



Leiden University
Medical Center

Confounding adjustment and estimating treatment effects

With models

Edouard Fu, PhD

Department of Clinical Epidemiology, LUMC



Content of lecture

1. What your dataset needs to look like
2. Fitting models for the weights
3. Fitting the outcome model

1. Dataset requirements

Dataset in longitudinal format

ID	Time	L ₀	L _k	A ₀	A _k	Y	C _k	IPTW	IPCW	C _{k-art}
1	0	0	0	0	0	0	0			
1	1	0	0	0	0	0	0			
1	2	0	1	0	1	0	0			
...			
1	59	0	1	0	1	0	0			
2	0	1	1	1	1	0	0			
2	1	1	1	1	1	0	0			
2	2	1	1	1	1	0	0			
...			
2	34	1	1	1	1	1	0			
3	0	0	0	0	0	0	0			
3	1	0	0	0	0	0	0			
3	2	0	0	0	0	0	1			

- ID: personal identifier
- Time: time (in months)
- L₀: baseline confounder
- L_k: time-varying confounder
- A₀: baseline treatment assignment
- A_k: time-varying treatment
- Y: all-cause mortality
- C_k: loss to follow-up
- IPTW: inverse probability of treatment weights
- IPCW: inverse probability of censoring weights
- C_{k-art}: artificial censoring

Temporality is key

ID	Time	L_0	L_k	A_0	A_k	Y_k	C_k	IPTW	IPCW	C_{k_art}
1	0	0	0	0	0	0	0			
1	1	0	0	0	0	0	0			
1	2	0	1	0	1	0	0			
...			
1	59	0	1	0	1	0	0			
<hr/>										
2	0	1	1	1	1	0	0			
2	1	1	1	1	1	0	0			
2	2	1	1	1	1	0	0			
...			
2	34	1	1	1	1	1	0			
<hr/>										
3	0	0	0	0	0	0	0			
3	1	0	0	0	0	0	0			
3	2	0	0	0	0	0	1			

Within each row, need to ensure temporality (L_k , A_k , Y)



2. Fitting weight models

IPTW (weights to adjust for time-varying confounding)

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	C _k	IPTW	IPCW	C _{k-art}
1	0	0	0	0	0	0	0			
1	1	0	0	0	0	0	0			
1	2	0	1	0	1	0	0			
...			
1	59	0	1	0	1	0	0			
2	0	1	1	1	1	0	0			
2	1	1	1	1	1	0	0			
2	2	1	1	1	1	0	0			
...			
2	34	1	1	1	1	1	0			
3	0	0	0	0	0	0	0			
3	1	0	0	0	0	0	0			
3	2	0	0	0	0	0	1			

Goal: A_k is not predicted anymore by the past (L_k) at each timepoint

How: Give everyone IPTW

$$W^A = \prod_{t=0}^{59} \frac{1}{Pr[A_k | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k]}$$

Fit the following pooled logistic model:

$$\text{logit}[pr(A_k = 1 | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1} = a, L_0, \bar{L}_k)] = \alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k$$

R code

```
# fit pooled logistic model
mod <- glm(A_k ~ Time + I(Time^2) + L_0 + L_k,
family = binomial(), data = dat)
```

$$\hat{f} = Pr(A_k = 1 | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k) = \frac{1}{1 + e^{-(\alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k)}}$$

$$mod = logit(\hat{f}) = \alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k$$

```
# predict
dat$probA.d <- predict(mod, type = 'response')
```

$$\hat{f}(A_k | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1} = \bar{a}_{k-1}, L_0 = l_0, \bar{L}_k = \bar{l}_k)$$

```
# calculate weight
dat$w <- ifelse(dat$A_k==1, (1/dat$probA.d),
(1/(1-dat$probA.d)))
```

$$\frac{1}{\hat{f}(A_k | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1} = \bar{a}_{k-1}, L_0 = l_0, \bar{L}_k = \bar{l}_k)}$$

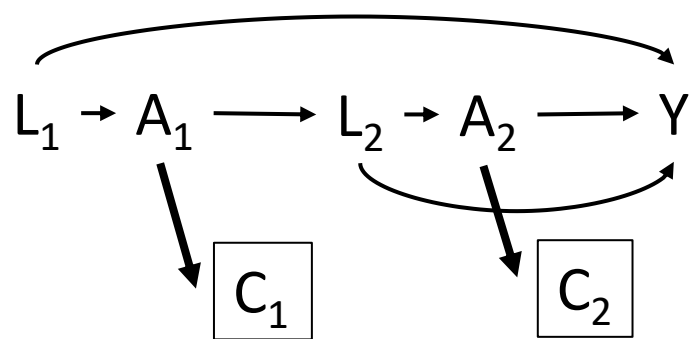
```
# calculate cumulative product of weights
dat$w_cum <- ave(dat$w, dat$id, FUN=function(x)
cumprod(x))
```

$$\hat{W}^A = \prod_{t=0}^{59} \frac{1}{\hat{f}(A_k | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1} = \bar{a}_{k-1}, L_0 = l_0, \bar{L}_k = \bar{l}_k)}$$

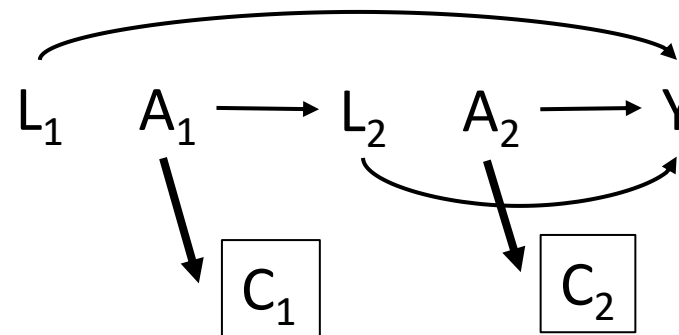
IPTW (weights to adjust for time-varying confounding)

ID	Time	L_0	L_k	A_0	A_k	Y_k	C_k	IPTW	IPCW	C_{k_art}
1	0	0	0	0	0	0	0	1.5		
1	1	0	0	0	0	0	0	2.2		
1	2	0	1	0	1	0	0	3.8		
...		
1	59	0	1	0	1	0	0	10.2		
2	0	1	1	1	1	0	0	1.3		
2	1	1	1	1	1	0	0	1.5		
2	2	1	1	1	1	0	0	2.6		
...		
2	34	1	1	1	1	1	0	5.4		
3	0	0	0	0	0	0	0	1.2		
3	1	0	0	0	0	0	0	2.0		
3	2	0	0	0	0	0	1	NA		

Weights ensure L_k no longer predicts A_k for every timepoint k

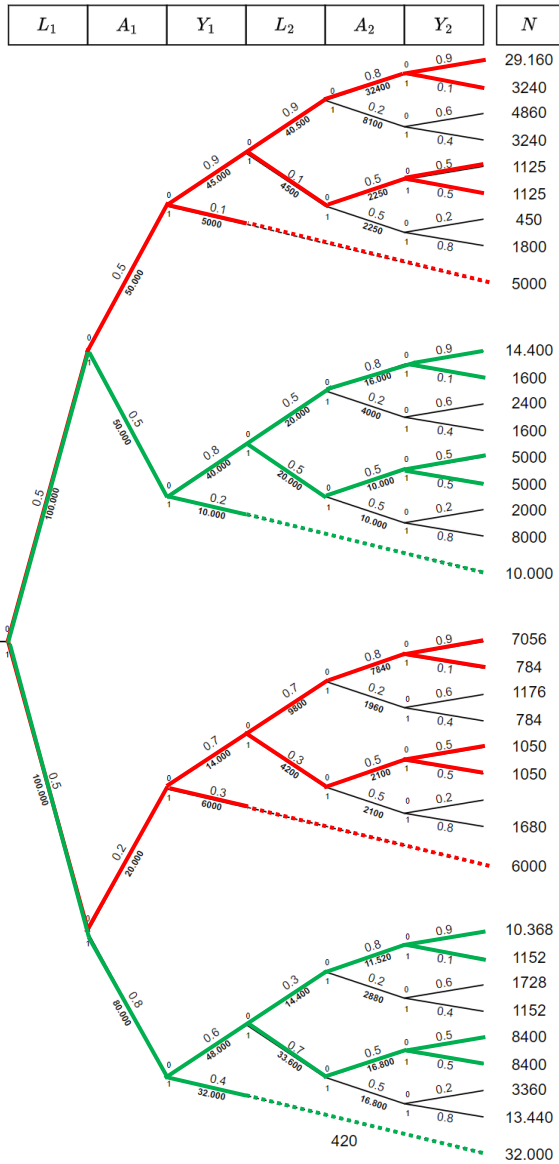


IPTW



“Association is causation”

Note that we are making a lot of assumptions!



```
mod <- glm(A_k ~ Time + I(Time^2) + L_0 + L_k, family = binomial(),
data = dat)
```

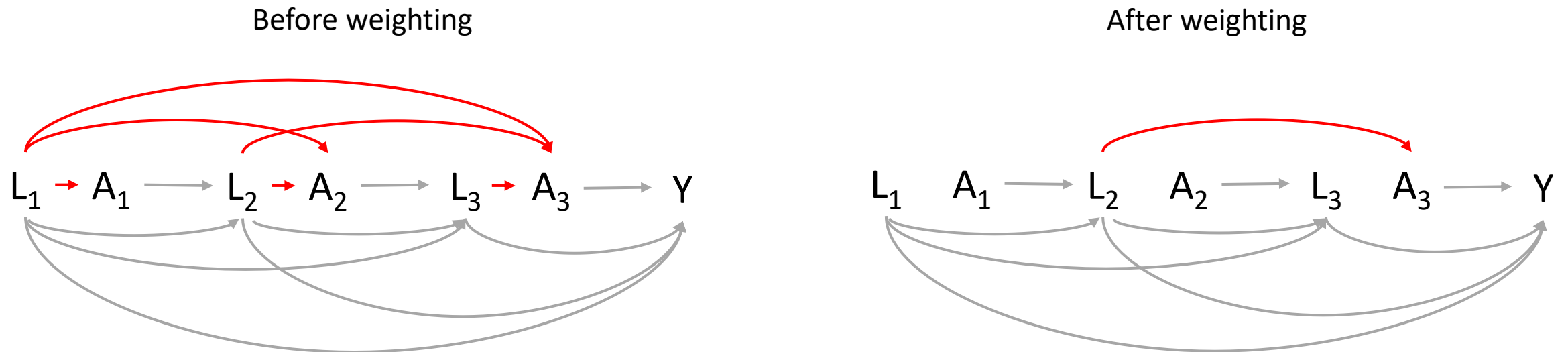
We fit one model on the entire tree... Is it realistic we can properly model the entire treatment process with one parametric model?

