1 \mid V = v, E = 2, X = x)P(E = 2 \mid V = v, X = x) \\
&\quad P(V = v \mid X = x)P(X = x)
\end{aligned}$$

Similarly,

$$\begin{aligned}
& P(M = 1, T = 0, V = v) \\
&= \sum_{x=0,1} P(T = 0 \mid E = 1)P(M = 1 \mid V = v, E = 1, X = x)P(E = 1 \mid V = v, X = x) \\
&\quad P(V = v \mid X = x)P(X = x)
\end{aligned}$$

Therefore,

$$\begin{aligned}
P(M = 1, T = 1, V = 1) &= \delta_{20}\gamma_{10}\alpha_0(1 - \pi) + \delta_{21}\gamma_{11}\alpha_1\pi \\
P(M = 1, T = 1, V = 0) &= \delta_{20}\gamma_{00}(1 - \alpha_0)(1 - \pi) + \delta_{21}\gamma_{01}(1 - \alpha_1)\pi \\
P(M = 1, T = 0, V = 1) &= \delta_{10}\beta_{10}\alpha_0(1 - \pi) + \delta_{11}\beta_{11}\alpha_1\pi \\
P(M = 1, T = 0, V = 0) &= \delta_{10}\beta_{00}(1 - \alpha_0)(1 - \pi) + \delta_{11}\beta_{01}(1 - \alpha_1)\pi
\end{aligned}$$

Therefore,

$$VE_A = 1 - \frac{[\delta_{20}\gamma_{10}\alpha_0(1-\pi) + \delta_{21}\gamma_{11}\alpha_1\pi][\delta_{10}\beta_{00}(1-\alpha_0)(1-\pi) + \delta_{11}\beta_{01}(1-\alpha_1)\pi]}{[\delta_{20}\gamma_{00}(1-\alpha_0)(1-\pi) + \delta_{21}\gamma_{01}(1-\alpha_1)\pi][\delta_{10}\beta_{10}\alpha_0(1-\pi) + \delta_{11}\beta_{11}\alpha_1\pi]}$$

The model based estimates from CC study is:

$$VE_B = 1 - OR_B \quad \text{where} \quad OR_B = \frac{P(C_B=1, V=1|B=1)P(C_B=0, V=0|B=1)}{P(C_B=1, V=0|B=1)P(C_B=0, V=1|B=1)}$$

Here  $\{C_B = 1\} = \{C_A = 1\} = \{M = 1, T = 1\}$ ,  $\{C_B = 0\}$  is a random subset of  $\{E = 0\}$  and  $\{B = 1\} = \{C_B = 1 \text{ or } C_B = 0\}$ .  $OR_B$  can be written as:

$$\begin{aligned} OR_B &= \frac{P(C_B = 1, V = 1)P(C_B = 0, V = 0)}{P(C_B = 1, V = 0)P(C_B = 0, V = 1)} \\ &= \frac{P(M = 1, T = 1, V = 1)P(E = 0, V = 0)}{P(M = 1, T = 1, V = 0)P(E = 0, V = 1)} \end{aligned}$$

As shown above,

$$P(M = 1, T = 1, V = 1) = \delta_{20}\gamma_{10}\alpha_0(1-\pi) + \delta_{21}\gamma_{11}\alpha_1\pi$$

$$P(M = 1, T = 1, V = 0) = \delta_{20}\gamma_{00}(1-\alpha_0)(1-\pi) + \delta_{21}\gamma_{01}(1-\alpha_1)\pi$$

For  $P(E = 0, V = v)$ ,

$$\begin{aligned} P(E = 0, V = v) &= \sum_{x=0,1} P(E = 0, V = v \mid X = x)P(X = x) \\ &= \sum_{x=0,1} P(E = 0 \mid V = v, X = x)P(V = v \mid X = x)P(X = x) \end{aligned}$$

