# STA 640 — Causal Inference

## Chapter 9 – Sequential/Longitudinal Treatments

Fan Li

Department of Statistical Science
Duke University

# Longitudinal observational studies

▶ All previous discussions focus on treatment at a single time (cross-section or panel settings)

▶ Common in real world situations, e.g. medical research, data (treatment, covariates and/or outcome) are repeatedly collected on subjects over a period – longitudinal studies

▶ Particularly interested in estimating the effect of a time-varying treatment on an outcome of interest measured at a later time

▶ Confounders can be time-varying, affected by past treatment and affecting future covariates and or outcomes

▶ Standard regression adjustment fails to give consistent estimators in the presence of time-varying confounders if those confounders are themselves affected by treatment

# Sequentially ignorable assignment

What type of studies we consider

- ▶ Treatments with multiple time points, where those treatments assignment is ignorable conditionally on the observed history.

- ▶ If we can justify the above assumption, this is a possible template for randomized experiments or observational studies

- ▶ Example 1: patients visiting doctors at different times.

- ▶ Example 2: workers exposed to hazards at the workplace (related to health worker survivor effect)

A common goal is to estimate the accumulative (over the study period) effect of the treatment on an outcome.
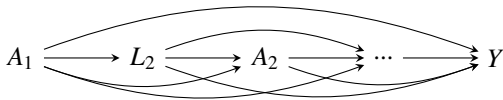
# Time-varying treatments: An example

- In many longitudinal medical studies, patients' treatment changes over time and is measured several times during the study, along with other time changing covariates.

- For example, type II diabetes patients recruited into a study comparing two antiglycaemic drugs may be followed up on several occasions, on each of which their HbA1c (a long-term measure of blood glucose level), blood pressure, cholesterol level, BMI, anaemia status, and others variables are measured

- Suppose wish to compare the effect of the two treatments on HbA1c 18 months after recruitment and on the risk (or hazard) of experiencing a cardiac event in the 18 months following recruitment.

# Time-varying treatments: The challenge

$$A_1 \longrightarrow L_2 \rightrightarrows A_2 \rightrightarrows \cdots \longrightarrow Y$$

- Study allows for the dose and type of treatment to be changed according to the current (and past) values of HbA1c and other covariates.

- A high HbA1c likely lead to increasing the dose of the current drug; but high HbA1c is also thought to lead to an increased risk of a cardiac event, making HbA1c at a particular time a *confounder of the relationship between subsequent treatment and the outcome*

- Because HbA1c varies over time (in a way that cannot be foreseen at baseline), it is called a *time-varying confounder*

# Time-varying treatments: The challenge

(Daniel et al. 2013, SIM)

- To estimate the causal effect of treatment on risk of cardiac event, it seems necessary to control for HbA1c in the analysis

- Not only does HbA1c affect treatment but also **the reverse is true**!!!

  - An effective antiglycaemic drug lowers HbA1c, and thus current value of the treatment variable has a causal effect on future values of HbA1c

- This means controlling for HbA1c is problematic, because future measurements of HbA1c lie on the causal pathway between past treatment and the outcome

  - conditioning on HbA1c blocks some of the effect of the treatment and, in addition, conditioning on a consequence of treatment risks inducing collider-stratification bias

# Notation with a toy example

► We will switch notation from previous lectures to be compatible with the literature in longitudinal treatment

► Toy example: patients with cancer, visiting doctor at two time points $t_1, t_2$.

► $a_t (t = 1, 2)$ : possible treatment at time $t$

► $A_{i,t}$: the observed treatment at time $t$ ($Z_{i,t}$ in previous notation system)

► $L_i^{obs}$: observed cancer progression at time 2

► $Y_i(a_1, a_2)$: potential outcome at time 3

► $Y_i = Y_i(A_{i,1}, A_{i,2})$: observed outcomes

# Toy example with $T = 2$

| month | action | potential outcome | observed value |
|-------|--------|-------------------|----------------|
| 1 | give treatment $a_1$ (1=high) | | $A_{i,1}$ |
| 2 | (i) measure cancer progression | | $L_i^{obs}$ |
| | (ii) give treatment $a_2$ (1=high) | | $A_{i,2}$ |
| 3 | measure cancer progression | $Y_i(a_1, a_2)$ | $Y_i$ |

▶ For each individual, there are a total of 4 potential outcomes $\{Y(1, 1), Y(1, 0), Y(0, 1), Y(0, 0)\}$, but only one will be observed
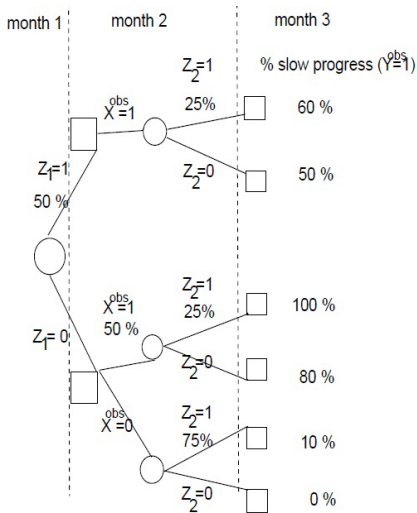
Figure: Example with treatment at two time points and sequentially ignorable assignment (Here $Z$ is the assignment, and $X^{obs}$ is the intermediate variable)

# Causal Estimand

- Typical target estimand - marginal causal effects due to treatment sequence:

$$\tau_{a_1 a_2, a_1' a_2'} = \mathbb{E}[Y_i(a_1, a_2) - Y_i(a_1', a_2')],$$

for all $(a_1, a_2) \neq (a_1', a_2') \in \{0, 1\}^2$.

- For example, compare cancer progression $Y$ between always taking high dose $\Pr(Y(1, 1) = 1)$ and always taking low dose $\Pr(Y(0, 0) = 1)$.

- Note that we only control $(a_1, a_2)$, that is why the potential outcomes $Y_i()$ are only a function of $a_1, a_2$ and not also of $L$.

# Problems with standard adjustment

Two "standard" approaches to estimate
$\Pr(Y(1,1) = 1) - \Pr(Y(0,0) = 1)$:

▶ Approach 1. "Do not condition on progression $L^{obs}$ because it is an intermediate outcome":

$$\Pr(Y = 1 | A_1 = 1, A_2 = 1) - \Pr(Y = 1 | A_1 = 0, A_2 = 0)$$
$$= \quad 60\% - 60\% = 0$$

▶ Approach 2. "Condition on intermediate progression $L^{obs}$ because it was used in deciding treatment $A_2$":

$$\Pr(Y = 1 | A_1 = 1, A_2 = 1, L^{obs} = 1) - \Pr(Y = 1 | A_1 = 0, A_2 = 0, L^{obs} = 1)$$
$$= \quad 60\% - 80\% = -20\%$$

# Problems with standard adjustment

Two "standard" approaches to estimate
$\Pr(Y(1, 1) = 1) - \Pr(Y(0, 0) = 1)$:

▶ Approach 1. "Do not condition on progression $L^{obs}$ because it is an intermediate outcome":
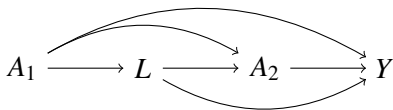
$$\Pr(Y = 1|A_1 = 1, A_2 = 1) - \Pr(Y = 1|A_1 = 0, A_2 = 0)$$
$$= \quad 60\% - 60\% = 0$$

▶ Approach 2. "Condition on intermediate progression $L^{obs}$ because it was used in deciding treatment $Z_2$":

$$\Pr(Y = 1|A_1 = 1, A_2 = 1, L^{obs} = 1) - \Pr(Y = 1|A_1 = 0, A_2 = 0, L^{obs} = 1)$$
$$= \quad 60\% - 80\% = -20\%$$

Both approaches are incorrect for the goal - adjusting for $L^{obs}$ alone is not enough

# Assumption 1: Positivity/Overlap



$$\tau_{a_1 a_2, a_1' a_2'} = \mathbb{E}[Y_i(a_1, a_2) - Y_i(a_1', a_2')]$$

► At every time point, units have positive probability to receive all levels of the treatment
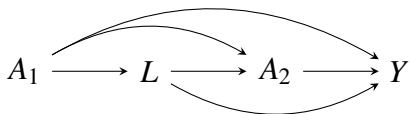
$$P(A_1 = a_1) > 0$$
$$P(A_2 = a_2 | A_1 = a_1, L = l) > 0$$

for all $a_1, a_2, l$.