- Could model time more flexibly (e.g. restricted cubic spline)
- Could add interactions (between time and confounders)
- Could fit separate models for each treatment group
- Could fit separate model for each timepoint

Model misspecification (bias-variance trade-off)

Assumptions about recency of confounders

If we only put most recent time-varying confounder value (+ baseline confounder) in our weighting model



Misspecified model leads to remaining red arrow after weighting, so residual confounding!
(even if all time-varying confounders are measured)

Weights can become very large

Solutions

1. Truncate the weights at the n^{th} percentile (e.g. 99th) or at a certain value

```
dat$w.trunc <- ifelse(dat$w>10, 10, dat$w)
```

2. Use stabilized weights

$$SW_{v1}^A = \prod_{t=0}^{59} \frac{\Pr[A_k | \bar{C}_k = \bar{0}, \bar{A}_{k-1}]}{\Pr[A_k | \bar{C}_k = \bar{0}, \bar{D}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k]}$$

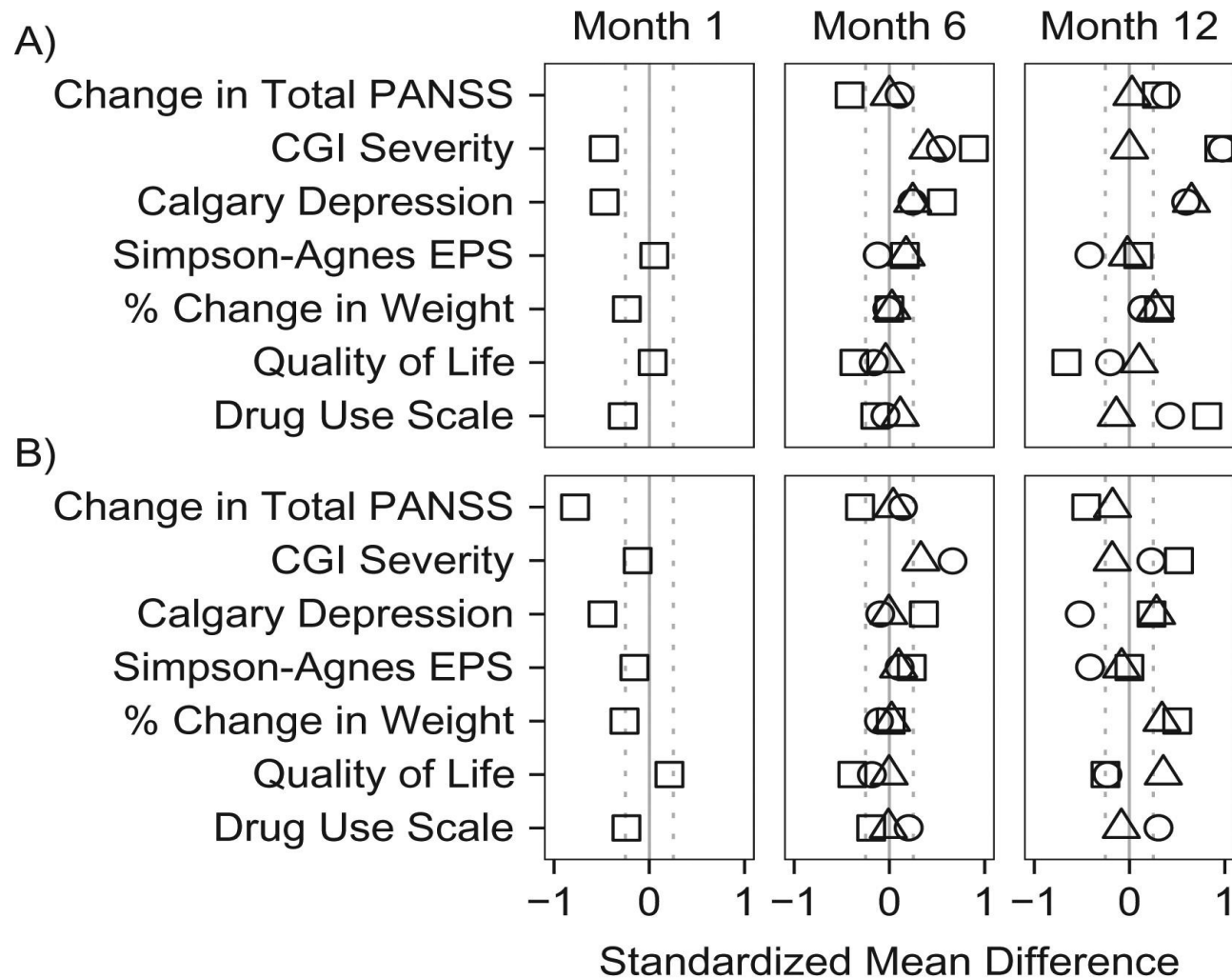
$$SW_{v2}^A = \prod_{t=0}^{59} \frac{\Pr[A_k | \bar{C}_k = \bar{0}, \bar{A}_{k-1}, \mathbf{L}_0]}{\Pr[A_k | \bar{C}_k = \bar{0}, \bar{D}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k]}$$

Fit two pooled logistic models:

Numerator: $\text{logit}[\text{pr}(A_k = 1 | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, = a, \bar{L}_k)] = \alpha_{0t} (+ \alpha_1^T L_0)$

Denominator: $\text{logit}[\text{pr}(A_k = 1 | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, = a, \bar{L}_k)] = \alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k$

Checking covariate balance at each timepoint



JW Jackson, *Am J Epidemiol* (2019), Diagnosing Covariate Balance Across Levels of Right-Censoring Before and After Application of Inverse-Probability-of-Censoring Weights

Repeat same process for IPCW

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	C _k	IPTW	IPCW	C _{k-art}
1	0	0	0	0	0	0	0	1.5	1.1	
1	1	0	0	0	0	0	0	2.2	1.3	
1	2	0	1	0	1	0	0	3.8	1.4	
...	
1	59	0	1	0	1	0	0	10.2	1.8	
2	0	1	1	1	1	0	0	1.3	1.2	
2	1	1	1	1	1	0	0	1.5	1.5	
2	2	1	1	1	1	0	0	2.6	1.8	
...	
2	34	1	1	1	1	1	0	5.4	2.0	
3	0	0	0	0	0	0	0	1.2	1.3	
3	1	0	0	0	0	0	0	2.0	2.7	
3	2	0	0	0	0	0	1	NA	NA	

$$W^C = \prod_{t=0}^{59} \frac{1}{Pr[C_k | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k]}$$

Fit the following pooled logistic model:

$$\text{logit}[pr(C_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1} = a, L_0, \bar{L}_k)] = \alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k + \alpha_3 A_{k-1}$$

3. Fitting outcome models

Artificial censoring

ID	Time	L_0	L_k	A_0	A_k	Y_k	C_k	IPTW	IPCW	C_{k_art}
1	0	0	0	0	0	0	0	1.5	1.1	0
1	1	0	0	0	0	0	0	2.2	1.3	0
1	2	0	1	0	1	0	0	3.8	1.4	1
...
1	59	0	1	0	1	0	0	10.2	1.8	1
2	0	1	1	1	1	0	0	1.3	1.2	0
2	1	1	1	1	1	0	0	1.5	1.5	0
2	2	1	1	1	1	0	0	2.6	1.8	0
...
2	34	1	1	1	1	1	0	5.4	2.0	0
3	0	0	0	0	0	0	0	1.2	1.3	0
3	1	0	0	0	0	0	0	2.0	2.7	0
3	2	0	0	0	0	0	1	NA	NA	0

Determine artificial censoring based on assigned strategy:

“Start treatment and always use”
vs. “Never start treatment”

Fit the outcome model

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	C _k	IPTW	IPCW	C _{k-art}
1	0	0	0	0	0	0	0	1.5	1.1	0
1	1	0	0	0	0	0	0	2.2	1.3	0
1	2	0	1	0	1	0	0	3.8	1.4	1
...
1	59	0	1	0	1	0	0	10.2	1.8	1
2	0	1	1	1	1	0	0	1.3	1.2	0
2	1	1	1	1	1	0	0	1.5	1.5	0
2	2	1	1	1	1	0	0	2.6	1.8	0
...
2	34	1	1	1	1	1	0	5.4	2.0	0
3	0	0	0	0	0	0	0	1.2	1.3	0
3	1	0	0	0	0	0	0	2.0	2.7	0
3	2	0	0	0	0	0	1	NA	NA	0

Fit the following **weighted** pooled logistic model:

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{C}_{k-1}(\text{art}) = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_0)] = \alpha_{0t} + \alpha_1 A_0$$

Then, the marginal ln(HR) for treatment is given by α_1 (under the assumption that outcome incidence is <10% in each time interval)

If baseline confounders were used in the numerator of the stabilized weights, then they have to be added to the outcome model:

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{C}_{k-1}(\text{art}) = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_0, \mathbf{L}_0)] = \alpha_{0t} + \alpha_1 A_0 + \alpha_2^T \mathbf{L}_0$$

R code

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{C}_{k-1}(\text{art}) = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_0)] = \alpha_{0t} + \alpha_1 A_0$$

```
# fit outcome model
```

```
outcome_mod <- glm(Y_k ~ Time + I(Time^2) + A_0 + L_0,  
family = binomial(), weight = IPTW*IPCW,  
data = subset(dat, C_k==0 & C_k_art == 0))
```

```
# obtain hazard ratio
```

```
exp(coef(outcome_mod))
```

Assessing effect modification by baseline variable

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{C}_{k-1}(\text{art}) = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_0, V)] = \alpha_{0t} + \alpha_1 A_0 + \alpha_2 V + \alpha_3 A_0 V$$

fit outcome model

```
outcome_mod <- glm(Y_k ~ Time + I(Time^2) + A_0 + V + A_0:V,  
family = binomial(), weight = IPTW*IPCW,  
data = subset(dat, C_k==0 & C_k_art == 0))
```

