

Causal Inference in Complex Longitudinal Settings

Tutorial on Personalized Medicine, Treatment Regimes, Reinforcement Learning, and Causal Inference

Jon Huang | Singapore Institute for Clinical Sciences | 12 February 2019

Schedule

Time	Topic
09:00 - 09:15 AM	Defining a (longitudinal) causal question
09:15 - 09:30 AM	Implementing the g-formula
09:30 - 10:00 AM	Case Study 1 (smoking cessation)

BREAK (10:00 - 10:30 AM)

Time	Topic
10:30 - 10:55 AM	Case Study 2 (HIV/ART)
10:55 - 11:15 AM	Case Study 3 (aspirin)
11:15 - 11:30 AM	Other considerations, methods, Q&A

Logistics

Logistics

- Opportunity to follow along with R code

Logistics

- Opportunity to follow along with R code
- Will walk through steps to access

Logistics

- **Opportunity to follow along with R code**
- **Will walk through steps to access**
- **Presentation and code are available**

Logistics

- Opportunity to follow along with R code
- Will walk through steps to access
- Presentation and code are available
- Engage and ask questions!

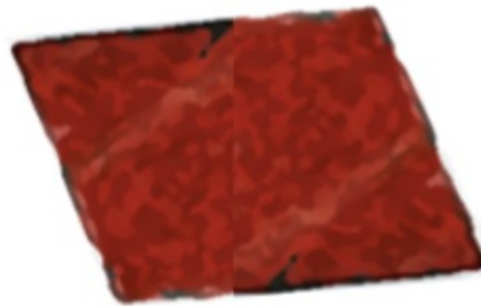
What's the question?

IT'S PAYBACK TIME!

Check out how many calories you need to burn from all the goodies that you ate during Chinese New Year!

1 slice
301
cal

BAK KWA



To burn off **1** slice

30
mins of
Running

1 slice
237
cal

KUEH LAPIS



To burn off **1** slice

35
mins of
Swimming

DEFINING A (LONGITUDINAL) CAUSAL QUESTION

What is a cause? (Review)

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

Causal effects can be defined by difference in "potential outcomes" (PO):

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

Causal effects can be defined by difference in "potential outcomes" (PO):

- For $X \in (0, 1)$, individual i has two POs:

$$Y_i^{x=1} \text{ or } Y_i^{x=0}$$

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

Causal effects can be defined by difference in "potential outcomes" (PO):

- For $X \in (0, 1)$, individual i has two POs:

$$Y_i^{x=1} \text{ or } Y_i^{x=0}$$

- If all POs could be observed \rightarrow causal effects (contrasts) e.g.

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

Causal effects can be defined by difference in "potential outcomes" (PO):

- For $X \in (0, 1)$, individual i has two POs:

$$Y_i^{x=1} \text{ or } Y_i^{x=0}$$

- If all POs could be observed \rightarrow causal effects (contrasts) e.g.

$$ACE_{X \rightarrow Y} = E[Y^{x=1}] - E[Y^{x=0}]$$

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

Causal effects can be defined by difference in "potential outcomes" (PO):

- For $X \in (0, 1)$, individual i has two POs:

$$Y_i^{x=1} \text{ or } Y_i^{x=0}$$

- If all POs could be observed \rightarrow causal effects (contrasts) e.g.

$$ACE_{X \rightarrow Y} = E[Y^{x=1}] - E[Y^{x=0}]$$

- Observational data can be considered incomplete subset of POs

What is a cause? (Review)

In counterfactual terms:

- A **cause** produces an outcome that is different from what **would have been** observed in its absence

Causal effects can be defined by difference in "potential outcomes" (PO):

- For $X \in (0, 1)$, individual i has two POs:

$$Y_i^{x=1} \text{ or } Y_i^{x=0}$$

- If all POs could be observed \rightarrow causal effects (contrasts) e.g.

$$ACE_{X \rightarrow Y} = E[Y^{x=1}] - E[Y^{x=0}]$$

- Observational data can be considered incomplete subset of POs
- Confounding can thus be defined in terms of bias in set of observed POs.

Translating observed statistical parameters to causal quantities

For confounders C , binary exposure X , continuous outcome Y :

Translating observed statistical parameters to causal quantities

For confounders C , binary exposure X , continuous outcome Y :

Exchangeability:

$$Y^x \perp X = x | C$$

Translating observed statistical parameters to causal quantities

For confounders C , binary exposure X , continuous outcome Y :

Exchangeability:

$$Y^x \perp X = x | C$$

Positivity:

$$0 < Pr(A = a | C = c) < 1$$

Translating observed statistical parameters to causal quantities

For confounders C , binary exposure X , continuous outcome Y :

Exchangeability:

$$Y^x \perp X = x | C$$

Positivity:

$$0 < Pr(A = a | C = c) < 1$$

Consistency:

$$\text{if } X = x \text{ then } Y^x = Y$$

Translating observed statistical parameters to causal quantities

For confounders C , binary exposure X , continuous outcome Y :

Exchangeability:

$$Y^x \perp X = x | C$$

Positivity:

$$0 < Pr(A = a | C = c) < 1$$

Consistency:

$$\text{if } X = x \text{ then } Y^x = Y$$

Causal quantity from observed quantities:

$$E(Y | X = x, C) = E(Y^x | C)$$

How to translate a question to an *estimable* causal contrast?

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

- Must know enough about the outcome / disease process

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

- Must know enough about the outcome / disease process
- What else influences it?

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

- Must know enough about the outcome / disease process
- What else influences it?
- Are these influences related to exposure?

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

- Must know enough about the outcome / disease process
- What else influences it?
- Are these influences related to exposure?

Do we have enough observational units at each level of exposure (and covariates) to model POs? (*positivity*)

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

- Must know enough about the outcome / disease process
- What else influences it?
- Are these influences related to exposure?

Do we have enough observational units at each level of exposure (and covariates) to model POs? (*positivity*)

Does the effect of exposure vary based on how that level was attained? (*consistency*)

How to translate a question to an *estimable* causal contrast?

Can you justify PO distribution is (conditionally) independent from observed exposure level? (correct model specification and *exchangeability*)

- Must know enough about the outcome / disease process
- What else influences it?
- Are these influences related to exposure?

Do we have enough observational units at each level of exposure (and covariates) to model POs? (*positivity*)

Does the effect of exposure vary based on how that level was attained? (*consistency*)

- Critical, varies by target population, and often least considered aspect

Example:

General: "What is the effect of smoking on body weight?"

Example:

General: "What is the effect of smoking on body weight?"

Better:

- "... among individuals who ever smoked vs. never?"
- "... among individuals who smoked 10 vs. 11 cigarettes per day?"
- "... among individuals who smoked daily for 10 years stopping smoking?"