Because  $P(E = 0 \mid V = v, X = x) = 1 - P(E = 1 \mid V = v, X = x) - P(E = 2 \mid V = v, X = x)$ , therefore,

$$P(E = 0, V = 0) = (1 - \beta_{00} - \gamma_{00})(1 - \alpha_0)(1 - \pi) + (1 - \beta_{01} - \gamma_{01})(1 - \alpha_1)\pi$$

$$P(E = 0, V = 1) = (1 - \beta_{10} - \gamma_{10})\alpha_0(1 - \pi) + (1 - \beta_{11} - \gamma_{11})\alpha_1\pi$$

Therefore,

$$VE_B = 1 - \frac{[\delta_{20}\gamma_{10}\alpha_0(1-\pi) + \delta_{21}\gamma_{11}\alpha_1\pi][(1-\beta_{00}-\gamma_{00})(1-\alpha_0)(1-\pi) + (1-\beta_{01}-\gamma_{01})(1-\alpha_1)\pi]}{[\delta_{20}\gamma_{00}(1-\alpha_0)(1-\pi) + \delta_{21}\gamma_{01}(1-\alpha_1)\pi][(1-\beta_{10}-\gamma_{10})\alpha_0(1-\pi) + (1-\beta_{11}-\gamma_{11})\alpha_1\pi]}$$

# Bibliography

- Aalen, O. O., Farewell, V. T., de Angelis, D., Day, N. E., and Nöel Gill, O. (1997), “A Markov model for HIV disease progression including the effect of HIV diagnosis and treatment: application to AIDS prediction in England and Wales,” *Statistics in medicine*, 16, 2191–2210.
- Aitkin, M. (1991), “Posterior bayes factors,” *Journal of the Royal Statistical Society. Series B (Methodological)*, 111–142.
- Allen, D., Ammann, A., Bailey, H., Arms, A., Baker, C., Berger, R., Birkhead, G., Boland, M., Colman, C., Permanente, K., et al. (1999), “Revised recommendations for HIV screening of pregnant women,” *Atlanta*.
- Bae, K. and Mallick, B. K. (2004), “Gene selection using a two-level hierarchical Bayesian model,” *Bioinformatics*, 20, 3423–3430.
- Becker, N. G., Lewis, J. J., Li, Z., and McDonald, A. (2003), “Age-specific back-projection of HIV diagnosis data,” *Statistics in medicine*, 22, 2177–2190.
- Bellocco, R. and Marschner, I. C. (2000), “Joint analysis of HIV and AIDS surveillance data in back-calculation,” *Statistics in medicine*, 19, 297–311.



- Belongia, E. A., Kieke, B. A., Donahue, J. G., Greenlee, R. T., Balish, A., Foust, A., Lindstrom, S., and Shay, D. K. (2009), “Effectiveness of inactivated influenza vaccines varied substantially with antigenic match from the 2004–2005 season to the 2006–2007 season,” *Journal of Infectious Diseases*, 199, 159–167.
- Beran, J., Ambrozaitis, A., Laiskonis, A., Mickuviene, N., Bacart, P., Calozet, Y., Demanet, E., Heijmans, S., Van Belle, P., Weber, F., et al. (2009), “Intradermal influenza vaccination of healthy adults using a new microinjection system: a 3-year randomised controlled safety and immunogenicity trial,” *BMC medicine*, 7, 13.
- Birrell, P. J., Gill, O. N., Delpech, V. C., Brown, A. E., Desai, S., Chadborn, T. R., Rice, B. D., and De Angelis, D. (2013), “HIV incidence in men who have sex with men in England and Wales 2001–10: a nationwide population study,” *The Lancet infectious diseases*.
- Blaxhult, A. and Svensson, A. (1992), “Assessing the extent of the HIV epidemic in Sweden, using information on the extent to which people who develop AIDS are already known to be HIV infected,” *International journal of epidemiology*, 21, 784–791.
- Branson, B. M., Handsfield, H. H., Lampe, M. A., Janssen, R. S., Taylor, A. W., Lyss, S. B., and Clark, J. E. (2006), “Revised recommendations for HIV testing of adults, adolescents, and pregnant women in health-care settings.” *MMWR. Recommendations and reports: Morbidity and mortality weekly report. Recommendations and reports/Centers for Disease Control and Prevention*, 55, 1–17.