95% confidence intervals

Need to account for use of IPTW/IPCW (and perhaps repeated use of same individual through sequential trials or cloning)

Solutions:

1. Robust standard error (e.g. survey package in R)

```
outcome_mod <- svyglm(Y_k ~ Time + I(Time^2) + A_0 + L_0,  
family = binomial(), design = svydesign(id = ~id, weights = ~IPTW*IPCW,  
data = subset(dat, C_k==0 & C_k_art == 0))  
exp(confint(outcome_mod))
```

2. Nonparametric bootstrap



Leiden University
Medical Center

Questions

e.l.fu@lumc.nl



Additional topics

- Dose-response models
- Constructing inverse probability weighted survival curves
- Competing risks
- Implementing clone-censor-weight
- Implementing sequential trials

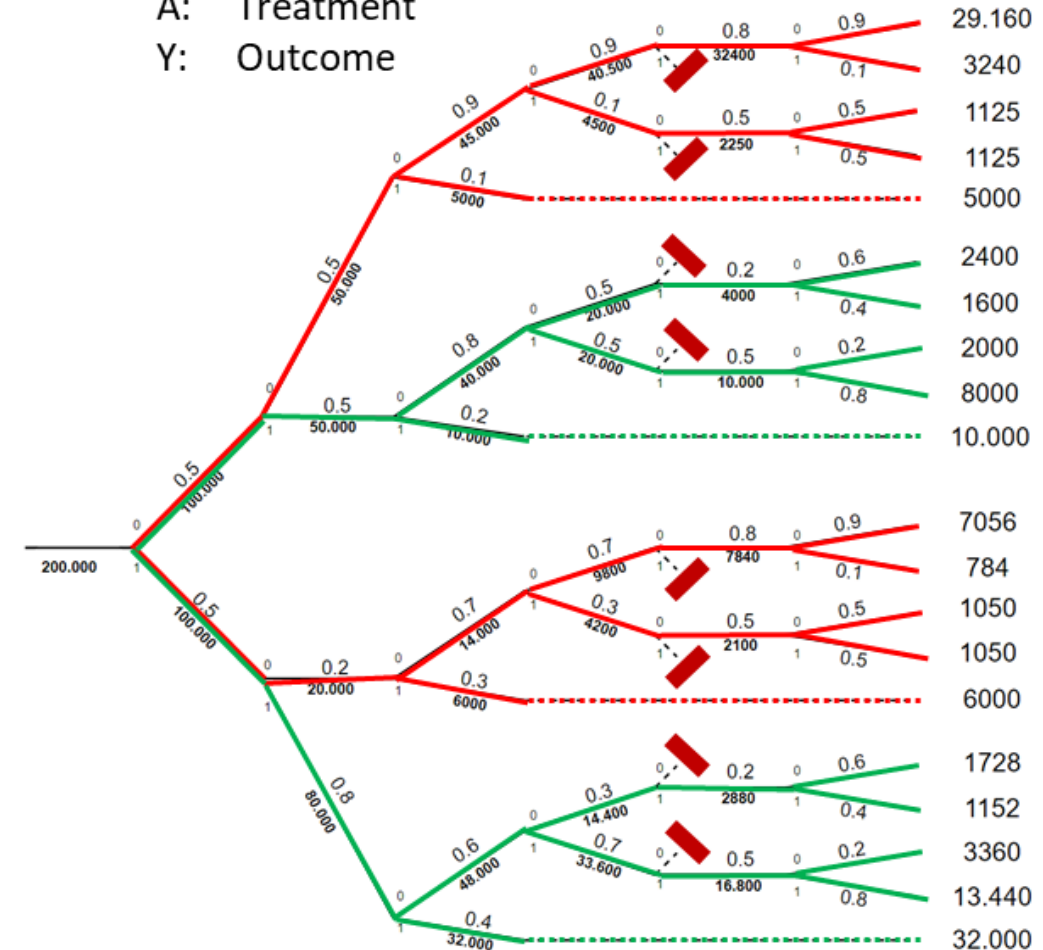
4. Dose-response models

Fitting a dose-response model instead of censoring

ID	Time	L_0	L_k	A_0	A_k	Y_k	C_k	IPTW	IPCW	C_{k_art}
1	0	0	0	0	0	0	0	1.5	1.1	0
1	1	0	0	0	0	0	0	2.2	1.3	0
1	2	0	1	0	1	0	0	3.8	1.4	1
...
1	59	0	1	0	1	0	0	10.2	1.8	1
2	0	1	1	1	1	0	0	1.3	1.2	0
2	1	1	1	1	1	0	0	1.5	1.5	0
2	2	1	1	1	1	0	0	2.6	1.8	0
...
2	34	1	1	1	1	1	0	5.4	2.0	0
3	0	0	0	0	0	0	0	1.2	1.3	0
3	1	0	0	0	0	0	0	2.0	2.7	0
3	2	0	0	0	0	0	1	NA	NA	0

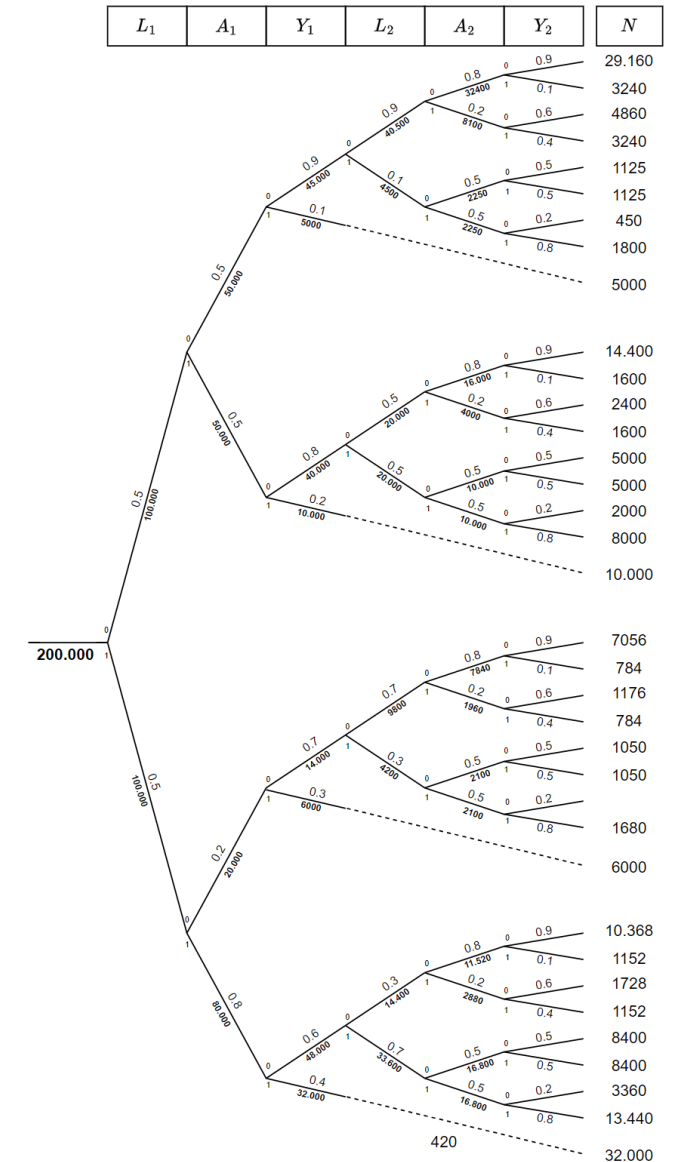
L_1	A_1	Y_1	L_2	A_2	Y_2	N
-------	-------	-------	-------	-------	-------	-----

L: Confounder
A: Treatment
Y: Outcome



Fitting a dose-response model instead of censoring

ID	Time	L_0	L_k	A_0	A_k	Y_k	C_k	IPTW	IPCW	C_{k_art}	A_{tot}
1	0	0	0	0	0	0	0	1.5	1.1	0	0
1	1	0	0	0	0	0	0	2.2	1.3	0	0
1	2	0	1	0	1	0	0	3.8	1.4	1	1
...
1	59	0	1	0	1	0	0	10.2	1.8	1	23
2	0	1	1	1	1	0	0	1.3	1.2	0	1
2	1	1	1	1	1	0	0	1.5	1.5	0	2
2	2	1	1	1	1	0	0	2.6	1.8	0	3
...
2	34	1	1	1	1	1	0	5.4	2.0	0	33
3	0	0	0	0	0	0	0	1.2	1.3	0	0
3	1	0	0	0	0	0	0	2.0	2.7	0	0
3	2	0	0	0	0	0	1	NA	NA	0	0



Censoring vs. dose-response model

Artificial censoring approach

Fit the following **weighted** pooled logistic model:

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{C}_{k-1}(\text{art}) = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_0)] = \alpha_{0t} + \alpha_1 A_0$$

HR for always treat vs. never treat:

$$e^{\alpha_1}$$

Dose-response approach

Fit the following **weighted** pooled logistic model:

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_k)] = \gamma_{0t} + \gamma_1 A_{tot} + \gamma_2 (A_{tot})^2$$

HR for each additional month of treatment:

$$e^{\gamma_1 A_{tot} + \gamma_2 (A_{tot})^2}$$

HR for always treat vs. never treat:

$$e^{\gamma_1 * 60 + \gamma_2 * 60^2}$$

Different dose-response models

Total duration of treatment

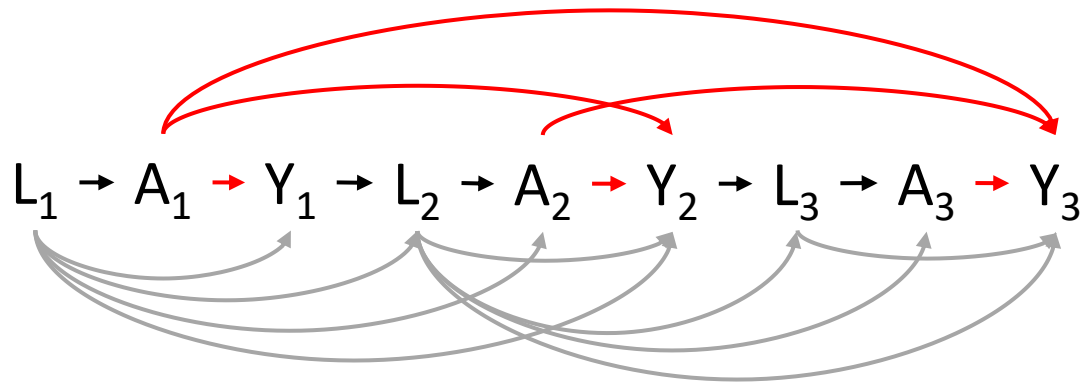
$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_k)] = \gamma_{0t} + \gamma_1 \sum_{k=0}^t A_k + \gamma_2 \left(\sum_{k=0}^t A_k \right)^2$$