Example:

General: "What is the effect of smoking on body weight?"

Better:

- "... among individuals who ever smoked vs. never?"
- "... among individuals who smoked 10 vs. 11 cigarettes per day?"
- "... among individuals who smoked daily for 10 years stopping smoking?"

Each question represents a different target population contrast (*consistency*)

Example:

General: "What is the effect of smoking on body weight?"

Better:

- "... among individuals who ever smoked vs. never?"
- "... among individuals who smoked 10 vs. 11 cigarettes per day?"
- "... among individuals who smoked daily for 10 years stopping smoking?"

Each question represents a different target population contrast (*consistency*)

Each may have different set of relevant confounders (*exchangeability*)

Example:

General: "What is the effect of smoking on body weight?"

Better:

- "... among individuals who ever smoked vs. never?"
- "... among individuals who smoked 10 vs. 11 cigarettes per day?"
- "... among individuals who smoked daily for 10 years stopping smoking?"

Each question represents a different target population contrast (*consistency*)

Each may have different set of relevant confounders (*exchangeability*)

Some may not be estimable from available data (*non-positivity*)

What makes the causal question *longitudinal*?

What makes the causal question *longitudinal*?

Longitudinal data = repeated observations over time

What makes the causal question *longitudinal*?

Longitudinal data = repeated observations over time

1. Exposure changes over time

What makes the causal question *longitudinal*?

Longitudinal data = repeated observations over time

1. Exposure changes over time
2. Covariates change over time

What makes the causal question *longitudinal*?

Longitudinal data = repeated observations over time

1. Exposure changes over time
2. Covariates change over time
3. Outcome occurs more than once

What makes the causal question *longitudinal*?

Longitudinal data = repeated observations over time

1. Exposure changes over time
2. Covariates change over time
3. Outcome occurs more than once
4. Any/all of the above

What makes the causal question *longitudinal*?

Longitudinal data = repeated observations over time

1. Exposure changes over time
2. Covariates change over time
3. Outcome occurs more than once
4. Any/all of the above
5. AND you believe such variations are important to your causal question!

Longitudinal example:

"What is the effect of smoking on CVD risk, *independent* of its effect on body weight?"

Longitudinal example:

"What is the effect of smoking on CVD risk, *independent* of its effect on body weight?"

Does "independent" mean...

- Intervening to fix everyone to the same BMI?
- Intervening to prevent anyone from gaining or losing weight after treatment?
- In another (magical?) world where smoking does not affect BMI?
- Body weight doesn't affect total smoking duration / intensity?

Longitudinal example:

"What is the effect of smoking on CVD risk, *independent* of its effect on body weight?"

Does "independent" mean...

- Intervening to fix everyone to the same BMI?
- Intervening to prevent anyone from gaining or losing weight after treatment?
- In another (magical?) world where smoking does not affect BMI?
- Body weight doesn't affect total smoking duration / intensity?

How these questions are answered → target estimand

Longitudinal example:

"What is the effect of smoking on CVD risk, *independent* of its effect on body weight?"

Does "independent" mean...

- Intervening to fix everyone to the same BMI?
- Intervening to prevent anyone from gaining or losing weight after treatment?
- In another (magical?) world where smoking does not affect BMI?
- Body weight doesn't affect total smoking duration / intensity?

How these questions are answered → target estimand AND determines whether it can be estimated with given data

Longitudinal example:

"What is the effect of smoking on CVD risk, *independent* of its effect on body weight?"

Does "independent" mean...

- Intervening to fix everyone to the same BMI?
- Intervening to prevent anyone from gaining or losing weight after treatment?
- In another (magical?) world where smoking does not affect BMI?
- Body weight doesn't affect total smoking duration / intensity?

How these questions are answered → target estimand AND determines whether it can be estimated with given data

Draw the DAG!

What makes a longitudinal question *complex*?



Not complex:

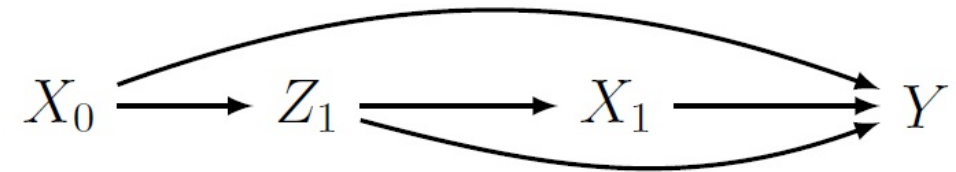
Individual and joint effects can be estimated by conventional methods.

What makes a longitudinal question *complex*?



Not complex:

Individual and joint effects can be estimated by conventional methods.



Complex:

Exposure at $t = 1$ is confounded by a consequence of X_0 (Z_1 is a time-varying confounder; exposure-covariate feedback).

Complex longitudinal examples:

Complex longitudinal examples:

"What is the effect of taking anti-retroviral therapy (ART) on CD4 counts, taking into account the effect of changes in HIV viral load on subsequent treatment?"

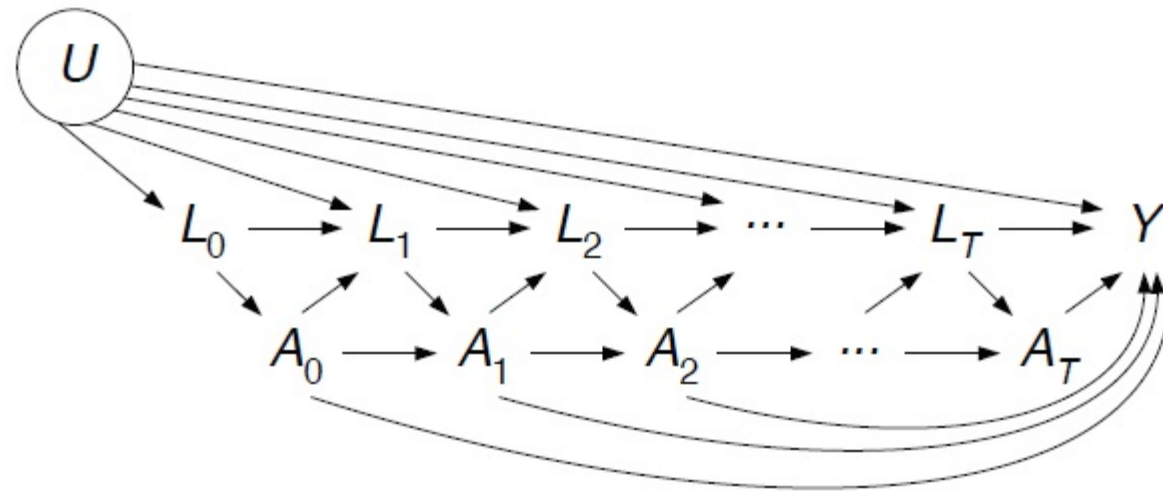
Complex longitudinal examples:

"What is the effect of taking anti-retroviral therapy (ART) on CD4 counts, taking into account the effect of changes in HIV viral load on subsequent treatment?"

