

LMTLE: A Brief Introduction

Ashley I Naimi

April 2023

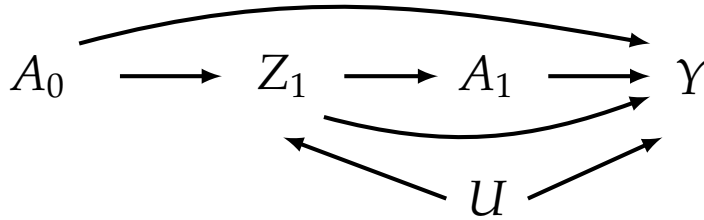
Contents

1	Introduction to Complex Longitudinal Data	2
2	IPW and g Computation for CLD	3
2.1	Inverse Probability Weighting	4

1 Introduction to Complex Longitudinal Data

Complex longitudinal data is becoming more common in a number of sector. These data are different from more traditional “longitudinal data” that one encounters in classical statistics courses.¹ Complex longitudinal data requires the presence of two features: first, repeated exposure, confounder, and (potentially) outcome measures; second, there has to be time-dependent feedback between these exposure, confounder, and (potentially) outcome variables.

This Figure demonstrates these basic conditions:



¹ For example, courses where you would learn how to estimate models with more sophisticated correlation structures, such as generalized estimating equations or mixed effects models

Figure 1: Causal diagram representing the structure from which the simple simulated data were generated.

This Figure is a simplified version, and several details can be added. For example, baseline covariates will always be present; there may be a Z_0 variable measured; there may be more than one time-dependent confounder; and the outcome may be measured at multiple time points, and may also serve as a time-dependent confounder.

Generally, the presence of a causal relation between Z_1 and A_1 suggests that Z_1 is a confounding variable. However, we cannot simply adjust for Z_1 in a standard regression model, since this would (i) block part of the effect of interest from $A_0 \rightarrow Z_1 \rightarrow Y$, and (ii) induce collider stratification bias through $A_0 \rightarrow Z_1 \leftarrow U \rightarrow Y$.

Because of this, we have to use specialized methods to estimate average treatment effects with complex longitudinal data.

2 IPW and g Computation for CLD

Perhaps the two most common techniques are inverse probability weighting or g computation (aka the parametric g formula) in such settings.

For example, if we're interested in the following ATE:

$$\psi = E(Y^{\bar{a}=1} - Y^{\bar{a}=0})$$

where \bar{a} denotes the entire history of the exposure measurement from the start to the end of follow-up for each person. To make our illustrations concrete, let's use a simple simulated dataset with longitudinal information on an exposure, several time-dependent confounders, and a time-to-event outcome.² Here's what the data look like:

² These data were simulated from Jessica Young's algorithm, modified by Erica Moodie. Details can be found in [Young et al. \(2010\)](#) and [Moodie et al. \(2014\)](#)

```
a <- read_csv(here("data", "2023_04_21-time-dependent.csv")) %>%
  group_by(ID) %>%
  mutate(exposure_lag = lag(exposure, n = 1L,
    default = 0), c1_lag = lag(c1, n = 1L,
    default = 0), c2_lag = lag(c2, n = 1L,
    default = 0)) %>%
  ungroup()

a
```

```
## # A tibble: 3,435 x 9
##       ID   Int exposure    c1    c2    Y exposure_lag c1_lag c2_lag
##   <dbl> <dbl>   <dbl> <dbl> <dbl> <dbl>   <dbl>  <dbl>  <dbl>
## 1     1     1     0     0     1     0     0     0     0
## 2     1     2     1     1     1     0     0     0     1
## 3     1     3     1     1     0     0     1     1     1
## 4     1     4     1     1     1     1     1     1     0
## 5     2     1     0     0     1     0     0     0     0
## 6     2     2     1     0     1     0     0     0     1
## 7     2     3     0     1     1     0     1     0     1
## 8     2     4     1     0     1     0     0     1     1
## 9     3     1     1     1     1     0     0     0     0
```

```
## 10      3      2      1      1      1      0      1      1      1
## # i 3,425 more rows
```

```
## look at proportion of outcome at
## each time point
a %>%
  group_by(Int) %>%
  summarise(mY = mean(Y))
```

```
## # A tibble: 4 x 2
##   Int    mY
##   <dbl> <dbl>
## 1     1 0.220
## 2     2 0.254
## 3     3 0.253
## 4     4 0.227
```

Let's talk a little about this data structure.