Average duration of treatment

$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_k)] = \delta_{0t} + \delta_1 \left(\frac{1}{t} \sum_{k=0}^t A_k \right) + \delta_2 \left(\frac{1}{t} \sum_{k=0}^t A_k \right)^2$$

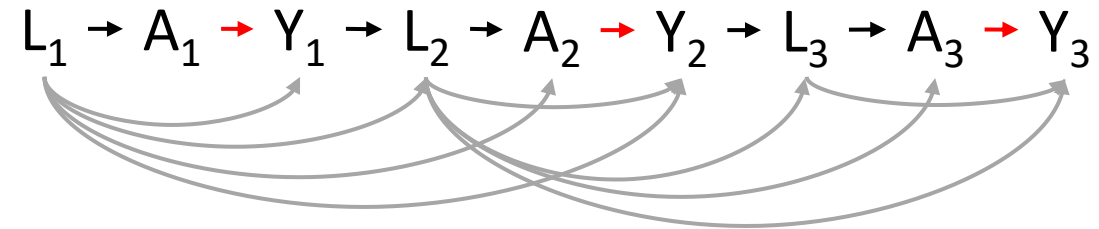
Sometimes dose-response model not needed

Hazard at each timepoint k depends on cumulative treatment history



$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_k)] = \gamma_{0t} + \gamma_1 \sum_{k=0}^t A_k + \gamma_2 \left(\sum_{k=0}^t A_k \right)^2$$

Hazard at each timepoint k only depends on most recent treatment



$$\text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_k)] = \beta_{0t} + \beta_1 A_k$$

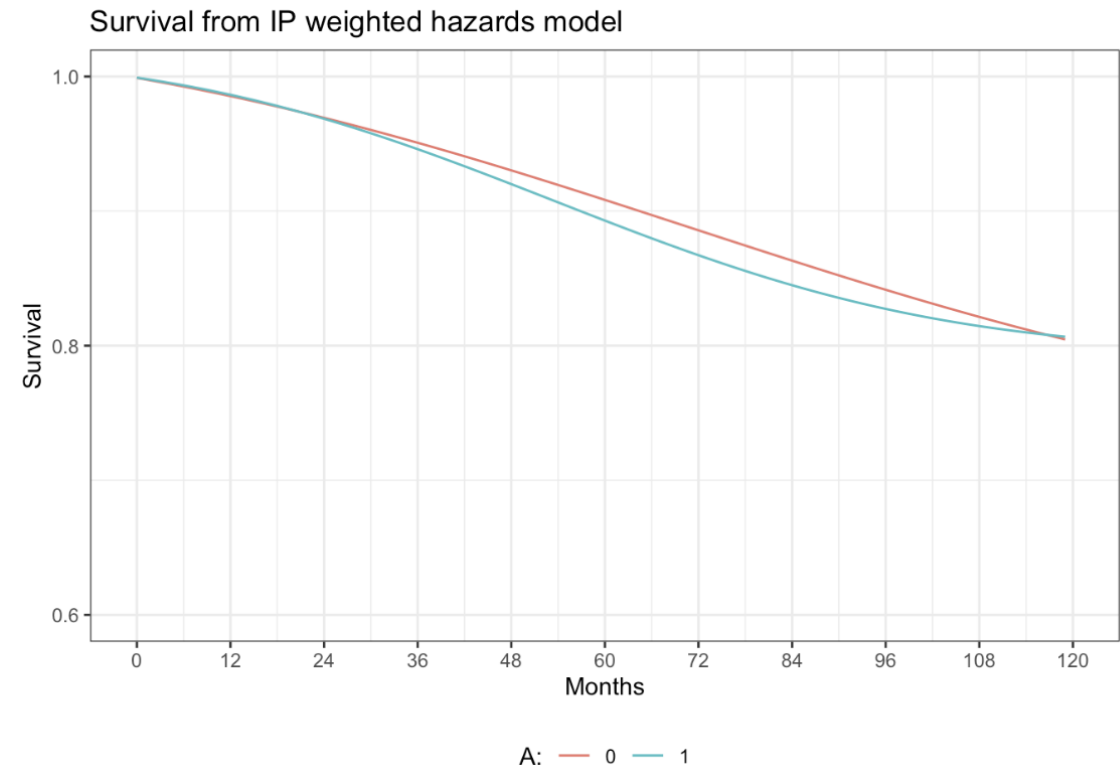
Useful references

- Danaei G, Rodríguez LA, Cantero OF, Logan R, Hernán MA. Observational data for comparative effectiveness research: an emulation of randomised trials of statins and primary prevention of coronary heart disease. *Stat Methods Med Res*. 2013 Feb;22(1):70-96. doi: 10.1177/0962280211403603. Epub 2011 Oct 19. PMID: 22016461; PMCID: PMC3613145.
- Toh S, Hernán MA. Causal inference from longitudinal studies with baseline randomization. *Int J Biostat*. 2008 Oct 19;4(1):Article 22. doi: 10.2202/1557-4679.1117. PMID: 20231914; PMCID: PMC2835458.

5. Parametric estimation of weighted survival curves

Making survival curves

ID	Time	L_0	L_k	A_0	A_k	Y_k
1	0	0	0	0	0	0
1	1	0	0	0	0	0
1	2	0	1	0	1	0
...
1	59	0	1	0	1	0
2	0	1	1	1	1	0
2	1	1	1	1	1	0
2	2	1	1	1	1	0
...
2	34	1	1	1	1	1
3	0	0	0	0	0	0
3	1	0	0	0	0	0
3	2	0	0	0	0	0



How is survival calculated?

ID	Time	L_0	L_k	A_0	A_k	Y_k
1	0	0	0	0	0	0
1	1	0	0	0	0	0
1	2	0	1	0	1	0
...
1	59	0	1	0	1	0
2	0	1	1	1	1	0
2	1	1	1	1	1	0
2	2	1	1	1	1	0
...
2	34	1	1	1	1	1
3	0	0	0	0	0	0
3	1	0	0	0	0	0
3	2	0	0	0	0	0

Survival

$$\Pr[Y_k = 0] = \prod_{m=1}^k \Pr[Y_m = 0 | Y_{m-1} = 0]$$

$$\begin{aligned} \Pr[Y_2 = 0] &= \Pr[Y_2 = 0 | Y_1 = 0] * \Pr[Y_1 = 0] \\ &= 0.95 * 0.90 = 0.855 \end{aligned}$$

Hazard

$$\Pr[Y_k = 1 | Y_{k-1} = 0]$$

$$\begin{aligned} \Pr[Y_2 = 1 | Y_1 = 0] &= \frac{\text{no. of deaths during interval 2}}{\text{no. of people alive during interval 2}} \\ &= 0.05 \end{aligned}$$

$$\Pr[Y_2 = 0 | Y_1 = 0] = 1 - \Pr[Y_2 = 1 | Y_1 = 0] = 1 - 0.05 = 0.95$$

Calculating survival from hazards

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	C _k	IPTW	IPCW	C _{k_art}
1	0	0	0	0	0	0	0	1.5	1.1	0
1	1	0	0	0	0	0	0	2.2	1.3	0
1	2	0	1	0	1	0	0	3.8	1.4	1
...
1	59	0	1	0	1	0	0	10.2	1.8	1
<hr/>										
2	0	1	1	1	1	0	0	1.3	1.2	0
2	1	1	1	1	1	0	0	1.5	1.5	0
2	2	1	1	1	1	0	0	2.6	1.8	0
...
2	34	1	1	1	1	1	0	5.4	2.0	0
<hr/>										
3	0	0	0	0	0	0	0	1.2	1.3	0
3	1	0	0	0	0	0	0	2.0	2.7	0
3	2	0	0	0	0	0	1	NA	NA	0

Survival from hazard

$$\begin{aligned} \Pr[Y_k = 0] &= \prod_{m=1}^k \Pr[Y_m = 0 | Y_{m-1} = 0] \\ &= \prod_{m=1}^k (1 - \Pr[Y_m = 1 | Y_{m-1} = 0]) \end{aligned}$$

Estimating hazards from a weighted logistic model

$$\text{logit}[\text{pr}(Y_{k+1} = 1 | Y_k = 0, C_k = 0, C_{k_art} = 0, A_0)] = \alpha_{0,k} + \alpha_1 A_0 + \alpha_2 A_0 * k + \alpha_3 A_0 * k^2$$

$$\text{where } \alpha_{0,k} = \alpha_0 + \alpha_4 * k + \alpha_5 * k^2$$

Use model to predict hazards at each timepoint

Time	Time ²	A ₀	h _k	S _k	S _{k_cum}
0	0	1			
1	1	1			
2	4	1			
...		...			
59	3481	1			

Dataset 1: Prediction under always treatment

$$\text{logit}[pr(Y_{k+1} = 1 | Y_k = 0, C_k = 0, C_{k_art} = 0, A_0 = 1)] = \alpha_{0,k} + \alpha_1 + \alpha_2 * k + \alpha_3 * k^2$$

Time	Time ²	A ₀	h _k	S _k	S _{k_cum}
0	0	0			
1	1	0			
2	4	0			
...		...			
59	3481	0			

Dataset 2: Prediction under never treatment

$$\text{logit}[pr(Y_{k+1} = 1 | Y_k = 0, C_k = 0, C_{k_art} = 0, A_0 = 1)] = \alpha_{0,k}$$