"What is the per-protocol effect in a trial of aspirin intake on pregnancy loss, taking into account side effects, non-compliance, and study withdrawal?"

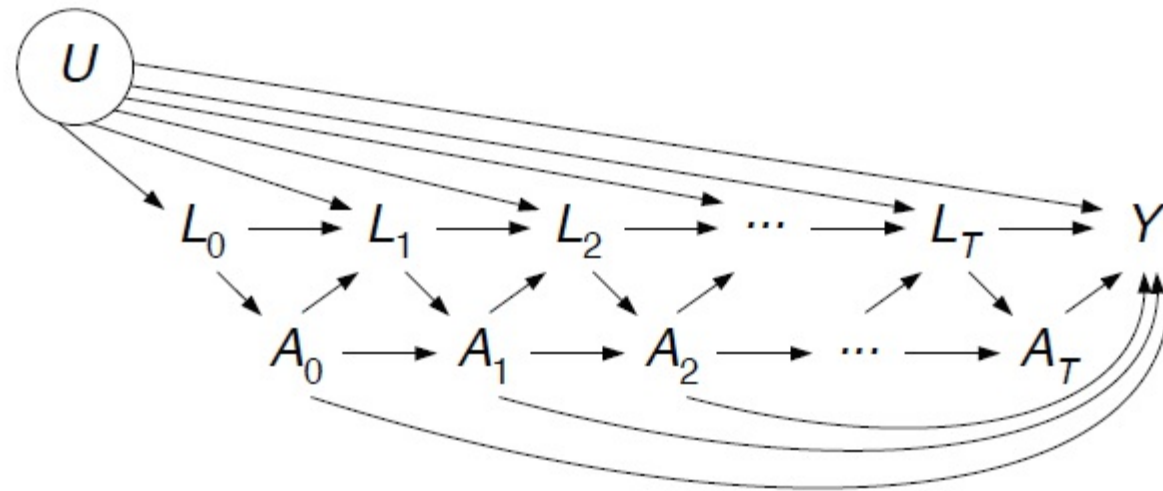
Estimating effects with time-varying treatment and confounding

Estimating effects with time-varying treatment and confounding



(Daniel, et. al. *The Stata Journal* 2011.)

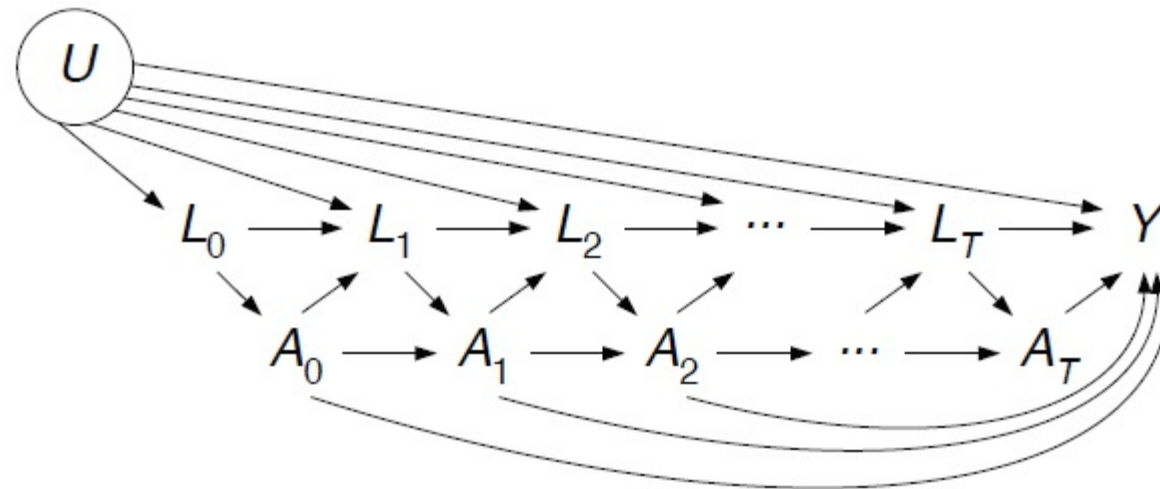
Estimating effects with time-varying treatment and confounding



(Daniel, *et. al. The Stata Journal* 2011.)

How would you estimate the causal of effect of A_0 ?

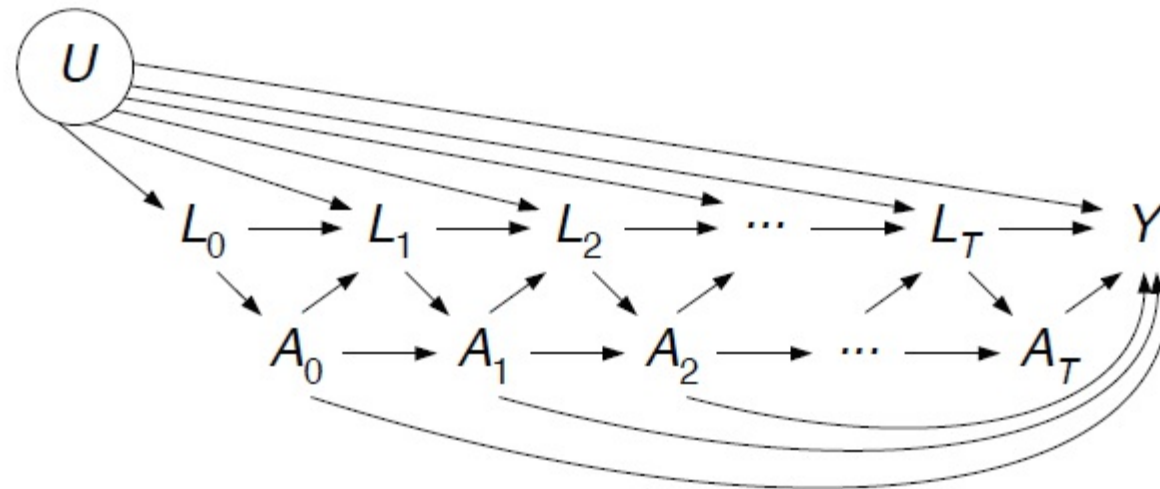
Estimating effects with time-varying treatment and confounding



(Daniel, *et. al. The Stata Journal* 2011.)

How would you estimate the causal effect of A_0 ? How about A_1 ?

