A Study of Non-collapsibility of the Odds Ratio via Marginal Structural Models and Logistic Regression Models

Menglan Pang

Degree of Master of Science

Department of Epidemiology, Biostatistics and Occupational Health

McGill University Montreal, Quebec, Canada

June, 2012

A thesis submitted to McGill University in partial fulfillment of the requirements of the degree of Master of Science

© Menglan Pang 2012

DEDICATION

This document is dedicated to the graduate students of McGill University.

ACKNOWLEDGMENTS

First and foremost, I would like to express my deepest gratitude to my supervisors, Dr. Robert Platt and Dr. Jay Kaufman, who have supported me throughout my graduate study. I thank them for reviewing my thesis and providing countless pieces of advice. This thesis could not be accomplished without their patient guidance.

I also thank the department of Epidemiology, Biostatistics and Occupational Health for a friendly academic atmosphere. I appreciated the opportunity of presenting my thesis work in the causal inference group. I thank the people who listened to me, offered great comments and suggestions on my thesis.

I would like to thank my parents for their love and care. They educated me to have a positive attitude to face difficulties. They encouraged me whenever I needed and supported me on every decision I made. Without them, I would not have achieved this goal.

Finally, I thank all my friends who have helped and supported me through my time in Montreal.

iii

TABLE OF CONTENTS

DEDICATIONii
ACKNOWLEDGMENTS iii
TABLE OF CONTENTSiv
LIST OF TABLES
LIST OF FIGURES
ABSTRACTix
RÉSUMÉxi
1 Introduction1
2 Literature Review
2.1 Definition of Non-collapsibility
2.1.1 Collapsibility in contingency tables5
2.1.2 Non-collapsibility in the context of regression
2.2 Sufficient conditions for collapsibility of RD, RR and OR
2.3 Non-collapsibility of the odds ratio9
2.3.1 Non-collapsibility and confounding9
2.3.2 Other problems of non-collapsibility of the OR11
2.4 Estimate the marginal (causal) OR in the presence of confounding12
2.4.1 G-computation12
2.4.2 IPW estimation in propensity score methods and Marginal structural models13
2.5 Measuring the non-collapsibility of the OR14
3 Introduction of the non-collapsibility effect on the odds ratio17
3.1 Non-collapsibility without confounding17
3.2 Confounding without non-collapsibility18

3.3 The measure of the non-collapsibility effect	20
4 The Non-collapsibility effect in a point-exposure study	21
4.1 An analytical approach	21
4.1.1 Introduction	21
4.1.2 Method	22
4.1.2.1 Notation and related parameters	22
4.1.2.2 Homogenous odds ratio across L	23
4.1.2.3 Estimation of the odds ratios	24
4.1.3 Result	25
4.1.3.1 Measure of non-collapsibility of the OR	25
4.1.3.2 Conditions for the absence of the non-collapsibility effect	26
4.1.3.3 Decomposition of the total discrepancy	30
4.1.3.4 Comparison of non-collapsibility effect with other formulae	32
4.1.4 Discussion	33
4.2 A graphical approach	34
4.2.1 Introduction	34
4.2.2 Method: relationship between non-collapsibility and the effect of A and L on Y	35
4.2.3 Results	36
4.2.3.1 Scenario without confounding	36
4.2.3.2 Scenario with confounding	43
4.2.4 Discussion	45
5 The non-collapsibility in the presence of time-varying confounding	47
5.1 Introduction	47
5.2 Method	47
5.2.1 Time-varying confounding scenario	47

5.2.1.1 Data Generation	49
5.2.1.2 Data Analysis	50
5.2.1.3 Results	51
5.2.2 Non-collapsibility with Z	55
5.2.2.1 Method	55
5.2.2.2 Result	58
5.3 Discussion	61
6 Conclusions and Discussion	63
Appendix A	68
Bibliography	72

LIST OF TABLES

Table 3-1: Non-collapsibility without confounding	. 18
Table 3-2: Confounding without non-collapsibility	. 19
Table 4-1: Data from the point-exposure study	. 22
Table 4-2: Decomposition of the total discrepancy	. 30
Table 4-3: Data correspond to the first row of table 4-2	. 30
Table 4-4: Data correspond to the second row of table 4-2	. 31
Table 4-5: Data correspond to the third row of table 4-2	. 31
Table 4-6: Some example results of comparison of non-collapsibility effect with other formula	e 33
Table 4-7: Non-collapsibility effect with different conditional outcome probabilities	. 39
Table 5-1: Estimates in time-varying confounding scenario.	. 51
Table 5-2: Comparison between the simulation and equation (13) for the total effect of A0 on V	Y
through L	. 54
Table 5-3: Comparison of the results from model 7 and equation (14)	. 55
Table 5-4: Estimates in time-varying confounding scenario with Z	. 58
Table 5-5: Comparison of non-collapsibility effect with respect to A ₀ effect	. 59
Table 5-6: Comparison of non-collapsibility effect with respect to A effect (1)	. 59
Table 5-7: Comparison of non-collapsibility effect with respect to A effect (2)	61

LIST OF FIGURES

Figure 4-1: Causal diagrams for the point-exposure studies	22
Figure 4-2: Non-collapsibility effect vs. the baseline risk with no confounding	36
Figure 4-3: Non-collapsibility effect vs. the marginal risk (1)	39
Figure 4-4: Non-collapsibility effect vs. the marginal risk (2)	40
Figure 4-5: Non-collapsibility effect vs. the marginal risk (3)	41
Figure 4-6: Non-collapsibility effect vs. the marginal risk (4)	42
Figure 4-7: Non-collapsibility effect vs. the marginal risk (5)	42
Figure 4-8: Non-collapsibility effect vs. the baseline risk with confounding	43
Figure 4-9: Non-collapsibility effect vs. the marginal risk with confounding	45
Figure 5-1: Time-varying confounding scenario	48
Figure 5-2: Time-varying confounding scenario with baseline variable Z (Scenario 1)	49
Figure 5-3: Time-varying confounding scenario with baseline variable Z (Scenario 2)	. 60

ABSTRACT

Background: It has been noted in epidemiology and biostatistics that when the odds ratio (OR) is used to measure the causal effect of a treatment or exposure, there is a discrepancy between the marginal OR and the conditional OR even in the absence of confounding. This is known as non-collapsibility of the OR. It is sometimes described (incorrectly) as a bias in the estimated treatment effect from a logistic regression model if an important covariate is omitted.

Objectives: Distinguish confounding bias from non-collapsibility and measure the non-collapsibility effect on the OR in different scenarios.

Methods: We used marginal structural models and standard logistic regression to measure the non-collapsibility effect and confounding bias. An analytic approach is proposed to assess the non-collapsibility effect in a point-exposure study. This approach can be used to verify the conditions for the absence of non-collapsibility and to examine the phenomenon of confounding without non-collapsibility. A graphical approach is employed to show the relationship between the noncollapsibility effect and the baseline risk or the marginal outcome probability, and it reveals the non-collapsibility behaviour with a range of different exposure effects and different covariate effects. In order to explore the non-collapsibility effect of the OR in the presence of time-varying confounding, an observational cohort study was simulated. **Results and Conclusion:** The total difference between the conditional and crude effects can be decomposed into a sum of the non-collapsibility effect and the confounding bias. We provide a general formula for expressing the noncollapsibility effect under different scenarios. Our analytic approach provided similar results to related formulae in the literature. Various interesting observations about non-collapsibility can be made from the different scenarios with or without confounding using the graphical approach. Somewhat surprisingly, the effect of the covariate plays a more important role in the noncollapsibility effect than does the effect of the exposure. In the presence of timevarying confounding, the non-collapsibility is comparable to the effect in the point-exposure study.

RÉSUMÉ

Contexte : Il a été observé en épidémiologie et en biostatistique que lorsque le "odds ratio" (OR) est utilisé pour mesurer l'effet causal d'un traitement ou d'une exposition, il y a une différence entre l'OR marginal et l'OR conditionnel et ce, même s'il y a absence de biais de confusion. Ceci est décrit comme le noncollapsibilité de l'OR. Il est parfois incorrectement décrit comme un biais dans l'effet estimé du traitement à partir d'un modèle de régression logistique, si une covariante importante est exclue.

Objectifs : Distinguer le biais provenant du biais de confusion du noncollapsibilité et mesurer l'effet du non-collapsibilité sur l'OR dans plusieurs scénarios.

Méthode : On a utilisé des modèles structuraux marginaux et la régression logistique ajustée pour mesurer l'effet du non-collapsibilité dans une étude d'exposition par points. Cette approche peut être utilisée pour vérifier les conditions de l'absence de non-collapsibilité et pour examiner le phénomène de biais de confusion sans non-collapsibilité. Une approche graphique est employée pour démontrer la relation entre le non-collapsibilité et le risque de base ou la probabilité du résultat marginal; ceci révèle le comportement de non-collapsibilité avec une étendue d'effets d'exposition et de covariance différents. De manière à explorer l'effet de non-collapsibilité de l'OR en présence de biais de confusion variant en fonction du temps, une étude d'observation de cohorte a été simulée.

Résultats et Conclusion : La différence entre les effets conditionnels et bruts peut être décomposée dans la somme de l'effet de non-collapsibilité et du biais de confusion. Nous suggérons une formule générale pour exprimer l'effet du noncollapsibilité dans plusieurs scénarios différents. Notre approche analytique expose des résultats similaires à d'autres étant trouvés avec des formules présentes dans la littérature. Plusieurs observations intéressantes sur le noncollapsibilité peuvent être faites à partir de différents scénarios, avec ou sans biais de confusion, en utilisant notre approche graphique. De manière surprenante, l'effet d'une covariable joue un plus grand rôle dans le non-collapsibilité que l'effet de l'exposition. En présence de biais de confusion reliée au temps, l'effet du non-collapsibilité est comparable à l'effet de l'étude d'exposition par point.

1 Introduction

Non-collapsibility of the odds ratio (OR) is the phenomenon that when estimating the exposure outcome association with the OR, collapsing over the other covariate(s), the conditional OR does not necessarily equal the marginal OR even in the absence of confounding and effect modification. It is also known as the discrepancy between estimates of the treatment effect in a logistic regression model if a necessary covariate is omitted (Gail et al. 1984). The non-collapsibility of the OR derives from the fact that when the expected of outcome is modeled as a log odds of exposure, the marginal effect cannot be expressed as a weighted average of conditional effects.

It is widely realized in epidemiologic research that the OR is not generally collapsible. Conditional and marginal ORs can be different in the absence of confounding and there could be confounding bias even when conditional and marginal ORs are equal (Greenland et al., 1999). Yet it remains puzzling exactly why and how the phenomena occur. This will be better understood if the non-collapsibility effect is defined and explored. A simple approach to quantify the magnitude of confounding is to compare the estimates with and without adjusting for the covariate(s) (Flanders and Khoury, 1990; Miettinen, 1972), but it is not appropriate for the OR due to the problem of noncollapsibility (Boivin and Wacholder 1985). Janes et al. showed that this simple quantification of confounding actually consists of two components: the true confounding bias and the nonlinearity (non-collapsibility) effect. They proposed a better measure of confounding that does not include the non-collapsibility effect (Janes et al. 2010). The magnitude of confounding was defined as the discrepancy between a marginal exposure effect and a crude exposure effect. It was estimated in two different ways, one with the standardization and the other involved inverse probability weights. Both of the methods need additional estimates which are obtained from logistic regression models. Furthermore, the nonlinearity effect was measured as the discrepancy between the marginal exposure effect and a conditional exposure effect. It was estimated by referring to the formula (equation 8) via approximation in the paper of Neuhaus et al. (1991). Hence, there is no explicit form of the two components given in the paper.

In this thesis, similar to the method developed by Janes et al., we propose an approach to measure the non-collapsibility effect by comparing the marginal OR (OR_m) from the marginal structural model (MSM) and the conditional OR (OR_c) from the standard logistic regression model (SLRM). The confounding bias is measured by comparing the marginal OR and the crude OR (OR_{crude}).

 OR_m corresponds to $e^{\alpha'}$ estimated by the MSM:

 $g[E(Y | A = a)] = \mu' + \alpha' A$ weighting each subject by the inverse of the probability of receiving his/her exposure conditional on L.

 OR_c corresponds to e^{α} estimated from the logistic regression model:

 $g[E(Y | A = a, L = l)] = \mu + \alpha A + \beta L$, where g is the logit link function.

 OR_{crude} corresponds to e^{α^*} estimated from the logistic regression model:

$$g[E(Y \mid A = a)] = \mu^* + \alpha^* A$$

In our approach, we measure the non-collapsibility effect and confounding bias analytically instead of by estimation and approximation. We will present the two components explicitly as the functions of a range of parameters, i.e. the probabilities of outcome conditional on the exposure and the covariate, the probability of the exposure conditional on the covariate, and the prevalence of the covariate. Given the values of the parameters, we can obtain the accurate true value of confounding bias and the noncollapsibility effect analytically and then investigate the relationship between the true value of non-collapsibility effect and the parameters. A decomposition of the total discrepancy between the conditional and the crude OR as the combination of the true confounding bias and the non-collapsibility effect is illustrated below by examples. Moreover, it can clearly explain the phenomena of non-collapsibility without confounding and confounding without non-collapsibility.

Some conditions for the collapsibility of the OR have been remarked on already in the literature. If the true effect of exposure is not null, OR is collapsible when the covariate is independent of exposure conditional on outcome, or is independent of the outcome conditional on exposure (Wermuth, 1987; Shapiro, 1982; Geng and Li, 2002; Hernán et al, 2011; Greenland and Pearl, 2011). It was mentioned that the marginal OR is attenuated by a factor which depends on parameters of the distribution of the covariate (Ritz and Spiegelman 2004); the magnitude of the non-collapsibility effect depends on the effects of the exposure and covariate on outcome and on the variance of the covariate

(Groenwold and Moons et al., 2011). Beyond that, some formulae were also developed to measure the discrepancy between the marginal OR and conditional OR in randomized trials or clustered data (Samuels, 1981; Gail, 1984; Neuhuas, 1991). However, to our knowledge, none of them investigated the relationship between the non-collapsibility effect and all the other related parameters in a general scenario with or without confounding. A point-exposure study is designed in this thesis to measure the noncollapsibility effect. A formula for the collapsibility effect on the OR will be expressed as the function of all parameters; figures under different scenarios and parameter settings will be presented to demonstrate the behavior of the non-collapsibility of the OR. When extending to scenarios in the presence of time-varying confounding, it remains difficult to derive a formula for the non-collapsibility effect. Thus simulations will be employed to explore the collapsibility behavior in the scenarios with a baseline covariate and compare to a similar point-exposure study.

As the marginal OR and the conditional OR have different interpretations, it is important for scientists to take care to consider the causal structure and the population they are interested in and then determine the appropriate measure. One should also be cautious about whether to include the covariates in the model. When comparing the marginal OR and the conditional OR, the confounding bias should be taken account of, as well as the non-collapsibility effect.

