

Title of the dissertation

A THESIS PRESENTED
BY
MORGAN F. BREITMEYER
TO
THE DEPARTMENTS OF STATISTICS AND MATHEMATICS

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
BACHELOR OF THE ARTS
IN THE SUBJECT OF
STATISTICS AND MATHEMATICS

HARVARD UNIVERSITY
CAMBRIDGE, MASSACHUSETTS
APRIL 2017

©2017 – MORGAN F. BREITMEYER
ALL RIGHTS RESERVED.

Title of the dissertation

ABSTRACT

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetur erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit, vel lacinia enim metus eu nunc.

Contents

0	INTRODUCTION	I
1	BACKGROUND	2
1.1	Causal Effects	3
1.2	Parametric G-formula	6
1.3	G-estimation	6
1.4	Doubly Robust Estimation	6
2	METHODS	7
3	SIMULATION DISCUSSION	8
4	CONCLUSION	9
	APPENDIX A APPENDIX	10
	REFERENCES	II

THIS IS THE DEDICATION.

Acknowledgments

LOREM IPSUM DOLOR SIT AMET, consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

0

Introduction

1

Background

THERE'S SOMETHING TO BE SAID for having a good opening line. Morbi commodo, ipsum sed pharetra gravida, orci $x = 1/\alpha$ magna rhoncus neque, id pulvinar odio lorem non turpis^{2,4}.

1.1 CAUSAL EFFECTS

A traditional understanding of causation comes from the field of medicine, where researchers can perform a controlled experiment to prove causation. This type of study contains two sample groups, one which receives no treatment (the placebo group) and one which receives the treatment (the treatment group). By comparing the outcome of these two groups, the researchers can demonstrate whether the outcome for patients receiving treatment differs significantly from the controls.

To translate this idea into statistical terms, some notation must be introduced. The random variable A represents the treatment status, where a value of 1 indicates treated and a value of 0 indicates untreated. The random variable Y is the outcome variable, often with a value of 0 indicating survival and a value of 1 indicating death. These interpretations of A and Y correspond to the above understanding of causation studies, but for various causal inference studies, the form of Y and in particular can change depending on the question of interest. For example, Y can be a continuous variable, such as the weight difference of an individual in a weight loss trial or the change in HDL levels in a cholesterol study.

To study the causal effect of A , the desired value is the difference in Y under the varying conditions of A . Notationally, this is the difference between $Y^{a=1}$, the outcome under treatment, and $Y^{a=0}$, the outcome under no treatment. A causal effect can be seen on an individual level if $Y_i^{a=1} \neq Y_i^{a=0}$ for individual i . By considering how each individual's responses to varying treatments differ,

causation (or lack thereof) can easily be determined using paired differences of the form

$$Y_i^{a=1} - Y_i^{a=0}$$

These differences would be tested against the null hypothesis of zero difference in outcome for varying treatments.

However, certain difficulties arise using this method. In many studies, it is impossible to have scenarios of both treatment and no treatment for the same individual, particularly if a potential outcome is death. Typically, individuals either have $Y_i^{a=1}$ or $Y_i^{a=0}$, but not both, making it impossible to calculate the paired differences. Therefore, a controlled double blinded experiment is often performed, where each individual is randomly assigned treatment or placebo. In these studies, the statistic of interest is the average causal effect in the population,

$$\mathbb{E}[Y^{a=1}] - \mathbb{E}[Y^{a=0}]$$

Mathematically, this is equivalent to

$$\mathbb{E}[Y^{a=1} - Y^{a=0}]$$

because the average of differences is equal to the difference of averages³. Note, that this is not the same as calculating the mean of paired differences as if each individual had received both treatments at different times to calculate individual causal effects. Rather, the difference in the means of the placebo and treatment groups is being calculated to estimate average causal effect across the popula-

tion.

1.1.1 IP WEIGHTING

Many of the concerns discussed above can be addressed using the method of IP weighting, which simulated a psuedo-population in which every individual has two data inputs, that of treatment and that of no treatment. The method by which this is done is by considering a covariate of the data, L , another value which is known before treatment is assigned. The idea of conditional exchangeability can be used to see that $Y^a \perp\!\!\!\perp A \mid L$ in the original population because A is independent of L ³. The pseudo-population can be calculated with the following for each of the possible A and L combinations

$$n \cdot \Pr[Y = y \mid A = a, L = 1] \cdot \Pr[A = a \mid L = 1] \cdot \Pr[L = 1] \cdot \frac{1}{\Pr[A = a \mid L = 1]}$$

where the last term here is the IP weight. This form can be used to solve for the standardized mean as follows,

$$E[Y^a] = \sum_l n \cdot \Pr[Y = y \mid A = a, L = 1] \cdot \Pr[A = a \mid L = 1] \cdot \Pr[L = 1] \cdot \frac{1}{\Pr[A = a \mid L = 1]} \quad (1.1)$$

$$= \sum_l n \cdot \Pr[Y = y \mid A = a, L = 1] \cdot \Pr[L = 1] \quad (1.2)$$

$$= \sum_l E[Y \mid A = a, L = 1] \Pr[L = 1] \quad (1.3)$$

1.1.2 ASSUMPTIONS

1.2 PARAMETRIC G-FORMULA

1.3 G-ESTIMATION

1.4 DOUBLY ROBUST ESTIMATION

2

Methods

LOREM IPSUM DOLOR SIT AMET,

3

Simulation Discussion

LOREM IPSUM DOLOR SIT AMET,

4

Conclusion



Appendix

References

- [1] Bang, H. & Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4), 962–973.
- [2] Eigen, M. (1971). Selforganization of matter and the evolution of biological macromolecules. *Naturwissenschaften*, 58(10), 465–523.
- [3] Hernan, M. A. & Robins, J. M. (2016). *Causal Inference*. Chapman & Hall/CRC.
- [4] Knuth, D. E. (1968). Semantics of context-free languages. *Mathematical Systems Theory*, 2(2), 127–145.
- [5] Lodi, S., Phillips, A., Logan, R., Olson, A., Costagliola, D., Abgrall, S., van Sighem, A., Reiss, P., Miró, J. M., Ferrer, E., et al. (2015). Comparative effectiveness of immediate antiretroviral therapy versus cd4-based initiation in hiv-positive individuals in high-income countries: observational cohort study. *The Lancet HIV*, 2(8), e335–e343.
- [6] Wright, J. D. (2015). *International encyclopedia of the social and behavioral sciences*.
- [7] Young, J. G., Cain, L. E., Robins, J. M., O'Reilly, E. J., & Hernán, M. A. (2011). Comparative effectiveness of dynamic treatment regimes: an application of the parametric g-formula. *Statistics in biosciences*, 3(1), 119–143.