- Burney, L. E. (1960), “Influenza immunization: statement,” *Public health reports*, 75, 944.
- Carrillo-Santistevé, P., Ciancio, B. C., Nicoll, A., and Lopalco, P. L. (2012), “The importance of influenza prevention for public health,” *Human vaccines & immunotherapeutics*, 8, 89–95.
- CDC (2012), “Monitoring selected national HIV prevention and care objectives by using HIV surveillance data – United States and 6 US dependent areas 2010,” *HIV surveillance supplemental report*, 17.
- (2013), “Prevention and control of seasonal influenza with vaccines. Recommendations of the Advisory Committee on Immunization Practices—United States, 2013–2014.” *MMWR. Recommendations and reports: Morbidity and mortality weekly report. Recommendations and reports/Centers for Disease Control*, 62, 1.
- CDC et al. (1987), “Public Health Service guidelines for counseling and antibody testing to prevent HIV infection and AIDS.” *MMWR. Morbidity and mortality weekly report*, 36, 509.
- Chau, P., Yip, P. S., and Cui, J. S. (2003), “Reconstructing the incidence of human immunodeficiency virus (HIV) in Hong Kong by using data from HIV positive tests and diagnoses of acquired immune deficiency syndrome,” *Journal of the Royal Statistical Society: series c (applied statistics)*, 52, 237–248.
- Cowling, B. J., Fang, V. J., Nishiura, H., Chan, K.-H., Ng, S., Ip, D. K., Chiu, S. S., Leung, G. M., and Peiris, J. M. (2012), “Increased risk of noninfluenza respiratory

virus infections associated with receipt of inactivated influenza vaccine,” *Clinical infectious diseases*, 54, 1778–1783.

Cui, J. and Becker, N. G. (2000), “Estimating HIV incidence using dates of both HIV and AIDS diagnoses,” *Statistics in medicine*, 19, 1165–1177.

De Serres, G., Skowronski, D., Wu, X., and Ambrose, C. (2013), “The test-negative design: validity, accuracy and precision of vaccine efficacy estimates compared to the gold standard of randomised placebo-controlled clinical trials,” *Euro Surveill*, 18.

Edwards, K. M., Dupont, W. D., Westrich, M. K., Plummer, W. D., Palmer, P. S., and Wright, P. F. (1994), “A randomized controlled trial of cold-adapted and inactivated vaccines for the prevention of influenza A disease,” *Journal of Infectious Diseases*, 169, 68–76.

Farewell, V., Aalen, O., De Angelis, D., and Mrc, N. D. (1994), “Estimation of the rate of diagnosis of HIV infection in HIV infected individuals,” *Biometrika*, 81, 287–294.

Ferdinands, J. M. and Shay, D. K. (2012), “Magnitude of potential biases in a simulated case-control study of the effectiveness of influenza vaccination,” *Clinical infectious diseases*, 54, 25–32.

Fielding, J. E., Grant, K. A., Garcia, K., and Kelly, H. A. (2011), “Effectiveness of seasonal influenza vaccine against pandemic (H1N1) 2009 virus, Australia, 2010,” *Emerging infectious diseases*, 17, 1181.