Estimating effects with time-varying treatment and confounding



(Daniel, et. al. *The Stata Journal* 2011.)

How would you estimate the causal effect of A_0 ? How about A_1 ? And overall A ?

Extension to longitudinal data

For time-varying confounders L and binary exposure A , continuous outcome Y , also consider:

$$L_T = l_t \dots, l_1, l_0 \text{ and } A_T = a_t \dots, a_1, a_0 \text{ where } A_1\{1,1\} = a_1 = 1, a_0 = 1$$

Extension to longitudinal data

For time-varying confounders L and binary exposure A , continuous outcome Y , also consider:

$$L_T = l_t \dots, l_1, l_0 \text{ and } A_T = a_t \dots, a_1, a_0 \text{ where } A_1 \{1,1\} = a_1 = 1, a_0 = 1$$

Sequential exchangeability:

$$Y^{a_1} \perp A_1 = a_1 | L_1, A_0 = a_0$$

Extension to longitudinal data

For time-varying confounders L and binary exposure A , continuous outcome Y , also consider:

$$L_T = l_t \dots, l_1, l_0 \text{ and } A_T = a_t \dots, a_1, a_0 \text{ where } A_1 \{1,1\} = a_1 = 1, a_0 = 1$$

Sequential exchangeability:

$$Y^{a_1} \perp A_1 = a_1 | L_1, A_0 = a_0$$

Positivity:

$$0 < Pr(A_1 = a_1 | A_0 = a_0, L_1 = l_1) < 1$$

Extension to longitudinal data

For time-varying confounders L and binary exposure A , continuous outcome Y , also consider:

$$L_T = l_t \dots, l_1, l_0 \text{ and } A_T = a_t \dots, a_1, a_0 \text{ where } A_1 \{1,1\} = a_1 = 1, a_0 = 1$$

Sequential exchangeability:

$$Y^{a_1} \perp A_1 = a_1 | L_1, A_0 = a_0$$

Positivity:

$$0 < Pr(A_1 = a_1 | A_0 = a_0, L_1 = l_1) < 1$$

Potential quantities of interest:

$$E(Y | A_1 = a_1, A_0, L_1) = E(Y_1^a | A_0, L_1)$$

Extension to longitudinal data

For time-varying confounders L and binary exposure A , continuous outcome Y , also consider:

$$L_T = l_t \dots, l_1, l_0 \text{ and } A_T = a_t \dots, a_1, a_0 \text{ where } A_1 \{1,1\} = a_1 = 1, a_0 = 1$$

Sequential exchangeability:

$$Y^{a_1} \perp A_1 = a_1 | L_1, A_0 = a_0$$

Positivity:

$$0 < Pr(A_1 = a_1 | A_0 = a_0, L_1 = l_1) < 1$$

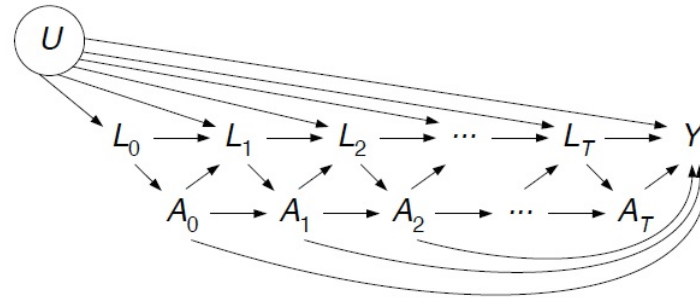
Potential quantities of interest:

$$E(Y | A_1 = a_1, A_0, L_1) = E(Y_1^a | A_0, L_1)$$

or,

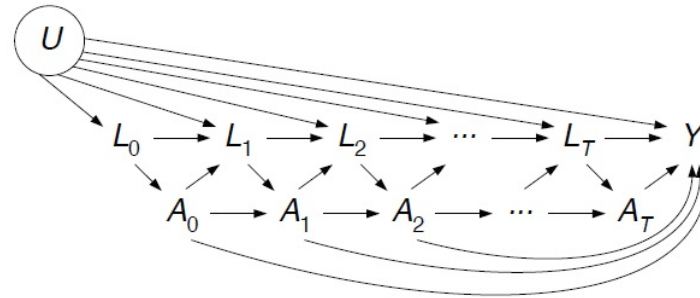
$$E(Y | A_1 = 1, A_0 = 1, L_1) = E(Y^{a_1} | L_1)$$

BUT!



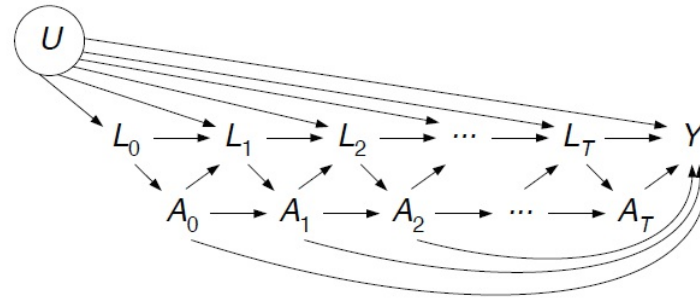
—
To estimate joint effect of continuous treatment $A_1 = \{1,1\}$; confounding by L_1

BUT!



—
To estimate joint effect of continuous treatment $A_1 = \{1,1\}$; confounding by L_1
e.g. health deteriorates; \uparrow treatment; poorer *observed* outcomes when $A_1 = 1$

BUT!

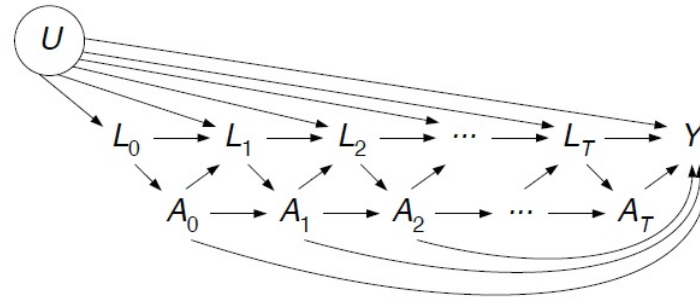


To estimate joint effect of continuous treatment $A_1 = \{1,1\}$; confounding by L_1

e.g. health deteriorates; \uparrow treatment; poorer *observed* outcomes when $A_1 = 1$

HOWEVER, standard regression on L_1 induces biasing pathway between A_0 and U and eliminates part of A_0 effect

BUT!



