

Disclaimer: This talk is not about TMLE.

Disclaimer: This talk is *not* about TMLE.

Aim: Provide you with the idea of sequential regression.

Disclaimer: This talk is not about TMLE.

Aim: Provide you with the idea of sequential regression.

→ Maximum likelihood based G-computation estimation

Notation

Consider the random vector, $\mathbf{X} = (L, A, Y)$. Let P_X denote the joint probability measure. We assume that $P_X \in \mathcal{P}$ where \mathcal{P} is an otherwise unspecified model which is dominated by a σ -finite measure, μ . With the notation,

$$Q_{Y}(dy \mid a, I) = P_{Y|A,L}(dy \mid A = a, L = I)$$

$$G_{A}(da \mid I) = P_{A|L}(da \mid L = I),$$

$$Q_{L}(dI) = P_{L}(dI),$$

the joint distribution P_X factorizes

$$P_X(dx) = Q_Y(dy \mid a, I) G_A(da \mid I) Q_L(dI).$$

Notation

For any measurable function f of the data we use the notation for the expected value of f under P_X :

$$\mathsf{P}_X f = \int f(x) \, \mathsf{P}_X(\mathsf{d} x).$$

Using the factorization,

$$\mathsf{P}_X(\mathsf{d} x) = Q_Y(\mathsf{d} y \mid \mathsf{a}, \mathsf{I}) \, G_A(\mathsf{d} \mathsf{a} \mid \mathsf{I}) \, Q_L(\mathsf{d} \mathsf{I}),$$

we extend this compact operator notation to conditional distributions and write

$$\mathsf{P}_X f = Q_L G_A Q_Y f = \iiint f(x) Q_Y (\mathsf{d} y \mid a, l) G_A (\mathsf{d} a \mid l) Q_L (\mathsf{d} l).$$

UNIVERSITY OF COPENHAGEN

Longitudinal data structures

For K time points, $\{0 = t_0 < t_1 < \cdots < t_K\}$, we consider longitudinal data given by

$$\mathbf{X} = (L(0), A(0), Y(1), L(1), A(1), \dots, Y(K)),$$

s.t. Y(k) is the value of a stochastic outcome process Y(t) at time t_k .

The state space of the process Y is $\{0,1\}$, i.e.,

Y(t) = 0: event has not yet occurred at time tY(t) = 1: event has occurred in the time interval [0, t]

Ex. 1: One time point, K = 1, no censoring

Consider the data, $\mathbf{X} = (L(0), A(0), Y(1)) \sim P_X \in \mathcal{P}$.

Our target parameter is the intervention-specific ATE defined by

$$\psi(\tilde{\mathsf{P}}) = \tilde{\mathsf{P}}(\mathsf{Y}^{(1)}(1) = 1) - \tilde{\mathsf{P}}(\mathsf{Y}^{(0)}(1) = 1),$$

where $(L(0), Y^{(1)}(1), Y^{(0)}(1)) \sim \tilde{P} \in \tilde{P}$ is the counterfactual data under interventions, where the treatment is set to either one or zero.

Want: To identify the target parameter through observed data, i.e.,

$$\psi(\tilde{\mathsf{P}}) = \theta(\mathsf{P}_X)$$
 for some functional, θ , of P_X

Ex. 1: One time point, K=1, no censoring

Want: To identify the target parameter through observed data.

$$\begin{split} \tilde{\mathsf{P}}(Y^{(1)}(1) = 1) &= \int \tilde{\mathsf{P}}(Y^{(1)}(1) = 1 \mid A(0) = 1, L(0) = I) Q_{L(0)}(\mathsf{d}I) \\ &\stackrel{!}{=} \int \mathsf{P}(Y(1) = 1 \mid A(0) = 1, L(0) = I) Q_{L(0)}(\mathsf{d}I) \\ &= \int \mathsf{P}(Y(1) = 1 \mid A(0) = a, L(0) = I) G_{A(0)}^*(\mathsf{d}a) Q_{L(0)}(\mathsf{d}I) \\ &= \int Q_{Y(1)}(1 \mid a, I) G_{A(0)}^*(\mathsf{d}a) Q_{L(0)}(\mathsf{d}I) \\ &= Q_{L(0)} G_{A(0)}^* Q_{Y(1)} f, \end{split} \tag{G-formula}$$

for $f(\mathbf{X}) = \mathbb{1}\{Y(1) = 1\}$. This leads to the estimator,

$$\hat{Q}_{L(0)}\hat{G}_{A(0)}^*\hat{Q}_{Y(1)}f = \frac{1}{n}\sum_{i=1}^n \hat{Q}_{Y(1)}(1\mid 1, L_i(0)).$$