- Figueiredo, M. A. and Nowak, R. D. (2003), “An EM algorithm for wavelet-based image restoration,” *Image Processing, IEEE Transactions on*, 12, 906–916.
- Fiore, A. E., Uyeki, T. M., Broder, K., Finelli, L., Euler, G. L., Singleton, J. A., Iskander, J. K., Wortley, P. M., Shay, D. K., Bresee, J. S., et al. (2010), “Prevention and control of influenza with vaccines: recommendations of the Advisory Committee on Immunization Practices (ACIP), 2010.” *MMWR. Recommendations and Reports: Morbidity and mortality weekly report. Recommendations and reports/Centers for Disease Control*, 59, 1.
- Foppa, I. M., Haber, M., Ferdinands, J. M., and Shay, D. K. (2013), “The case test-negative design for studies of the effectiveness of seasonal influenza vaccine,” *Vaccine*.
- Frey, S., Vesikari, T., Szymczakiewicz-Multanowska, A., Lattanzi, M., Izu, A., Groth, N., and Holmes, S. (2010), “Clinical Efficacy of Cell CultureDerived and Egg-Derived Inactivated Subunit Influenza Vaccines in Healthy Adults,” *Clinical Infectious Diseases*, 51, 997–1004.
- Gelman, A., Meng, X.-L., and Stern, H. (1996), “Posterior predictive assessment of model fitness via realized discrepancies,” *Statistica Sinica*, 6, 733–760.
- Gelman, A. and Rubin, D. B. (1992), “Inference from iterative simulation using multiple sequences,” *Statistical science*, 457–472.
- Genkin, A., Lewis, D. D., and Madigan, D. (2007), “Large-scale Bayesian logistic regression for text categorization,” *Technometrics*, 49, 291–304.

- Gilks, W. R., Best, N., and Tan, K. (1995), “Adaptive rejection Metropolis sampling within Gibbs sampling,” *Applied Statistics*, 455–472.
- Green, T. A. (1998), “Using surveillance data to monitor trends in the AIDS epidemic,” *Statistics in medicine*, 17, 143–154.
- Haber, M., An, Q., Foppa, M. I., Shay, K. D., Ferdinands, M. J., and Orenstein, A. W. (2014), “A probability model for evaluating the bias and precision of influenza vaccine effectiveness estimates from case-control studies,” *Epidemiology and Infectious Disease*.
- Hak, E., Nordin, J., Wei, F., Mullooly, J., Poblete, S., Strikas, R., and Nichol, K. L. (2002), “Influence of high-risk medical conditions on the effectiveness of influenza vaccination among elderly members of 3 large managed-care organizations,” *Clinical Infectious Diseases*, 35, 370–377.
- Hall, H. I., Song, R., Rhodes, P., Prejean, J., An, Q., Lee, L. M., Karon, J., Brookmeyer, R., Kaplan, E. H., McKenna, M. T., et al. (2008), “Estimation of HIV incidence in the United States,” *Jama*, 300, 520–529.
- Hardelid, P., Fleming, D., McMenemy, J., Andrews, N., Robertson, C., Sebastian-Pillai, P., Ellis, J., Carman, W., Wreghitt, T., Watson, J., et al. (2011), “Effectiveness of pandemic and seasonal influenza vaccine in preventing pandemic influenza A (H1N1) 2009 infection in England and Scotland 2009-2010,” *Euro Surveill*, 16, 19763.
- Jackson, L. A., Gaglani, M. J., Keyserling, H. L., Balser, J., Bouveret, N., Fries, L., and Treanor, J. J. (2010), “Safety, efficacy, and immunogenicity of an inactivated

influenza vaccine in healthy adults: a randomized, placebo-controlled trial over two influenza seasons,” *BMC infectious diseases*, 10, 71.

Jackson, L. A., Jackson, M. L., Nelson, J. C., Neuzil, K. M., and Weiss, N. S. (2006a), “Evidence of bias in estimates of influenza vaccine effectiveness in seniors,” *International Journal of Epidemiology*, 35, 337–344.

Jackson, L. A., Nelson, J. C., Benson, P., Neuzil, K. M., Reid, R. J., Psaty, B. M., Heckbert, S. R., Larson, E. B., and Weiss, N. S. (2006b), “Functional status is a confounder of the association of influenza vaccine and risk of all cause mortality in seniors,” *International journal of epidemiology*, 35, 345–352.