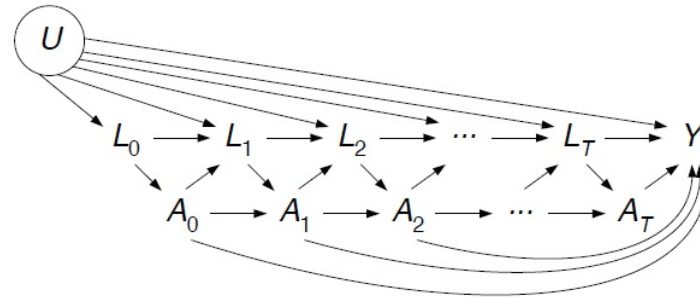
To estimate joint effect of continuous treatment $A_1 = \{1,1\}$; confounding by L_1

e.g. health deteriorates; \uparrow treatment; poorer *observed* outcomes when $A_1 = 1$

HOWEVER, standard regression on L_1 induces biasing pathway between A_0 and U and eliminates part of A_0 effect

"Adjusting for the future" \rightarrow biased post-treatment strata (collider stratification bias)

BUT!



To estimate joint effect of continuous treatment $A_1 = \{1,1\}$; confounding by L_1

e.g. health deteriorates; \uparrow treatment; poorer *observed* outcomes when $A_1 = 1$

HOWEVER, standard regression on L_1 induces biasing pathway between A_0 and U and eliminates part of A_0 effect

"Adjusting for the future" \rightarrow biased post-treatment strata (collider stratification bias)

Solution?

IMPLEMENTING THE G-FORMULA

Effect estimation under exposure-confounder feedback

True causal effect of A_0 and $A_1 = 0$;

$Pr(L_1 | A_0 = 1) = 0.5, 0.75$ untreated; $Pr(A_1 | L_1 = 1) = 0.8, 0.4$ otherwise

$$E[Y^{a_0=1}] = 76(48/160) + 76(32/160) + 44(16/160) + 44(64/160) = 60$$

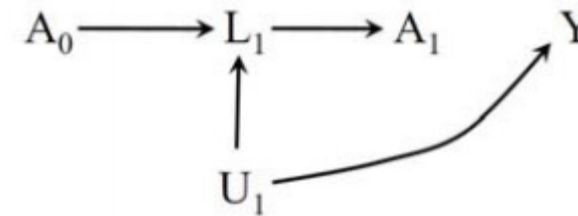
$$E[Y^{a_0=0}] = 84(24/160) + 84(16/160) + 52(24/160) + 52(96/160) = 60$$

$$E[Y^{a_0=1}] - E[Y^{a_0=0}] = 0$$

Also,

$$E[Y^{a_1=1} | A_0, L_1] - E[Y^{a_1=0} | A_0, L_1] = 0$$

HOWEVER, joint effect...



N	A_0	L_1	A_1	Mean Y
2400	0	0	0	84
1600	0	0	1	84
2400	0	1	0	52
9600	0	1	1	52
4800	1	0	0	76
3200	1	0	1	76
1600	1	1	0	44
6400	1	1	1	44

Hernán & Robins (2019)

Effect under continuous treatment

Unadjusted (for L_1):

Effect under continuous treatment

Unadjusted (for L_1):

$$E[Y|A_0 = 1, A_1 = 1] = 76(32/96) + 44(64/96) = 54.7 \quad E[Y|A_0 = 0, A_1 = 0] = 84(24/48) + 52(24/48) = 68$$

$$E[Y|A_0 = 1, A_1 = 1] - E[Y|A_0 = 0, A_1 = 0] = \mathbf{-13.3}$$

Effect under continuous treatment

Unadjusted (for L_1):

$$E[Y|A_0 = 1, A_1 = 1] = 76(32/96) + 44(64/96) = 54.7 \quad E[Y|A_0 = 0, A_1 = 0] = 84(24/48) + 52(24/48) = 68$$

$$E[Y|A_0 = 1, A_1 = 1] - E[Y|A_0 = 0, A_1 = 0] = \mathbf{-13.3}$$

Adjust for L_1 by stratification:

Effect under continuous treatment

Unadjusted (for L_1):

$$E[Y|A_0 = 1, A_1 = 1] = 76(32/96) + 44(64/96) = 54.7 \quad E[Y|A_0 = 0, A_1 = 0] = 84(24/48) + 52(24/48) = 68$$

$$E[Y|A_0 = 1, A_1 = 1] - E[Y|A_0 = 0, A_1 = 0] = \mathbf{-13.3}$$

Adjust for L_1 by stratification:

Within $L_1 = 0$:

$$E[Y|A_0 = 1, A_1 = 1, L_1 = 0] - E[Y|A_0 = 0, A_1 = 0, L_1 = 0] = 76 - 84 = \mathbf{-8}$$

(Subjects where $A_0 = 0, L_1 = 0$ healthier than a random draw of subjects)

Within $L_1 = 1$:

$$E[Y|A_0 = 1, A_1 = 1, L_1 = 1] - E[Y|A_0 = 0, A_1 = 0, L_1 = 1] = 44 - 52 = \mathbf{-8}$$

(Subjects where $A_0 = 1, L_1 = 1$ sicker than a random draw of subjects)

Stratification is the issue!

Standardization to estimate counterfactual means

For point exposure and time-fixed confounders:

$$E(Y) = \sum_A \sum_L E(Y | A, L)P(A | L)P(L)$$

Standardization to estimate counterfactual means

For point exposure and time-fixed confounders:

$$E(Y) = \sum_A \sum_L E(Y | A, L)P(A | L)P(L)$$

If our model is correct, and causal assumptions hold:

$$E[Y^a] = \sum_l E[Y | A = a, L = l] * f(l), \text{ where } f(l) = Pr(L = l)$$

i.e. weighted conditional mean of Y standardized to observed covariate distribution

Standardization to estimate counterfactual means

For point exposure and time-fixed confounders:

$$E(Y) = \sum_A \sum_L E(Y | A, L)P(A | L)P(L)$$

If our model is correct, and causal assumptions hold:

$$E[Y^a] = \sum_l E[Y | A = a, L = l] * f(l), \text{ where } f(l) = Pr(L = l)$$

i.e. weighted conditional mean of Y standardized to observed covariate distribution

Extention to time-varying exposures and confounders:

$$E[Y^{a_0, a_1}] = \sum_{l_1} E[Y | A_0 = a_0, A_1 = a_1, L_1 = l_1] * f(l_1 | a_0), \text{ where } f(l_1 | a_0) = Pr(L_1 = l | A_0 = a)$$

i.e., weighted conditional mean of Y standardized to exposure and covariate history as they would have been observed under sequential randomization

Applied to the example:

N	A_0	L_1	A_1	Mean Y
2400	0	0	0	84
1600	0	0	1	84
2400	0	1	0	52
9600	0	1	1	52
4800	1	0	0	76
3200	1	0	1	76
1600	1	1	0	44
6400	1	1	1	44