Ex. 2: Multiple time points, K=2, with censoring

Consider data, $\mathbf{X} = (L(0), A(0), C(1), Y(1), L(1), A(1), C(2), Y(2)),$ where C censoring variable taking values $\{0,1\}$ with C(t)=0 to mean that the event has been censored. Our target parameter is given by

$$\psi(\tilde{\mathsf{P}}) = \tilde{\mathsf{P}}(\mathsf{Y}^{(1)}(2) = 1) - \tilde{\mathsf{P}}(\mathsf{Y}^{(0)}(2) = 1).$$

Want: Identifiability, i.e., $\psi(\hat{P}) = \theta(P_X)$. Need some more notation!

Ex. 2: Multiple time points, K=2, with censoring

Consider data, $\mathbf{X} = (L(0), A(0), C(1), Y(1), L(1), A(1), C(2), Y(2)),$ where C censoring variable taking values $\{0,1\}$ with C(t)=0 to mean that the event has been censored. Our target parameter is given by

$$\psi(\tilde{\mathsf{P}}) = \tilde{\mathsf{P}}(\mathsf{Y}^{(1)}(2) = 1) - \tilde{\mathsf{P}}(\mathsf{Y}^{(0)}(2) = 1).$$

Want: Identifiability, i.e., $\psi(\tilde{P}) = \theta(P_X)$. Need some more notation!

History:
$$\bar{A}(1) = (A(0), A(1))$$

Factorization: $dP_{\mathbf{X}} = Q_{Y(2)} G_{C(2)} G_{A(1)} Q_{L(1)} Q_{Y(1)} G_{C(1)} G_{A(0)} Q_{L(0)}$

Then for some function of the data, $f(\mathbf{X})$, we can write

$$P_{\mathbf{X}}f = \int f(x)dP_{\mathbf{X}}(x) = Q_{L(0)}G_{A(0)}G_{C(1)}Q_{Y(1)}Q_{L(1)}G_{A(1)}G_{C(2)}Q_{Y(2)}f.$$

Ex. 2: Multiple time points, K=2, with censoring

Consider the counterfactual distribution defined by

$$\begin{split} &G_{\bar{A}(1)}^*(\mathsf{d}(\mathsf{a}_0,\mathsf{a}_1)) = \!\! 1 \{ \mathsf{a}_0 = 1, \mathsf{a}_1 = 1 \} \\ &G_{\bar{C}(2)}^*(\mathsf{d}(\mathsf{c}_1,\mathsf{c}_2)) = \!\! 1 \{ \mathsf{c}_1 = 1, \mathsf{c}_2 = 1 \}. \end{split}$$

The G-formula in this case is obtained by

$$\begin{split} \tilde{\mathsf{P}}(Y^{(1)}(2) = 1) & \stackrel{!}{=} \int \mathsf{P}(Y(2) = 1 \mid \bar{C}(2) = c, \bar{A}(1) = a, \bar{L}(1) = I, Y(1) = y) \\ & \times G^*_{\bar{C}(2)} G^*_{\bar{A}(1)} Q_{L(1)} Q_{Y(1)} G^*_{C(1)} G^*_{A(0)} Q_{L(0)} \\ & = Q_{L(0)} G^*_{A(0)} G^*_{C(1)} Q_{Y(1)} Q_{L(1)} G^*_{\bar{A}(1)} G^*_{\bar{C}(2)} Q_{Y(2)} f, \end{split} \tag{G-formula}$$

for
$$f(\mathbf{X}) = \mathbb{1}\{Y(2) = 1\}.$$

G-formula:
$$Q_{L(0)}G_{A(0)}^*G_{C(1)}^*Q_{Y(1)}Q_{L(1)}G_{\bar{A}(1)}^*G_{\bar{C}(2)}^*Q_{Y(2)}f$$

Step 1. Regress Y(2) on past covariates, $\bar{A}(1)$ and $\bar{L}(1)$, for subjects at risk

$$\,\leadsto\,\,G_{\bar{C}(2)}^*Q_{Y(2)}f=\mathsf{P}(Y(2)=1\mid\bar{C}(2)=1,Y(1)=y,\bar{A}(1),\bar{L}(1))$$

Deterministic info: $Y(1) = 1 \Rightarrow Y(2) = 1$ and C(1) = 0 or $C(2) = 0 \Rightarrow Y(2) = 0$

head(d)