Jackson, M. L. and Nelson, J. C. (2013), “The test-negative design for estimating influenza vaccine effectiveness,” *Vaccine*.

Janjua, N. Z., Skowronski, D. M., De Serres, G., Dickinson, J., Crowcroft, N. S., Taylor, M., Winter, A.-L., Hottes, T. S., Fonseca, K., Charest, H., et al. (2012), “Estimates of Influenza Vaccine Effectiveness for 2007–2008 From Canada’s Sentinel Surveillance System: Cross-Protection Against Major and Minor Variants,” *Journal of Infectious Diseases*, 205, 1858–1868.

Janssen, R., Onorato, I., Valdiserri, R., Durham, T., Nichols, W., Seiler, E., and Jaffe, H. (2003), “Advancing HIV prevention: new strategies for a changing epidemic—United States, 2003.” *MMWR. Morbidity and mortality weekly report*, 52, 329.

Kass, R. E. and Raftery, A. E. (1995), “Bayes factors,” *Journal of the american statistical association*, 90, 773–795.

- Kelly, H., Carville, K., Grant, K., Jacoby, P., Tran, T., and Barr, I. (2009), “Estimation of influenza vaccine effectiveness from routine surveillance data,” *PLoS One*, 4, e5079.
- Kissling, E., Valenciano, M., Falcao, J., Larrauri, A., Widgren, K., Pitigoi, D., Oroszi, B., Nunes, B., Savulescu, C., Mazick, A., et al. (2009), “I-MOVE towards monitoring seasonal and pandemic influenza vaccine effectiveness: lessons learnt from a pilot multi-centric case-control study in Europe, 2008-9,” .
- Kyung, M., Gill, J., Ghosh, M., and Casella, G. (2010), “Penalized regression, standard errors, and Bayesian lassos,” *Bayesian Analysis*, 5, 369–411.
- Larson, H. J. and Heymann, D. L. (2010), “Public health response to influenza A (H1N1) as an opportunity to build public trust,” *JAMA: The Journal of the American Medical Association*, 303, 271–272.
- Long, P. (1964), “RECOMMENDATIONS FOR INFLUENZA IMMUNIZATION AND CONTROL: 1964-1965.” *Medical times*, 92, 1203.
- Longini, I. M., Clark, W. S., Byers, R. H., Ward, J. W., Darrow, W. W., Lemp, G. F., and Hethcote, H. W. (1989), “Statistical analysis of the stages of HIV infection using a Markov model,” *Statistics in medicine*, 8, 831–843.
- Longini, I. M., Clark, W. S., Gardner, L. I., and Brundage, J. F. (1991), “The dynamics of CD4+ T-lymphocyte decline in HIV-infected individuals: a Markov modeling approach,” *Journal of Acquired Immune Deficiency Syndromes*, 4, 1141–1147.

- Madhi, S. A., Maskew, M., Koen, A., Kuwanda, L., Besselaar, T. G., Naidoo, D., Cohen, C., Valette, M., Cutland, C. L., and Sanne, I. (2011), “Trivalent inactivated influenza vaccine in African adults infected with human immunodeficient virus: double blind, randomized clinical trial of efficacy, immunogenicity, and safety,” *Clinical Infectious Diseases*, 52, 128–137.
- Mallitt, K.-A., Wilson, D. P., McDonald, A., and Wand, H. (2012), “HIV incidence trends vary between jurisdictions in Australia: an extended back-projection analysis of men who have sex with men,” *Sexual Health*, 9, 138–143.
- Marschner, I. C. (1994), “Using time of first positive HIV test and other auxiliary data in back-projection of AIDS incidence,” *Statistics in medicine*, 13, 1959–1974.
- Meiklejon, G. (1994), “Commission on influenza,” *The histories of the commissions. Falls Church, VA: The Borden Institute, Office of the Surgeon General, Department of the Army.*
- Molinari, N.-A. M., Ortega-Sanchez, I. R., Messonnier, M. L., Thompson, W. W., Wortley, P. M., Weintraub, E., and Bridges, C. B. (2007), “The annual impact of seasonal influenza in the US: measuring disease burden and costs,” *Vaccine*, 25, 5086–5096.
- Moore, R. D. and Chaisson, R. E. (1999), “Natural history of HIV infection in the era of combination antiretroviral therapy,” *Aids*, 13, 1933–1942.
- O’Brien, S. J. and Hendrickson, S. L. (2013), “Host genomic influences on HIV/AIDS,” *Genome biology*, 14, 201.