2.1 Inverse Probability Weighting

Let's start by constructing stabilized IP weights to estimate the ATE in these data:

```
# numerator
num <- glm(exposure ~ factor(Int), data = a,
  family = binomial("logit"))$fitted.values

# denominator
den <- glm(exposure ~ factor(Int) + exposure_lag +
  c1 + c2 + c1_lag + c2_lag, data = a,
  family = binomial("logit"))$fitted.values

a <- a %>%
  mutate(sw_ = num/den) %>%
  group_by(ID) %>%
```

```

mutate(sw = cumprod(sw_)) %>%
ungroup() %>%
select(-sw_)

a %>%
  group_by(Int) %>%
  summarise(meanSW = mean(sw), maxSW = max(sw))

## # A tibble: 4 x 3
##       Int meanSW maxSW
##   <dbl> <dbl> <dbl>
## 1     1     1.02  1.60
## 2     2     1.12  5.21
## 3     3     1.38 22.2
## 4     4     1.93 72.5

```

We can then use these weights to fit an IP weighted MSM:

```
library(lmtest)
```

```

## Loading required package: zoo

##
## Attaching package: 'zoo'

## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric

```

```

library(sandwich)

modMSM <- glm(Y ~ factor(Int) + exposure,
  data = a, weights = sw, family = binomial("logit"))