4

2 Literature Review

2.1 Definition of Non-collapsibility

2.1.1 Collapsibility in contingency tables

A 2x2xK contingency table is commonly constructed in an epidemiological study when the exposure (A) and outcome (Y) are binary and are stratified on a third variable (L) with K levels. A given measure of the exposure-outcome association is said to be collapsible, if collapsing over L, the value obtained from the marginal table can be expressed as a weighted average of the stratum-specific values. The measure is said to be strictly collapsible if the marginal value and the K stratum-specific values are all equal (Whittemore, 1978; Ducharme and Lepage, 1986; Greenland and Mickey, 1988; Greenland, Robins and Pearl, 1999 ; Newman, 2001). In other words, it is strictly collapsible if and only if the measure is averageable and the stratum-specific values are homogeneous.

Non-collapsibility is the phenomenon that the measure violates the collapsibility condition. The most notable non-collapsibility phenomenon attracts many researchers' attention when the marginal value does not equal the stratum-specific values even when the latter are constant and there is no confounding.

5

2.1.2 Non-collapsibility in the context of regression

In regression contexts, the measure is strictly collapsible over the third covariate (L) in a generalized linear model if the estimate is constant when L is included and omitted (Clogg, Petkova and Shihadeh, 1992). The measure is implicitly assumed to be homogenous across L from models conditional on L. Non-collapsibility is present if the estimates are distinct in the regression models with and without adjusting for L.

Non-collapsibility has been also considered in the clustered and longitudinal data setting. The random effect of the cluster corresponds to the covariate L in the preceding regression models. It has been noticed, for example, that the cluster-specific effect from a mixed-effect logistic model can be different from the population averaged effect from a marginal regression model (Neuhaus, Kalbfleisch and Hauck, 1991; Ritz and Spiegelman, 2004).

2.2 Sufficient conditions for collapsibility of RD, RR and OR

Many examples show that in one study, collapsibility can differ for the risk difference (RD), risk ratio (RR) and odds ratio (OR) (Greenland et al., 1999; Newman, 2001; Greenland and Morgenstern, 2001; Groenwold and Moons et al., 2011). That is to say collapsibility depends on the chosen measure. Newman in his book *Biostatistical Methods in Epidemiology* systematically investigates the conditions that are sufficient to guarantee the collapsibility of measures (Newman, 2011, pp 46-53).

All three measures are collapsible when the risk of outcome is the same across strata of L in the unexposed population. It requires that the covariate is independent of the outcome

conditional on the exposure, i.e. $L \perp Y \mid A = 0$. When this condition is satisfied and the stratum-specific values are homogenous, all the three measures are strictly collapsible.

Another condition guaranteeing the collapsibility of RD and RR is that the distributions of L are the same in the exposed and unexposed population. This condition says that L is independent of A and can be related or unrelated to Y. When L is independent of both A and Y, it is a factor irrelevant to the study, thus the measures are collapsible over L naturally. When L is a risk factor for Y but marginally independent of A, the RD and RR are strictly collapsible if they are homogenous across L. However, this is not necessarily true for the OR. The marginal OR can be distinct from the stratum-specific ORs under this condition. This is exactly the condition that the randomized trial tries to create to balance the distributions of the risk factor in the unexposed and unexposed group. The non-collapsibility of the OR can be present even in the absence of confounding and effect modification. This particular phenomenon "non-collapsibility without confounding" was observed by Mietten and Cook and then discussed deeply in the literature (Mietten and Cook, 1981; Greenland, 1987; Greenland et al., 1999; Hernán et al, 2011).

A particular condition for OR to be collapsible (but not for RD and RR) is that when L is independent of the exposure among the individuals in the population who do not have the outcome, i.e. $L \perp A \mid Y = 0$. When this condition is satisfied and the ORs are homogenous across L, the OR is strictly collapsible. This is the phenomenon of "confounding without non-collapsibility" that occurs when L is associated with both A and Y, but independent of A given Y. Conditions for collapsibility discussed by many other researchers all coincide with Newman's (Wermuth, 1987; Shapiro, 1982; Geng and Li, 2002; Hernán et al, 2011; Greenland and Pearl, 2011). It should also be noticed that non-confounding also does not guarantee the collapsibility of person-time RD and RR (i.e. rate differences and ratios), because the distribution of person-time can be altered over time by exposure and the other risk factors even if the exposed and unexposed cohort are exchangeable at the beginning of the study (Greenland, 1996).

The non-collapsibility phenomenon has also been investigated with respect to the collider and the mediator of the exposure and outcome besides the potential confounding variable (Hernán, Clayton and Keiding, 2011). Moreover, it was extended to the descendant and the ancestor of the covariate and a set of variables jointly or conditionally (Greenland and Pearl, 2011).

In the regression context, collapsibility has been also widely researched in randomized trials and clustered and longitudinal data. A variety of models were applied to evaluate if the marginal effect is equivalent with the conditional effect (Gail et al., 1984; Ritz et al., 2004). Measures from most of the generalized linear models are collapsible, for example with identity or log link. But it turns out to be non-collapsible in logistic regression (i.e., a logit link). This is accordant with the preceding discussion that after randomization, the RD (measured from an additive linear model) and RR (measured from an additive model on the log scale) are collapsible but the OR (measured from a multiplicative model on the log scale) is not. The OR is non-collapsible unless the effect of exposure is zero, or the covariate does not vary or it is not associated with the outcome (Gail et al., 1984).

Several tests of collapsibility and strict collapsibility have been established (Whittemore, 1978; Asmussen and Edwards, 1983; Ducharme and LePage, 1986; Greenland and Mickey, 1988), though we will not focus on these tests in this thesis.

2.3 Non-collapsibility of the odds ratio

As mentioned before, the RD and RR are collapsible in the absence of confounding, while the OR is a non-collapsible measure in this case. This phenomenon is quite counterintuitive, and thus arouses great interest in many researchers. We will focus on non-collapsibility of the OR in the remainder of the literature review and the thesis. Jensen's inequality provides theoretical justification for non-collapsibility, showing that the marginal OR is shifted towards the null compared to the conditional OR (Samuels, 1981). The non-collapsibility of the OR just reflects the fact that the marginal OR cannot be expressed as a weighted average of conditional OR values (Greenland, 1987; Greenland and Morgenstern, 2001).

2.3.1 Non-collapsibility and confounding

There is a collapsibility-based definition of confounding, but with a risk to detect confounding even when the covariate is not associated with the exposure. For example, the marginal effect can be said to be biased when the risk factor is perfectly balanced among the groups (Gail et al., 1984). Furthermore, in some statistics literature, noncollapsibility has been treated as confounding bias incorrectly when comparing logistic regression with or without a baseline covariate (Boivin and Wacholder, 1985; Becher, 1992). Similarly, a covariate is defined to be a non-confounder if the marginal OR equals the conditional OR (Guo and Geng, 1995). Non-collapsibility has been a source of confusion when defining confounding (Grayson, 1987; Greenland et al., 1989). Due to non-collapsibility of the OR, it is not appropriate to define the discrepancy between the marginal OR and the conditional OR as a bias. They are estimates of two different parameters and both of them are valid. The marginal OR estimates the average effect of exposure in the entire population, whereas the conditional OR estimates the effect of exposure across strata (Greenland, 1987).

Another definition of a confounding variable is a variable that associated with both exposure and outcome. It introduces bias in effect estimation when one analyzes data with a crude statistical model that ignores the confounding variable. In the counterfactual approach, "confounding is present if our substitute population imperfectly represents what our target would have been like under the counterfactual condition" (Maldonado and Greenland 2002). It follows that confounding and non-collapsibility are different concepts (Miettinen and Cook, 1981; Greenland and Robins, 1986). Collapsibility is based on the observable distribution alone, while non-confounding is a property of potential outcome distributions.

Randomization is often considered to be a gold standard for control of confounding. Successful randomization ensures exchangeability of the exposed and unexposed populations, therefore confounding is absent. However non-collapsibility of the OR is present, since the crude estimate is not equal to a weighted average of stratum-specific estimates even after the randomization. We emphasize that the counterfactual definition of confounding depends on the target population and focuses on causal inference, while collapsibility depends on the selected parameter and has no definitive implication for causality or confounding (Greenland et al., 1999; Greenland and Morgenstern, 2001).

2.3.2 Other problems of non-collapsibility of the OR

In a multi-dimensional study, some use stratification based on the propensity score to measure the effect of exposure across strata. The non-collapsibility of the OR implies that the weighted average of the stratum-specific OR values is not representative of the marginal OR in the entire population. This is perhaps one of the reasons why Austin showed that in some settings, propensity score methods are biased for the marginal (causal) OR (Austin, 2007), and that estimates from propensity score methods are further from the null than the marginal OR (Stampf and Graf et al., 2010; Forbes and Shortreed ,2008; Zhang 2009).

Non-collapsibility of the OR also explains some findings in randomized trials. Using logistic regression for adjustment will result in estimates that are further from null and that will have an increased standard error compared with unadjusted analysis (Breslow and Day, 1987; Wickramaratne and Holford, 1989; Robinson and Jewell, 1991; Steyerberg 2000; Hernandez 2004). One is actually comparing two standard errors for two different parameters (the marginal effect and the conditional effect). The increased standard error after adjustment reveals the fact that the standard error of the conditional OR (OR_c) exceeds that of the marginal OR (OR_m). On the other hand, the RD and RR are not affected by the non-collapsibility problem in randomized trials. When using a linear regression or a generalized liner model with log link, adjusting baseline covariates

11

will result in the constant estimate. It gives a decreased standard error due to a reduction of residual variance after adjustment (Steyerberg, 2009).

Non-collapsibility can also be used to illustrate why after adding null constant to null observed data in contingency tables, the modified table can show evidence against the null (Greenland 2010).

2.4 Estimate the marginal (causal) OR in the presence of confounding

Logistic regression model is a generalized linear model with a logit link when the outcome is from a Bernoulli distribution. It provides an unbiased estimate of the conditional effect of exposure (OR_c) when L is included in the model. The effect of exposure is assumed to be constant across the strata; OR_c corresponds to the stratum-specific ORs from the contingency table. When the covariate is not associated with the exposure, the logistic regression model excluding the covariate yields a marginal exposure effect which is the average effect in the entire population. However, when the covariate is a confounding variable, omitting the covariate provides a biased marginal effect in the population. Due to the non-collapsibility of the OR, the weighted average of the stratum-specific ORs is also biased for the marginal OR (OR_m). Other approaches have been proposed to estimate the OR_m in the presence of confounding.

2.4.1 G-computation

Based on the counterfactual concept, G-computation or a simple imputation-type method was proposed to measure the OR_m (Petersen and Wang et al., 2006; Snowden and Rose

et al, 2011; Stampf and Graf et al., 2010; Zhang 2009). Y(A=1) denotes a subject's actual outcome among the exposed group, and Y(A=0) denote a subject's actual outcome among the unexposed group. Y(a=1) and Y(a=0) denote a subject's potential outcome if exposed and unexposed. Y(a=1) = Y(A=1) is only observed among the individuals who are exposed, and Y(a=0) is missing in that group. Whereas, for the individual who is unexposed, the potential outcome Y(a=0) = Y(A=0) is observed, and Y(a=1) is missing. The marginal (casual) OR can be estimated in the entire population by the mean of the probabilities of the potential outcome for each individual exposed and unexposed (Equation 1). The probabilities are predicted from the full logistic regression model (Equation 2) when the individual is under exposed and unexposed respectively.

(1) $\log(OR_m) = \operatorname{logit}(\hat{P}[Y(a=1)=1]) - \operatorname{logit}(\hat{P}[Y(a=0)=1])$ where \hat{P} is the mean of the probability of the potential outcome of each individual.

(2) logit
$$(P(Y=1) | A, L) = \beta_0 + \beta_1 A + \beta_2 L$$

2.4.2 IPW estimation in propensity score methods and Marginal structural models

In point-exposure studies, inverse probability weighted (IPW) estimation using the propensity score can be used to measure the OR_m in the presence of confounding. It can be estimated from a marginal logistic regression by weighting each treated individual with the inverse of the propensity score (probability of receiving treatment given the confounding variables) and weighting each untreated individual with the inverse of one minus their propensity score (Hernán and Robins, 2006; Forbes and Shortreed, 2008).

In a time-varying confounding scenario where L can be simultaneously a confounder and an intermediate variable, the marginal structural model (MSM) can be used to estimate the OR_m . The parameters of a MSM are estimated by the inverse-probability-of-treatment weighted (IPTW) estimators. Each subject is assigned a weight, which is the product of the inverse of receiving treatment conditional on the past confounder history at each time point. By IPTW, one creates a pseudo-population where the treatment is not confounded by the covariates. A standard marginal logistic regression is performed in the pseudopopulation, providing an unbiased marginal treatment effect (Robins, Hernán and Brumback, 2000).

2.5 Measuring the non-collapsibility of the OR

It has been described that in randomized trials and clustered data analysis, the discrepancy between OR_m and OR_c is related to the effect of L on Y and the distribution of L in the population (Ritz and Spiegelman 2004, Groenwold and Moons et al., 2011). Furthermore, various quantitative measures of non-collapsibility have been proposed.

In a randomized trial, the magnitude of discrepancy between a model with or without covariates was derived by Gail et al. (1984). Assume that the conditional expectation of Y given treatment (A) and covariate (L) satisfy the full regression model $E(Y \mid A, L) = h(\mu + \alpha A + \beta L), \text{ where } h(.) \text{ is a known function. It was shown that if the covariate is a random variable with expectation zero and variance <math>\sigma^2$, the discrepancy can be approximated by (3):

(3)
$$\frac{1}{4}\beta'\Omega\beta\{\frac{h''(\mu+\alpha)}{h'(\mu+\alpha)}-\frac{h''(\mu-\alpha)}{h'(\mu-\alpha)}\}$$
 where Ω is the covariance of the omitted

covariates, β is the effect of the covariates on the outcome, and μ is the intercept and α is the conditional treatment effect from the full model.

In the clustered data setting, when comparing the marginal model and the cluster-specific model, the discrepancy was approximated by (4):

(4)
$$-\beta \frac{Var(p)}{E(p)E(q)}$$
 where β is the cluster-specific effect, and logit(p) is the distribution

of the effect of the cluster variable (Neuhaus et al., 1991; Janes et al., 2010).

Samuels showed the relationship between the OR_m and the OR_c when the covariate and treatment are unconditionally independent (Samuels, 1981). The discrepancy was computed as (5) :

(5)
$$p^{-1}(1-p)\{\sum h(p_i)r_i\}\{1-OR_c\sum h(p_i)r_i\}^{-1}$$

where
$$r_i = P(L=i)$$
, $p_i = P(Y=1 | A=0, L=i)$, $p = \sum p_i r_i$ and
 $h(p_i) = p_i \{1 + p_i (OR_C - 1)\}^{-1}$.

By applying Jensen's inequality to the convex function, it follows that the OR_m will be always closer to 1 compared to the OR_c . The same conclusion can be also reached by the (3) and (4). However these measures of collapsibility of the OR are all proposed under the scenario where the covariate is not a confounding variable. In my thesis, we will develop a general formula to measure the non-collapsibility effect in the scenario with or without confounding.