- Orenstein, E. W., De Serres, G., Haber, M. J., Shay, D. K., Bridges, C. B., Gargiullo, P., and Orenstein, W. A. (2007), “Methodologic issues regarding the use of three observational study designs to assess influenza vaccine effectiveness,” *International journal of epidemiology*, 36, 623–631.
- Osterholm, M. T., Kelley, N. S., Sommer, A., and Belongia, E. A. (2012), “Efficacy and effectiveness of influenza vaccines: a systematic review and meta-analysis,” *The Lancet infectious diseases*, 12, 36–44.
- Park, T. and Casella, G. (2008), “The bayesian lasso,” *Journal of the American Statistical Association*, 103, 681–686.
- Punyacharoensin, N. and Viwatwongkasem, C. (2009), “Trends in three decades of HIV/AIDS epidemic in Thailand by nonparametric backcalculation method,” *AIDS*, 23, 1143–1152.
- Skowronski, D., Gilbert, M., Tweed, S., Petric, M., Li, Y., Mak, A., et al. (2005), “Effectiveness of vaccine against medical consultation due to laboratory-confirmed influenza: results from a sentinel physician pilot project in British Columbia, 2004–2005,” *Can Commun Dis Rep*, 31, 181–191.
- Skowronski, D., Janjua, N., De Serres, G., Dickinson, J., Winter, A., Mahmud, S., Sabaiduc, S., Gubbay, J., Charest, H., Petric, M., et al. (2013), “Interim estimates of influenza vaccine effectiveness in 2012/13 from Canadas sentinel surveillance network, January 2013,” *Euro Surveill*, 18.
- Skowronski, D., Masaro, C., Kwindt, T., Mak, A., Petric, M., Li, Y., Sebastian, R., Chong, M., Tam, T., and De Serres, G. (2007), “Estimating vaccine effectiveness

against laboratory-confirmed influenza using a sentinel physician network: results from the 2005–2006 season of dual A and B vaccine mismatch in Canada,” *Vaccine*, 25, 2842–2851.

Skowronski, D. M., De Serres, G., Crowcroft, N. S., Janjua, N. Z., Boulianne, N., Hottes, T. S., Rosella, L. C., Dickinson, J. A., Gilca, R., Sethi, P., et al. (2010), “Association between the 2008–09 seasonal influenza vaccine and pandemic H1N1 illness during spring–summer 2009: four observational studies from Canada,” *PLoS medicine*, 7, e1000258.

Skowronski, D. M., De Serres, G., Dickinson, J., Petric, M., Mak, A., Fonseca, K., Kwindt, T. L., Chan, T., Bastien, N., Charest, H., et al. (2009), “Component-specific effectiveness of trivalent influenza vaccine as monitored through a sentinel surveillance network in Canada, 2006–2007,” *Journal of Infectious Diseases*, 199, 168–179.

Skowronski, D. M., Janjua, N. Z., De Serres, G., Hottes, T. S., Dickinson, J. A., Crowcroft, N., Kwindt, T. L., Tang, P., Charest, H., Fonseca, K., et al. (2011), “Effectiveness of AS03 adjuvanted pandemic H1N1 vaccine: case-control evaluation based on sentinel surveillance system in Canada, autumn 2009,” *BMJ: British Medical Journal*, 342.