R code (1/2)

fit of weighted hazards model

```
outcome_mod <- glm(Y_k==1 ~ Time + Timesq + A_0 + I(A_0*Time) + I(A_0*Timesq),  
family = binomial(), weight = IPTW*IPCW,  
data = subset(dat, C_k==0 & C_k_art == 0))
```

creation of “treated” and “untreated” empty datasets

```
dat_notreat <- data.frame(cbind(0, seq(0, 59), (seq(0, 59))^2))  
dat_treat <- data.frame(cbind(1, seq(0, 59), (seq(0, 59))^2))
```

```
colnames(dat_notreat) <- c("A_0", "Time", "Timesq")  
colnames(dat_treat) <- c("A_0", "Time", "Timesq")
```

Time	Time ²	A ₀	h _k	S _k	S _{k_cum}
0	0	1			
1	1	1			
2	4	1			
...		...			
59	3481	1			

R code (2/2)

Calculating hazard in each person-month

```
dat_notreat$h_k <- predict(outcome_mod, dat_notreat, type="response")
```

```
dat_treat$h_k <- predict(outcome_mod, dat_treat, type="response")
```

Calculating survival in each person-month

```
dat_notreat$S_k <- 1-dat_notreat$h_k
```

```
dat_treat$S_k <- 1- dat_treat$h_k
```

Calculating cumulative survival

```
dat_notreat$S_k_cum <- cumprod(dat_notreat$ S_k)
```

```
dat_treat$S_k_cum <- cumprod(dat_treat$S_k)
```

Time	Time ²	A ₀	h _k	S _k	S _{k_cum}
0	0	1			
1	1	1			
2	4	1			
...		...			
59	3481	1			

Useful references

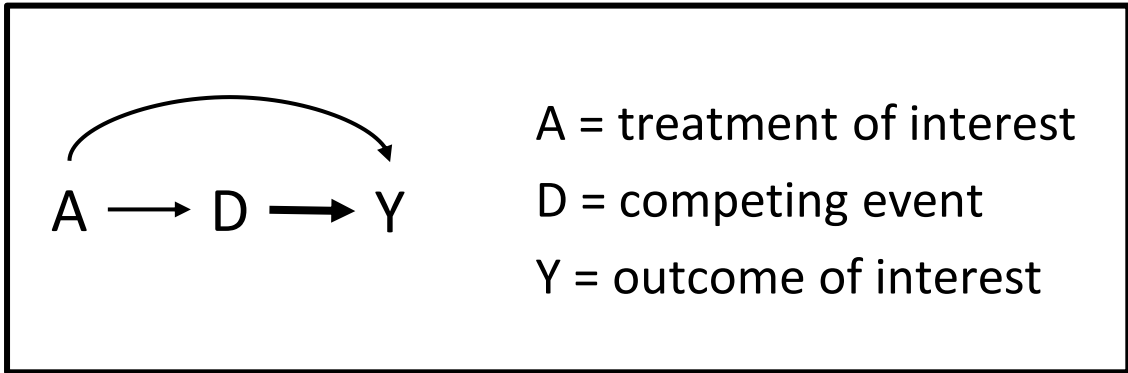
- Hernán MA, Robins JM (2020). Causal Inference: What If. Boca Raton: Chapman & Hall/CRC. Chapter 17 Causal survival analysis
- <https://remlapmot.github.io/cibookex-r/causal-survival-analysis.html> (R code)
- Cole SR, Hernán MA. Adjusted survival curves with inverse probability weights. Comput Methods Programs Biomed. 2004 Jul;75(1):45-9. doi: 10.1016/j.cmpb.2003.10.004. PMID: 15158046.

6. Competing risks

What is a competing event?

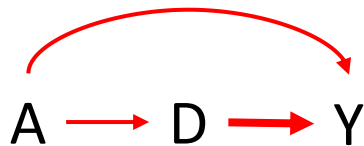
- A competing (risk) event is any event that makes it impossible for the event of interest to occur
- E.g., if interested in the effect of SGLT-2 inhibitors vs. placebo on dialysis, then death is a competing event
- Similarly applies to randomized trials and observational studies

How to handle competing events?



1. Total effect of treatment

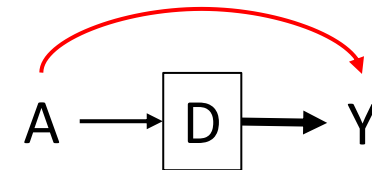
$$\Pr[Y^{a=1} = 1] \text{ vs. } \Pr[Y^{a=0} = 1]$$



“What is the total effect of treatment on the outcome, part of which may be mediated by the competing event?”

2. Controlled direct effect of treatment

$$\Pr[Y^{a=1,d=0} = 1] \text{ vs. } \Pr[Y^{a=0,d=0} = 1]$$

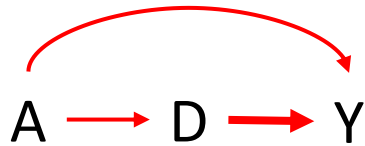


“What is the direct effect of treatment on the outcome, in a world where we eliminate the competing event?”

Total effect

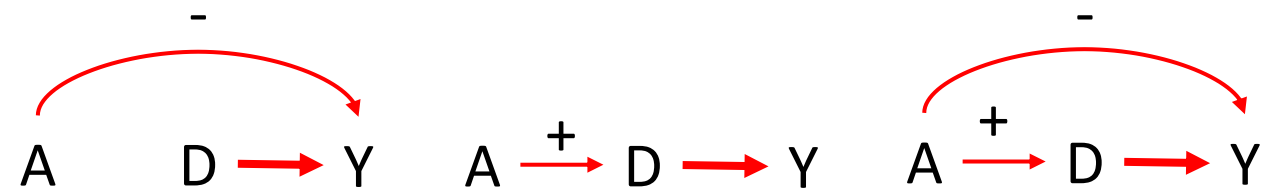
1. Total effect of treatment

$\Pr[Y^{a=1} = 1] \text{ vs. } \Pr[Y^{a=0} = 1]$



“What is the total effect of treatment on the outcome, part of which may be mediated by the competing event?”

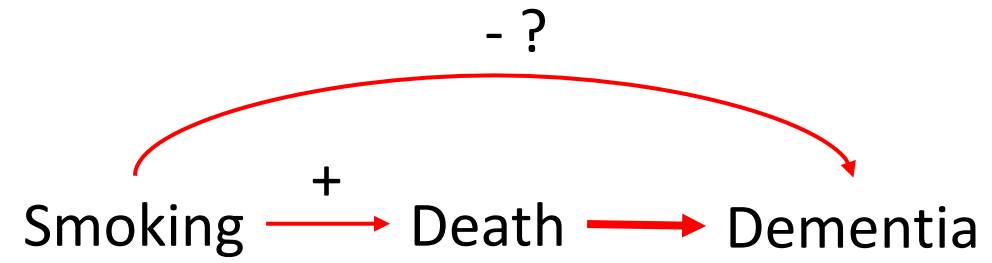
- Can be easily identified in a perfect randomized trial
- However, does not answer question about mechanism: if we find $\Pr[Y^{a=1} = 1] < \Pr[Y^{a=0} = 1]$, is this due to treatment A lowering Y, due to A increasing D (thereby preventing A), or a combination of both?



Most extreme example

- We conduct a RCT testing a new pill vs. placebo on the 5-year risk of dialysis
 - Assume that the trial is perfect (infinite sample size, perfect adherence, no loss to follow-up etc)
- After completing the trial, we find $\Pr[Y = 1|A = 1] = 0$ and $\Pr[Y = 1|A = 0] = 0.4$
- We conclude that the new pill is very effective in preventing dialysis
- However, the pill is poisonous and kills those that ingest it within 1 minute
- Are we still interested in the total effect?

Less extreme example



How to estimate the total effect

Dataset to fit the **outcome model**

ID	Time	L_0	L_k	A_0	A_k	Y_k	IPTW	D_k
1	0	0	0	0	0	0	1.5	0
1	1	0	0	0	0	0	2.2	0
1	2	0	1	0	0	0	3.8	0
1	3	0	1	0	0	0	4.2	1
1	4	0	NA	0	NA	0	4.2	1
1	5	0	NA	0	NA	0	4.2	1
...
1	59	0	NA	0	1	0	4.2	1

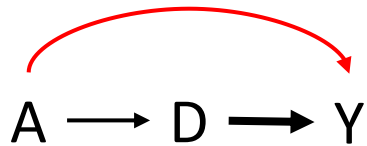
Dataset to fit the **weight model**

ID	Time	L_0	L_k	A_0	A_k	Y_k	IPTW	D_k
1	0	0	0	0	0	0	1.5	0
1	1	0	0	0	0	0	2.2	0
1	2	0	1	0	0	0	3.8	0
1	3	0	1	0	0	0	4.2	1

Controlled direct effect

2. Controlled direct effect of treatment

$$\Pr[Y^{a=1,d=0} = 1] \text{ vs. } \Pr[Y^{a=0,d=0} = 1]$$



“What is the direct effect of treatment on the outcome, in a world where we eliminate the competing event?”