```

```
## Warning in eval(family$initialize): non-integer #successes in a binomial glm!
```

```
coeftest(modMSM, vcov. = vcovHC(modMSM, type = "HC3"))

##
## z test of coefficients:
##
##              Estimate Std. Error  z value  Pr(>|z|)
## (Intercept)  -1.77091    0.11198 -15.8152 < 2.2e-16 ***
## factor(Int)2   0.01620    0.10716   0.1512 0.8798338
## factor(Int)3  -0.10762    0.13177  -0.8167 0.4140977
## factor(Int)4  -0.61707    0.17900  -3.4473 0.0005662 ***
## exposure      0.71606    0.11950   5.9919 2.074e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
install.packages("ltml", repos = "http://lib.stat.cmu.edu/R/CRAN",
  dependencies = T)
```

```
##
## The downloaded binary packages are in
## /var/folders/zm/rqfq5xs0fs86qs2mcxk6q0r0000gr/T/Rtmpkgs1D0/downloaded_packages
```

```
library(ltml)
```

```
# read in data again, to simplify
```

```
a <- read_csv(here("data", "2023_04_21-time-dependent.csv"))
```

```
## Rows: 3435 Columns: 6
## -- Column specification -----
## Delimiter: ","
## dbl (6): ID, Int, exposure, c1, c2, Y
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
# convert data from long to wide
```

```
a %>%
```

```
  print(n = 3)
```

```
## # A tibble: 3,435 x 6
```

```
##       ID   Int exposure    c1    c2    Y
```

```
##    <dbl> <dbl>    <dbl> <dbl> <dbl> <dbl>
```

```
## 1     1     1        0     0     1     0
```

```
## 2     1     2        1     1     1     0
```

```
## 3     1     3        1     1     0     0
```

```
## # i 3,432 more rows
```

```
# TO DEAL WITH BASELINE CONFOUNDERS,
```

```
# KEEP THEM IN DATA
```

```
b <- a %>%
```

```
  pivot_wider(names_from = Int, values_from = c(exposure,
    c1, c2, Y)) %>%
```

```
  mutate(Y_1 = if_else(is.na(Y_1), 1, Y_1),
```

```
    Y_2 = if_else(is.na(Y_2), 1, Y_2),
```

```
    Y_3 = if_else(is.na(Y_3), 1, Y_3),
```

```
    Y_4 = if_else(is.na(Y_4), 1, Y_4)) %>%
```

```
  select(exposure_1, c1_1, c2_1, Y_1, exposure_2,
```

```
    c1_2, c2_2, Y_2, exposure_3, c1_3,
```

```
    c2_3, Y_3, exposure_4, c1_4, c2_4,
```

```
    Y_4)
```

```
b
```

```
## # A tibble: 1,228 x 16
```

```
##    exposure_1 c1_1 c2_1 Y_1 exposure_2 c1_2 c2_2 Y_2 exposure_3 c1_3
```

```
##          <dbl> <dbl> <dbl> <dbl>          <dbl> <dbl> <dbl> <dbl>          <dbl> <dbl>
```

```
## 1           0     0     1     0           1     1     1     0           1     1
```

```
## 2           0     0     1     0           1     0     1     0           0     1
```

```
## 3           1     1     1     0           1     1     1     0           1     1
```

```
## 4           0     1     1     0           0     1     1     0           1     1
```

```
## 5      1      1      1      0      1      0      1      1      NA      NA
## 6      1      1      1      0      1      1      1      1      NA      NA
## 7      1      0      1      0      0      0      1      1      NA      NA
## 8      1      1      0      0      1      1      1      0      1      1
## 9      0      0      0      0      1      1      1      0      1      1
## 10     1      1      1      1      NA      NA      NA      1      NA      NA

## # i 1,218 more rows
## # i 6 more variables: c2_3 <dbl>, Y_3 <dbl>, exposure_4 <dbl>, c1_4 <dbl>,
## #   c2_4 <dbl>, Y_4 <dbl>
```

```
# ltmle

# super learner library
sl.lib <- c("SL.mean",
           "SL.glm",
           "SL.ranger")

# ltmle
result <- ltmle(b,
               Anodes=c(paste0("exposure_",1:4)),
               Lnodes=c("c1_1", "c2_1", "c1_2", "c2_2",
                        "c1_3", "c2_3", "c1_4", "c2_4"),
               Ynodes=c("Y_1","Y_2","Y_3","Y_4"),
               survivalOutcome = TRUE,
               SL.library = list(Q = sl.lib, g = sl.lib),
               abar=list(treatment = c(1, 1, 1, 1),
                        control = c(0, 0, 0, 0)),
               # estimate.time = T, need to comment out to run
               stratify = T)
```

```
## Loading required namespace: SuperLearner
```

```
## Qform not specified, using defaults:
```

```
## formula for c1_1:
```

```
## Q.kplus1 ~ 1
```



```

## formula for c1_2:

## Q.kplus1 ~ c1_1 + c2_1

## formula for c1_3:

## Q.kplus1 ~ c1_1 + c2_1 + c1_2 + c2_2

## formula for c1_4:

## Q.kplus1 ~ c1_1 + c2_1 + c1_2 + c2_2 + c1_3 + c2_3

##

## gform not specified, using defaults:

## formula for exposure_1:

## exposure_1 ~ 1

## formula for exposure_2:

## exposure_2 ~ c1_1 + c2_1

## formula for exposure_3:

## exposure_3 ~ c1_1 + c2_1 + c1_2 + c2_2

## formula for exposure_4:

## exposure_4 ~ c1_1 + c2_1 + c1_2 + c2_2 + c1_3 + c2_3

##

## Loading required package: nnls

## Loading required namespace: ranger

## Timing estimate unavailable

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