# Assumption 1: Positivity/Overlap

▶ Positivity states that conditional on covariate history, the probability of receiving each treatment sequence is bounded away from zero and one

▶ The above illustration is only for two time periods ($T = 2$), and therefore the covariate history includes only $L$ right after treatment $A_1$

▶ Positivity is less likely to be satisfied

  ▶ for large values of $T$

  ▶ when $A$ includes more than two values (multiple treatments)

# Assumption 2: Sequential Ignorability (Robins, 1986)



$$A_1 \longrightarrow L \longrightarrow A_2 \longrightarrow Y$$

Let $L_t^{obs}$ (often shorthanded to $L_t$) denote the time-varying confounders, including both time-varying covariates and intermediate outcome at time $t$.

Let $\bar{a}_t = (a_1, a_2, ..., a_t)$, $\bar{A}_t = (A_1, \ldots, A_t)$

▶ Sequential ignorability: treatment at time $t$ is randomized with probabilities depending on the observed past, *including covariates, intermediate outcomes*, that is, at any time $t$ :

$$\{Y_i(\bar{a}_t), \forall \bar{a}_t\} \perp A_{i,t} \mid H_{i,t},$$

where $H_{i,t} = (A_1, ..., A_{t-1}; L_1, ...L_{t-1})$ is the observed history

▶ In the previous simple example, $H_1$=nothing, and $H_2 = (A_{i,1}, L_i^{obs})$

# Identifiability

- $E[Y(a_1, a_2)]$ is a function of potential outcomes
- Idenitifiability of $E[Y(a_1, a_2)]$ means that it can be written in terms of observed data

$$E[Y(a_1, a_2)]$$
$$= E[Y(a_1, a_2)|A_1 = a_1] \tag{1}$$
$$= E[Y(A_1, a_2)|A_1 = a_1] \tag{2}$$
$$= \sum_{l=0,1} E[Y(A_1, a_2)|A_1 = a_1, L^{obs} = l]P(L^{obs} = l|A_1 = a_1) \tag{3}$$
$$= \sum_{l=0,1} E[Y(A_1, a_2)|A_1 = a_1, L^{obs} = l, A_2 = a_2]P(L^{obs} = l|A_1 = a_1) \tag{4}$$
$$= \sum_{l=0,1} E[Y(A_1, A_2)|A_1 = a_1, L^{obs} = l, A_2 = a_2]P(L^{obs} = l|A_1 = a_1) \tag{5}$$

- (1), (4) from sequential ignorability
- (2), (5) from consistency (SUTVA); (3) from law of total probability

# g-computation

(Robins, 1986)

▶ Causal effects are identified under the assumption of sequential
ignorability leading to the g-computation

$$
\begin{aligned}
\Pr(Y(0,0)=1) &= \Pr(Y(0,0)=1|A_1=0) \\
&= \sum_{L^{obs}=0,1} \Pr(Y(0,0)=1|A_1=0, L^{obs}) \Pr(L^{obs}|A_1=0) \\
&= \sum_{L^{obs}=0,1} \Pr(Y(0,0)=1|A_1=0, L^{obs}, A_2=0) \Pr(L^{obs}|A_1=0) \\
&= \sum_{L^{obs}=0,1} \Pr(Y^{obs}=1|A_1=0, L^{obs}, A_2=0) \Pr(L^{obs}|A_1=0) \\
&= 0\%(50\%) + 80\%(50\%) = 40\%
\end{aligned}
$$

# g-computation

▶ Similarly, we can estimate $\Pr(Y_i(1,1) = 1) = 60\%$
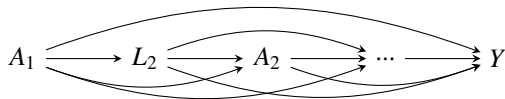
▶ Therefore, the causal effect is given by

$$\Pr(Y_i(1,1) = 1) - \Pr(Y_i(0,0) = 1) = 20\%$$

▶ Use analogous arguments, we can also estimate
$\Pr(Y_i(0,1) = 1) = 55\%$ and $\Pr(Y_i(1,0) = 1) = 50\%$.

▶ This procedure can be generated to longitudinal treatments with any $T$
time points: g-computation

# g-computation

▶ For $T$ time points, let $\bar{a}_t = (a_1, a_2, ..., a_t)$, $\bar{A}_t = (A_1, \ldots, A_t)$



$$\Pr(Y(\bar{a}_T)) =$$
$$\sum_{L_2^{obs}, ..., L_T^{obs}} \Pr(Y^{obs} \mid L_2^{obs}, ..., L_T^{obs}, \overline{A}_T = \bar{a}_T)$$
$$\times \Pr(L_2^{obs} | A_1 = a_1) \times \cdots \times$$
$$\times \Pr(L_T^{obs} | A_1 = a_1, L_2^{obs}, A_2 = a_2, \ldots, L_{T-1}^{obs}, A_{T-1} = a_{T-1}).$$

▶ This is the basic g-computation formula, or g-formula
▶ Can pose models for all distributions in the RHS and estimate $\Pr(\bar{a})$ –
   in essence, this is an outcome modeling approach

# g-computation

- To operationalize the g-formula, a key component is to specify models for all components

    - outcome regression $\Pr(Y^{obs} \mid L_2^{obs}, ..., L_T^{obs}, \overline{A}_T = \overline{a}_T)$

    - models for time-varying confounders
    $\Pr(L_t^{obs}|A_1 = a_1, L_2^{obs}, A_2 = a_2, \ldots, L_{t-1}^{obs}, A_{t-1} = a_{t-1}) \ \forall t$

- Model for time-varying confounder can be complex

    - involves a large number of time points $T$

    - involves many covariates, some of which are continuous

    - may further factor
    $\Pr(L_t^{obs}|A_1 = a_1, L_2^{obs}, A_2 = a_2, \ldots, L_{t-1}^{obs}, A_{t-1} = a_{t-1}) \ \forall t$ with a series of conditional distributions

    - variable selection with longitudinal treatments still an open question

# Parametric g-formula

Specify models for the joint density of time-varying confounders, treatments, and outcomes over time via parametric modeling (Keil et al. 2014)

Step 1: Fit models for each component in the g-computation formula

(1.1) Fit a pooled (over persons and time) model for the conditional distribution of each confounder $L_t$ at time $t$ as function of $t$, past treatment, and covariate history, for example, with a single binary $L_t$

$$\text{logit}\{P(L_t = 1)\} = \beta_0 + \beta_1 t + g_1\{\bar{A}_{t-1}; \beta_2\} + g_2\{\bar{L}_t; \beta_3\}$$