3 Introduction of the non-collapsibility effect on the odds ratio

Non-collapsibility of the odds ratio (OR) is the phenomenon that occurs when estimating the exposure-outcome association by the OR, collapsing over the other covariate(s), the conditional ORs do not equal the marginal OR even in the absence of confounding and effect measure modification. As noted, Jensen's inequality provides theoretical justification for non-collapsibility; the marginal OR is shifted towards the null compared to the conditional OR. Non-collapsibility of the OR also implies that the marginal OR cannot be expressed as a weighted average of the conditional ORs.

In the literature review, two interesting phenomena were mentioned, namely "noncollapsibility without confounding" and "confounding without non-collapsibility". We next demonstrate these by the following numerical examples. Hypothetical cohort studies are used in order to understand the phenomena in 2×2 tables. To avoid issues related to random error, we assume that the entire population has been recruited into the cohort. Probability refers to proportions in the source population in this section.

3.1 Non-collapsibility without confounding

Table 3-1 adapted from Greenland and Morgenstern (2001), shows data on exposure A, outcome Y, and covariate L. No confounding is observed in the data, as

P(A = 1 | L = 1) = P(A = 1 | L = 0) = 0.5. In the stratum L=0, $OR_{L=0} = \frac{20 \times 40}{10 \times 30} = 2.667$

and in the stratum L=1, $OR_{L=1} = \frac{40 \times 20}{30 \times 10} = 2.667$. The conditional ORs are equal across

L, indicating that there is no modification of the effect of A on Y by L. Collapsing over L, intuitively one might expect the crude OR to be consistent with the common stratified OR of 2.667. However, because of non-collapsibility, the crude odds ratio

$$OR_{crude} = \frac{(20+40) \times (40+20)}{(10+30) \times (30+10)} = 2.25$$
 is distinct from the homogenous conditional OR

2.667. This could be interpreted (incorrectly) as confounding bias even if L is perfectly balanced among A (Boivin and Wacholder, 1985; Becher, 1992).

L=0			L=1		
	A=0	A=1		A=0	A=1
Y=0	20	10	Y=0	40	30
Y=1	30	40	Y=1	10	20
	$OR_{L=0} = \frac{20 \times 40}{10 \times 30} = 2.667$			$OR_{L=1} = \frac{40 \times 20}{30 \times 10} = 2.667$	

Table 3-1: Non-collapsibility without confounding

3.2 Confounding without non-collapsibility

In contrast, collapsibility of OR does not necessarily imply non-confounding. Table 3-2 shows data to illustrate the constant crude OR with the homogenous conditional OR in the presence of confounding. Confounding can be observed in the data, as

$$P(A=1|L=1) = 0.432$$
, and $P(A=1|L=0) = 0.578$. The conditional ORs are

homogenous,
$$OR_{L=0} = \frac{24 \times 100}{15 \times 60} = 2.667$$
, and $OR_{L=1} = \frac{80 \times 20}{50 \times 12} = 2.667$. But

$$OR_{crude} = \frac{(24+80) \times (100+20)}{(15+50) \times (60+12)} = 2.667$$
 even in the presence of confounding.

L=0			L=1		
	A=0	A=1		A=0	A=1
Y=0	24	15	Y=0	80	50
Y=1	60	100	Y=1	12	20
	$OR_{L=0} = \frac{24 \times 100}{15 \times 60} = 2.667$			$OR_{L=1} = \frac{80 \times 20}{50 \times 12} = 2.667$	

Table 3-2: Confounding without non-collapsibility

It is clearly shown in the preceding numerical examples that non-collapsibility of the OR cannot be used as a criterion for detecting confounding. In order to distinguish the non-collapsibility of the OR from confounding, we introduce the concept of the *non-collapsibility effect*. It is a measure of the discrepancy between the conditional OR and the marginal OR after taking account of the confounding bias. Though it is not a real effect in the causal sense, we will use this terminology to refer to the disparity between two different effect estimators.

When we have data with confounding, ignoring confounding by simply collapsing the table over L will induce confounding bias. The discrepancy between the crude OR and the conditional OR actually consists of two components: the non-collapsibility effect and the confounding bias. Janes et al (2010) demonstrated that the total discrepancy is equal to the sum of the confounding bias plus the non-collapsibility effect. Under certain circumstances, the confounding bias and the non-collapsibility effect can cancel each other out thus lead to the equivalent crude OR and conditional OR. This explains the possibility of observing confounding without non-collapsibility. Condition for the occurrence of this phenomenon is that L is independent of A conditional on Y (Greenland and Robins, 1999, Hernán et al, 2011).

3.3 The measure of the non-collapsibility effect

As demonstrated above, when collapsing over L, one needs to consider the noncollapsibility effect and confounding bias separately. We propose an approach to measure the two distinct components. The total discrepancy while collapsing the table over L is defined as the difference between the crude odds ratio (OR_{crude}) and the conditional odds ratio (OR_c). For a measure of the non-collapsibility effect, we need to compare a marginal model and a conditional model after adjusting for confounding. Odds ratios estimated with the marginal structural model (OR_m) have good properties (Hernán and Robins, 2006). It is a marginal measure of the effect of exposure and there is no confounding bias since the confounding is adjusted through the weights. Odds ratios estimated with the standard logistic regression model (SLRM) are conditional (stratumspecific) odds ratios (OR_c), and are also adjusted for confounding. Therefore, we can compare OR_m and OR_c to measure the non-collapsibility effect, and compare OR_{crude} and OR_m to measure the confounding bias. As Janes et al. did in their work (2010), if the measures are taken on the log scale, the decomposition can be written as

$$\log(OR_{crude}) - \log(OR_{c}) = [\log(OR_{m}) - \log(OR_{c})] + [\log(OR_{crude}) - \log(OR_{m})]$$
(6)

The decomposition will be shown in more examples in the following section.

4 The Non-collapsibility effect in a point-exposure study

In this chapter, we are going to describe two different approaches to study the noncollapsibility effect in a point-exposure study – an analytical approach and a graphical approach. The analytical approach is developed to derive the true value of the noncollapsibility effect as a function of the related parameters, while the graphical approach is used to present how the non-collapsibility effect is affected by the related parameters.

4.1 An analytical approach

4.1.1 Introduction

After realizing the non-collapsibility of the OR and the decomposition of the total discrepancy, it is also of great interest to measure the non-collapsibility effect in different scenarios. Samuels provided a formula to compute this in a setting where A is independent of L (Samuels, 1981), while Gail and Neuhaus provided an approximation of the non-collapsibility effect in a randomized trial and in clustered data respectively (Gail, 1984; Neuhuas, 1991). In this section, we intend to develop a formula for measuring the non-collapsibility effect in a general scenario. This formula can be used to assess the conditions of collapsibility of OR, show the decomposition of the total discrepancy explicitly and compare with the other related formulae in the literature.

Starting from the simplest scenario, figure 4-1 represents casual diagrams for pointexposure studies with exposure A, outcome Y, and covariate L. L is a baseline covariate on the left of figure 4-1, while L is a confounding variable on the right of figure 4-1.



Figure 4-1: Causal diagrams for the point-exposure studies

Table 4-1 shows the data in the point-exposure study with all individuals recruited from the entire population. Probability refers to proportions in the source population. We focus on the setting where all the variables are dichotomous with two levels, 0 and 1.

L=0				L=1	
	A=0	A=1		A=0	A=1
Y=0	а	b	Y=0	е	f
Y=1	С	d	Y=1	g	h

Table 4-1: Data from the point-exposure study

4.1.2 Method

4.1.2.1 Notation and related parameters

We have assumed a deterministic procedure, the observed count and expected value for each cell are just different ways of referring to the same quantity, as well as the true OR and the estimated OR. We use the notation and terminology of probability for convenience.

The outcome Y_i for each individual is a random variable following the Bernoulli distribution with parameter $P = P(Y_i = 1 | A = a_i, L = l_i)$. The probabilities of outcome conditional on A and L, the probability of A conditional on L, and the prevalence of L are the related parameters. Table 4-1 can be constructed from these parameters, and OR_m ,

 OR_{crude} , OR_{c} can be computed analytically.

Denote the conditional outcome probabilities by $P_{00} = P(Y = 1 | A = 0, L = 0)$,

$$P_{10} = P(Y = 1 | A = 1, L = 0), P_{01} = P(Y = 1 | A = 0, L = 1), \text{ and } P_{11} = P(Y = 1 | A = 1, L = 1).$$

Denote the conditional exposure probabilities by $P_{A0} = P(A = 1 | L = 0)$ and

 $P_{A1} = P(A = 1 | L = 1)$], and denote the prevalence of L by P_L . Denote $q_{00} = 1 - P_{00}$,

 $q_{A0} = 1 - P_{A0}$, and the analogous notations for q_{10} , q_{01} , q_{11} , and q_{A1} .

4.1.2.2 Homogenous odds ratio across L

Constraints were needed in order to guarantee two homogeneous conditional ORs across

strata of L. In the stratum L=0, $OR_{L=0} = \frac{(1 - P_{00}) \times P_{10}}{P_{00} \times (1 - P_{10})}$ and in the stratum L=1,

 $OR_{L=1} = \frac{(1 - P_{01}) \times P_{11}}{P_{01} \times (1 - P_{11})}. OR_{L=0} = OR_{L=1} \text{ indicates the effect of A on Y is identical across L.}$

Denote the common effect of A on Y across L by OR_c , then P_{10} can be written as a

function of P_{00} and OR_C : $P_{10} = \frac{P_{00} \times OR_C}{1 - P_{00} + P_{00} \times OR_C}$ (7)

The similar constraint can be made for P_{11} ,

$$P_{11} = \frac{P_{01} \times OR_C}{1 - P_{01} + P_{01} \times OR_C} \quad (8)$$

After fixing OR_c and constraining P_{10} and P_{11} by (7) and (8), we can ensure the common odds ratio OR_c across L.

4.1.2.3 Estimation of the odds ratios

 OR_c , OR_{crude} , OR_m can be computed analytically from table 4-1, and also can be estimated from the SLRM and the marginal structural model (MSM). As demonstrated in the introduction,

 OR_c corresponds to e^{α} estimated from the logistic regression model:

 $g[E(Y | A = a, L = l)] = \mu + \alpha A + \beta L$, where g is the logit link function.

 OR_{crude} corresponds to e^{α^*} estimated from the logistic regression model:

 $g[E(Y \mid A = a)] = \mu^* + \alpha^* A$

 OR_m corresponds to $e^{\alpha'}$ estimated by the MSM:

 $g[E(Y | A = a)] = \mu' + \alpha' A$ weighting each subject by the inverse of the probability of receiving his/her exposure conditional on L.
4.1.3 Result

4.1.3.1 Measure of non-collapsibility of the OR

The expected value or observed value in each cell in Table 4-1 can be computed using the conditional outcome probabilities, the conditional exposure probabilities and the prevalence of L. For example,

$$a = P(Y = 0 | A = 0, L = 0) \times P(A = 0 | L = 0) \times P(L = 0) * N = q_{00} \times q_{A0} \times (1 - P_L) \times N,$$

where N is the sample size. The same calculations can be done for the observed count in the other cells.

It is easy to show that:

$$OR_{crude} = \frac{[q_{00} \times q_{A0} \times (1 - P_L) + q_{01} \times q_{A1} \times P_L] \times [P_{10} \times P_{A0} \times (1 - P_L) + P_{11} \times P_{A1} \times P_L]}{[q_{10} \times P_{A0} \times (1 - P_L) + q_{11} \times P_{A1} \times P_L] \times [P_{00} \times q_{A0} \times (1 - P_L) + P_{01} \times q_{A1} \times P_L]}$$

$$OR_{C} = \frac{P_{10} \times q_{00}}{P_{00} \times q_{10}} = \frac{P_{11} \times q_{01}}{P_{01} \times q_{11}}$$

$$OR_m = \frac{(a \times w_a + e \times w_e) \times (d \times w_d + h \times w_h)}{(b \times w_b + f \times w_f) \times (c \times w_c + g \times w_g)} \text{ where } w_a = w_c = \frac{1}{P(A = 0 \mid L = 0)} = \frac{1}{q_{A0}};$$

$$w_b = w_d = \frac{1}{P(A=1 \mid L=0)} = \frac{1}{P_{A0}}; \ w_e = w_g = \frac{1}{P(A=0 \mid L=1)} = \frac{1}{q_{A1}};$$
$$w_f = w_h = \frac{1}{P(A=1 \mid L=1)} = \frac{1}{P_{A1}};$$

 OR_m can be written as:

$$OR_{m} = \frac{[q_{00} \times q_{A0} \times (1-P_{L}) \times \frac{1}{q_{A0}} + q_{01} \times q_{A1} \times P_{L} \times \frac{1}{q_{A1}}] \times [P_{10} \times P_{A0} \times (1-P_{L}) \times \frac{1}{P_{A0}} + P_{11} \times P_{A1} \times P_{L} \times \frac{1}{P_{A1}}]}{[q_{10} \times P_{A0} \times (1-P_{L}) \times \frac{1}{P_{A0}} + q_{11} \times P_{A1} \times P_{L} \times \frac{1}{P_{A1}}] \times [P_{00} \times q_{A0} \times (1-P_{L}) \times \frac{1}{q_{A0}} + P_{01} \times q_{A1} \times P_{L} \times \frac{1}{q_{A1}}]}$$
$$= \frac{[q_{00} \times (1-P_{L}) + q_{01} \times P_{L}] \times [P_{10} \times (1-P_{L}) + P_{11} \times P_{L}]}{[q_{10} \times (1-P_{L}) + q_{11} \times P_{L}] \times [P_{00} \times (1-P_{L}) + P_{01} \times P_{L}]}$$

As illustrated by the equation (6), $log(\frac{OR_m}{OR_c})$ was used to measure the magnitude of the

non-collapsibility effect.

4.1.3.2 Conditions for the absence of the non-collapsibility effect

A historically important contribution to this literature was the work of Gail et al. (1984). Corollary 1 in the paper demonstrated that there is no non-collapsibility effect if L does not vary, or $\alpha = 0$ or if $\beta = 0$. These conclusions can be verified with our analytical approach.

1. No non-collapsibility effect if L does not vary

L does not vary implies the prevalence of L is 0 or 1. If the prevalence of L is 0, then

$$OR_m = \frac{q_{00} \times P_{10}}{q_{10} \times P_{00}} = OR_c; \text{ If the prevalence of L is 1, then } OR_m = \frac{q_{01} \times P_{11}}{q_{11} \times P_{01}} = OR_c;$$
$$\log(\frac{OR_m}{OR_c}) = 0.$$

In this case, we just have one stratum of L, either 0 or 1. There is no L covariate variation. Logically, and trivially, therefore, there is no non-collapsibility effect.

2. No non-collapsibility effect if $\alpha = 0$

 $\alpha = 0$ indicates that the effect of A on Y conditional on L is 0 ($Y \perp A \mid L$), it follows that $P_{00} = P_{10}$ and $P_{01} = P_{11}$ then

$$OR_{m} = \frac{[q_{10} \times (1 - P_{L}) + q_{11} \times P_{L}] \times [P_{00} \times (1 - P_{L}) + P_{01} \times P_{L}]}{[q_{10} \times (1 - P_{L}) + q_{11} \times P_{L}] \times [P_{00} \times (1 - P_{L}) + P_{01} \times P_{L}]} = 1; \quad OR_{C} = 1, \text{ and}$$

 $log(\frac{OR_m}{OR_c}) = 0$. The non-collapsibility effect is 0 if $\alpha = 0$ or the conditional effect of A

on Y is 0.