Applied to the example:

N	A_0	L_1	A_1	Mean Y
2400	0	0	0	84
1600	0	0	1	84
2400	0	1	0	52
9600	0	1	1	52
4800	1	0	0	76
3200	1	0	1	76
1600	1	1	0	44
6400	1	1	1	44

$$E[Y^{a_0=1, a_1=1}] =$$

$$E[Y|A_0 = 1, A_1 = 1, L_1 = 0] * Pr[L_1 = 0 | A_0 = 1] +$$
$$E[Y|A_0 = 1, A_1 = 1, L_1 = 1] * Pr[L_1 = 1 | A_0 = 1]$$

$$= 76 (80/160) + 44 (80/160) = \mathbf{60}$$

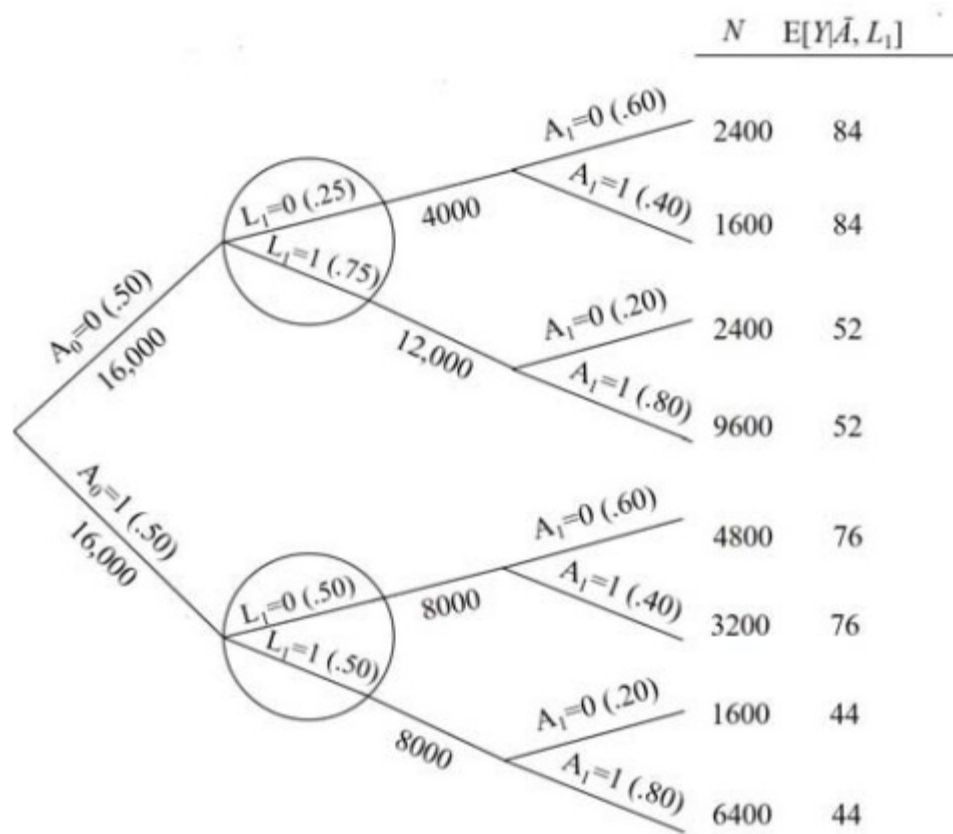
$$E[Y^{a_0=0, a_1=0}] =$$

$$E[Y|A_0 = 0, A_1 = 0, L_1 = 0] * Pr[L_1 = 0 | A_0 = 0] +$$
$$E[Y|A_0 = 0, A_1 = 0, L_1 = 1] * Pr[L_1 = 1 | A_0 = 0]$$

$$= 84 (40/160) + 52 (120/160) = \mathbf{60}$$

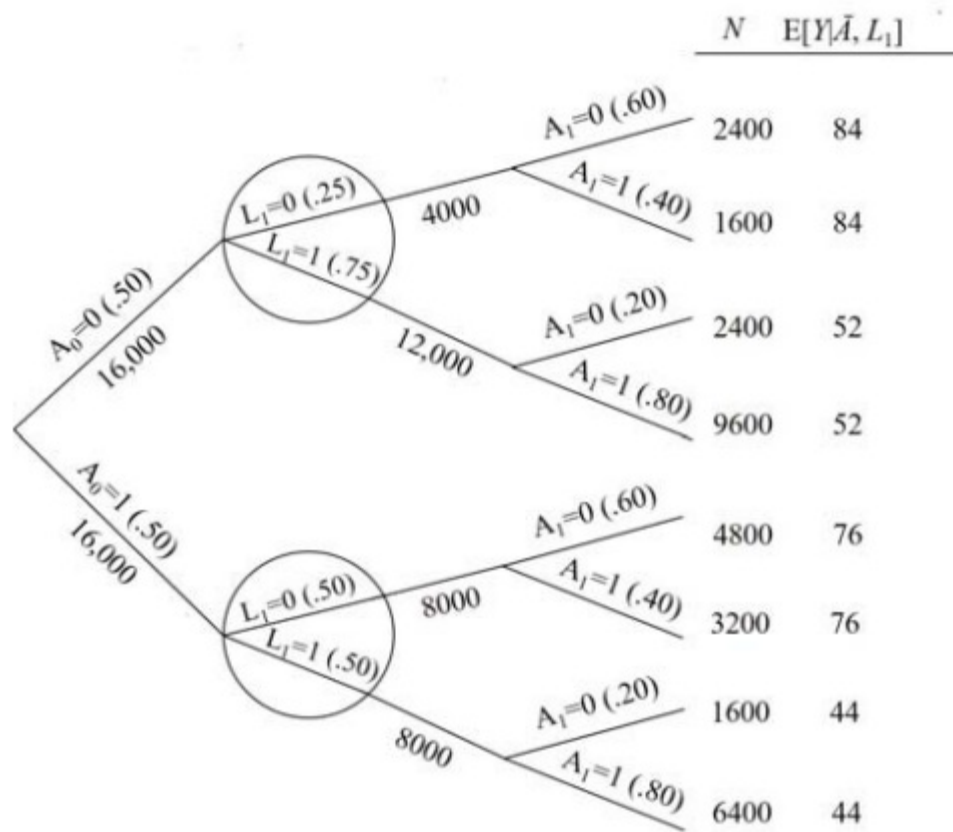
$$E[Y^{a_0=1, a_1=1}] - E[Y^{a_0=0, a_1=0}] = \mathbf{0}$$

Considered as a simulation:

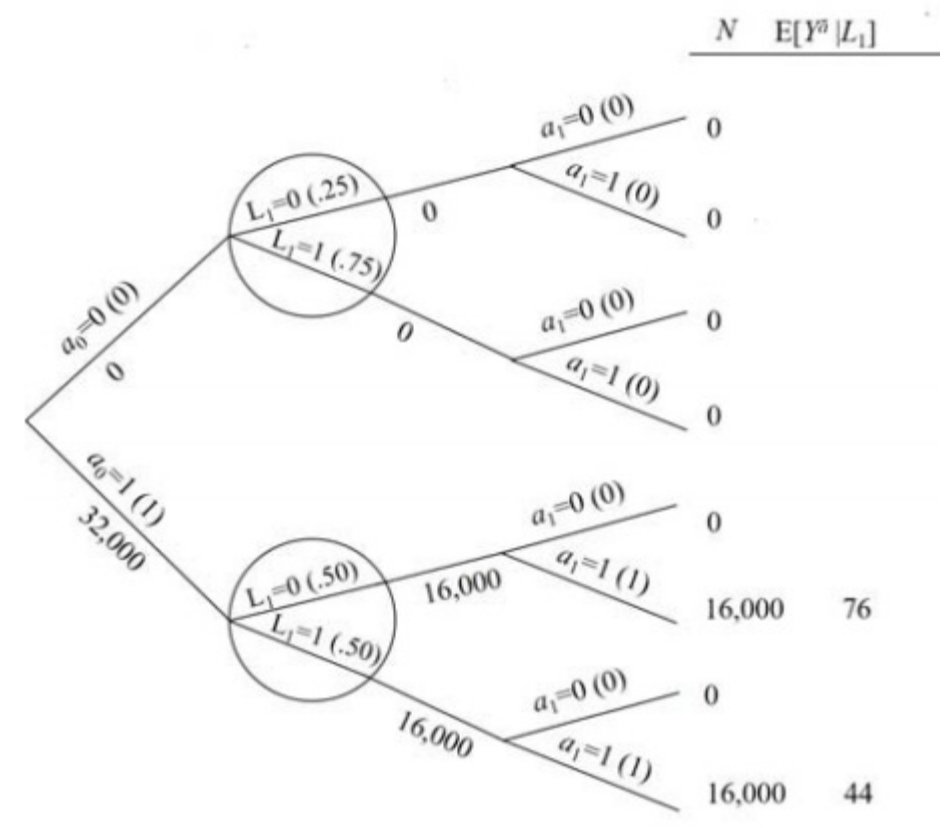


Observed Data

Considered as a simulation:



Observed Data



Simulated

Parametric g-formula

$$\sum_{\bar{l}} E[Y | \bar{A} = \bar{a}, \bar{L} = \bar{l}] * \prod_{k=0}^K f(l_k | \bar{a}_{k-1}, \bar{l}_{k-1})$$


where, e.g. conditional means estimated by linear, and distribution of discrete covariates by logistic, regressions

"plug-in g-formula" or "parametric g-formula"

can be generalized further to densities

Estimation Steps In Practice

- Draw a DAG representing causal structure
- Model temporal / structural relationships (Q-model) implied by DAG
- Sample baseline values simulate POs using model
- Repeat setting exposures to levels reflecting desired contrast
- Bootstrap to obtain SE



INTERLUDE:
RStudio Cloud Setup

1. bit.ly/CI_tutorial_cloud





Email

Password

Log in

— or —

 Log in with Google

 Log in with GitHub

Forgot your password? [Sign up](#)

By clicking log in, you agree to the RStudio.cloud terms of use.

1. Create Account

R Studio

— or —

1. Create Account

R Studio

Welcome Jon!

To get started, please provide a name for your account:

Account names can contain letters, numbers, and hyphens, but can't start or end with a hyphen.

Create Account

2. Navigate to project

☰ NUS-IMS CI Tutorial

Projects

Members

Info



Jon Huang

Welcome to NUS-IMS CI Tutorial

12 February 2019

"Introduction to Causal Inference in Complex Longitudinal Settings" Session

Jonathan Y Huang, MPH, PhD

NUS-IMS "Tutorial on Personalized Medicine, Treatment Regimes, Reinforcement Learning, and Causal Inference"

Lecture notes and practice code and data will be accessible here.


Navigate to "bit.ly/CI_tutorial_cloud" to access!

2. Navigate to project

☰ NUS-IMS CI Tutorial **Projects** Members Info

All Projects New Project ▾

START Longitudinal Causal Inference Demo

 Jon Huang

Presentation notes, code, and example data.

Created Jan 15, 2019 2:00 PM

3. Load files

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The top right corner shows the R version as 3.5.2 and the user as Jon Huang. The main window is divided into several panes:

- Console:** Displays the R startup message: "R version 3.5.2 (2018-12-20) -- 'Eggshell Igloo' Copyright (C) 2018 The R Foundation for Statistical Computing Platform: x86_64-pc-linux-gnu (64-bit)". It also includes the standard R license text and instructions for help and quitting.
- Environment:** Shows "Global Environment" and "Environment is empty".
- Files:** A file browser window showing the contents of the "project" directory. The files listed are:

Name	Size	Modified
..		
.Rhistory	0 B	Feb 7, 2019, 1:06 PM
20190212-Longitudinal_CI_Tutorial_X.html	12.6 KB	Feb 6, 2019, 10:34 AM
20190212-Longitudinal_CI_Tutorial_X.Rmd	10.9 KB	Feb 6, 2019, 10:34 AM
20190212-Longitudinal_CI_Tutorial_X_files		
background-blur-darker.jpg	1.5 MB	Feb 6, 2019, 10:32 AM
complex.png	10 KB	Feb 6, 2019, 10:32 AM
IMSAbstractTemplate-JYH.pdf	47.1 KB	Feb 6, 2019, 10:30 AM
not_complex.png	7.6 KB	Feb 6, 2019, 10:32 AM
pearl_cause.png	16.5 KB	Feb 6, 2019, 10:32 AM
project.Rproj	205 B	Feb 7, 2019, 1:06 PM
SICS_logo_small.jpg	6.5 KB	Feb 6, 2019, 10:36 AM

Two ways to follow along:

☰ NUS-IMS CI Tutorial / Longitudinal Causal Inference Demo ⚙️ ⋮ JH Jon Huang

The screenshot shows the RStudio interface with the following components:

- Source Editor:** Contains R code for loading NHANES data. Lines 642-643 show a chunk of code starting with ````{r, message=F}`. A red box highlights the `Run Current Chunk` button next to line 643.
- Environment:** Shows 'Global Environment' and 'Environment is empty'.
- Files:** Shows a file browser with a list of files including `20190212-Longitudinal_CI_Tutorial_X.Rmd`, `project.Rproj`, and various image files.
- Console:** Shows the output of the R Markdown rendering process, including the command `ackslash+smart --output 20190212-Longitudinal_CI_Tutorial_X.html` and the output `Output created: 20190212-Longitudinal_CI_Tutorial_X.html`.