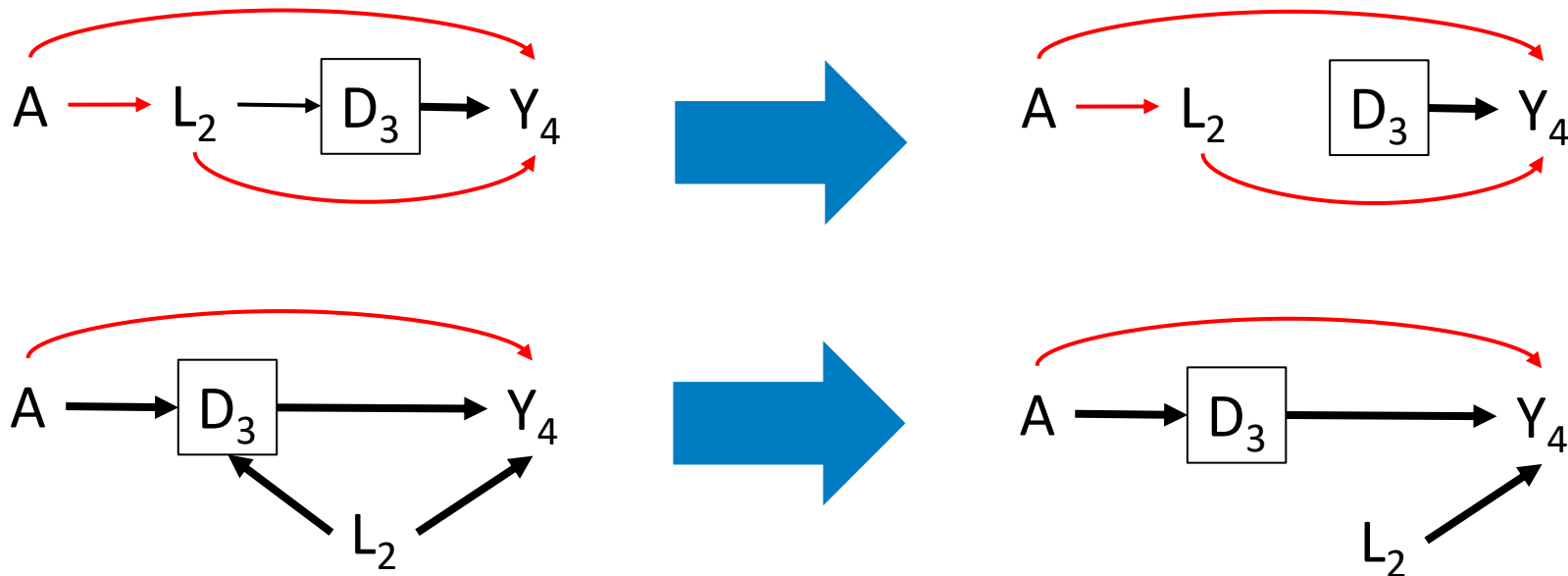
- Helps to elucidate mechanisms
- However, also difficult to interpret: “a world where we eliminate the competing event” → How are we going to eliminate this in the real world? What is this potential intervention?
- Additional assumptions are required to identify $Y^{a,d=0}$: $Y^{a,d=0} \perp A$ and $Y^{a,d=0} \perp D$

Controlled direct effect

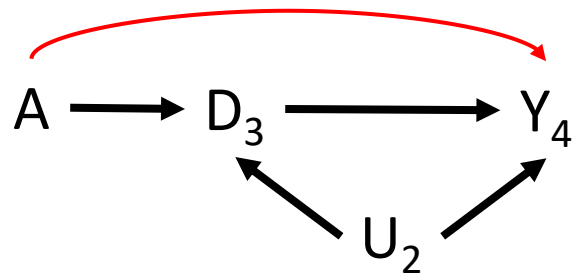
- Competing event is considered a censoring event: value of $Y^{a=1,d=0}$ is unknown after competing event occurs
- “A censoring event is any event occurring in the study that ensures the values of all future counterfactual outcomes under treatment level a that are of interest are unknown/missing, even for an individual who actually received treatment level a .”
- Thus, if you censor for competing events you are implicitly targeting the CDE
- We try to simulate what would have happened, had the competing event not occurred
 - Intuitively, we upweight people without the competing event who have similar characteristics as those with the competing event

Assumptions for censoring

- Unbiased estimation requires absence of backdoor paths between A and Y_4 and no backdoor paths between D_3 and Y_4 (data shown below are from a randomized controlled trial)
- Use IPCW to remove arrow between L_2 and D_3



Violation of assumptions



If there are unmeasured common causes of D_3 and Y_4 , then we cannot validly estimate the controlled direct effect

How to estimate the controlled direct effect

Dataset format

ID	Time	L_0	L_k	A_0	A_k	Y_k	IPTW	D_k	IPCW
1	0	0	0	0	0	0	1.5	0	1.1
1	1	0	0	0	0	0	2.2	0	1.3
1	2	0	1	0	0	0	3.8	0	1.8
1	3	0	1	0	0	0	4.2	1	2.1

Step 1: Fit weight model for censoring due to competing event

$$W^D = \prod_{t=0}^{59} \frac{1}{Pr[D_k | \bar{D}_{k-1} = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k]}$$

Step 2: Use this model to calculate IPCW

Step 3: Fit outcome model adding these additional IPCW (on top of IPTW and IPCW for loss-to-follow-up)

Useful references

Introductory:

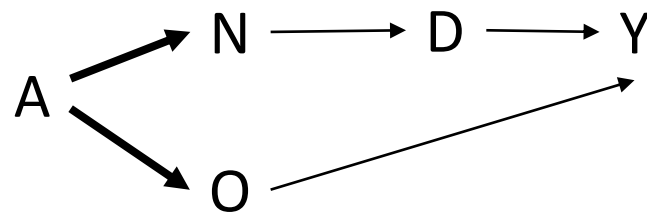
- Rojas-Saunero LP, Young JG, Didelez V, Ikram MA, Swanson SA. Considering Questions Before Methods in Dementia Research With Competing Events and Causal Goals. *Am J Epidemiol*. 2023 Aug 4;192(8):1415-1423. doi: 10.1093/aje/kwad090. PMID: 37139580; PMCID: PMC10403306.
- Mansournia MA, Nazemipour M, Etminan M. A practical guide to handling competing events in etiologic time-to-event studies. *Glob Epidemiol*. 2022 Jul 11;4:100080. doi: 10.1016/j.gloepi.2022.100080. PMID: 37637022; PMCID: PMC10446108.

Technical:

- Young JG, Stensrud MJ, Tchetgen Tchetgen EJ, Hernán MA. A causal framework for classical statistical estimands in failure-time settings with competing events. *Stat Med*. 2020 Apr 15;39(8):1199-1236. doi: 10.1002/sim.8471. Epub 2020 Jan 27. PMID: 31985089; PMCID: PMC7811594.

Separable effects

- Decompose medication into two separable components (N and O): one only affecting competing event death, the other component affecting only the outcome of interest
 - E.g. with the poisonous pill, one component (ACE) directly reduces risk of dialysis, whereas the other component (K+) leads to cardiac arrest and death
- The effect of this new medication is our separable direct effect: $E[Y^{n=0,o=1}] - E[Y^{n=0,o=0}]$
- Using data from a trial of the original medication to try to emulate the trial of a hypothetical yet-to-exist treatment



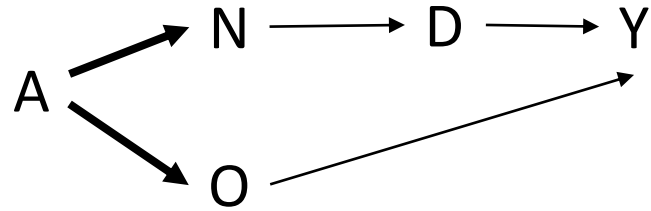
$$Y^{a=1} = Y^{n=1,o=1}$$

$$Y^{a=0} = Y^{n=0,o=0}$$

$Y^{n=0,o=1}$ can be identified by the mediation formula and is equivalent to $Y^{a=1, Ma=0}$

- Define separable direct effects/indirect effects in potential outcomes notation
- We can use information on A, D and Y to identify the separable effects of N and O**
 - Assumptions: (i) no unmeasured common causes of mediator D and outcome Y and (ii) no direct effects of component O on mediator D and of component N on outcome Y
- This is an interventionist way of thinking

G-formula for identification of $Y^{n=0,o=1}$



- In our randomized trial, where we randomize to $A=1$ and $A=0$, we can readily identify $Y^{n=1,o=1}$ (because $Y^{n=1,o=1} = Y^{a=1}$) and we can also readily identify $Y^{n=0,o=0}$ (because $Y^{n=0,o=0} = Y^{a=0}$)
- However, nobody in our population has $Y^{n=0,o=1}$, but we need this quantity since we are interested in the causal effect $Y^{n=0,o=1}_{Y^{a=0}}$
- If data on N and O were available, then we could identify $E[Y^{n=0,o=1}]$ with

$$E[Y^{n=0,o=1}] = \sum_m E[Y|O = 1, M = m] \Pr[M = m, N = 0]$$

- However, we don't have data about N and O . Nevertheless, $O = 1$ iff $A = 1$, and $N = 0$ iff $A = 0$, so we can replace M and N by A ! There is a deterministic relationship between A and N/O

$$E[Y^{n=0,o=1}] = \sum_m E[Y|A = \mathbf{1}, M = m] \Pr[M = m, A = \mathbf{0}]$$

7. Clone-censor-weight implementation

Clone-censor-weight algorithm

ID	Time	L_0	L_k	A_0	A_k	Y_k	IPTW	C_{k_art}
1	0	0	0	NA	0	0		
1	1	0	0	NA	0	0		
1	2	0	1	NA	1	0		
...		
1	59	0	1	NA	1	0		
2	0	1	1	NA	1	0		
2	1	1	1	NA	1	0		
2	2	1	1	NA	1	0		
...		
2	34	1	1	NA	1	1		
3	0	0	0	NA	0	0		
3	1	0	0	NA	0	0		
3	2	0	0	NA	0	0		

Step 0. **Fit the following pooled logistic model** on the dataset before cloning and censoring:

$$\text{logit}[\text{pr}(A_k = 1 | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1} = a, \bar{L}_k)] \\ = \alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k$$

(We already know how to do this)

Clone-censor-weight algorithm

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	IPTW	C _{k_art}
1	0	0	0	NA	0	0		
1	1	0	0	NA	0	0		
1	2	0	1	NA	1	0		
...		
1	59	0	1	NA	1	0		
2	0	1	1	NA	1	0		
2	1	1	1	NA	1	0		
2	2	1	1	NA	1	0		
...		
2	34	1	1	NA	1	1		
3	0	0	0	NA	0	0		
3	1	0	0	NA	0	0		
3	2	0	0	NA	0	0		

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	IPTW	C _{k_art}
1	0	0	0	0	0	0		
1	1	0	0	0	0	0		
1	2	0	1	0	1	0		
...		
1	59	0	1	0	1	0		
2	0	1	1	0	1	0		
2	1	1	1	0	1	0		
2	2	1	1	0	1	0		
...		
2	34	1	1	0	1	1		
3	0	0	0	0	0	0		
3	1	0	0	0	0	0		
3	2	0	0	0	0	0		

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	IPTW	C _{k_art}
1	0	0	0	1	0	0		
1	1	0	0	1	0	0		
1	2	0	1	1	1	0		
...		
1	59	0	1	1	1	0		
2	0	1	1	1	1	0		
2	1	1	1	1	1	0		
2	2	1	1	1	1	0		
...		
2	34	1	1	1	1	1		
3	0	0	0	1	0	0		
3	1	0	0	1	0	0		
3	2	0	0	1	0	0		