3. No non-collapsibility effect if $\beta = 0$

 $\beta = 0$ implies that the effect of L on Y conditional on A is 0 ($Y \perp L \mid A$), it follows $P_{00} = P_{01}$ and $P_{10} = P_{11}$, then

$$OR_{m} = \frac{[q_{10} \times (1 - P_{L}) + q_{10} \times P_{L}] \times [P_{11} \times (1 - P_{L}) + P_{11} \times P_{L}]}{[q_{11} \times (1 - P_{L}) + q_{11} \times P_{L}] \times [P_{01} \times (1 - P_{L}) + P_{01} \times P_{L}]} = \frac{q_{01} \times P_{11}}{q_{11} \times P_{01}} = OR_{C}, \text{ and}$$

$$log(\frac{OR_m}{OR_c}) = 0$$
. The non-collapsibility effect is 0 if $\beta = 0$ or L independent of Y

conditional on A.

4. Confounding without non-collapsibility

The analytical approach can also be used to illustrate the situation where confounding is present without non-collapsibility. It has been noted that the crude OR equals the

conditional OR if L is independent of A conditional on Y (Wermuth, 1987; Newman 2011). $L \perp A \mid Y$ indicates that

$$P(A=1 \mid L=0, Y=0) = P(A=1 \mid L=1, Y=0)$$
(9)

$$P(A=1 \mid L=0, Y=1) = P(A=1 \mid L=1, Y=1)$$
(10)

$$P(A = 0 \mid L = 1, Y = 0) = P(A = 0 \mid L = 0, Y = 0)$$
(11)

$$P(A = 0 | L = 1, Y = 1) = P(A = 0 | L = 0, Y = 1)$$
(12)

Denote $P_{Y_1} = P(Y = 1 | L = 1)$ and $P_{Y_0} = P(Y = 1 | L = 0)$. The conditional probabilities can be written as following:

$$P(A = 1 | L = 0, Y = 0) = \frac{P(A = 1, L = 0, Y = 0)}{P(Y = 0 | L = 0) \times P_L} = \frac{P(Y = 0 | A = 1, L = 0) \times P(A = 1 | L = 0) \times P_L}{P(Y = 0 | L = 0) \times P_L}$$
$$= \frac{q_{10} \times P_{A0}}{q_{Y0}}$$

$$P(A=1 | L=1, Y=0) = \frac{q_{11} \times P_{A1}}{q_{Y1}};$$

$$P(A=1 | L=0, Y=1) = \frac{P_{10} \times P_{A0}}{P_{Y0}};$$

$$P(A=1 | L=1, Y=1) = \frac{P_{11} \times P_{A1}}{P_{Y1}};$$

$$P(A=0 \mid L=0, Y=0) = \frac{q_{00} \times q_{A0}}{q_{Y0}};$$

$$P(A=0 | L=1, Y=0) = \frac{q_{01} \times q_{A1}}{q_{Y1}};$$

$$P(A=0 | L=0, Y=1) = \frac{P_{00} \times q_{A0}}{P_{Y0}};$$

$$P(A = 0 | L = 1, Y = 1) = \frac{P_{01} \times q_{A1}}{P_{Y1}}.$$

By (9)-(12), we have
$$\frac{q_{10} \times P_{A0}}{q_{Y0}} = \frac{q_{11} \times P_{A1}}{q_{Y1}}; \frac{P_{10} \times P_{A0}}{P_{Y0}} = \frac{P_{11} \times P_{A1}}{P_{Y1}}; \frac{q_{00} \times q_{A0}}{q_{Y0}} = \frac{q_{01} \times q_{A1}}{q_{Y1}};$$

$$\frac{P_{00} \times q_{A0}}{P_{Y0}} = \frac{P_{01} \times q_{A1}}{P_{Y1}}; \text{ then, } q_{00} = \frac{q_{01} \times q_{A1} \times q_{Y0}}{q_{Y1} \times q_{A0}}; P_{10} = \frac{P_{11} \times P_{A1} \times P_{Y0}}{P_{Y1} \times P_{A0}};$$

$$q_{10} = \frac{q_{11} \times P_{A1} \times q_{Y0}}{q_{Y1} \times P_{A0}}; P_{00} = \frac{P_{01} \times q_{A1} \times P_{Y0}}{P_{Y1} \times q_{A0}}.$$
 Consequently, substituting the above

equations into the calculation of the crude OR, OR_{crude} can be computed as:

$$\begin{split} OR_{crude} &= \frac{\left[q_{00} \times q_{A0} \times (1 - P_L) + q_{01} \times q_{A1} \times P_L\right] \times \left[P_{10} \times P_{A0} \times (1 - P_L) + P_{11} \times P_{A1} \times P_L\right]}{\left[q_{10} \times P_{A0} \times (1 - P_L) + q_{11} \times P_{A1} \times P_L\right] \times \left[P_{00} \times q_{A0} \times (1 - P_L) + P_{01} \times q_{A1} \times P_L\right]}{\left[\frac{q_{01} \times q_{A1} \times q_{Y0}}{q_{Y1}} \times (1 - P_L) + q_{01} \times q_{A1} \times P_L\right] \times \frac{P_{11} \times P_{A1} \times P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_{11} \times P_{A1} \times P_L\right]} \\ &= \frac{q_{01} \times q_{A1} \times \left[\frac{q_{Y0}}{q_{Y1}} \times (1 - P_L) + q_{11} \times P_{A1} \times P_L\right] \times \left[\frac{P_{01} \times q_{A1} \times P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_{01} \times q_{A1} \times P_L\right]}{\left[q_{11} \times P_{A1} \times \left[\frac{q_{Y0}}{q_{Y1}} \times (1 - P_L) + P_L\right] \times P_{11} \times P_{A1} \times \left[\frac{P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_L\right]} \times \left[\frac{q_{01} \times q_{A1} \times P_L}{P_{Y1}}\right]}{\left[q_{11} \times P_{A1} \times \left[\frac{q_{Y0}}{q_{Y1}} \times (1 - P_L) + P_L\right] \times P_{01} \times q_{A1} \times \left[\frac{P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_L\right]} \\ &= \frac{q_{01} \times q_{A1} \times \left[\frac{q_{Y0}}{q_{Y1}} \times (1 - P_L) + P_L\right] \times P_{01} \times q_{A1} \times \left[\frac{P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_L\right]}{\left[q_{11} \times P_{A1} \times \left[\frac{q_{Y0}}{q_{Y1}} \times (1 - P_L) + P_L\right] \times P_{01} \times q_{A1} \times \left[\frac{P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_L\right]} \\ &= \frac{q_{01} \times P_{A1}}{q_{11} \times P_{A1}} \times \left[\frac{q_{Y0}}{q_{Y1}} \times (1 - P_L) + P_L\right] \times P_{01} \times q_{A1} \times \left[\frac{P_{Y0}}{P_{Y1}} \times (1 - P_L) + P_L\right]} \\ &= \frac{q_{01} \times P_{A1}}{q_{11} \times P_{01}} = OR_C \end{split}$$

Under the circumstance that the covariate is independent of exposure conditional on the outcome, the total discrepancy equals 0, $log(\frac{OR_m}{OR_c}) = 0$.

4.1.3.3 Decomposition of the total discrepancy

Table 4-2 shows the decomposition of the total discrepancy by the analytical approach. From the equation (6), the total discrepancy was computed by $\log(OR_{crude}) - \log(OR_{c})$, the non-collapsibility effect was computed by $\log(OR_m) - \log(OR_c)$, and the confounding bias was computed by $\log(OR_{crude}) - \log(OR_m)$.

P(Y = 1 A = 0, L = 0)	P(Y = 1 A = 0, L = 1)	OR _c	OR _m	OR _{crude}	<i>P</i> (<i>L</i> = 1)	P_{A1}	P_{A0}	Non- collapsibility effect	Confounding bias	Total discrepancy
0.2	0.6	2.667	2.253	2.253	0.45	0.500	0.500	-0.169	0	-0.169
0.2	0.6	2.667	2.253	2.667	0.45	0.556	0.455	-0.169	0.169	0
0.2	0.9	2.667	1.764	2.302	0.45	0.556	0.455	-0.413	0.266	-0.147

Table 4-2: Decomposition of the total discrepancy

Table 4-3: Data correspond to the first row of table 4-2

	L=0		L=1	
	A=0	A=1	A=0	A=1
Y=0	220	165	90	45
Y=1	55	110	135	180
	275	275	225	225

	L=0		L=1	
	A=0	A=1	A=0	A=1
Y=0	240	150	80	50
Y=1	60	100	120	200
	300	250	200	250

Table 4-4: Data correspond to the second row of table 4-2

Table 4-5: Data correspond to the third row of table 4-2

	L=0		L=1	
	A=0	A=1	A=0	A=1
Y=0	240	150	20	10
Y=1	60	100	180	240
	300	250	200	250

Table 4-3, table 4-4 and table 4-5 present data in each row of table 4-2 with sample size N=1000. All the conditional ORs are homogeneous and equal to 2.667. Table 4-3 shows the data that $P_{A1} = P_{A0}$, A is independent of L, indicating no confounding. The first row of table 4-2 shows the confounding bias equals 0, and the total discrepancy equals the non-collapsibility effect. Table 4-4 shows data in which A is independent of L conditional on Y, which is the condition that the total discrepancy is 0. $P_{A1} \neq P_{A0}$ in table 4-4 implies confounding, and the second row of table 4-2 shows that the non-collapsibility effect and confounding bias cancel each other out, thus leading to an equivalent crude OR and conditional OR in the presence of confounding. But this is not generally true. Table 4-5 shows data in which the non-collapsibility effect and the confounding bias both exist. The non-collapsibility effect and the confounding bias do not cancel each other out in the third row of table 4-2. It is clearly evident in table 4-2

that the total discrepancy is the sum of the non-collapsibility effect and the confounding bias.

4.1.3.4 Comparison of non-collapsibility effect with other formulae

Other formulae to measure the non-collapsibility effect can also be found in the literature, i.e. formula by Samuels (1981), equation (8) by Neuhaus et al. (1991) and equation (2.9) by Gail et al (1984). As discussed in the literature review, Samuels developed the formula in a model in which A and L are unconditionally independent. Neuhaus considered non-collapsibility in the context of cluster-specific and population-averaged approaches for correlated binary data. In our setting, there are just two clusters, and L is the variable that differentiates clusters and acts as the random effect in the context of Neuhaus' work. The non-collapsibility effect was measured by $-\beta \frac{Var(p)}{E(p)E(q)}$, where logit(p) takes the value $\mu + \beta$ with probability P_L and value μ with probability 1- P_L , and q = 1 - p. Gail et al. assumed that the covariate is continuous with expectation 0 and variance σ^2 , while in our

setting, the covariate is a binary variable. Hence we needed to modify Gail's formula to compare the results. It can be shown in the appendix that when the covariate is binary, the non-collapsibility effect measured by Gail's formula is

$$\frac{1}{4}\beta^2 E(L^2)\left\{\frac{h''(\mu+\alpha)}{h'(\mu+\alpha)}-\frac{h''(\mu-\alpha)}{h'(\mu-\alpha)}\right\}.$$

Table 4-6 presents the comparison of results by our analytical approach and Samuels, Neuhaus and Gail formulae under various scenarios. The non-collapsibility effect by our analytical approach is equivalent with the formula by Samuels, and is fairly close with the formula by Neuhaus. All of the four approaches give us the same qualitative findings.

P_{A0}	P_{A1}	$P(Y = 1 \mid A = 0, L = 0)$	P(Y = 1 A = 0, L = 1)	OR_{C}	P(L = 1)	$\log(\frac{OR_m}{OR_c})$	Samuels	Neuhaus	Gail
0.5	0.5	0.1	0.5	5	0.3	-0.3438	-0.3438	-0.3151	-0.243
0.5	0.5	0.4	0.5	2	0.3	-0.0058	-0.0058	-0.0059	-0.0079
0.5	0.5	0.2	0.4	1	0.45	0	0	0	0
0.5	0.5	0.5	0.8	0.25	0.2	0.0992	0.0992	0.081	0.1153
0.5	0.5	0.5	0.6	0.5	0.7	0.0058	0.0058	0.0059	0.01918
0.5	0.5	0.75	0.5	0.2	0.6	0.1054	0.1054	0.1006	0.2037
0.45	0.55	0.20	0.60	2.67	0.45	-0.169	-0.169	-0.165	-0.227
0.5	0.2	0.40	0.50	1	0.4	0	0	0	0
0.3	0.6	0.80	0.50	0.2	0.40	0.163	0.163	0.160	0.195

Table 4-6: Some example results of comparison of non-collapsibility effect with other formulae

4.1.4 Discussion

It is not appropriate to use OR to detect confounding bias conventionally due to the noncollapsibility effect. However, one can use our analytical approach to measure the confounding bias and non-collapsibility effect separately if all confounders are measured. If a simpler approach is preferred, one can also just avoid using the OR to define the confounding but use a collapsible measure such as the RR or RD instead, since the conventional way of detecting non-confounding agrees with the collapsibility of the RR and RD. As we know, when the outcome is rare, the OR approximates the RR, and thus can be also used to detect confounding validly.

We provide the same results as Samuels, mainly because the two formulae are closely related. Though Samuels developed the formula in the model that A is independent of L, after extending it to a more general case, it agrees with our results. Gail's results do not agree as closely with the others, presumably because it is only the first term of a Taylor

expansion, and thus makes the approximation less precise. Although they provide the same or similar results, all of the other three approaches concentrate on the situation where A is independent of L, while our approach gives a general formula to measure the non-collapsibility effect with or without confounding.

4.2 A graphical approach

4.2.1 Introduction

There are some graphical approaches for the non-collapsibility of OR as well. For example, Shapiro (1982) used a coordinate system to represent the odds and revealed the non-collapsibility of OR and the conditions of collapsibility of the OR. Some conclusions were made, stating that the non-collapsibility depends on the effects of the exposure and covariate on outcome and on the variance of the covariate (Ritz and Spiegelman 2004; Groenwold and Moons et al., 2011).

The magnitude of confounding usually requires that values for three types of parameters be specified: the prevalence of the covariate in the population or subpopulation, the association between the exposure and the covariate and the effect of the covariate on the outcome. Flanders and Khoury (1990) illustrated the relationship between the magnitude of confounding and these parameters graphically, and derived limits of the magnitude if some of the relevant parameters cannot be specified.

Inspired by graphical approaches such as Flanders and Khoury (1990) to represent relations between the parameters related to the magnitude of confounding, we will investigate the relationship between the non-collapsibility effects with the related parameters in different scenarios by using figures.

4.2.2 Method: relationship between non-collapsibility and the effect of A and L on Y

In practice, one might be also interested in the relationship between the non-collapsibility effect and the effect of A (OR_c) and the effect of L (OR_L) on Y.