G-formula:
$$Q_{L(0)}G_{A(0)}^*G_{C(1)}^*Q_{Y(1)}Q_{L(1)}G_{\bar{A}(1)}^*G_{\bar{C}(2)}^*Q_{Y(2)}f$$

Step 1. Regress Y(2) on past covariates, $\bar{A}(1)$ and $\bar{L}(1)$, for subjects at risk

$$\longrightarrow G_{\bar{C}(2)}^* Q_{Y(2)} f = \mathsf{P}(Y(2) = 1 \mid \bar{C}(2) = 1, Y(1) = y, \bar{A}(1), \bar{L}(1))
Y(1) + (1 - Y(1)) \mathsf{P}(Y(2) = 1 \mid \bar{C}(2) = 1, Y(1) = 0, \bar{A}(1), \bar{L}(1))$$

Deterministic info: $Y(1) = 1 \Rightarrow Y(2) = 1$ and C(1) = 0 or $C(2) = 0 \Rightarrow Y(2) = 0$

head(d)

G-formula:
$$Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f$$

Step 2. Predict according to counterfactual distribution

$$\,\leadsto\,\, G^*_{\bar{A}(1)}\,G^*_{\bar{C}(2)}\,Q_{Y(2)}f=\mathsf{P}(\,Y(2)=1\mid\bar{C}(2)=1,\,Y(1)=y,\bar{A}(1)=1,\bar{L}(1))$$

$$d1 \leftarrow copy(d); d1[,A_0 := 1]; d1[,A_1 := 1]$$

head(d1)

$$GQhat2_0 \leftarrow predict(fit1, newdata = d1, type = "response") d[, GQhat2 := (Y_1 + (1 - Y_1)*GQhat2_0)] # 1 if Y_1 = 1$$

$$\text{G-formula:} \quad Q_{L(0)} \, G_{A(0)}^* \, G_{C(1)}^* \, Q_{Y(1)} \, Q_{L(1)} \, G_{\overline{A}(1)}^* \, G_{\overline{C}(2)}^* \, Q_{Y(2)} f$$

Step 2. Predict according to counterfactual distribution

$$\Rightarrow G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f = \mathsf{P}(Y(2) = 1 \mid \bar{C}(2) = 1, Y(1), \bar{A}(1) = 1, \bar{L}(1)) = \bar{\mathbf{Q}}_{L(2)}^{d,2}$$

 $d1 \leftarrow copy(d); d1[,A_0 := 1]; d1[,A_1 := 1]$ head(d1)

$$\label{eq:GQhat2_0} \begin{split} \mathsf{GQhat2_0} &<- \ \mathsf{predict} \big(\, \mathsf{fit1} \, , \, \, \mathsf{newdata} \, = \, \mathsf{d1} \, , \, \, \mathsf{type} \, = \, \texttt{"response"} \big) \\ \mathsf{d} \big[\, , \mathsf{GQhat2} \, := \, \big(\mathsf{Y_1} \, + \, \big(1 \, - \, \mathsf{Y_1} \big) * \mathsf{GQhat2_0} \big) \big] \, \, \# \, 1 \, \, \, \mathsf{if} \, \, \, \mathsf{Y_1} \, = \, 1 \end{split}$$

G-formula:
$$Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f$$

Step 3. Regress $G^*_{\overline{A}(1)}G^*_{\overline{C}(2)}Q_{Y(2)}f$ on A(0) and L(0) for subjects at risk

$$\Rightarrow G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f$$

$$= \mathbb{E} \left[G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f \mid C(1) = 1, A(0), L(0) \right]$$

Note: $G^*_{\bar{A}(1)}G^*_{\bar{C}(2)}Q_{Y(2)}f\in(0,1)$ not binary. We use quasi-binomial logistic regression with an extra dispersion parameter to describe additional variation in data – in R with glm, coefficient estimates are the same, but the SE's differ.