Step 1. **Duplicate the dataset**, and assign each individual to each of the strategies he is compatible with (cloning)

Clone-censor-weight algorithm

ID	Time	L_0	L_k	A_0	A_k	Y_k	IPTW	C_{k_art}
1	0	0	0	0	0	0		0
1	1	0	0	0	0	0		0
1	2	0	1	0	1	0		1
...
1	59	0	1	0	1	0		1
2	0	1	1	0	1	0		0
2	1	1	1	0	1	0		0
2	2	1	1	0	1	0		0
...
2	34	1	1	0	1	1		0
3	0	0	0	0	0	0		0
3	1	0	0	0	0	0		0
3	2	0	0	0	0	0		0

Step 2. **Artificially censor** if and when the individual no longer follows his assigned strategy. Next, remove the rows that are artificially censored

(here, illustrated on one of the cloned datasets)

Clone-censor-weight algorithm

ID	Time	L_0	L_k	A_0	A_k	Y_k	IPTW	C_{k_art}
1	0	0	0	0	0	0	1.5	0
1	1	0	0	0	0	0	2.2	0
1	2	0	1	0	1	0	3.8	1
...
1	59	0	1	0	1	0		1
2	0	1	1	0	1	0	1.3	0
2	1	1	1	0	1	0	1.5	0
2	2	1	1	0	1	0	2.6	0
...
2	34	1	1	0	1	1	5.4	0
3	0	0	0	0	0	0	1.2	0
3	1	0	0	0	0	0	2.0	0
3	2	0	0	0	0	0	NA	0

Step 3. **Calculate the IPTW** (we could also call them IPCW) using the model we previously fit on the remaining rows

Clone-censor-weight algorithm

ID	Time	L ₀	L _k	A ₀	A _k	Y _k	IPTW	C _{k_art}
1	0	0	0	0	0	0	1.5	0
1	1	0	0	0	0	0	2.2	0
1	2	0	1	0	1	0	3.8	1
...
1	59	0	1	0	1	0		1
2	0	1	1	0	1	0	1.3	0
2	1	1	1	0	1	0	1.5	0
2	2	1	1	0	1	0	2.6	0
...
2	34	1	1	0	1	1	5.4	0
3	0	0	0	0	0	0	1.2	0
3	1	0	0	0	0	0	2.0	0
3	2	0	0	0	0	0	NA	0

Step 4. Fit the **weighted** outcome model using pooled logistic regression:

$$\begin{aligned} & \text{logit}[\text{pr}(Y_k = 1 | \bar{C}_{k-1} = \bar{0}, \bar{C}_{k-1}(\text{art}) = \bar{0}, \bar{Y}_{k-1} = \bar{0}, A_0)] \\ & = \alpha_{0t} + \alpha_1 A_0 \end{aligned}$$

Alternative implementation

Implementation 1:

Step 0: Fit weight model on dataset before cloning/censoring

Step 1: Clone/duplicate dataset and assign to strategies

Step 2: Artificially censor

Step 3: Calculate weights

Step 4: Fit outcome model

Both approaches are equivalent non-parametrically

Implementation 2:

Step 1: Clone/duplicate dataset and assign to strategies

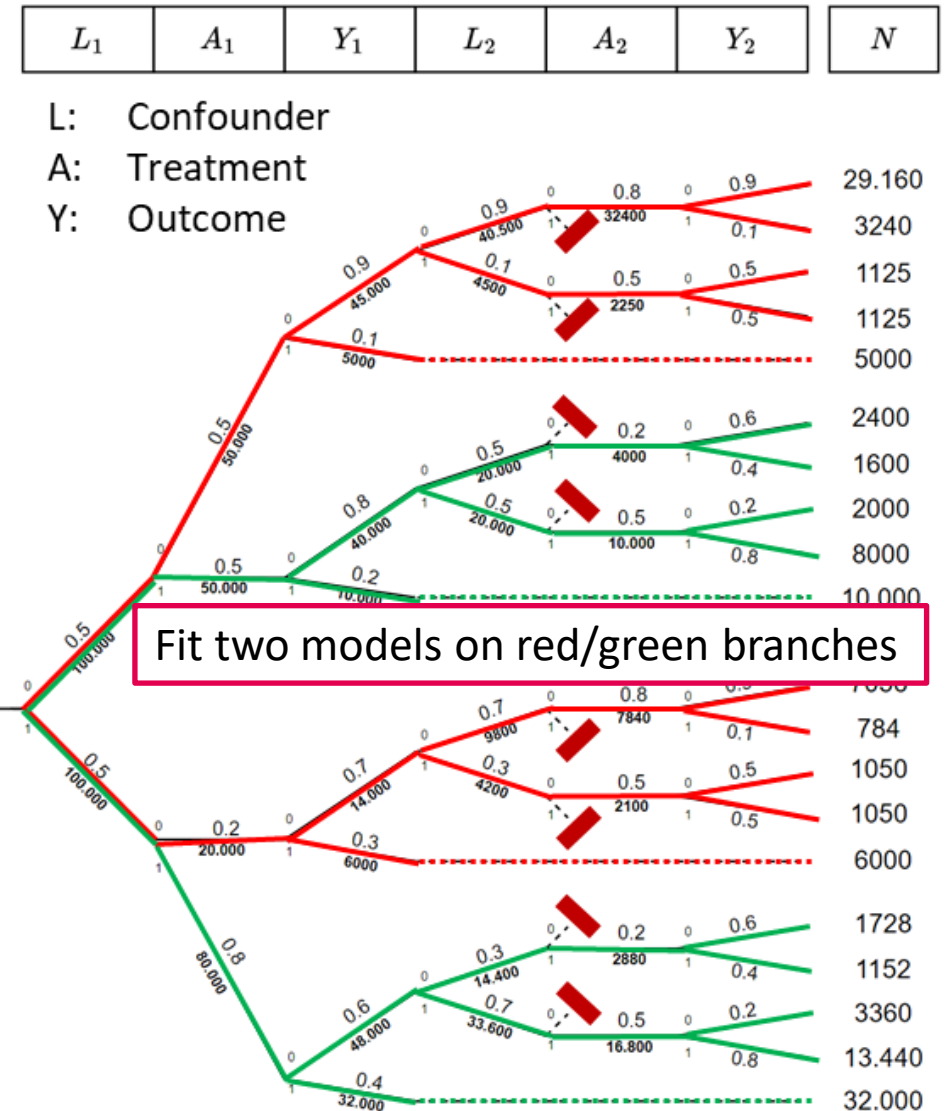
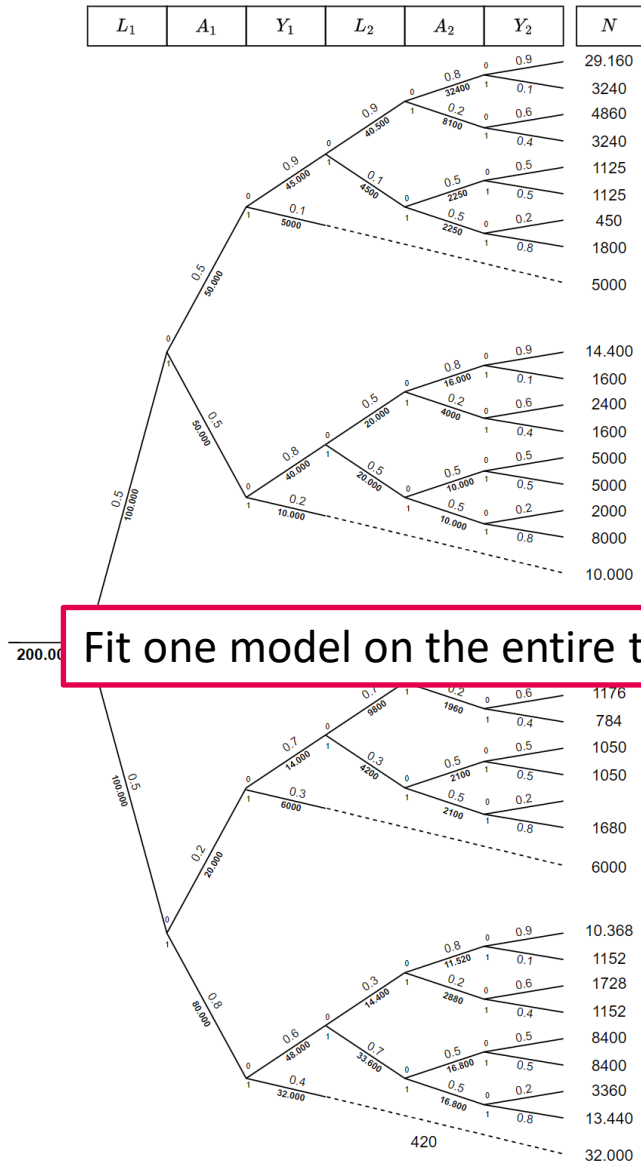
Step 2: Artificially censor

Step 3a: Estimate weight models (separately for each cloned dataset)

Step 3b: Calculate weights

Step 4: Fit outcome model

Difference between implementations



8. Sequential trial implementation

Target trial specification

	Specified target trial
Eligibility criteria	<ul style="list-style-type: none">• 55-84 years• No history of coronary heart disease, stroke, peripheral vascular disease, heart failure, schizophrenia, dementia• 2 years of continuous recording in database• January 2000-November 2006• No previous use of statins
Treatment strategies	<ol style="list-style-type: none">1. Start statins and always use2. Never start statins