- ▶ Denote $\bar{A}_t = (A_1, \ldots, A_t)$ and $\bar{L}_t = (L_1, \ldots, L_t)$
- ▶ $g_1\{\bar{A}_{t-1}; \beta_2\} = g_1\{A_{t-1}; \beta_2\}$ (concurrent); $g_1\{\bar{A}_{t-1}; \beta_2\} = g_1\{\sum_{k=1}^{t-1} A_k; \beta_2\}$ (cumulative)
- ▶ For multivariate $L_t$, either use a series of conditional models or specify a model for each component of $L_t$

# Parametric g-formula

Step 1: Fit models for each component in the g-computation formula

(1.2) Fit a pooled (logistic) model for the (binary) outcome $Y_T$ as a function of past treatment, and confounder history, for example,

$$\text{logit}\{P(Y_T = 1)\} = \eta_0 + h_1\{\bar{A}_T; \eta_2\} + h_2\{\bar{L}_T; \eta_3\}$$

▶ Again, this is a time-averaged model (or a set of time-specific models)

▶ Here we illustrate ideas based on relatively simple models – they are likely oversimplifications for realistic settings but convenient choices

# Parametric g-formula

Step 2: Approximate the sum (or integral) by performing Monte Carlo simulation for $S$ number of times based on the intervention sequence (regimen) of interest. For each $t \geq 2$

(2.1) Simulate time-varying confounders from the fitted models in Step (1.1) using previously simulated confounders and assigned treatment values. The assignment treatment values $\bar{a}_T = (a_1, \ldots, a_T)$ will be set according to the target estimand of interest

- ► For $t = 2$, simulate $L_2$ from $\text{logit}\{P(L_2 = 1)\} = \hat{\beta}_0 + 2\hat{\beta}_1 + g_1\{a_1; \hat{\beta}_2\}$

- ► For $t = 3$, simulate $L_3$ from
  $\text{logit}\{P(L_3 = 1)\} = \hat{\beta}_0 + 3\hat{\beta}_1 + g_1\{a_1, a_2; \hat{\beta}_2\} + g_2\{\hat{L}_2; \hat{\beta}_3\}$

- ► and so on . . .

# Parametric g-formula

(2.2) Based on all simulated confounders and the treatment assigned, compute the average potential outcomes for all patients using the outcome model fitted in Step (2.1)

- ▶ standard averaging, but used in the more complex longitudinal settings
- ▶ obtain $\hat{\mathbb{E}}^s[Y(\bar{a}_T)]$ based on the $s$th simulation

Step 3: Calculate the average of the estimated potential outcomes over all generated simulation, obtain

$$\widehat{\mathbb{E}}[Y(\bar{a}_T)] = \frac{1}{S} \sum_{s=1}^{S} \hat{\mathbb{E}}^s[Y(\bar{a}_T)]$$

With a large $S$, Step 3 tries to minimize the simulation error

# Compatibility

- Parametric g-estimators can require many modeling assumptions

- Depending on their functional forms, it is possible that the parametric models

$$\left\{ \mu_t(\bar{a}_T, \bar{L}_t; \delta_t) : t = T, \ldots, 1 \right\}$$

  are mutually incompatible, i.e. no joint distribution satisfies all $K$ simultaneously

- One could hope for small bias of each model does not add up

- Compatibility itself is not a practical drawback

  - Even for parametric models that are mutually compatible, the models are practically (although not logically) certain to be misspecified (Bang and Robins, 2005)

- Main point is misspecification of outcome models can bring bias

- Highlight the need for flexible outcome models

# g-null paradox

- g-null paradox means "model misspecification leads to hypothesis tests that inevitably reject the null hypothesis as sample size increases, even when the causal null hypothesis is true" (Robins, 2003)

- "Worst" consequence of model incompatibility: choices of model form can rule out the null hypothesis *a priori* because no parameter values in the model parameter space are consistent with the causal null

- Cannot be easily assessed by frequentist approaches, but can be somewhat relieved by expanding the model space with flexible modeling

- In the Bayesian setting (Bayesian g-formula), g-null paradox can be assessed by examining whether the prior predictive distribution of the potential outcomes rules out the g-null hypothesis (Keil et al. 2018)

# Dimension reduction and propensity score

- ▶ When all conditional distributions in g-computation are correctly specified, g-computation leads to the most efficient estimates with the smallest large sample variances

- ▶ However, dimension of variables increases exponentially with $T$, due to time-varying covariates

- ▶ With medium to large $T$, model building and model checking in g-computation can be very demanding

- ▶ Dimension reduction is crucial. Propensity score again plays a central role to achieve dimension reduction: weighting or outcome regression. Matching is less suitable.

- ▶ Positivity/overlap can be checked in terms of the propensity score, instead of directly on covariates (ignored by the g-computation estimator)

# Dimension reduction and propensity score

▶ Define the propensity score at time $t$ given the observed history as:

$$e_{it} = \Pr(A_{it} = 1 \mid \boldsymbol{H}_{it}), \quad i = 1, ..., T,$$

where $\boldsymbol{H}_{it} = \{\bar{L}_{it}, \bar{A}_{i,t-1}\}$ is the observed history for unit $i$ up to time $t$

▶ Under SI, easy to show SI holds for the longitudinal propensity scores: For a given $t$, and for all $\bar{a}_t$

$$\{Y_i(\bar{a}_t)\} \perp A_{i,t} \mid e_{i,1}, A_{i,1}, ..., e_{i,t-1}, A_{i,t-1} \tag{6}$$

▶ Equation (6) imply that instead of adjusting for the history of covariates, we can adjust for the history of propensity scores - substantially reduce the covariate dimension in modeling

▶ Two approaches: weighting (marginal structural models (MSM)) and regression on PS history

# Longitudinal treatments: Recap of notation

- $N$ subjects

- $A_0$: baseline treatment at baseline time $\tau_0$

- $L_0$: vector of baseline covariates

  Suppose we have $T + 1$ subsequent follow-up visits at times $\tau_1, \ldots, \tau_{T+1}$

- $A_t$: treatment at visit $t$ during interval $[\tau_t, \tau_{t+1})$

- $L_t$: covariates measured just before $A_t$ and remain unchanges during interval $[\tau_t, \tau_{t+1})$

# Longitudinal treatments: Recap of notation

- $L_0$ can contain time-fixed covariates (age and gender etc, often denoted as $X$) and other baseline measure of <span style="color:red">tima-varying covariates</span>

- An outcome $Y_i$ is measured at the final visit $T + 1$

- Assume $A_t$ is binary for simplicity, but generalization to multiple treatments is possible

- Goal: estimate the causal effect of the <span style="color:red">time-varying treatment</span> on the outcome <span style="color:red">in the combined population</span>, using the observational data (the sequence of $A_t$ is non-randomized)

- Question: what is the ideal randomized trial that we wish to mimic?