Skowronski, D. M., Janjua, N. Z., De Serres, G., Winter, A.-L., Dickinson, J. A., Gardy, J. L., Gubbay, J., Fonseca, K., Charest, H., Crowcroft, N. S., et al. (2012), “A sentinel platform to evaluate influenza vaccine effectiveness and new variant circulation, Canada 2010–2011 season,” *Clinical infectious diseases*, 55, 332–342.

- Song, R., Hall, H. I., and Frey, R. (2005), “Uncertainties associated with incidence estimates of HIV/AIDS diagnoses adjusted for reporting delay and risk redistribution,” *Statistics in medicine*, 24, 453–464.
- Sweeting, M. J., De Angelis, D., and Aalen, O. O. (2005), “Bayesian back-calculation using a multi-state model with application to HIV,” *Statistics in medicine*, 24, 3991–4007.
- Taffe, P. and May, M. (2008), “A joint back calculation model for the imputation of the date of HIV infection in a prevalent cohort,” *Statistics in medicine*, 27, 4835–4853.
- Thompson, M., Shay, D., Zhou, H., Bridges, C., Cheng, P., Burns, E., Bresee, J., and Cox, N. (2010), “Estimates of Deaths Associated With Seasonal Influenza-United States, 1976-2007 (Reprinted from MMWR, vol 59, pg 1057-1062, 2010),” .
- Thompson, W. W., Shay, D. K., Weintraub, E., Brammer, L., Bridges, C. B., Cox, N. J., and Fukuda, K. (2004), “Influenza-associated hospitalizations in the United States,” *JAMA: the journal of the American Medical Association*, 292, 1333–1340.
- Thompson, W. W., Shay, D. K., Weintraub, E., Brammer, L., Cox, N., Anderson, L. J., and Fukuda, K. (2003), “Mortality associated with influenza and respiratory syncytial virus in the United States,” *JAMA: the journal of the American Medical Association*, 289, 179–186.
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., and Knight, K. (2005), “Sparsity and smoothness via the fused lasso,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 91–108.

- Treanor, J. J., Talbot, H. K., Ohmit, S. E., Coleman, L. A., Thompson, M. G., Cheng, P.-Y., Petrie, J. G., Lofthus, G., Meece, J. K., Williams, J. V., et al. (2012), “Effectiveness of seasonal influenza vaccines in the United States during a season with circulation of all three vaccine strains,” *Clinical Infectious Diseases*, 55, 951–959.
- Valenciano, M., Kissling, E., Cohen, J.-M., Oroszi, B., Barret, A.-S., Rizzo, C., Nunes, B., Pitigoi, D., Cámara, A. L., Mosnier, A., et al. (2011), “Estimates of pandemic influenza vaccine effectiveness in Europe, 2009–2010: results of Influenza Monitoring Vaccine Effectiveness in Europe (I-MOVE) multicentre case-control study,” *PLoS medicine*, 8, e1000388.
- Wand, H., Wilson, D., Yan, P., Gonnermann, A., McDonald, A., Kaldor, J., and Law, M. (2009), “Characterizing trends in HIV infection among men who have sex with men in Australia by birth cohorts: results from a modified back-projection method,” *Journal of the International AIDS Society*, 12, 19.
- Wand, H., Yan, P., Wilson, D., McDonald, A., Middleton, M., Kaldor, J., and Law, M. (2010), “Increasing HIV transmission through male homosexual and heterosexual contact in Australia: results from an extended back-projection approach,” *HIV medicine*, 11, 395–403.
- Ward, J., Janssen, R., and Jaffe, H. (1993), “Recommendations for HIV testing services for inpatients and outpatients in acute-care hospital settings.” *MMWR. Morbidity and Mortality Weekly Report*, 42, 1–6.

WHO Media Centre, W. H. O. (2013), “HIV/AIDS fact sheets,”  
*<http://www.who.int/mediacentre/factsheets/fs360/en/>*.

Yan, P., F, Z., and Wand, H. (2011), “Using HIV diagnostic data to estimate HIV incidence: method and simulation,” *Statistical Communications in Infectious Diseases*, 3, 6.