It can be shown in the appendix that P_{01} can be expressed as a function of the baseline risk P_{00} , the effect of A and L on Y, the prevalence of L, and the exposure probability conditional on L. Therefore, we have the relationship between the non-collapsibility effect and all the other parameters. Again, $\log(\frac{OR_m}{OR_c})$ was used to measure the non-

collapsibility effect. Since there is no simple form of the non-collapsibility effect function of the relevant parameters, a graphical approach was implemented. The non-collapsibility behavior can be explored in different scenarios by plotting the effect with the baseline risk and different combinations of values of the other parameters. A range of different exposure effects and different covariate effects was specified. A strong association was defined by OR of 5 or 0.2, while a moderate association was defined by OR of 2 or 0.5. The value >1 represents a harmful effect, while the value <1 represents a protective effect.

4.2.3 Results

4.2.3.1 Scenario without confounding

Figure 4-2 shows the relationship between the non-collapsibility effect and the baseline risk (P_{00}). It shows the scenario in which there is no confounding, L is a baseline covariate, A is randomly assigned with probability 0.5 regardless of L (P(A = 1 | L = 1) = P(A = 1 | L = 0) = 0.5), and the prevalence of L is 0.5.



Figure 4-2: Non-collapsibility effect vs. the baseline risk with no confounding Figures for other values of the prevalence of L are similar to Figure 4-2. But the noncollapsibility effect tends to be milder than that at $P_L = 0.5$, which implies that one will observe a relatively larger non-collapsibility effect when the prevalence of L is closer to 0.5.

In Figure 4-2, the solid lines present the non-collapsibility effect when the effects of A on Y are positive, while the dashed lines present the non-collapsibility effect when the

effects of A are negative. It is shown that all the dashed lines are above 0, and all the solid lines are below 0. This corresponds to the well-known fact that the marginal OR is always shifted towards the null compared to the conditional OR (Samuels, 1981). Therefore, the non-collapsibility effect is negative when the A effect is harmful, and the non-collapsibility effect is positive when the A effect is protective.

The non-collapsibility effect was observed to be symmetric in Figure 4-2. Each pair of solid line and dashed line with the same color shows the non-collapsibility effect with the opposite A and L effect with the same magnitude. The shapes of each pair of lines are the same, which indicates that the magnitude of the non-collapsibility effect is the same but

in a different direction if
$$OR_C(1) = \frac{1}{OR_C}(2)$$
, $OR_L(1) = \frac{1}{OR_L}(2)$, and $P_{00}(1) = 1 - P_{00}(2)$.

For example, the non-collapsibility effect is the same but in different direction for $OR_C(1) = 5$, $OR_L(1) = 5$, $P_{00}(1) = 0.2$ and $OR_C(1) = 0.2$, $OR_L(1) = 0.2$, $P_{00}(1) = 0.8$.

We can observe that the black lines are further from 0 compared to the red lines, indicating that after fixing the effect of A on Y, the non-collapsibility effect is smaller as the effect of L on Y gets smaller. Whereas, the black lines are further from 0 compared to the light blue lines, indicating that after fixing the effect of L on Y, the non-collapsibility effect is smaller as the effect of A on Y gets smaller. If we compare the solid light blue line ($OR_C = 2$, $OR_L = 5$) to the solid red line ($OR_C = 5$, $OR_L = 2$), it is observed that the light blue line is further form 0 compared to the red one, which indicates the novel finding that the effect of L plays a more important role in the non-collapsibility effect compared to the effect of A. The non-collapsibility behavior is quite modest when the effect of L and the effect of A are both small, by observing that the yellow lines and the purple lines are very close to the x axis.

It can also be observed that as the baseline risk goes to 0, the non-collapsibility effect disappears. It is well known that when the risk is small the OR can be used to estimate the RR, and there is no non-collapsibility effect in a non-founding scenario when the measure of association is RR. However, it is less well known that as the baseline risk goes to 1, the non-collapsibility effect also tends to disappear.

When both A and L effects are protective, the most extreme non-collapsibility effect is present when P_{00} is around 0.9. This is quite sensible, since both A and L have negative effects, the other risks are smaller than the baseline risk. Only if the baseline risk is high can the other risks have moderate values. When the baseline risk is small, the other risks are even smaller, which makes the non-collapsibility effect trivial (vice versa when A and L are both harmful). We infer that when the four conditional outcome probabilities are all moderate, we can have a relatively large non-collapsibility effect. However when the four conditional outcome probabilities are all very small or very large, the non-collapsibility effect with different conditional outcome probabilities as well as the marginal outcome probabilities (P_{γ}) after fixing the A and L effect. The third and the fourth row show that the non-collapsibility effects are relatively small when the outcome probabilities are all large or small.

38

A effect	L effect	P_{00}	P ₁₀	<i>P</i> ₀₁	<i>P</i> ₁₁	P_{Y}	$\log(\frac{OR_m}{OR_C})$
5	5	0.1	0.357	0.414	0.78	0.413	-0.274
5	5	0.5	0.833	0.852	0.966	0.788	-0.149
5	5	0.8	0.952	0.954	0.991	0.924	-0.048
5	0.2	0.1	0.357	0.02	0.092	0.142	-0.096
5	0.2	0.5	0.833	0.134	0.437	0.476	-0.288
5	0.2	0.8	0.952	0.402	0.77	0.731	-0.191

Table 4-7: Non-collapsibility effect with different conditional outcome probabilities

The relationship between the marginal outcome probability (P_{y}) and the non-

collapsibility effect is shown in Figure 4-3.



P.L=0.5, no confounding, A assigned randomly with p=0.5

Figure 4-3: Non-collapsibility effect vs. the marginal risk (1)

Figure 4-3 shows the scenario in which L is a baseline covariate, the prevalence of L is 0.5, and A is randomly assigned regardless of L with probability 0.5. In this scenario the non-collapsibility effect is symmetric about $P_Y = 0.5$ and the x axis. There are pairs of lines overlapping, for example the blue lines and the black lines.

The non-collapsibility effect is the same for two identical A effects if the magnitude of the L effect is the same, regardless of the direction. For example, the non-collapsibility effect is the same for $OR_C(1) = 5$, $OR_L(1) = 0.5$ and $OR_C(2) = 5$, $OR_L(2) = 2$. For any two A effects with the same magnitudes but different directions (e.g. $OR_C(1) = 0.5$ vs. $OR_C(2) = 2$), the magnitude of the non-collapsibility effect will be the same but in the other direction. This is true for any fixed magnitude of the L effect, regardless of the direction of the L effect. For example, the magnitude of the non-collapsibility effect is the same for $OR_C(1) = 5$, $OR_L(1) = 2$; $OR_C(2) = 0.2$, $OR_L(2) = 2$; and

$$OR_{C}(3) = 0.2, OR_{L}(3) = 0.5$$

Similarly to what we have illustrated above, after fixing one of the A or L effects, the non-collapsibility effect is smaller as the magnitude of the other effect gets smaller. Again we find the somewhat surprising result that the effect of L plays a more important role in the non-collapsibility effect compared to the effect of A.



Figure 4-4: Non-collapsibility effect vs. the marginal risk (2)





Figure 4-5: Non-collapsibility effect vs. the marginal risk (3)

Figure 4-4 and Figure 4-5 show the scenario in which the prevalence of L is 0.3 or 0.7. There are no lines overlapping in the figures, but the symmetric property remains. It seems that when there is no confounding, after fixing the A effect, the prevalence of L determines whether the non-collapsibility has the same behavior for $OR_L(1) = \frac{1}{OR_L}(2)$. When the A effect is fixed, the non-collapsibility effect is symmetric about $P_Y = 0.5$ if the L effects are in the opposite direction with the same magnitude. That is to say, the non-collapsibility effect is the same for $OR_L(1) = \frac{1}{OR_L}(2)$ and $P_Y(1) = 1 - P_Y(2)$ after

fixing the A effect. And we still have the observation that if $OR_C(1) = \frac{1}{OR_C}(2)$,

$$OR_L(1) = \frac{1}{OR_L}(2)$$
, and $P_Y(1) = 1 - P_Y(2)$, the non-collapsibility effect has the same

magnitude in the opposite direction.

P.L=0.5, no confounding, A assigned randomly with p=0.67



Figure 4-6: Non-collapsibility effect vs. the marginal risk (4)



P.L=0.5, no confounding, A assigned randomly with p=0.33

Figure 4-7: Non-collapsibility effect vs. the marginal risk (5)

Figure 4-6 and Figure 4-7 show scenarios in which A is randomly assigned (i.e., unconfounded) with probability not equal to 0.5. The non-collapsibility effect is not symmetric about $P_Y = 0.5$ in these figures. It seems that when there is no confounding, after fixing the A effect, the prevalence of A determines whether the non-collapsibility effect is symmetric about $P_{\gamma} = 0.5$ when the L effects are in the opposite direction with the same magnitude. But there are still pairs of lines overlapping. If the prevalence of L is 0.5, then the non-collapsibility effect is identical for the same L effect regardless of its direction after fixing the A effect.





Figure 4-8: Non-collapsibility effect vs. the baseline risk with confounding Figure 4-8 shows the relationship between the non-collapsibility effect and the baseline risk in different combinations of A and L effects and the effect of L on A (denoted by OR_{LA}). The non-collapsibility effect is again symmetric. It is the same but in a different direction for $OR_{c}(1) = \frac{1}{OR_{c}}(2)$, $OR_{L}(1) = \frac{1}{OR_{L}}(2)$, $OR_{LA}(1) = OR_{LA}(2)$, and

 $P_{00}(1) = 1 - P_{00}(2)$.

An interesting and novel finding can be observed that if the A effect and L effect on Y are both in the same direction, the non-collapsibility effect is larger if the effect of L on A is negative than when the effect of L on A is positive with the same magnitude. And the more the effect of L on A tends to be negative, the larger the non-collapsibility effect is. For example, the non-collapsibility effect is decreasing for $OR_C(1) = 0.2$,

$$OR_{L}(1) = 0.2, OR_{LA}(1) = 0.2; OR_{C}(2) = 0.2, OR_{L}(2) = 0.2, OR_{LA}(2) = 0.5;$$

$$OR_{C}(3) = 0.2, OR_{L}(3) = 0.2, OR_{LA}(3) = 2$$
; and $OR_{C}(4) = 0.2$,

 $OR_L(4) = 0.2$, $OR_{LA}(4) = 5$. When the A effect and L effect on Y are in different directions, the non-collapsibility effect is larger if the effect of L on A is positive than that when the effect of L on A is negative with the same magnitude. And the more the effect of L on A tends to be positive, the larger the non-collapsibility effect is. For example, the non-collapsibility effect is decreasing for $OR_C(1) = 5$,

 $OR_{L}(1) = 0.2, OR_{LA}(1) = 5; OR_{C}(2) = 5, OR_{L}(2) = 0.2, OR_{LA}(2) = 2; OR_{C}(3) = 5,$ $OR_{L}(3) = 0.2, OR_{LA}(3) = 0.5; OR_{C}(4) = 5, OR_{L}(4) = 0.2, OR_{LA}(4) = 0.2.$

P.L=0.5, with confounding



Figure 4-9: Non-collapsibility effect vs. the marginal risk with confounding

Figure 4-9 shows the relationship between the non-collapsibility effect with the marginal outcome probability in different combinations of A and L effects and the effect of L on A. All the observations from Figure 4-8 appear in Figure 4-9. In addition, the non-

collapsibility effect is the same if $P_Y(1) = 1 - P_Y(2)$, $OR_L(1) = \frac{1}{OR_L}(2)$, and

 $OR_{LA}(1) = \frac{1}{OR_{LA}}(2)$ after fixing the A effect. For example, the non-collapsibility is the same for $OR_C(1) = 5$, $OR_L(1) = 5$, $OR_{LA}(1) = 0.2$, and $P_Y(1) = 0.2$, and $OR_C(2) = 5$, $OR_L(2) = 0.2$, $OR_{LA}(2) = 5$, and $P_Y(2) = 0.8$.

4.2.4 Discussion

It is logical to think that the L effect and A effect do not play an identical role in the noncollapsibility effect. The A effect is the effect of exposure and is assumed to be homogenous across the strata, whereas L is the variable that we collapse over. In each scenario, it indicates that the L effect does have a larger influence on the noncollapsibility effect.

We presents the relationship between the non-collapsibility effect and the baseline risk or marginal outcome probability in a two dimensional figure while fixing the values of the other parameters. Since there are four additional parameters, it is much more difficult to present non-collapsibility with all the parameters (every possible value) in one figure. Further investigation is needed about how to make the figure and present the noncollapsibility effect in a more vivid and perceptually intuitive way. 5 The non-collapsibility in the presence of time-varying confounding

5.1 Introduction

In some observational studies, exposure may vary over time, and a confounding variable that is also affected by the previous exposure can be also encountered frequently. If the exposure and other factors of the individuals are measured multiple times during the follow-up, in the presence of time-varying confounding, will the non-collapsibility effect change? It appears difficult to use the analytical approach in this relatively complex scenario. Observational cohort studies were simulated to explore the non-collapsibility effect in these scenarios.

5.2 Method

5.2.1 Time-varying confounding scenario

To make sure everything is well understood, we started from the simplest scenario in the presence of time-varying confounding that is most like our point-exposure study. Figure 5-1 shows the scenario we want to explore first, assuming there are two time points, no unmeasured confounding and no loss to follow-up. Let A_0 denote the exposure at the first time point, let A_1 denote the exposure at the second time point and let Y denote the outcome. L is a time-varying confounding variable in the Figure 5-1. It is an intermediate variable when we measure the A_0 effect, so we must not adjust for L. Meanwhile it is a common cause of the exposure and the outcome that confounds the A_1 effect on Y,

therefore we must adjust for L to obtain the unbiased estimate of the A₁ effect. MSM can be employed to assess the exposure effect in this longitudinal setting with inverseprobability-of-treat weighted (IPTW) estimators and it provides a marginal exposure effect. Nevertheless, it is always a biased estimation of the cumulative treatment effect of A₀ and A₁ with the SLRM. There is therefore no good way to research the noncollapsibility effect with respect to L. But we can add another baseline variable Z which has an effect on the outcome but independent of any other variables. This makes the diagram (Figure 5-2) be analogous to the one on the left of Figure 4-1 in spite of the timevarying confounding L. Therefore we can study the non-collapsibility effect of Z and compare it to the point-exposure study.

In order to understand the models and interpret the results better, we first assume the direct effects of A_0 and A_1 on Y are both null. The simulation and the analysis were firstly implemented in the scenario without the Z variable.