# Marginal structural model (MSM)

▶ Under sequential ignorability and longitudinal positivity/overlap, we have used the g-computation to estimate the causal estimand $\mathbb{E}[Y(\bar{a})]$

▶ Even if the treatment is binary, there are $2^{T+1}$ values of $\bar{A} = \bar{a}$, where we use $\bar{A}$ to denote the history of treatment or treatment path

▶ Causal effects are characterized by $2^{T+1}$ average potential outcomes corresponding to each treatment path:

$$\left\{ \mathbb{E}[Y(\bar{a}) : \bar{a} \in \bar{\mathcal{A}}] \right\}$$

▶ As $T$ increases, the high-dimensional nature of this characterisation leads to difficulties both with estimation (due to an insufficient number of subjects following any given trajectory) and with interpretation (due to too many potential comparisons)

# Marginal structural model

▶ Some simplification is necessary, for example, one can use the following GLM with inverse link function $h$

$$\mathbb{E}[Y(\bar{a})] = h(\bar{a}; \gamma)$$

▶ For example, we can posit $h(\bar{a}; \gamma) = \gamma_+ \gamma_2 \mathrm{cum}(\bar{a})$, where $\mathrm{cum}(\bar{a}) = \int_0^{T+1} a(t)dt = \sum_{t=0}^{T} a(t)$ is the cumulative treatment

▶ Can further adjust for baseline covariates (effect modifiers), for example with $V \in L_0$

$$\mathbb{E}[Y(\bar{a})] = h\left(\gamma_1 + \gamma_2 \sum_{t=0}^{T} a(t) + \gamma_3 V\right)$$

▶ We call such models marginal structural models (MSM)
  ▶ models for some aspect of the conditional distribution of the counterfactuals given baseline covariates
  ▶ always marginal with respect to post-baseline confounders

# MSM versus Associational Model

- ▶ MSM is a structural model, different from associational models, which are

$$\mathbb{E}(Y|\bar{A} = \bar{a}) = h(\bar{a}; \alpha)$$

- ▶ The associational models are concerned with only observed outcomes, and therefore $\alpha \neq \gamma$ if there is time-varying confounding (when would they be equal?)

- ▶ Directly fitting the associational model to the observed data leads to bias - remember to the simple cancer progression example?

- ▶ Can use weighted estimation of the associational model to remove time-varying confounding and recover $\gamma \Rightarrow$ IPW estimation of MSM

# Generalizing IPW (Horvitz-Thompson) estimator

▶ Define the <span style="color:red">propensity score</span> at time $t$ given the observed history as:

$$e_t = P(A_t = 1 \mid \bar{A}_{t-1}, \bar{L}_t), \quad i = 1, ..., T,$$

▶ Obtain the <span style="color:red">stabilised inverse probability weights</span> for each individual

$$SW = \frac{\prod_{t=0}^{T} P(A_t = A_t^{obs} \mid \bar{A}_{t-1})}{\prod_{t=0}^{T} P(A_t = A_t^{obs} \mid \bar{A}_{t-1}, \bar{L}_t)}$$

with $A_{-1} = \emptyset$

▶ Replace the denominator with 1 leads to the unstablized weights

▶ If the MSM further conditions on baseline covariates $V$, the stabilised inverse probability weights can be further modified as

$$SW\text{-}V = \frac{\prod_{t=0}^{T} P(A_t = A_t^{obs} \mid \bar{A}_{t-1}, V)}{\prod_{t=0}^{T} P(A_t = A_t^{obs} \mid \bar{A}_{t-1}, \bar{L}_t)}$$

# IPW estimation of MSM

- ▶ For $T = 2$
- ▶ Specify a model for the outcome, $f(y, \bar{a}, ; \gamma)$ with score function
  $S(y, a_0, a_1; \gamma) = \frac{\partial}{\partial \gamma} \log f(y, a_0, a_1; \gamma)$.
- ▶ For example, for a binary outcome $Y$ with two time points, two possible models:
  - ▶ $\text{logit}\{\Pr(Y_i = 1 | A_0, A_1)\} = \gamma_1 + \gamma_2 A_0 + \gamma_3 A_1 + \gamma_4 A_0 A_1$
  - ▶ $\text{logit}\{\Pr(Y_i = 1 | A_0, A_1)\} = \gamma_1 + \gamma_2 (A_0 + A_1)$.
- ▶ Solving for the following estimating equation

$$\sum_{i=1}^{N} \text{SW}_i \times S(Y_i^{obs}, A_{i0}, A_{i1}; \gamma) = 0, \qquad (7)$$

  gives consistent estimates of the parameters $\gamma$
- ▶ Eq (7) is solved via maximizing the weighted likelihood

# Estimating weights

- As usual, the weights are unknown, and therefore need to be estimated from data

- Consistency of $\hat{\gamma}$ depends on correct estimation of propensity score weights

- Estimate the propensity score at each time, and normalizing the inverse propensity score weights by the unconditional probability of being assigned to the observed treatment at each $t$

$$\text{SW}_i = \frac{\prod_{t=0}^{T} P(A_{i,t} = A_{i,t}^{obs} \mid \bar{A}_{i,t-1}; \hat{\varphi})}{\prod_{t=0}^{T} P(A_{i,t} = A_{i,t}^{obs} \mid \bar{A}_{i,t-1}, \bar{L}_{i,t}; \hat{\beta})} = \frac{\prod_{t=0}^{T} P(A_{i,t} = A_{i,t}^{obs} \mid \bar{A}_{i,t-1}; \hat{\varphi})}{\prod_{t=0}^{T} \hat{e}_{i,t}^{A_{i,t}^{obs}} \{1 - \hat{e}_{i,t}\}^{1 - A_{i,t}^{obs}}}$$

where $\hat{\phi}$ and $\hat{\beta}$ the MLE of the respective propensity score models

# Estimating weights

- $\Pr(A_t = A_t^{obs} | A_{i,t-1}, ..., A_{i,0})$ can be estimated by the proportion of subjects in cell of $A_{i,t-1}^{obs}, ..., A_{i,0}^{obs}$ in the study sample with $A_{i,t} = A_{i,t}^{obs}$

- When $T$ is large, this may lead to many zero weights due to empty cells

- Instead estimate both denominators and numerators from models. For example, let

$$\text{logit}\{\Pr(A_{i,1} = 1 | A_0)\} = \phi_1 + \phi_2 A_0.$$

$$\text{logit}\{\Pr(A_{i,1} = 1 | \bar{L}_1, A_0)\} = \beta_1 + \beta_2 A_0 + \beta_3^T L_1 + \beta_4^T A_0 L_1.$$

- The last model assumes a "Markov" type condition so that $A_{i,1} \perp L_0 | A_0, L_1$ – usually the lagged covariates and treatment may be most predictive (rationale for a pooled model when $T$ is large)

- Important to check overlap of propensity scores at each time (what about balance?)

# MSM: procedure

Step 1  Build an outcome model: $\Pr(Y(a_0, a_1)) = \Pr(Y|A_1, A_2)$ under "randomization"

Step 2  Build a propensity score model for each time: $\Pr(A_1|A_0, \bar{L}_1)$ and $\Pr(A_0|L_0)$; also build model for $\Pr(A_2|A_1)$ and $\Pr(A_1)$ (for stabilized weights) – this can be replaced with pooled models

Step 3  Estimate the propensity scores at each time, check overlap and remove units in the non-overlap region

Step 4  Calculate the stabilized weights for each unit at each time point

Step 5  Estimate the parameters of the outcome model by maximizing the weighted likelihood (weighted regression)

Case study: Hernan, Robin, Brumback (2000). outcome is survival, use the marginal structural Cox model

# An illustrative example

Daniel et al. 2013

- Let $\gamma_0 = -0.5$, $\gamma_1 = -0.75$ and $\gamma_{01} = 0.2$, and generate data from the true MSM, and induce time-varying confounding through $L_1$

$$\mathbb{E}\{Y(a_0, a_1)\} = \gamma_{int} + \gamma_0 a_0 + \gamma_1 a_1 + \gamma_{01} a_0 a_1$$

- This is a nonparametric MSM because it is saturated

- Fit a naive associational model without adjusting for $L_1$

$$\mathbb{E}[Y|A_0, A_1] = \alpha_{int} + \alpha_0 A_0 + \alpha_1 A_1 + \alpha_{01} A_0 A_1$$

- Fit an adjusted associational model

$$\mathbb{E}[Y|A_0, A_1, L_1] = \beta_{int} + \beta_0 A_0 + \beta_1 A_1 + \beta_{01} A_0 A_1 + \beta_l L_1$$

- Using what we have learned, which one of these models could estimate the causal parameters $\gamma$?