12

Sequential outcome regression

G-formula:
$$Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f$$

Step 4. Predict according to counterfactual distribution

$$\longrightarrow G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f
= \mathbb{E} \left[G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f \mid C(1) = 1, A(0) = 1, L(0) \right]$$

 $\mathsf{GQhat1} \leftarrow \mathsf{predict}(\mathsf{fit1}, \mathsf{newdata} = \mathsf{d1}, \mathsf{type} = \mathsf{"response"}))$

Step 5. Take sample average over L(0)

$$ightarrow Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f$$

mean (GQhat1)
[1] 0.03377855

12

UNIVERSITY OF COPENHAGEN

Sequential outcome regression

G-formula:
$$Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f$$

Step 4. Predict according to counterfactual distribution

$$\Rightarrow G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f
= \mathbb{E} \Big[\underbrace{G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f}_{\bar{Q}_{L(2)}^{d,2}} \mid C(1) = 1, A(0) = 1, L(0) \Big] = \bar{Q}_{L(1)}^{d,2}$$

 $\mathsf{GQhat1} \leftarrow \mathsf{predict}(\mathsf{fit1}, \mathsf{newdata} = \mathsf{d1}, \mathsf{type} = \mathsf{"response"}))$

Step 5. Take sample average over L(0)

$$ightharpoonup Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\overline{A}(1)}^* G_{\overline{C}(2)}^* Q_{Y(2)} f$$

mean (GQhat1)
[1] 0.03377855

12

Sequential outcome regression

G-formula:
$$Q_{L(0)}G_{A(0)}^*G_{C(1)}^*Q_{Y(1)}Q_{L(1)}G_{ar{A}(1)}^*G_{ar{C}(2)}^*Q_{Y(2)}f$$

Step 4. Predict according to counterfactual distribution

$$\Rightarrow G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f
= \mathbb{E} \left[\underbrace{G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f}_{\bar{Q}_{L(2)}^{d,2}} \mid C(1) = 1, A(0) = 1, L(0) \right] = \bar{Q}_{L(1)}^{d,2}$$

GQhat1 <- predict(fit1, newdata = d1, type = "response"))

Step 5. Take sample average over L(0)

$$\rightarrow Q_{L(0)} G_{A(0)}^* G_{C(1)}^* Q_{Y(1)} Q_{L(1)} G_{\bar{A}(1)}^* G_{\bar{C}(2)}^* Q_{Y(2)} f = \bar{Q}_{L(0)}^{d,2}$$

mean (GQhat1) [1] 0.03377855

Sequential outcome regression with the Ltmle package

Maximum likelihood based G-computation estimate with Ltmle

Step 0. Prepare data.

head(data)

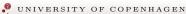
gform

Qform

Sequential outcome regression with the Ltmle package

Maximum likelihood based G-computation estimate with Ltmle

```
Use setting: gcomp = TRUE (default is FALSE)
fit_ltmle <- Ltmle(data = data,
                    Anodes = c("A_0", "A 1").
                    Cnodes = c("C 1","C 2").
                    Lnodes = c("L_0", "L_1"),
                    Ynodes = c("Y_0", "Y_1"),
                    survivalOutcome = TRUE,
                    Qform = Qform,
                    gform = gform,
                    abar = list (c(1,1),c(0,0)),
                    gcomp = TRUE,
                    SL.library = "glm")
```



Extension to K > 2 time points is straightforward

Extension to K > 2 time points is straightforward

Step 1+2. Regress Y(K) on past cov. (until time K-1) for subjects at risk and predict according to intervention rule $\leadsto G^*_{\bar{A}(K-1)} G^*_{\bar{C}(K)} Q_{Y(K)} f$

Step 3 + 4. Regress $G^*_{\bar{A}(K-1)}G^*_{\bar{C}(K)}Q_{Y(K)}f$ on past cov. (until time K-2) for subjects at risk and predict according to intervention rule

Step 2K + 1. Take sample average over L(0)

Extension to K>2 time points is straightforward

Extension to competing risk:

- * Suppose D(1) competing risk, e.g., a value of one means death
- * At risk: $d[Y_1 == 0 \& C_1 == 1 \& C_2 == 1 \& D_1 == 0]$
- * Deterministic info about Y(2), e.g., $D(1) = 1 \Rightarrow Y(2) = 0$
- * In Ltmle: deterministic.Q.function

Extension to K > 2 time points is straightforward

Extension to competing risk:

- st Suppose D(1) competing risk, e.g., a value of one means death
- * At risk: $d[Y_1 == 0 \& C_1 == 1 \& C_2 == 1 \& D_1 == 0]$
- * Deterministic info about Y(2), e.g., $D(1) = 1 \Rightarrow Y(2) = 0$
- * In Ltmle: deterministic.Q.function

References

- [1] Heejung Bang and James M Robins. Doubly robust estimation in missing data and causal inference models. Biometrics, 61(4):962–973, 2005. (Hard to read!)
- [2] Samuel D Lendle, Joshua Schwab, Maya L Petersen, and Mark J van der Laan. Itmle: an r package implementing targeted minimum loss-based estimation for longitudinal data. Journal of Statistical Software, 81:1–21, 2017. (Ltmle doc.)