Figure 5-1: Time-varying confounding scenario



Figure 5-2: Time-varying confounding scenario with baseline variable Z (Scenario 1) 5.2.1.1 Data Generation

We simulated an observational cohort study with N=100,000 subjects who were randomly assigned to be exposed at the first time point with probability 0.5, i.e. $P(A_0 = 1) = 0.5$ (A₀=1 indicates exposed and A₀=0 indicates unexposed). A₀ was generated from a Bernoulli distribution with $P(A_0 = 1) = 0.5$. L is a factor that is affected by A₀ and also has an effect on A₁, so L was generated from a Bernoulli distribution with $P(L=1) = \exp(\alpha_0 + \alpha_1 A_0)$, where $\exp(x) = \frac{\exp(x)}{1 + \exp(x)}$. α_1 should be the unbiased estimate of the effect of A₀ on L. Likewise, A₁ was generated from a Bernoulli distribution with $P(A_1 = 1) = \exp((\beta_0 + \beta_1 L))$, β_1 should be the unbiased estimate of the effect of L on A1. Since we assume that there is no direct effect of A0 and A1 on Y, Y was generated from a Bernoulli distribution with $P(Y = 1) = \exp((\gamma_0 + \gamma_1 L))$. α_0 , β_0 and γ_0 are the parameters which are not really important to us, hence we are willing to assume arbitrarily that $\alpha_0 = \beta_0 = \gamma_0 = 0$. We set α_1 , β_1 and γ_1 to be log(5), log(2), log(0.5) or log(0.2), and ran the simulation with all the combinations of those values, so as to illustrate the scenario overall with strong or moderate effects and with harmful or protective effects.

The datasets were generated according to the parameters settings described previously. For each setting, we generated 100 independent random samples. All simulations were performed in R version 2.14.1 running on a Linux platform.

5.2.1.2 Data Analysis

We proposed 7 models (Model 1 to Model 7) to analyze the data. Model 1 and model 2 provide the estimates of the effect of A_0 by adjusting L marginally and conditionally. As discussed above, one should not adjust for L which is on the causal pathway between A_0 and Y. It implies that model 3 which contains only A_0 in the model should provide the unbiased estimate of the total effect of A_0 on Y through L, while model 1 and model 2 yield biased estimates of the A_0 effect. We can obtain unbiased estimates of the marginal and conditional A_1 effects by adjusting L with MSM and SLRM from model 4 and model 5 respectively. Model 6 provides a crude and biased estimate of the A_1 effect. Model 7 is a MSM that can be used to estimate the cumulative effect of the exposure, where L was adjusted by the IPTW estimators. The cumulative effect of A_0 and A_1 is denoted by A which equals the value of A_0+A_1 deterministically.

Model 1: $\operatorname{logit}(P(Y = 1)) = \omega_0 + \omega_1 A_0$ Model 2: $\operatorname{logit}(P(Y = 1)) = \omega_0 + \omega_1 A_0 + \omega_2 L$ Model 3: $\operatorname{logit}(P(Y = 1)) = \omega_0 + \omega_1 A_0$ Model 4: $\operatorname{logit}(P(Y = 1)) = \omega_0 + \omega_1 A_1$ $\operatorname{model} 5: \operatorname{logit}(P(Y = 1)) = \omega_0 + \omega_1 A_1 + \omega_2 L$ Model 6: $\operatorname{logit}(P(Y = 1)) = \omega_0 + \omega_1 A_1$

Model 7: logit
$$(P(Y = 1)) = \omega_0 + \omega_1 A$$
 weights $= \frac{P(A_1 | A_0)}{P(A_1 | A_0, L)}$

5.2.1.3 Results

We repeated the simulation and analysis 100 times. The averages of the 100 log odds ratios are considered to be estimates of the exposure effects from those models. Table 5-1 presents the exponential of the means under different parameter settings. The results are estimates of the ORs.

e^{α_1}	e^{β_1}	e^{γ_1}	OR _{Model 1}	OR _{Model 2}	OR _{Model 3}	OR _{Model 4}	OR _{Model 5}	OR _{Model 6}	OR _{Model 7}
5	5	5	1.000	1.000	1.750	1.002	1.001	1.800	1.353
2	5	5	0.999	0.998	1.299	1.001	1.002	1.792	1.151
0.5	5	5	1.000	1.000	0.785	0.999	0999	1.649	0.881
0.2	5	5	0.999	0.999	0.624	1.002	1.003	1.549	0.786
5	2	5	0.999	0.999	1.747	0.999	0.998	1.290	1.326
5	0.5	5	1.000	1.000	1.751	0.999	0.999	0.773	1.333
5	0.2	5	0.999	0.999	1.750	1.003	1.003	0.557	1.375
5	5	2	1.001	1.001	1.263	0.999	0.998	1.288	1.136
5	5	0.5	1.001	1.001	0.791	1.000	1.000	0.773	0.879
5	5	0.2	1.001	1.001	0.573	1.001	1.001	0.557	0.741

Table 5-1: Estimates in time-varying confounding scenario

In each parameter setting, the estimates of the effect of A_0 on Y we obtained from model 1 and model 2 are null and those we obtained from model 3 are not. As A_0 affects L and L affects Y, even though there is no direct effect of A_0 on Y, there should be a total effect through L. As we can see from the results, A_0 and Y are conditionally independent. Model 1 and model 2 provide biased estimates by conditioning on the intermediate variable. Without adjusting for the intermediate variable, estimates from model 3 are the unbiased estimates of the total effect of A_0 on Y through L.

But what is the magnitude of the total effect? How is it related to the effect of A_0 on L and the effect of L on Y?

According to the data generation,

logit($P(L=1 | A_0)) = \alpha_0 + \alpha_1 A_0$, we have

 $P(L=1 | A_0 = 1) = \text{expit}(\alpha_0 + \alpha_1)$, and $P(L=1 | A_0 = 0) = \text{expit}(\alpha_0)$

Similarly,

logit($P(Y = 1 | L)) = \gamma_0 + \gamma_1 L$, then it follows that

 $P(Y = 1 | A_0 = 1) = \exp((\gamma_0 + \gamma_1 \exp((\alpha_0 + \alpha_1))))$ and

 $P(Y=1 \mid A_0=0) = \operatorname{expit}(\gamma_0 + \gamma_1 \operatorname{expit}(\alpha_0))$

$$OR_{A_{0},Y} = \frac{P(Y=1 | A_{0} = 1) \times P(Y=0 | A_{0} = 0)}{P(Y=0 | A_{0} = 1) \times P(Y=1 | A_{0} = 0)}$$

=
$$\frac{\operatorname{expit}(\gamma_{0} + \gamma_{1} \operatorname{expit}(\alpha_{0} + \alpha_{1}))}{1 - \operatorname{expit}(\gamma_{0} + \gamma_{1} \operatorname{expit}(\alpha_{0} + \alpha_{1}))} \times \frac{1 - \operatorname{expit}(\gamma_{0} + \gamma_{1} \operatorname{expit}(\alpha_{0}))}{\operatorname{expit}(\gamma_{0} + \gamma_{1} \operatorname{expit}(\alpha_{0}))}$$

After setting $\alpha_0 = \gamma_0 = 0$, it follows,

$$P(Y=1 \mid A_0=1) = \frac{(e^{\gamma_1})^{\exp(i(\alpha_1)}}{1+(e^{\gamma_1})^{\exp(i(\alpha_1)}} \text{ and } P(Y=1 \mid A_0=0) = \frac{(e^{\gamma_1})^{\exp(i(0)}}{1+(e^{\gamma_1})^{\exp(i(0)}} = \frac{(e^{\gamma_1})^{\frac{1}{2}}}{1+(e^{\gamma_1})^{\frac{1}{2}}}$$

$$OR_{A_{0},Y} = \frac{(e^{\gamma_{1}})^{\exp(\alpha_{1})}}{(e^{\gamma_{1}})^{\frac{1}{2}}} = (e^{\gamma_{1}})^{\exp(\alpha_{1})-\frac{1}{2}}$$

$$\log(OR_{A_0,Y}) = \log((e^{\gamma_1})^{\exp(\alpha_1) - \frac{1}{2}}) = \gamma_1(\exp(\alpha_1) - \frac{1}{2}) \quad (13)$$

By using equation (13), if we assume there is no random error, one can derive the true magnitude of the total effect of A₀ on Y through L as a function of α_1 and γ_1 . To compare the results from the simulations and equation (13), uncertainty of the simulation must be estimated. The empirical standard error can be used, calculated as the standard deviation from the 100 simulations (Burton, Altman and et al., 2006). 95% confidence intervals of the log(OR_{Model3}) were computed. A one-sample t-test was performed to compare the estimates and the true effects. Table5-2 shows that the confidence intervals of the log(OR_{Model3}) include the results from equation (13) and the p values are generally >0.05. After taking account of the uncertainty of the simulation, results are comparable to the ones calculated by the formula. Theoretically, as the sample size goes to infinity and the number of repetitions in the simulation becomes large, the simulation results should converge to the analytical results.

e^{α_1}	e^{β_1}	e^{γ_1}	$log(OR_{Model3})$	Results from equation(13)	$\frac{\text{CI for}}{\log(OR_{\text{Model3}})}$	P value
5	5	5	0.560	0.536	[0.529, 0.591]	0.135
2	5	5	0.261	0.268	[0.234, 0.289]	0.641
0.5	5	5	-0.242	-0.268	[-0.269, -0.215]	0.055
0.2	5	5	-0.471	-0.536	[-0.496, -0.445]	0.000
5	2	5	0.558	0.536	[0.531, 0.586]	0.113
5	0.5	5	0.560	0.536	[0.529, 0.591]	0.123
5	0.2	5	0.560	0.536	[0.529, 0.591]	0.132
5	5	2	0.234	0.231	[0.208, 0.259]	0.828
5	5	0.5	-0.235	-0.231	[-0.258, -0.211]	0.758
5	5	0.2	-0.557	-0.536	[-0.582, -0.533]	0.094

Table 5-2: Comparison between the simulation and equation (13) for the total effect of A0 on Y through L

When we focus on the effect of A_1 on Y, the results are consistent with the discussion above, showing there is no effect of A_1 on Y from model 4 and model 5 in each parameter setting. It again confirms the condition of the absence of non-collapsibility when the exposure effect is null. Model 6 ignoring L provides biased estimates caused by confounding.

Model 7 provides estimates of the cumulative effect of A_0 and A_1 on Y, where $A=A_0+A_1$ and takes value of 0, 1, 2. Assuming A acts as a numeric variable, the estimate we obtained is the effect of increasing one unit of A (i.e. A_0 or A_1) on Y, yet we don't know whether the increasing unit is due to A_0 or A_1 . Therefore we presume that the cumulative effect is the average of A_0 effect and A_1 effect. By the study design and learning from the results, A_1 has no effect on Y, and A_0 has a total effect through L, hence the cumulative effect is a half of A_0 effect and we should have the relationship below:

 $\log(OR_{\text{Model7}}) = \frac{\log(OR_{\text{Model3}})}{2}$ (14)

Table 5-3 shows the cumulative effect from the simulation by $OR_{Model 7}$ and from equation (14) by $OR_{Model 3}$. A two-sample t-test was performed to compare the estimates for each parameters setting. The 95% confidence intervals are overlapped for the estimates, and the p values are generally >0.05. The results are comparable after taking account of the uncertainty in the simulation.

e^{α_1}	e^{β_1}	e^{γ_1}	$\log(OR_{Model7})$	$\frac{\log(OR_{\text{Model3}})}{2}$	CI for $log(OR_{Model7})$	$\frac{\text{CI for}}{\frac{\log(OR_{\text{Model3}})}{2}}$	P value
5	5	5	0.302	0.280	[0.281, 0.323]	[0.264, 0.295]	0.092
2	5	5	0.141	0.131	[0.122, 0.160]	[0.117, 0.145]	0.400
0.5	5	5	-0.126	-0.121	[-0.147, -0.106]	[-0.134, -0.107]	0.655
0.2	5	5	-0.241	-0.235	[-0.262, -0.221]	[-0.248, -0.223]	0.636
5	2	5	0.282	0.279	[0.264, 0.301]	[0.265, 0.293]	0.791
5	0.5	5	0.288	0.280	[0.265, 0.311]	[0.265, 0.296]	0.591
5	0.2	5	0.319	0.280	[0.292, 0.346]	[0.264, 0.295]	0.014
5	5	2	0.127	0.117	[0.109, 0.146]	[0.104, 0.130]	0.353
5	5	0.5	0.129	-0.117	[-0.146, -0.111]	[-0.129, -0.105]	0.289
5	5	0.2	0.300	-0.279	[-0.319, -0.280]	[-0.291, -0.266]	0.071

Table 5-3: Comparison of the results from model 7 and equation (14)

5.2.2 Non-collapsibility with Z

5.2.2.1 Method

Now we add the baseline variable Z into the simulation study to explore the noncollapsibility effect in the time-varying confounding scenario (Figure 5-2). Z is a risk factor for Y but independent of any other variables. It was generated from a Bernoulli distribution with P(Z = 1) = 0.5. Y was generated from a Bernoulli distribution with $P(Y = 1) = \exp((\gamma_0 + \gamma_1 L + \gamma_2 Z))$. γ_2 was specified with the value of 0.2. The generation of A₀, L and A₁ is accordant with the previous data generation.

We proposed 12 models (Model 8 to Model 19) below to explore the non-collapsibility effect. Learning from the diagram and the results from the previous simulation, not adjusting for L yields unbiased estimates of the total effect of A_0 on Y through L. Model 10 which is the same as model 3 provides us the marginal effect of A_0 . Based on model 10, model 8 and model 9 adjust for Z marginally and conditionally. As we defined in

section 3.3, $\log(\frac{OR_{Model8}}{OR_{Model9}})$ can be measured as the non-collapsibility effect of Z with

respect to A_0 effect. Since there is no confounding bias, model 8 and model 10 provide the same estimates. As for the effect of A_1 on Y, we must adjust for L marginally or conditionally to obtain an unbiased estimates. Model 11 to model 13 adjust for L marginally and model 14 to model 16 adjust for L conditionally. After adjusting for L, those models provide null estimates of effect by adjusting for Z with MSM, SLRM or without any adjustment. There should be neither a non-collapsibility effect nor a confounding bias by Z after adjusting for L, as the effect of A1 on Y is null. To estimate the cumulative effect of A on Y, after adjusting L marginally by IPTW, Z is adjusted marginally or conditionally or left unadjusted in model 17 to model 19 respectively. Comparing model 17 and model 19, we should get the same estimates since there is no confounding bias. By comparing model 17 and model 18 we get the non-collapsibility of Z with respect to cumulative effect A.