# An illustrative data analysis

Daniel et al. 2013

**Table I.** The results of the naïve analyses of simulated dataset I, with and without adjusting for $L_1$, along with the true values of the parameters of (4).

| Parameter | Estimate | 95% CI | Parameter | Estimate | 95% CI | Parameter | True value |
|---|---|---|---|---|---|---|---|
| $\alpha_0$ | $-0.390$ | $(-0.453, -0.327)$ | $\beta_0$ | $-0.585$ | $(-0.658, -0.511)$ | $\gamma_0$ | $-0.5$ |
| $\alpha_1$ | $-0.806$ | $(-0.874, -0.738)$ | $\beta_1$ | $-0.746$ | $(-0.813, -0.678)$ | $\gamma_1$ | $-0.75$ |
| $\alpha_{01}$ | $0.096$ | $(0.002, 0.190)$ | $\beta_{01}$ | $0.258$ | $(0.160, 0.356)$ | $\gamma_{01}$ | $0.2$ |

▶ None! Same essence as the earlier simple (nonparametric) example

▶ Unadjusted associational model ignores time-varying confounding

▶ Adjusted associational model induces collider-stratification bias

▶ The correct approach is either g-computation or MSM

▶ Read more about this example in a tutorial by Daniel et al. (2015)

# An illustrative data analysis

Daniel et al. 2013

**Table III.** The results of the analysis of simulated dataset I, as analysed using the g-computation formula to obtain the parameters of the marginal structural model defined in Equation (4), with bootstrap standard errors; the 95% CIs are based on a normal approximation, using the bootstrap standard errors.

| Parameter | True value | G-computation estimate | Bootstrap SE | 95% CI | |
|---|---|---|---|---|---|
| $\gamma_{int}$ | 2.1 | 2.075 | 0.021 | 2.034 | 2.116 |
| $\gamma_0$ | −0.5 | −0.500 | 0.047 | −0.592 | −0.407 |
| $\gamma_1$ | −0.75 | −0.739 | 0.037 | −0.811 | −0.667 |
| $\gamma_{01}$ | 0.2 | 0.218 | 0.061 | 0.098 | 0.339 |

**Table V.** The results of the analysis of the simulated dataset I, as analysed using inverse probability weighting in the marginal structural model (4).

| Parameter | True value | IPW estimate | SE[‡] | 95% CI | |
|---|---|---|---|---|---|
| $\gamma_{int}$ | 2.1 | 2.075 | 0.021 | 2.034 | 2.116 |
| $\gamma_0$ | −0.5 | −0.500 | 0.048 | −0.593 | −0.406 |
| $\gamma_1$ | −0.75 | −0.739 | 0.039 | −0.815 | −0.663 |
| $\gamma_{01}$ | 0.2 | 0.218 | 0.065 | 0.091 | 0.345 |

[‡]This is the sandwich estimator of standard error, which takes into account the non-independence of pseudo-subjects as a result of weighting.

# Inference for MSM

- Bootstrap will provide valid variance and interval estimates using MSM

- Treating $SW_i$ as fixed, and use the robust sandwich variance to provide a conservative variance estimate

    - essentially a "survey" weighted GLM

    - In R, this is done by the svyglm() in the survey package

- Can invoke the M-estimation theory to provide a more accurate sandwich variance that takes into account the estimation of the stablized weights

    - can be cumbersome if the weights are not estimated via pooled models over person and time

# Loss to follow up

▶ Write $R_{i,t} = 1$ is subject $i$ observed at visit $t$, and zero otherwise; assume missing at random (MAR),

$$P(R_{i,t} = 1|\bar{A}_i, \bar{L}_i, R_{i,t-1} = 1) = P(R_{i,t} = 1|\bar{A}_{i,t-1}, \bar{L}_{i,t-1}, R_{i,t-1} = 1)$$

▶ The stablized IPW becomes

$$\text{SW} = \frac{\prod_{t=0}^{T} P(A_{i,t} = A_{i,t}^{obs} \mid \bar{A}_{i,t-1}, R_{i,t-1} = 1)P(R_{i,t} = 1|\bar{A}_{i,t-1}, R_{i,t-1} = 1)}{\prod_{t=0}^{T} P(A_{i,t} = A_{i,t}^{obs} \mid \bar{A}_{i,t-1}, \bar{L}_{i,t}, R_{i,t-1} = 1)P(R_{i,t} = 1|\bar{A}_{i,t-1}, \bar{L}_{i,t-1}, R_{i,t-1} = 1)}$$

▶ An advantage of IPW is that specialized software routines are not, in general, needed, because the models can be fitted using standard regression commands, incorporating weights (SAS, R and Stata)

# Doubly Robust MSM

(Bang and Robins, 2005 Biometrics)

- ▶ Simple weighting is not efficient, and lead to bias if the weights are incorrectly specified

- ▶ The ICE estimator (one form of g-formula) can be used to estimate MSM parameters as well, with a simple modification in the final sequential regression

  - ▶ recall ICE estimator imputes all potential outcomes (in an exhaustive fashion), and so in the final regression model, we can simply just fit a MSM to estimate $\gamma$

- ▶ Combine IPW and ICE to create a doubly-robust estimator for MSM

  - ▶ consistent for $\gamma$ is either series of propensity score models or outcome models are correct

  - ▶ more efficient than IPW alone by exploiting a series of conditional outcome models

# Doubly Robust MSM

Bang and Robins, 2005 Biometrics

- First estimate the propensity scores at each time, and create
  $\bar{\pi}_t(\hat{\beta}) = \prod_{k=0}^t P(A_t = A_t^{obs} | \bar{A}_{t-1}, \bar{L}_t; \hat{\beta})$

    - e.g. one can use a pooled logistic model over persons and time

- The DR MSM estimator uses the idea of "clever covariate" in the sequential regression steps (augmenting the ICE estimator by the "clever covariate")

    - the clever covariate at each time is defined as $\bar{\pi}_t^{-1}(\hat{\beta})$

    - when performing (canonical link) regression at each time, include the linear term $\bar{W}(\bar{A}_t) = \bar{\pi}_t^{-1}(\hat{\beta})$

- Connecting to Lecture 7 on DR estimator for point treatment with a clever covariate

# Operationalizing Doubly Robust MSM

- ▶ Recall the example with $T = 3(t = 0, 1, 2)$ time points as in the last lecture, and the observed data consists of $(L_0, A_0, L_1, A_1, L_2, A_2, Y)$

- ▶ Step 1: Compute the inverse probability weights at each time point to create the "clever covariate", e.g., using a pooled logistic regression for person-time data

$$\text{logit}\{P(A_t = A_t^{obs}|\bar{A}_{t-1}, \bar{L}_t; \hat{\beta})\} = \beta_{0,t} + \beta_1 A_{t-1} + \beta_2 L_t$$

the weight at time $t = 0, 1, 2$ for each unit then becomes

$$\hat{W}(\bar{A}_t) = \prod_{t=0}^{t} 1/\hat{P}(A_t = A_t^{obs}|\bar{A}_{t-1}, \bar{L}_t; \hat{\beta}) = \prod_{t=0}^{t} \hat{\pi}^{-1}(A_t|\bar{A}_{t-1}, \bar{L}_t)$$