summary(result)
```

```
## Estimator:  tmle
```

```
## Call:
```

```
## ltmle(data = b, Anodes = c(paste0("exposure_", 1:4)), Lnodes = c("c1_1",
```

```
##      "c2_1", "c1_2", "c2_2", "c1_3", "c2_3", "c1_4", "c2_4"),
```

```
##      Ynodes = c("Y_1", "Y_2", "Y_3", "Y_4"), survivalOutcome = TRUE,
```

```
##      abar = list(treatment = c(1, 1, 1, 1), control = c(0, 0, 0,
```

```
##      0)), stratify = T, SL.library = list(Q = sl.lib, g = sl.lib))
```

```
##
```

```
## Treatment Estimate:
```

```
##      Parameter Estimate:  0.68633
```

```
##      Estimated Std Err:  0.019478
```

```
##              p-value:  <2e-16
```

```
##      95% Conf Interval: (0.64816, 0.72451)
```

```
##
```

```
## Control Estimate:
```

```
##      Parameter Estimate:  0.4596
```

```
##      Estimated Std Err:  0.051186
##                p-value:  <2e-16
##      95% Conf Interval: (0.35928, 0.55992)
##
## Additive Treatment Effect:
##      Parameter Estimate:  0.22673
##      Estimated Std Err:  0.054766
##                p-value:  3.4727e-05
##      95% Conf Interval: (0.11939, 0.33407)
##
## Relative Risk:
##      Parameter Estimate:  1.4933
##      Est Std Err log(RR):  0.11493
##                p-value:  0.0004845
##      95% Conf Interval: (1.1921, 1.8706)
##
## Odds Ratio:
##      Parameter Estimate:  2.5728
##      Est Std Err log(OR):  0.22507
##                p-value:  2.686e-05
##      95% Conf Interval: (1.6551, 3.9993)
```

```
result$fit$g
```

```
## [[1]]
## [[1]]$exposure_1
##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 0.4061438 0.05827759  6.969124 5.19302e-12
##
## [[1]]$exposure_2
##           Risk      Coef
## SL.mean_All  0.09835384 0.002072906
## SL.glm_All   0.09883476 0.000000000
## SL.ranger_All 0.09788259 0.997927094
##
```

```
## [[1]]$exposure_3
##              Risk      Coef
## SL.mean_All    0.1342642 0.7904696
## SL.glm_All     0.1364648 0.2095304
## SL.ranger_All  0.1387118 0.0000000
##
## [[1]]$exposure_4
##              Risk      Coef
## SL.mean_All    0.1162729 0.4032324
## SL.glm_All     0.1173651 0.0000000
## SL.ranger_All  0.1155993 0.5967676
##
##
## [[2]]
## [[2]]$exposure_1
##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 0.4061438 0.05827759 6.969124 5.19302e-12
##
## [[2]]$exposure_2
##              Risk      Coef
## SL.mean_All    0.2355227 0.0000000
## SL.glm_All     0.2303725 0.6424337
## SL.ranger_All  0.2309712 0.3575663
##
## [[2]]$exposure_3
##              Risk      Coef
## SL.mean_All    0.2493752 0.6738447
## SL.glm_All     0.2581789 0.0000000
## SL.ranger_All  0.2560680 0.3261553
##
## [[2]]$exposure_4
##              Risk      Coef
## SL.mean_All    0.2577241 0.6450922
## SL.glm_All     0.2760526 0.3549078
## SL.ranger_All  0.2871431 0.0000000
```

```
result$fit$Q
```

```
## [[1]]
## [[1]]$c1_1
##           Estimate Std. Error  t value      Pr(>|t|)
## (Intercept) 0.7830343 0.03023081 25.90186 1.196625e-105
##
## [[1]]$c1_2
##           Risk      Coef
## SL.mean_All  0.06778137 0.1729875
## SL.glm_All   0.06718843 0.8270125
## SL.ranger_All 0.16142769 0.0000000
##
## [[1]]$c1_3
##           Risk      Coef
## SL.mean_All  0.1105063 0.5142731
## SL.glm_All   0.1106044 0.4857269
## SL.ranger_All 0.1433396 0.0000000
##
## [[1]]$c1_4
##           Risk Coef
## SL.mean_All  0.1811478 1
## SL.glm_All   0.1902496 0
## SL.ranger_All 0.1890942 0
##
##
## [[2]]
## [[2]]$c1_1
##           Estimate Std. Error  t value      Pr(>|t|)
## (Intercept) -0.1619542 0.04814336 -3.363999 0.0008284433
##
## [[2]]$c1_2
##           Risk Coef
## SL.mean_All  0.09869432 1
## SL.glm_All   0.10083616 0
```

```

## SL.ranger_All 0.12279127    0
##
## [[2]]$c1_3
##
## Risk Coef
## SL.mean_All 0.06287278    1
## SL.glm_All 0.07407730    0
## SL.ranger_All 0.08018747    0
##
## [[2]]$c1_4
##
## Risk Coef
## SL.mean_All 0.09408145    1
## SL.glm_All 0.25518219    0
## SL.ranger_All 0.11725773    0

```

References

- Erica E. M. Moodie, David A. Stephens, and Marina B. Klein. A marginal structural model for multiple-outcome survival data: assessing the impact of injection drug use on several causes of death in the canadian co-infection cohort. *Stat Med*, 33(8):1409–1425, 2014.
- J. G. Young, M. A. Hernán, S. Picciotto, and J. M. Robins. Relation between three classes of structural models for the effect of a time-varying exposure on survival. *Lifetime Data Anal*, 16(1):71–84, 2010.