Model 8: logit(
$$P(Y=1)$$
) = $\omega_0 + \omega_1 A_0$ weights = $\frac{P(A_0)}{P(A_0 \mid Z)}$

Model 9:	$logit (P(Y=1)) = \omega_0 + \omega_1 A_0 + \omega_2 Z$	
Model 10:	$logit(P(Y=1)) = \omega_0 + \omega_1 A_0$	
Model 11:	$logit(P(Y=1)) = \omega_0 + \omega_1 A_1$	weights = $\frac{P(A_1)}{P(A_1 \mid L, Z)}$
Model 12:	logit($P(Y=1)$) = $\omega_0 + \omega_1 A_1 + \omega_2 Z$	weights = $\frac{P(A_1)}{P(A_1 \mid L)}$
Model 13:	$logit(P(Y=1)) = \omega_0 + \omega_1 A_1$	weights = $\frac{P(A_1)}{P(A_1 \mid L)}$
Model 14:	logit($P(Y=1)$) = $\omega_0 + \omega_1 A_1 + \omega_2 L$	weights = $\frac{P(A_1)}{P(A_1 \mid Z)}$
Model 15:	logit($P(Y=1)$) = $\omega_0 + \omega_1 A_1 + \omega_2 L + \omega_3 Z$	
Model 16:	$logit(P(Y=1)) = \omega_0 + \omega_1 A_1 + \omega_2 L$	
Model 17:	$logit(P(Y=1)) = \omega_0 + \omega_1 A$	weights = $\frac{P(A_1 A_0)}{P(A_1 A_0, L, Z)}$
Model 18:	logit $(P(Y=1)) = \omega_0 + \omega_1 A + \omega_2 Z$	weights = $\frac{P(A_1 A_0)}{P(A_1 A_0, L)}$
Model 19:	$logit(P(Y=1)) = \omega_0 + \omega_1 A$	weights = $\frac{P(A_1 \mid A_0)}{P(A_1 \mid A_0, L)}$

In order to compare the non-collapsibility effect in the longitudinal study with the pointexposure study, we need to change Figure 5-2 slightly. The structure remains the same except the relationship between L and A₁ needs to be removed. A₁ was generated from a Bernoulli distribution with $P(A_1 = 1) = 0.5$. We keep L in the Figure in order to have the same randomness as what we have in the time-varying confounding case. For the pointexposure study, the non-collapsibility effect with respect to A₀ is again measured by comparing model 8 and model 9. The non-collapsibility effect with respect to A is measured by comparing model 20 and model 21. L was left unadjusted, since it is a mediator rather than a time-varying confounding variable. Model 20: logit(P(Y = 1)) = $\omega_0 + \omega_1 A$ weights = $\frac{P(A_1 | A_0)}{P(A_1 | A_0, Z)}$

Model 21: logit(P(Y = 1)) = $\omega_0 + \omega_1 A + \omega_2 Z$

To assess the comparability of the non-collapsibility effect in the time-varying confounding and the point-exposure study, 95% confidence intervals of the non-collapsibility effect were computed for each parameters setting. Two-sample t-test was performed to compare the non-collapsibility effect between the two scenarios.

5.2.2.2 Result

The results in the time-varying confounding scenario are presented in table 5-4. Estimates from model 11 to model 16 all show no effect, and therefore they are not listed.

-									
e^{α_1}	e^{β_1}	e^{γ_1}	e^{γ_2}	OR _{Model 8}	OR _{Model 9}	OR _{Model 10}	OR _{Model 17}	OR _{Model 18}	OR _{Model 19}
5	5	5	0.2	1.573	1.668	1.573	1.282	1.324	1.282
2	5	5	0.2	1.249	1.285	1.249	1.128	1.145	1.128
0.5	5	5	0.2	0.799	0.777	0.798	0.889	0.877	0.889
0.2	5	5	0.2	0.636	0.600	0.636	0.793	0.770	0.792
5	2	5	0.2	1.575	1.671	1.574	1.262	1.301	1.262
5	0.5	5	0.2	1.571	1.666	1.570	1.261	1.300	1.261
5	0.2	5	0.2	1.572	1.667	1.572	1.291	1.333	1.290
5	5	2	0.2	1.215	1.254	1.216	1.115	1.135	1.115
5	5	0.5	0.2	0.808	0.789	0.807	0.890	0.879	0.890
5	5	0.2	0.2	0.588	0.566	0.588	0.754	0.738	0.754

Table 5-4: Estimates in time-varying confounding scenario with Z

Table 5-5 and table 5-6 shows the non-collapsibility effect of Z with respect to A_0 effect and the cumulative effect (A) in different parameter settings measured from the timevarying confounding and point-exposure study. The 95% confidence intervals for the two scenarios are overlapping and the p values are >0.05. After adjusting for the time-varying confounding by IPTW, the non-collapsibility effect of the baseline variable Z with respect to the cumulative effect is comparable with that from the point-exposure study. So is the non-collapsibility of Z with respect to A_0 .

	0	ν.	- ^γ 2	time-varying	g confounding	point-exp	posure study	
e^{α_1}	e^{β_1}	e^{γ_1}	e^{γ_2}	$\log(\frac{OR_{\rm Model8}}{OR_{\rm Model9}})$	CI	$\log(rac{OR_{ m Model8}}{OR_{ m Model9}})$	CI	P value
5	5	5	0.2	-0.059	[-0.063, -0.055]	-0.059	[-0.062, -0.055]	0.986
2	5	5	0.2	-0.028	[-0.031, -0.025]	-0.028	[-0.031, -0.025]	0.951
0.5	5	5	0.2	0.028	[0.025, 0.031]	0.028	[0.025, 0.032]	0.912
0.2	5	5	0.2	0.059	[0.055, 0.062]	0.059	[0.055, 0.063]	0.950
5	2	5	0.2	-0.059	[-0.063, -0.056]	-0.059	[-0.062, -0.055]	0.806
5	0.5	5	0.2	-0.059	[-0.063, -0.055]	-0.059	[-0.063, -0.055]	0.931
5	0.2	5	0.2	-0.059	[-0.063, -0.055]	-0.059	[-0.063, -0.055]	0.966
5	5	2	0.2	-0.031	[-0.035, -0.027]	-0.031	[-0.035, -0.027]	0.968
5	5	0.5	0.2	0.023	[0.020, 0.026]	0.024	[0.021, 0.027]	0.919
5	5	0.2	0.2	0.037	[0.035, 0.040]	0.037	[0.034, 0.040]	0.862

Table 5-5: Comparison of non-collapsibility effect with respect to A₀ effect

Table 5-6: Comparison of non-collapsibility effect with respect to A effect (1)

e^{α_1}	e^{β_1}	e^{γ_1}	e^{γ_2}	time-varying confounding		point-exposure study		
				$\log(\frac{OR_{\rm Model17}}{OR_{\rm Model18}})$	CI	$\log(\frac{OR_{\text{Model20}}}{OR_{\text{Model21}}})$	CI	P value
5	5	5	0.2	-0.032	[-0.038, -0.026]	-0.029	[-0.034, -0.025]	0.463
2	5	5	0.2	-0.015	[-0.021, -0.009]	-0.014	[-0.019, -0.009]	0.678
0.5	5	5	0.2	0.014	[0.009, 0.019]	0.014	[0.008, 0.020]	0.931
0.2	5	5	0.2	0.030	[0.024, 0.035]	0.029	[0.024, 0.034]	0.847
5	2	5	0.2	-0.030	[-0.034, -0.026]	-0.029	[-0.034, -0.024]	0.700
5	0.5	5	0.2	-0.030	[-0.035, -0.025]	-0.029	[-0.034, -0.024]	0.770
5	0.2	5	0.2	-0.032	[-0.039, -0.026]	-0.029	[-0.033, -0.025]	0.358
5	5	2	0.2	-0.017	[-0.024,-0.011]	-0.016	[-0.020, -0.011]	0.707
5	5	0.5	0.2	0.013	[0.006, 0.019]	0.012	[0.006, 0.017]	0.810
5	5	0.2	0.2	0.021	[0.016, 0.026]	0.019	[0.014, 0.023]	0.464

We extend the simulation study to a more general scenario by removing the assumption that the cumulative effect of A on Y is null. Figure 5-3 shows the scenario where there is an effect of A₀ on A₁ and there is a cumulative effect of A on Y, where $A=A_0+A_1$ deterministically. Consequently, we generated A₁ from a Bernoulli distribution with $P(A_1 = 1) = \text{expit}(\beta_0 + \beta_1 L + \beta_2 A_0)$, and Y from a Bernoulli distribution with $P(Y = 1) = \text{expit}(\gamma_0 + \gamma_1 L + \gamma_2 Z + \gamma_3 A)$. The generation of A₀, L and Z is identical to the previous data generation. α_0 , β_0 and γ_0 were set to 0. The other parameters were set to be the values in table 5-7.



Figure 5-3: Time-varying confounding scenario with baseline variable Z (Scenario 2) Model 17 and model 18 are again utilized to measure the non-collapsibility effect of Z with respect to the cumulative effect of A on Y. A point-exposure study was simulated by removing the relationship between L, A_0 and A_1 in figure 5-3. The non-collapsibility effect is measured by comparing estimates from model 20 and model 21 in the pointexposure study. Results were presented in Table 5-7. The 95% confidence intervals of the non-collapsibility effect for time-varying confounding scenario and point-exposure
scenario are overlapping and the p values are >0.05, indicating that the non-collapsibility effect of Z with respect to the cumulative effect is comparable with point-exposure study.

	e^{β_1}	e^{β_2}	e^{γ_1}	e^{γ_2}	e^{γ_3}	time-varying confounding		point-exposure study		
e^{α_1}						$\log(\frac{OR_{\text{Modell 7}}}{OR_{\text{Modell 8}}})$	CI	$\log(\frac{OR_{\text{Model20}}}{OR_{\text{Model21}}})$	CI	P value
2	2	2	2	0.2	5	-0.166	[-0.175, -0.157]	-0.176	[-0.187, -0.165]	0.148
5	2	2	2	0.2	5	-0.171	[-0.180, -0.162]	-0.182	[-0.192, -0.173]	0.090
0.5	2	2	2	0.2	5	-0.159	[-0.169, -0.150]	-0.167	[-0.176, -0.157]	0.282
0.2	2	2	2	0.2	5	-0.157	[-0.165, -0.149]	-0.163	[-0.172, -0.154]	0.297
2	5	2	2	0.2	5	-0.163	[-0.173, -0.154]	-0.176	[-0.185, -0.166]	0.079
2	0.5	2	2	0.2	5	-0.175	[-0.184, -0.165]	-0.177	[-0.187, -0.167]	0.759
2	0.2	2	2	0.2	5	-0.181	[-0.192, -0.170]	-0.175	[-0.185, -0.166]	0.475
2	2	5	2	0.2	5	-0.164	[-0.174, -0.155]	-0.176	[-0.186, -0.167]	0.072
2	2	0.5	2	0.2	5	-0.178	[-0.188, -0.167]	-0.176	[-0.185, -0.166]	0.786
2	2	0.2	2	0.2	5	-0.187	[-0.185, -0.167]	-0.176	[-0.199, -0.176]	0.132
2	2	2	5	0.2	5	-0.124	[-0.134, -0.115]	-0.131	[-0.140, -0.123]	0.263
2	2	2	0.5	0.2	5	-0.175	[-0.182, -0.167]	-0.178	[-0.188, -0.169]	0.522
2	2	2	0.2	0.2	5	-0.137	[-0.145, -0.129]	-0.135	[-0.142, -0.128]	0.749
2	2	2	2	0.1	5	-0.327	[-0.339, -0.314]	-0.340	[-0.354, -0.326]	0.166
2	2	2	2	0.2	10	-0.220	[-0.233, -0.206]	-0.233	[-0.248, -0.219]	0.170
2	2	2	2	5	0.2	0.174	[0.166, 0.183]	0.178	[0.168, 0.188]	0.586

Table 5-7: Comparison of non-collapsibility effect with respect to A effect (2)

5.3 Discussion

Informed by the observations from the point-exposure study, we anticipate that when both the effects of Z and A are harmful, the risks will be high thus the magnitude of the non-collapsibility will become trivial. Therefore, in the process of exploring the noncollapsibility effect of Z in the time-varying confounding scenario, we mainly set the effect of Z on Y with the OR of 0.2, and the effect of A on Y with the OR of 5. From the figures, we can also predict the symmetric behavior of the non-collapsibility effect when the effect of Z and A are on the opposite direction with the same magnitude. One of the parameter settings in table 5-7 again confirms this. Moreover, the main objective in this section was to find out whether time-varying confounding has an influence on the non-collapsibility effect. We have already described the other basic features of the non-collapsibility effect from the point-exposure study. Thus it is not necessary to carry out the simulation with every combination of the parameters. This setting would be a good representative of the non-collapsibility effect in time-varying confounding scenarios.

The Z effect only goes into the study at the very last point on Y, and is independent of any other variables in the time-varying confounding scenario. Under the assumptions that there is no unmeasured confounding and no loss to follow-up, the results are expected to be comparable with the point-exposure study whenever the time-varying confounding was controlled appropriately.

One can also use more models to analyze the data, for example adjusting Z in both the numerator and the denominator of the weights in the conditional model to examine if the confidence interval will be decreased. However, we do not focus on the precision of the estimates in this thesis.

6 Conclusions and Discussion

In this thesis, first we attempted to distinguish the concept of non-collapsibility and confounding bias. Considering the conditions of the collapsibility of the RD, RR and OR measures, non-collapsibility and confounding are equivalent when the RD or the RR is used as the measure of association. However, when the OR is applied, the two concepts are not always concordant. One should always keep in mind that comparing the OR_c and OR_{crude} cannot be deemed as a good strategy to detecting confounding bias. It was shown that the total discrepancy between the OR_c and OR_{crude} is the combination of the true confounding bias and the non-collapsibility effect. We emphasis again that both OR_m and OR_{c} are legitimate effect estimates. They are unbiased estimators for two different parameters, i.e. marginal effect at the population level and conditional effect at the individual or clustered level. They just happen to be different due to the non-collapsibility effect. The choice of OR_m or OR_c should depend on the question under research. It is important for the investigator to consider which parameter is mainly of interest and choose the corresponding estimator.

The analytical approach provides a tool to compute the non-collapsibility effect in different scenarios. Given a dataset from a real life study, one can measure the non-collapsibility of OR using the formula provided in this thesis. It explicitly shows the decomposition of the total discrepancy and verification of various conditions of the collapsibility of the OR. Other formulae in the literature provided similar results in the

setting that A is independent of L. However, the formula in our novel approach provides the non-collapsibility effect in general scenarios. This allows one to assess the noncollapsibility effect and the confounding bias simultaneously if a potential confounding variable is excluded in the models.

The graphical approach gives us a straightforward impression of the non-collapsibility effect by visualizing it with figures. It provides a comprehensive view of the noncollapsibility effect with a variety of parameter values. Therefore, we can draw conclusions about the relationship between the non-collapsibility effect with one specified related factor given that the others are fixed. Many interesting observations from the figures were discussed, providing another aspect of understanding the noncollapsibility effect besides the magnitude.

The objective of the simulation study was to explore the non-collapsibility effect in the time-varying confounding scenario. The results are comparable under the settings with or without a direct cumulative effect on the outcome. Some scenarios and models that do not involve the non-collapsibility effect were also demonstrated in this part of the data analysis. Though they were not the aims of the study, it provided additional confirmation of model specification and adjustment of variables concerning mediation and confounding in the causal diagram. However, a few results in table 5-2 and table 5-3 imply that the simulation results do not agree the theoretical results under some settings. We suppose that it could happen by chance. Moreover, if the alpha level is set to 0.01 instead of 0.05, we would not reject the null hypothesis. We can conclude that the results are comparable in table 5-3.

Some assumptions and limitations of the analytical approach have to be acknowledged. Firstly, several assumptions were made in the analytical approach. MSM and SLRM were performed to measure the marginal effect and the conditional effect. They are unbiased estimates under the assumption that there is no unmeasured confounding, no misclassification and no missing data. Furthermore, we used the un-stabilized weighted estimators in the MSM. The assumption of the positivity of weights was made, and it would be problematic if the denominator of the weights approaches 0 or 1 due to sparse or even empty cells in the observed data. In addition, we've implemented a deterministic process to measure the non-collapsibility effect. We are studying the true value of the non-collapsibility effect in a source population and there is no sampling variability. However, random error and sampling variability are inevitable in practice. When one measures the non-collapsibility effect in finite samples, the precision of the estimate is also of interest in practice. Bootstrapping would be a good method to obtain the standard deviation.