# Operationalizing Doubly Robust MSM

▶ Step 2: At time $t = 2$, postulate the mean model to compute regression coefficient vector $\delta_2$ with ordinary least squares

$$Y = \delta_{2,0} + \delta_{2,1} cum(\bar{A}_2) + \delta'_{2,2} L_2 + \delta_{2,3} \hat{W}(\bar{A}_2) + \epsilon_2$$

For time $t = 1$, use the above fitted model to compute the pseudo-outcomes for each unit under observed history $\bar{A}_1$ but all possible treatment status at time $t = 2$. In other words, for each unit, we expand the data set with two rows per patient; one with outcome $\hat{Y}(\bar{A}_1, A_2 = 1)$ and $\hat{Y}(\bar{A}_1, A_2 = 0)$ and same covariate and treatment history up to time 1 otherwise

# Operationalizing Doubly Robust MSM

- **Step 3**: Define $Y^{(1)} = \hat{Y}(\bar{A}_1, a_2)$ for $a_2 \in \{0, 1\}$, then perform linear regression on this pseudo-outcome with sample size $2N$ to estimate $\delta_2$, allowing for the clever covariate

$$Y^{(1)} = \delta_{1,0} + \delta_{1,1} cum(\bar{A}_1) + \delta'_{1,2} L_1 + \delta_{1,3} \bar{W}(\bar{A}_1) + \epsilon_1$$

- Repeat the above outcome imputation step, but now for time $t = 0$, we use the above fitted model to compute pseudo-outcomes for each unit under observed treatment history $A_0$ but all possible treatment status for the future

# Operationalizing Doubly Robust MSM

▶ That is to say, we further expand the data set to have four rows per patient, with replicated histories $A_0$ and $L_0$, but different outcomes given by $\hat{Y}(A_0, A_1 = A_2 = 1)$, $\hat{Y}(A_0, A_1 = 0, A_2 = 1)$, $\hat{Y}(A_0, A_1 = 1, A_2 = 0)$ and $\hat{Y}(A_0, A_1 = A_2 = 0)$ (expanded data set with $4N$ sample size)

▶ Step 4: Define $Y^{(0)} = \hat{Y}(A_0, a_1, a_2)$ for $a_1, a_2 \in \{0, 1\}$, then perform linear regression on this pseudo-outcome with sample size $4N$ to estimate $\delta_0$

$$Y^{(0)} = \delta_{0,0} + \delta_{0,1} A_0 + \delta'_{0,2} L_0 + \delta_{0,3} \bar{W}(A_0) + \epsilon_0$$

# Operationalizing Doubly Robust MSM

▶ Step 4: As before, expanding the data set to obtain eight row per patient, with outcomes $\hat{Y}(A_0 = A_1 = A_2 = 1)$, $\hat{Y}(A_0 = 0, A_1 = A_2 = 1)$, $\hat{Y}(A_0 = A_1 = 0, A_2 = 1)$, $\hat{Y}(A_0 = 1, A_1 = 0, A_2 = 1)$, $\hat{Y}(A_0 = 0, A_1 = 1, A_2 = 0)$, $\hat{Y}(A_0 = A_1 = 1, A_2 = 0)$ and $\hat{Y}(A_0 = A_1 = A_2 = 0)$, $\hat{Y}(A_0 = 1, A_1 = A_2 = 0)$

▶ At step 4, we have already computed all eight potential outcomes for each unit, and therefore we can just compute the final average potential outcomes by simple averaging!

# Operationalizing Doubly Robust MSM

► For example,

$$\widehat{\mathbb{E}}[Y(0, 0, 0)] = \frac{1}{N} \sum_{i=1}^{N} \hat{Y}_i(0, 0, 0)$$

$$\widehat{\mathbb{E}}[Y(1, 1, 1)] = \frac{1}{N} \sum_{i=1}^{N} \hat{Y}_i(1, 1, 1) \ldots \ldots$$

► Given the interest is in the marginal structural model parameter with baseline covariate $V \in L_0$, then can perform a final regression based on all imputed potential outcomes, for example, run the model with $8N$ sample size

$$\hat{Y}(a_1, a_2, a_3) = \gamma_1 + \gamma_2 cum(\bar{A}_2) + \gamma_3 V + \epsilon^*$$

# Doubly Robust MSM

▶ A direct extension of doubly robust estimator with a cross-sectional, point treatment in Lecture 7

▶ Leads to consistent estimator of the structural regression parameters if

  ▶ either the sequence of propensity score models are correctly specified

  ▶ or the sequence of outcome regression models are correctly specified

  ▶ but not necessarily both

▶ When both sequences of models are correct, the resulting estimator for structural regression parameters are semiparametric efficient

▶ The diagnostic nature of the DR MSM estimator is also similar to the case with a point treatment

  ▶ has not been used much in practice, because no software + a bit more complicated construction. . .

# Covariate-balancing propensity scores

Imai and Ratkovic, 2015 JASA

- ▶ For IPW estimators, may have extreme weights under lack of overlap

- ▶ Lack of overlap is much more likely with large $T$

- ▶ The covariate-balancing propensity scores directly balance the mean covariates, and bypass the need to estimate a parametric propensity score model

  - ▶ usually more efficient and robust

- ▶ But the balancing conditions at each time $t$ are based on the entire weight history (across $T$ time points), and becomes much more complicated with large $T$

- ▶ Other calibration estimator exists (e.g. residual balancing for MSM; Zhou and Geoffrey, 2020)

# Strength of MSM

- Intuitive and relatively easy to explain – compared to g-formula, most closely related to standard methods

- easily extended to different types of outcome variable

- only require to specify models for the treatment assignment/propensity score and the MSM itself

  - conditional distributions of (1) outcome $Y$ given the covariates and (2) time-varying covariates given past covariates and treatments are left unspecified (unlike g-formula)

  - less prone to model misspecification than the g-formula (high-dimensional covariates)

- not prone to g-null paradox

# Limitations of MSM

- inverse weighting can be unstable and inefficient if there are extreme weights

    - stablized weights can help, but as we see in the cross-sectional setting, not a lot

    - prone to extreme weights, can use weight trimming or truncation (Cole and Hernan, AJE 2008), but similar issue as in one time point (sensitive to cutoff and ambiguous target population

    - possible to extend to target populations, e.g. other balancing weights such as overlap weights, but remain an open question

- possible interactions between treatment and time-varying covariates cannot be explored because the MSM is marginal with respect to the latter

# Balancing Weights for MSM?

- We have seen from prior lectures that IPW can be generalized to the family of balancing weights, among which overlap weights address the positivity issues from a design based perspective

    - focus on interpretable overlap population at equipoise

- This is in general much harder to operationalize with longitudinal treatments and confounding under the MSM framework

    - whether the tilling function depends only on baseline confounders?