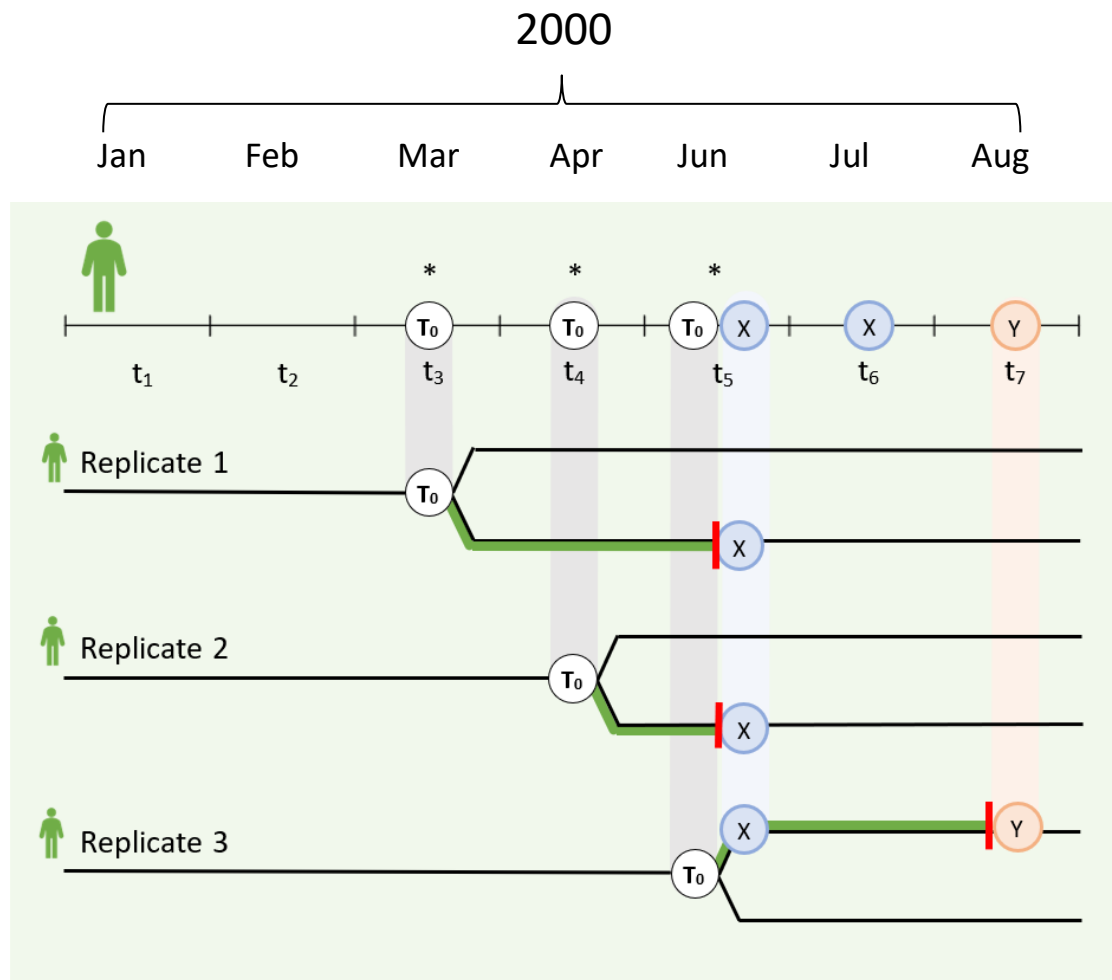
Emulating this trial

	Specified target trial
Eligibility criteria	<ul style="list-style-type: none">• 55-84 years• No history of coronary heart disease, stroke, peripheral vascular disease, heart failure, schizophrenia, dementia• 2 years of continuous recording in database• January 2000-November 2006• No previous use of statins
Treatment strategies	<ol style="list-style-type: none">1. Start statins and always use2. Never start statins

- First trial starts January 2000: check eligibility and do treatment assignment
 - Second trial starts February 2000: check eligibility and do treatment assignment
- Etc. etc. for a total of 83 trials

People can be eligible for multiple trials and hence have multiple time zeros

Sequential trial design



Data format (longitudinal history) – weight models

Table A2. Data for three hypothetical individuals

Individual	Trial	Eligible	Initiator	Current user	Baseline LDL(mmol/L)	LDL(mmol/L)	Month CHD	Month dead
1	12	1	0	0	2.47	2.47	0	14
1	13	1	0	0	2.47	2.47	0	14
2	24	1	0	0	2.77	2.77	26	26
2	25	1	1	1	2.77	2.77	26	26
2	26	0	0	1	2.77	2.84	26	26
3	43	1	1	1	2.88	2.88	0	0
3	44	0	0	1	2.88	2.88	0	0
3	45	0	0	0	2.88	2.77	0	0
3	⋮	⋮	⋮	⋮	2.88	⋮	0	0
3	67	0	0	0	2.88	2.71	0	0
3	68	1	0	0	2.88	2.71	0	0
3	69	1	1	1	2.88	2.73	0	0
3	70	0	0	1	2.88	2.73	0	0

$$W_{m+t}^A = \prod_{k=m}^{m+t} \frac{1}{Pr[A_k | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, L_0, \bar{L}_k]}$$

Fit the following logistic model: $logit[pr(A_k = 1 | \bar{C}_k = \bar{0}, \bar{Y}_{k-1} = \bar{0}, \bar{A}_{k-1}, = a, L_0, \bar{L}_k)] = \alpha_{0t} + \alpha_1^T L_0 + \alpha_2^T L_k$

Expand dataset and create replicates

Table A2. Data for three hypothetical individuals

Individual	Trial	Eligible	Initiator	Current user	Baseline LDL(mmol/L)	LDL(mmol/L)	Month CHD	Month dead
1	12	1	0	0	2.47	2.47	0	14
1	13	1	0	0	2.47	2.47	0	14
2	24	1	0	0	2.77	2.77	26	26
2	25	1	1	1	2.77	2.77	26	26
2	26	0	0	1	2.77	2.84	26	26
3	43	1	1	1	2.88	2.88	0	0
3	44	0	0	1	2.88	2.88	0	0
3	45	0	0	0	2.88	2.77	0	0
3	⋮	⋮	⋮	⋮	2.88	⋮	0	0
3	67	0	0	0	2.88	2.71	0	0
3	68	1	0	0	2.88	2.71	0	0
3	69	1	1	1	2.88	2.73	0	0
3	70	0	0	1	2.88	2.73	0	0



Table A3. The expanded dataset for the three hypothetical individuals in Table A2

Individual	Trial (m)	Follow-up month (t)	Eligible (E_m)	Initiator (A_m)	Current user (A_{m+t})	Baseline LDL (L_m)	Time-varying LDL (L_{m+t})	Event (D_{m+t})
1	12	0	1	0	0	2.47	2.47	0
1	12	1	1	0	0	2.47	2.47	0
1	13	0	1	0	0	2.47	2.47	0
2	24	0	1	0	0	2.77	2.77	0
2	24	1	1	0	1	2.77	2.77	0
2	24	2	1	0	1	2.77	2.84	1
2	25	0	1	1	1	2.77	2.77	0
2	25	1	1	1	1	2.77	2.84	1
3	43	0	1	1	1	2.88	2.88	0
3	43	1	1	1	1	2.88	2.88	0
3	43	2	1	1	0	2.88	2.77	0
3	43	⋮	⋮	⋮	⋮	⋮	⋮	⋮
3	43	26	1	1	1	2.88	2.73	0
3	43	27	1	1	1	2.88	2.73	.
3	68	0	1	0	0	2.71	2.71	0
3	68	1	1	0	1	2.71	2.73	0
3	68	2	1	0	1	2.71	2.73	.
3	69	0	1	1	1	2.73	2.73	0
3	69	1	1	1	1	2.73	2.73	.

Artificially censor & Calculate weights on expanded dataset

Table A3. The expanded dataset for the three hypothetical individuals in Table A2

Individual	Trial (m)	Follow-up month (t)	Eligible (E_m)	Initiator (A_m)	Current user (A_{m+t})	Baseline LDL (L_m)	Time-varying LDL (L_{m+t})	Event (D_{m+t})	C_k_art	IPTW
1	12	0	1	0	0	2.47	2.47	0		
1	12	1	1	0	0	2.47	2.47	0		
1	13	0	1	0	0	2.47	2.47	0		
2	24	0	1	0	0	2.77	2.77	0		
2	24	1	1	0	1	2.77	2.77	0		
2	24	2	1	0	1	2.77	2.84	1		
2	25	0	1	1	1	2.77	2.77	0		
2	25	1	1	1	1	2.77	2.84	1		
3	43	0	1	1	1	2.88	2.88	0		
3	43	1	1	1	1	2.88	2.88	0		
3	43	2	1	1	0	2.88	2.77	0		
3	43	⋮	⋮	⋮	⋮	⋮	⋮	⋮		
3	43	26	1	1	1	2.88	2.73	0		
3	43	27	1	1	1	2.88	2.73	.		
3	68	0	1	0	0	2.71	2.71	0		
3	68	1	1	0	1	2.71	2.73	0		
3	68	2	1	0	1	2.71	2.73	.		
3	69	0	1	1	1	2.73	2.73	0		
3	69	1	1	1	1	2.73	2.73	.		

Fit weighted outcome model

$$\text{logit}[\text{pr}(Y_{m+t+1} = 1 | \bar{C}_{m+t+1}(\text{art}) = \bar{0}, \bar{Y}_{m+t} = \bar{0}, A_m)] = \alpha_{0,m+t} + \alpha_1 A_m$$

$$\text{where } \alpha_{0,m+t} = \alpha_0 + \alpha_2 * m + \alpha_3 * m^2 + \alpha_4 * t + \alpha_5 * t^2$$

Table A3. The expanded dataset for the three hypothetical individuals in Table A2

Individual	Trial (m)	Follow-up month (t)	Eligible (E _m)	Initiator (A _m)	Current user (A _{m+t})	Baseline LDL (L _m)	Time-varying LDL (L _{m+t})	Event (D _{m+t})
1	12	0	1	0	0	2.47	2.47	0
1	12	1	1	0	0	2.47	2.47	0
1	13	0	1	0	0	2.47	2.47	0
2	24	0	1	0	0	2.77	2.77	0
2	24	1	1	0	1	2.77	2.77	0
2	24	2	1	0	1	2.77	2.84	1
2	25	0	1	1	1	2.77	2.77	0
2	25	1	1	1	1	2.77	2.84	1
3	43	0	1	1	1	2.88	2.88	0
3	43	1	1	1	1	2.88	2.88	0
3	43	2	1	1	0	2.88	2.77	0
3	43	⋮	⋮	⋮	⋮	⋮	⋮	⋮
3	43	26	1	1	1	2.88	2.73	0
3	43	27	1	1	1	2.88	2.73	.
3	68	0	1	0	0	2.71	2.71	0
3	68	1	1	0	1	2.71	2.73	0
3	68	2	1	0	1	2.71	2.73	.
3	69	0	1	1	1	2.73	2.73	0
3	69	1	1	1	1	2.73	2.73	.