Furthermore, in the point-exposure studies, we have explored the non-collapsibility effect in very simple scenarios, whereas reality is often much more complicated. For example, there can be more than two strata of L, and the non-collapsibility behavior could be different depending on how the strata were collapsed over. It is also possible that there are more than one baseline variable in one study, while we've only observed the noncollapsibility effect of one variable. When L is a vector, further study can be performed to investigate the non-collapsibility effect over some variables in L given that the others are conditionally or marginally adjusted. It could be interesting to know whether the noncollapsibility effect is additive with more than one baseline variable.

The graphical approach presented the analytical non-collapsibility effect with figures, however the derivation of the formula expressed by the other parameters is implicit. Along with the figures, a concise and understandable formula is needed to convey more precise information in a mathematical and statistical perspective.

For the time-varying confounding scenario, more simulations and more settings are required to better illustrate the non-collapsibility effect compared to the time-point exposure scenario. The simulation was performed 100 times for each parameters setting. More simulations will provide a better accuracy of the estimates. We demonstrated that the non-collapsibility effect is comparable between the time-exposure and time-varying confounding scenario. However, from the results in table 5-7, it appears that the results are more likely to be further from the null in point-exposure scenarios. There might be some still unknown mechanism or related factors that are driving this subtle behavior. More research is needed to demonstrate the reason for this unexpected phenomenon. Moreover, I focused on the relatively limited settings so far. For example, the intercepts of the models are set to be zero, and only a few values of the parameters were specified. Further research should aim at exploring more possible values of the parameters to make the comparison more complete.

In conclusion, the non-collapsibility problem of the OR should be neither ignored nor confused with confounding. The non-collapsibility effect depends on a variety of parameters, i.e. the baseline risk, the effect of the exposure, the prevalence and effect of the covariate. Particularly, the effect of the covariate plays a more important role than does the effect of the exposure, a result that has never been reported previously. Lastly, simulation results suggested that the non-collapsibility effect over a baseline variable in a time-varying confounding scenario is comparable to the time-exposure study if the time-varying confounder was adjusted appropriately in the model.

Appendix A

A.1 The approximation of non-collapsibility effect by Gail

In Gail's work, X is the covariate with expectation 0 and covariance σ^2 . Equation (2.9) shows the approximation of non-collapsibility effect over X. In the context of this thesis, L is a binary variable. We will derive a formula analogous to (2.9), assuming L is the covariate we are considering.

By (2.5),

$$\alpha^* = \frac{1}{2}h^{-1}(\xi_1) - \frac{1}{2}h^{-1}(\xi_2)$$

For a small β , by the second-order Taylor series,

$$\xi_1 = E\{h(\mu + \alpha + \beta L)\} = h(\mu + \alpha) + \beta E(L)h'(\mu + \alpha) + \frac{1}{2}\beta' E(L^2)\beta h''(\mu + \alpha)$$

Thus,

$$h^{-1}(\xi_1) = h^{-1}(h(\mu + \alpha) + \beta E(L)h'(\mu + \alpha) + \frac{1}{2}\beta' E(L^2)\beta h''(\mu + \alpha))$$

By the Taylor series expansion, we have

$$h^{-1}(\xi_1) = h^{-1}(h(\mu + \alpha)) + (h^{-1})'(h(\mu + \alpha))[\beta E(L)h'(\mu + \alpha) + \frac{1}{2}\beta' E(L^2)\beta h''(\mu + \alpha)]$$

From
$$(h^{-1})'(h(\mu + \alpha)) = \frac{d}{d\mu}(h^{-1}(h(\mu + \alpha))) = \frac{d}{d\mu}(\mu + \alpha) = 1$$
 and
$$\frac{d}{d\mu}(h^{-1}(h(\mu + \alpha))) = (h^{-1})'(h(\mu + \alpha)) \times h'(\mu + \alpha)$$

It follows that:

$$(h^{-1})'(h(\mu + \alpha)) = \frac{1}{h'(\mu + \alpha)}$$

Thus,

$$h^{-1}(\xi_1) = \mu + \alpha + \beta E(L) + \frac{1}{2}\beta' E(L^2)\beta h''(\mu + \alpha) / h'(\mu + \alpha)$$

Similarly,
$$h^{-1}(\xi_2) = \mu - \alpha + \beta E(L) + \frac{1}{2}\beta' E(L^2)\beta \frac{h''(\mu - \alpha)}{h'(\mu - \alpha)}$$

Submitting $h^{-1}(\xi_1)$ and $h^{-1}(\xi_2)$ into (2.5) gives

$$\alpha^* = \frac{1}{2}h^{-1}(\xi_1) - \frac{1}{2}h^{-1}(\xi_2) = \frac{1}{2}(2\alpha + \frac{1}{2}\beta' E(L^2)\beta\{\frac{h''(\mu+\alpha)}{h'(\mu+\alpha)} - \frac{h''(\mu-\alpha)}{h'(\mu-\alpha)}\}$$

Then the approximation of non-collapsibility effect is:

$$\alpha^* - \alpha = \frac{1}{4}\beta' E(L^2)\beta\{\frac{h''(\mu+\alpha)}{h'(\mu+\alpha)} - \frac{h''(\mu-\alpha)}{h'(\mu-\alpha)}\}$$

A.2 The expression of P_{01} as a function of the other parameters

 OR_L can be constructed in table 3.

$$OR_{L} = \frac{(a+b) \times (g+h)}{(c+d) \times (e+f)} = \frac{(q_{00} \times q_{A0} + q_{10} \times P_{A0}) \times (P_{01} \times q_{A1} + P_{11} \times P_{A1})}{(P_{00} \times q_{A0} + P_{10} \times P_{A0}) \times (q_{01} \times q_{A1} + q_{11} \times P_{A1})}$$

Firstly, we define

$$S = OR_L \times \frac{P_{00} \times q_{A0} + P_{10} \times P_{A0}}{q_{00} \times q_{A0} + q_{10} \times P_{A0}}$$

It follows that

$$S \times (q_{01} \times q_{A1} + q_{11} \times P_{A1}) = P_{01} \times q_{A1} + P_{11} \times P_{A1}$$

$$S \times [(1 - P_{01}) \times q_{A1} + (1 - P_{11}) \times P_{A1}) = P_{01} \times q_{A1} + P_{11} \times P_{A1}$$

$$S = (S+1) \times (P_{01} \times q_{A1} + P_{11} \times P_{A1})$$

$$\frac{S}{S+1} = P_{01} \times q_{A1} + P_{11} \times P_{A1}$$

$$\frac{S}{S+1} = P_{01} \times q_{A1} + \frac{P_{01} \times OR_C}{1 - P_{01} + P_{01} \times OR_C} \times P_{A1}$$

$$P_{01}q_{A1} + (OR_C - 1)P_{01}^2q_{A1} + P_{01}OR_CP_{A1} = \frac{S}{S+1} + \frac{S}{S+1}(OR_C - 1)P_{01}$$

$$(OR_{C}-1)P_{01}^{2}q_{A1} + P_{01}[q_{A1} + OR_{C}P_{A1} - \frac{S}{S+1}(OR_{C}-1)] - \frac{S}{S+1} = 0$$

Solving the above quadratic equation, we define

 $a = (OR_C - 1)q_{A1}$

$$b = q_{A1} + OR_C P_{A1} - \frac{S}{S+1}(OR_C - 1)$$

$$c = -\frac{S}{S+1}$$

 P_{01} can be expressed as

$$P_{01} = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

Bibliography

- Asmussen, S. and Edwards, D. (1983). Collapsibility and response variables in contingency tables. Biometrika, 70, 567-578.
- Austin, P. (2007). The performance of different propensity score methods for estimating marginal odds ratios. Statistics in Medicine, 26, 3078-3094.
- Becher, H. (1992). The concept of residual confounding in regression models and some applications. Statistics in Medicine, 11, 1747-1758.
- Boivin, J.F. and Wacholder, S. (1985). Conditions for confounding of the risk ratio and of the odds ratio. Am. J. Epidemiol.,121,152-158.
- Breslow, N.E. and Day, N.E. (1987). Statistical methods in cancer research, 2: The design and analysis of cohort studies. Lyon, France: IARC Scientific Publications.
- Clogg, C.C., Petkova, E. and Shihadeh, S. (1992). Statistical methods for analyzing collapsibility in regression models. J. Educ. Statist, 17, 51-74.
- Ducharme, G.R. and LePage, Y. (1986). Testing collapsibility in contingency tables. J. Roy. Statist. Soc. Ser. B, 48, 197-205.
- Flanders W.D. and Khoury M.J. (1990). Indirect assessment of confounding: graphic description and limits on effect of adjusting for covariates. Epidemiology, 1(3), 239-246.

- Forbes, A. and Shortreed, S. (2008). Letter to the editor Inverse probability weighted estimation of the marginal odds ratio: correspondence regarding 'the performance of different propensity score methods for estimating marginal odds ratios'.
 Statistics in Medicine, 27, 5556-5559.
- Gail, M.H., Wieand, S. and Piantadosi, S. (1984). Biased estimates of treatment effect in randomized experiments with nonlinear regressions and omitted covariates.Biometrika, 71, 431-444.
- Geng, Z. and Li, G. (2002). Conditions for non-confounding and collapsibility without knowledge of completely constructed causal diagrams. Board of the Foundation of the Scandinavian Journal of Statistics, 29, 169-181.
- Grayson, D.A. (1987). Confounding confounding. American Journal of Epidemiology, 126, 546–553.
- Greenland, S. and Robins, J.M. (1986). Identifiability, exchangeability, and epidemiological confounding. Internat. J. Epidemiol., 15, 413-419.
- Greenland, S. (1987). Interpretation and choice of effect measures in epidemiologic analyses. Amer. J. Epidemiology, 125, 761-768.
- Greenland, S. and Mickey, R.M. (1988). Closed-form and dually consistent methods for inference on collapsibility in 2 * 2*K and 2* J* K tables. J. Roy. Statist. Soc. Ser. C, 37, 335-343.

- Greenland, S., Morgenstern, H., Poole, C., and Robins, J.M. (1989). RE: Confounding confounding (letter). American Journal of Epidemiology, 129, 1086–1089. (Newman)
- Greedland, S. (1996). Absence of confounding does not correspond to collapsibility of the rate ratio or rate difference. Epidemiology, 7, 498-501.
- Greenland, S. and Robins, J.M. and Pearl, J. (1999). Confounding and collapsibility in causal inference. Statistical Science, 14, 1, 29-46.
- Greenland, S. and Morgenstern, H. (2001). Confounding in health research. Annu. Rev. Public Health, 22, 189-212.
- Greenland, S. (2010). Simpson's paradox from adding constants in contingency tables as an example of Bayesian noncollapsibility. American Statistical Association, 64, 340-344.
- Greenland, S. and Pearl, J. (2011). Adjustments and their consequences-collapsibility analysis using graphical models. International Statistical Review, 79, 3, 401-426.
- Groenwold, R.H.H. and Moons, K.G.M. et al. (2011). Reporting of treatment effects from randomized trials: A plea for multivariable risk ratios. Contemporary Clinical Trials, 32, 399-402.
- Guo, J. and Geng, Z. (1995). Collapsibility of logistic regression coefficients. J. Roy Statist. Soc. Ser. B, 57, 263-267.

- Hernán, M.A. And Robins, J.M. (2006). Estimating causal effects from epidemiological data. J. Epidemiol. Community Health, 60, 578-586.
- Hernán, M.A., Clayton, D. and Keiding, N. (2011). The simpson's paradox unraveled. Int. J. Epidemiol., 1-6.
- Hernandez, A.V., Steyerberg, E.W. and Habbema, J.D. (2004). Covariate adjustment in randomized controlled trials with dichotomous outcomes increases statistical power and reduces sample size requirements. J Clin. Epidemiol, 57, 454–60.
- Janes, H., Dominici, F. and Zeger, S. (2010). On quantifying the magnitude of confounding. Biostatistics, 11, 572-582.
- Maldonado G. and Greenland S. (2002). Estimating causal effects. International Journal of Epidemiology, 31, 422-429.
- Mietten O.S. (1972). Components of the crude risk ratio. Amer. J. Epidemiol, 96, 168-172.
- Mietten, O.S. and Cook, E.F. (1981). Confounding: essence and detection. American Journal of Epidemiology, 114, 593-603.
- Neuhaus, J.M., Kalbfleisch, J.D. and Hauck, W.W. (1991). A comparison of clusterspecific and population-averaged approaches for analyzing correlated binary data. Internat. Statist. Rev. 59, 25-35.

Newman, S.C. (2001). Biostatistical methods in epidemiology. John Wiley & Sons, Inc.

- Petersen M.L., Wang Y., van der Laan M.J., and David R.B. (2006). Assessing the effectiveness of antiretroviral adherence interventions. Using marginal structural models to replicate the findings of randomized controlled trials. J Acquir Immune Defic Syndr., 43(suppl 1), S96-S103.
- Ritz, J. and Spiegelman, D. (2004). Equivalence of conditional and marginal regression models for clustered and longitudinal data. Statistical Methods in Medical Research, 13, 309-323.
- Robins, J.M., Hernán, M.A., and Brumback, B. (2000). Marginal structural models and causal inference in epidemiology. Epidemiology, 11, 550-560.
- Robinson, L.D. and Jewell, N.P. (1991). Some surprising results about covariate adjustment in logistic regression. Int. Statist. Rev. 59, 227-240.
- Samuels, M.L. (1981). Matching and design efficiency in epidemiological studies. Biometrika, 68, 577–588.
- Shapiro, S.H. (1982). Collapsing contingency tables-approach. J Am Stat Assoc, 36, 43-46.
- Snowden J.M., Rose S., and Mortimer K.M. (2011). Implementation of g-computation on a simulated data set: Demonstration of a causal inference technique. American Journal of Epidemiology, 173(7), 731-738.

- Stampf, S., Graf, E. et. al (2010). Estimators and confidence intervals for the marginal odds ratio using logistic regression and propensity score stratification. Statistics in Medicine, 29, 760-769.
- Steyerberg, E,W., Bossuyt, P.M, and Lee, K.L. (2000). Clinical trials in acute myocardial infarction: should we adjust for baseline characteristics? Am Heart J, 139, 745– 751.
- Steyerberg, E,W. (2009). Clinical prediction models: a practical approach to development, validation, and updating. New York, NY: Springer.
- Wermuth, N. (1987). Parametric collapsibility and lack of moderating effects in contingency tables with a dichotomous response variable. J. Roy. Statist. Soc. Ser. B, 49, 353-364.
- Whittemore, A.S. (1978). Collapsing multidimensional contingency tables. J. Roy. Statist. Soc. Ser. B, 40, 328-340.
- Wickramaratne, P.J. and Holford, T.R. (1989). Confounding in epidemiologic studies. Response. Biometrics, 45, 1319-1322.
- Zhang, Z. (2009). Estimating a marginal causal odds ratio subject to confounding. Communications in Statistics-Theory and Methods, 38, 309-321.