    - if the tilling function depends on time-varying confounders, the counterfactual time-varying covariates are not fully observable

    - how to conceptualize the overlap target population with time-varying treatment patterns

    - open methodological questions to address

# Regression on longitudinal propensity score (RLPS)

- ▶ MSM uses the longitudinal propensity score as inverse weight
- ▶ An alternative approach is through regressing the history of the longitudinal propensity scores.
- ▶ Achy-Brou et al. (2010) showed that, similar to the g-computation:

$$
\Pr(Y_i(\bar{a}_T))
$$
$$
= \sum_{e_1,\ldots,e_T} \Pr(Y_i^{obs} \mid e_{i,1}, A_{i,1} = a_1, \ldots, e_{i,T}, A_{i,T} = a_T)
$$
$$
\times \Pr(e_{i,T} | e_{i,1}, A_{i,1} = a_1, \ldots, e_{i,T-1}, A_{i,T-1} = a_{T-1})
$$
$$
\ldots \times \Pr(e_{i,2} | e_{i,1}, A_{i,1} = a_1) \Pr(e_{i,1}). \tag{8}
$$

- ▶ Therefore, given models for the RHS of the equation, we can estimate the target quantities $\Pr(Y(\bar{a}_t))$

# Regression on longitudinal propensity score (RLPS)

Here instead of conditioning on all covariates, we use the propensity scores as the single predictor at each time

Pros:

► Simpler to specify model and conduct model checking

► If the models are corrected specified, regression estimators are more efficient (smaller variance in large samples) than the weighting estimators

► Regression estimators are not as sensitive to extreme weights as weighting estimators

Cons:

► Not as efficient as g-computation; and need to model distributions of longitudinal PS

# RLPS: case study

Achy-Brou, Frangakis and Griswold (2009, Biometrics)

- **Units**: patients with diabetes who took active treatments.
- **Time points**: $T = 3$; Sample size: $n = 131,714$.
- **Treatments at one time**: $a_1$ Insulin; $a_2$ Exenatide; $a_3$ both; $a_4$ other.
- **Outcome**: (1) hospitalization rate; (2) total health care cost.
- **Covariates**: 3 baseline and 18 time-varying ones.

- **Goal**: predict and compare patient outcomes if all patients had been assigned to "Insulin-Insulin-Insulin" (In3), "Exenatide-Exenatide-Exenatide" (Ex3), "Other-Other-Other" (O3) longitudinal treatments, adjusting for time-varying confounding.

## RLPS: case study

Step 1 For each time point $t = 1, ..., 3$ and each treatment $k = 1, ...K$, estimate the propensity score models: $\Pr(A_{i,t} = k | H_{i,t}) = e_{i,t,k}$.

Step 2 At each time point $t$ and for each treatment $k$: Stratify subjects into five blocks by the quintiles of the propensity scores. Check covariates balance within each block. Remove units in the non-overlapping region of the propensity scores.

Step 3 Fit regression models to estimate the probabilities $\Pr(e_t | e_{t-1}, A_{t-1}, ..., e_1, A_1)$.

Step 4 Fit a regression model for $\Pr(Y | e_T, A_T, e_{t-1}, A_{t-1}, ..., e_1, A_1)$ and estimate the model parameters by, say, MLE

Step 5 Estimate the average potential outcome $\mathbb{E}(Y(\bar{a}_T))$ by (8), with the parameter estimates obtained in Steps 3 and 4.

# Bayesian Inference of Longitudinal Treatments

▶ Similar to single-time treatments, Bayesian inference considers the observed values of the four quantities to be realizations of random variables and the unobserved values to be unobserved random variables

▶ Use $T = 2$ to illustrate: there are six potential outcomes for each units:

$$\mathbb{V} \equiv (L(0), L(1), Y(0, 0), Y(1, 0), Y(0, 1), Y(1, 1)).$$

▶ For each unit $i$, we observe one out of two intermediate potential outcomes at time 1, $L_i^{obs} = L_i(A_{i1})$, and one out of four potential outcomes at time 2, $Y_i^{obs} = Y_i(A_{i1}, A_{i2})$

▶ Potential outcomes under unassigned treatment sequences are missing: $L_i^{mis} = L_i(1 - A_{i1})$ and $\mathbf{Y}_i^{mis} = \{Y_i(1 - A_{i1}, A_{i2}), Y_i(A_{i1}, 1 - A_{i2}); Y_i(1 - A_{i1}, 1 - A_{i2})\}$.

# Bayesian Inference: Joint Outcome Modeling

- ▶ Goal: simulate the posterior predictive distributions of the missing potential outcomes
- ▶ Assuming
  - ▶ Sequential ignorability
  - ▶ Exchangeability
  - ▶ Prior independence of parameters in the models for outcome and for the assignment mechanism
- ▶ Then the posterior predictive distribution of the missing potential outcomes is

$$
\Pr(\mathbf{Y}^{mis}, \boldsymbol{L}^{mis} \mid \mathbf{Y}^{obs}, \boldsymbol{L}^{obs}, \boldsymbol{A}_1, \boldsymbol{A}_2)
$$

$$
= \frac{\Pr(\mathbb{V}) \Pr(\boldsymbol{A}_1, \boldsymbol{A}_2 \mid \mathbb{V})}{\iint \Pr(\mathbb{V}) \Pr(\boldsymbol{A}_2, \boldsymbol{A}_1 \mid \mathbb{V},) d\mathbf{Y}^{mis} d\boldsymbol{L}^{mis}}
$$

$$
\propto \int \prod_i \Pr(Y_i(0,0), Y_i(1,0), Y_i(0,1), Y_i(1,1), L_{i1}(0), L_{i1}(1) \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}
$$

# Bayesian Inference: Joint Outcome Modeling

► The most straightforward Bayesian approach for time-varying treatments requires specifying a joint outcome model of all variables (intermediate variables, final potential outcomes and treatment assignment and possibly conditional on baseline covariates)

► And then derive posterior predictive distributions of the missing potential outcomes $\mathbf{Y}_i^{mis}, {}_i^{mis}$ given the observed data

► When $T > 2$, the joint outcome modeling approach may require specifying a very complex high-dimensional model, raising inferential challenges

  ► simplification via g-computation

# Bayesian Inference: g-computation

(Gustafson, 2015, Biometrics; Keil et al. 2018)

- ▶ Bayesian version of g-computation: build a Bayesian model for the outcome given history at each time period, and then use g-formula to combine the posterior draws from these models

- ▶ A simple case of $T = 2$ with binary outcome and treatment

  - ▶ Use a Bayesian saturated binary regression model (one parameter for each cell in the contingency table): $A_1$, $L|A_1$, $Y|(A_1, L, A_2)$

  - ▶ Independent uniform prior for each of the parameters, posterior is independently Beta distribution

  - ▶ Results are very similar to the Bayesian marginal structural model (later), but arguably simpler implementation, and fully Bayesian

# Dynamic treatment regimes

- We have focused on estimating the population effects of a static treatment regimes

- Static treatment regimes assign treatment based only on baseline covariates

- Dynamic treatment regimes assign treatment based on time-varying covariates and baseline covariates

- Much research has been done in the frequentist paradigm, e.g. Murphy (2003) proposed frequentist semiparametric plug-in methods to find the dynamic treatment regime that maximizes the expected final outcome $\mathbb{E}[Y_i(\boldsymbol{a})]$ asymptotically

- From the Bayesian perspective, the procedure is straightforward as that for static treatment regimes - a decision theory perspective (Zajonc, 2012)

# Dynamic treatment regimes: Notations

▶ Notations: baseline covariates $L_0$; treatment at time $t$, $A_t$; intermediate outcome at time 1 (time-varying covariates) $L_1$; final outcome at time 2, $Y$.

▶ A <span style="color:red">dynamic treatment regime</span> (DTR) $\delta$ is a pair of decision functions $\delta_1 : \mathcal{X}_0 \to \{0, 1\}$ and $\delta_2 : \mathcal{X}_0 \times \{0, 1\} \times \mathcal{X}_1 \to \{0, 1\}$ that assign units with observed covariates $(a_1, l_1, l_0)$ to a treatment sequence:

$$\boldsymbol{\delta} \equiv (\delta_1(l_0), \delta_2(l_0, a_1, l_1))$$

▶ More generally, decision functions can be conditional probability distributions over treatments

▶ Potential outcomes indexed by decision functions instead of treatment sequences: $L_{i1}(\delta_1) \equiv L_{i1}(\delta_1(L_{i0}))$ and

$$Y_i(\boldsymbol{\delta}) \equiv Y_i(\delta_1(L_{i0}), \delta_2(L_{i0}, \delta_1(L_{i0}), L_{i1}(\delta_1))).$$

# Average treatment regime effects

- Let $\delta'$ be the reference or placebo treatment

- Average treatment regime effect of treatment regime $\delta$ is

$$\tau(\delta, \delta') = \mathbb{E}[Y_i(\delta) - Y_i(\delta')]$$

- Or the improvement from the status-quo (observed treatment):

$$\tau(\delta) = \mathbb{E}[Y_i(\delta) - Y_i]$$

# Optimal dynamic treatment regimes: Bayesian perspective

- $\mathcal{D}$: the specific class of treatment regimes under consideration

- $u_\delta(\cdot)$: a utility function over outcomes (defined by investigators or policymakers)

- $\boldsymbol{D} = (\boldsymbol{L}_0, \boldsymbol{A}_1, \boldsymbol{L}_1, \boldsymbol{A}_2, \mathbf{Y})$: all observed data from a sample

- The outcome of interest: the <span style="color:red">posterior expected utility</span> for regime $\boldsymbol{\delta}$ given prior $p(\boldsymbol{\theta})$:

$$U(\boldsymbol{\delta}, p \mid \boldsymbol{D})$$
$$\propto \iint u_{\boldsymbol{\delta}}(\tilde{Y}(\boldsymbol{\delta})) \Pr(\tilde{Y}(\boldsymbol{\delta})|\boldsymbol{\theta}) \prod_i^N \Pr(\boldsymbol{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} d\tilde{Y}(\boldsymbol{\delta}),$$

- General strategy: integrate over all sources of uncertainty by maximizing using the posterior predictive distribution

# Optimal dynamic treatment regimes: Bayesian perspective

▶ $\prod_i^N \Pr(\mathbb{V} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$ is the likelihood function of the observed data

▶ Expected utility averages over the posterior distribution of the unknown outcome $\tilde{Y}(\boldsymbol{\delta})$ conditional on the observed data $\boldsymbol{D}$ and incorporate uncertainty in the parameters $\boldsymbol{\delta}$

▶ The optimal treatment regime selects the maximizing rule, conditional on the observed data and prior,

$$\boldsymbol{\delta}^*(p, \boldsymbol{D}) = \arg \max_{\boldsymbol{\delta} \in \mathcal{D}} U(\boldsymbol{\delta}, p | \boldsymbol{D})$$

▶ Given the optimal treatment regime, we can calculate the average improvement over the status-quo: $\mathbb{E}[Y_i(\boldsymbol{\delta}^*) - Y_i]$

# Choosing the utility function

- Zajonc (2012) proposed a mean-variance utility function:

$$U(\boldsymbol{\delta}, p \mid \lambda, \boldsymbol{D}) = \lambda_1 \, \mathbb{E}[\tilde{Y}(\boldsymbol{\delta})|\boldsymbol{D}] + \lambda_2 \mathbb{V}[\tilde{Y}(\boldsymbol{\delta})|\boldsymbol{D}]$$

  Varying $\lambda_1$ and $\lambda_2$ reflects mean-variance preference

- In single time treatment, Dehejia (2005) also considered the Bayesian decision theory perspective.

- Slight different goals of analyzing static and dynamic treatment rules, but rely on same set of assumptions

- DTRs helps better understand causal mechanisms to improve policy and practice across different disciplines.

# References

Robins J. M. (1986). A new approach to causal inference in mortality studies with sustained exposure periods - Application to control of the healthy worker survivor effect. Mathematical Modelling, 7, 1393-1512.

Robins, J. M. (1999). Association, causation, and marginal structural models. Synthese, 151-179.

Keil, A. P., Edwards, J. K., Richardson, D. R., Naimi, A. I., Cole, S. R. (2014). The parametric G-formula for time-to-event data: towards intuition with a worked example. Epidemiology (Cambridge, Mass.), 25(6), 889.

McGrath, S., Lin, V., Zhang, Z., Petito, L. C., Logan, R. W., Hernan, M. A., Young, J. G. (2020). gfoRmula: An R Package for Estimating the Effects of Sustained Treatment Strategies via the Parametric g-formula. Patterns, 100008.

Robins, J. M., Hernan, M. A., Brumback, B. (2000). Marginal structural models and causal inference in epidemiology. Epidemiology, 560-560.

Hernan, M. A, Brumback, B., Robins, J. M. (2000). Marginal structural models to estimate the causal effect of zidovudine on the survival of HIV-positive men. Epidemiology, 561-570.

# References

Cole, S. R., Hernan, M. A. (2008). Constructing inverse probability weights for marginal structural models. American journal of epidemiology, 168(6), 656-664.

Daniel, R. M., Cousens, S. N., De Stavola, B. L., Kenward, M. G., Sterne, J. A. C. (2013). Methods for dealing with time-dependent confounding. Statistics in medicine, 32(9), 1584-1618.

Naimi, A. I., Cole, S. R., Kennedy, E. H. (2017). An introduction to g methods. International journal of epidemiology, 46(2), 756-762.

Robins, J. M. (1997). Causal inference from complex longitudinal data. In Latent Variable Modeling and Applications to Causality. Lecture Notes in Statistics, M. Berkane (eds), vol 120, New York, NY: Springer.

Bang, H., Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. Biometrics, 61(4), 962-973.

Imai, K., Ratkovic, M. (2015). Robust estimation of inverse probability weights for marginal structural models. Journal of the American Statistical Association, 110(511), 1013-1023.

Achy-Brou, A. C., Frangakis, C. E., Griswold, M. (2010). Estimating treatment effects of longitudinal designs using regression models on propensity scores. Biometrics, 66(3), 824-833.

# References

Shinohara RT, Narayan AK, Hong K, et al. Estimating parsimonious models of longitudinal causal effects using regressions on propensity scores. Statistics in Medicine 2013; 32(22): 3829-3837.

Keil, A. P., Daza, E. J., Engel, S. M., Buckley, J. P., Edwards, J. K. (2018). A Bayesian approach to the g-formula. Statistical methods in medical research, 27(10), 3183-3204.

Gustafson P. (2015). Discussion of the article by Saarela et. al. on Bayesian estimation of marginal structural models. Biometrics 71, 291-293

Robins JM. General methodological considerations. J Econometrics. 2003; 112(1):89-106

Hernan, M.A. and Robins, J.M. (2020) Causal Inference: What If. Boca Raton, FL: Chapman & Hall CRC.

Zajonc, T. (2012). Bayesian inference for dynamic treatment regimes: Mobility, equity, and efficiency in student tracking. Journal of the American Statistical Association, 107(497), 80-92.

Murphy, S. A. (2003). Optimal dynamic treatment regimes. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 65(2), 331-355.