1. Easy Way: Scroll to line 700. Click "Play" buttons.

Two ways to follow along:

☰ NUS-IMS CI Tutorial / Longitudinal Causal Inference Demo ⚙️ ⋮ JH Jon Huang

The screenshot shows the RStudio interface with the following components:

- Editor:** Contains R code for loading and processing NHANES data. A tooltip "Run Current Chunk" is visible over the code.
- Environment:** Shows "Global Environment" and "Environment is empty".
- Files:** A file explorer showing a list of files in a "project" directory, including R scripts, images, and HTML files.
- Console:** Shows the output of the R Markdown rendering process, including the path to the generated HTML file. The output line "Output created: 20190212-Longitudinal CI Tutorial X.html" is highlighted with a red box.

2. For fast typers: Read code. Duplicate in Console.

WORKED EXAMPLES

CASE STUDY #1

Smoking cessation on weight gain (NHEFS)

Quitting smoking on body weight change (NHEFS)

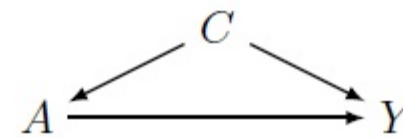
Question of interest: What is the effect of quitting smoking (A) in a population of adult U.S. smokers (1971-1982) on risk of weight gain (Y)?

Quitting smoking on body weight change (NHEFS)

Question of interest: What is the effect of quitting smoking (A) in a population of adult U.S. smokers (1971-1982) on risk of weight gain (Y)?

Simplifying assumptions:

- ONLY factors confounding (C) the relationship (*i.e.* biasing distribution POs): *sex, age, race, income, marital status, schooling, physical activity, heart failure, high BP, asthma, diabetes, chronic bronchitis / emphysema*
- No informative missingness, and all C measured prior to A , time-invariant, and dichotomous such that causal



relationships faithfully represented by:

IMPORT DATA

Load NHEFS Data (NHANES 1 Epidemiologic Follow-Up Study)

```
aa <- read_csv(paste0(here(), "/nhefs.csv"))
# original sample size
nrow(aa)
```

```
## [1] 1746
```

Simplify data and create a weight gain indicator `delta` as outcome

```
a <- aa %>%
  select(seqn, qsmk, smkintensity82_71, smokeintensity, active, exercise, wt82_71,
         sbp, dbp, hbp, hf, ht, hbpm, sex, age, hf, race, income, marital, school,
         asthma, bronch, diabetes) %>%
  mutate(hbp_71 = hbp) %>% na.omit()

a$delta <- as.numeric(a$wt82_71 > 0)

a <- a %>% select(delta, qsmk, sex, age, race, income, marital, school, active,
                hf, hbpm, asthma, bronch, smokeintensity, exercise,
                diabetes, hbp_71)
```

WARNING: Complete case not recommended in general practice

To reduce residual confounding, continuous or interval measures should not be casually dichotomized. However, we will do so only for the purposes of this exercise:

```
a$smokeintensity <- as.numeric(a$smokeintensity>median(a$smokeintensity))
a$age <- as.numeric(a$age>median(a$age))
a$exercise <- as.numeric(a$exercise>0)
a$income <- as.numeric(a$income>median(a$income))
a$marital <- as.numeric(a$marital>median(a$marital))
a$school <- as.numeric(a$school>median(a$school))
a$active <- as.numeric(a$active>0)
a$hbpmed <- as.numeric(a$hbpmed>0)
a$smokeintensity <- as.numeric(a$smokeintensity>median(a$smokeintensity))
a$exercise <- as.numeric(a$exercise>0)
a$diabetes <- as.numeric(a$diabetes>0)
a$hbp_71 <- as.numeric(a$hbp_71>0)
```

DESCRIBE DATA

Describe population and covariates, e.g.:

```
aa %>% ggplot() + geom_histogram(aes(age))
```

Describe basic exposure-outcome relationship:

```
a %>% group_by(qsmk) %>%  
  summarize(gain = sum(delta),  
            no_gain = sum(!delta),  
            mu = round((gain)/(gain + no_gain), 2))
```

```
## # A tibble: 2 x 4  
##   qsmk gain no_gain mu  
##   <dbl> <dbl> <int> <dbl>  
## 1     0   724   390 0.65  
## 2     1   268    94 0.74
```

TRANSLATE DAG TO MODELS

Causal quantity of interest:

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

DAG gives us:

Variable	Modelled by
Y	$E(Y A, C) = \alpha_0 + \alpha_1 A + \alpha_2 C$
A	$P(A = 1 C) = \text{expit}(\beta_0 + \beta_1 C)$

* $\text{expit}(a) = 1/[1+\exp(-a)]$

Causal assumptions allow us to transform:

$$E(Y) = \sum_A \sum_C E(Y | A, C)P(A | C)P(C)$$

into:

$$E(Y^a) = \sum_C E(Y | A = a, C)P(C)$$

when assigning a value $A = a$.

FIT MODELS

Fit models for A and Y conditional on covariates:

```
model_A <- glm(qsmk ~ sex+age+race+income+marital+school+active+hf+hbpmed+asthma+  
              bronch+smokeintensity+exercise+diabetes+hbp_71,  
              data=a,family=binomial("logit"))  
summary(model_A)  
  
model_Y <- glm(delta ~ qsmk+sex+age+race+income+marital+school+active+hf+hbpmed+  
                  asthma+bronch+smokeintensity+exercise+diabetes+hbp_71,  
                  data=a,family=binomial("logit"))  
summary(model_Y)
```

Note: Can survey some of the estimates for face validity. But should avoid over-interpretation, due to the joint conditioning ("Table 2 Fallacy").

RESAMPLE FOR SIMULATION / DIAGNOSTICS

Monte Carlo resampling (with replacement). Usually, N * some factor (1 - 1000).

```
# resample data
index <- sample(1:nrow(a),size=1e4,replace=T)
MC <- a[index,]
nrow(MC)
```

```
## [1] 10000
```

```
MC$qsmk<-NULL
# predict exposure
pA <- predict(model_A,newdata=MC,type="response")
```

The variable pA is the predicted exposure under model.

Predicted prevalence of exposure: 0.2452247 versus observed prevalence: 0.2452575.

SIMULATE DISCRETE EXPOSURE STATES

Transform predicted variable ($0 < p_A < 1$) to binary exposure status (q_A) by comparing to a uniform random value:

```
u <- runif(1e4)
qA <- as.numeric(pA>u)
head(qA)
```

```
## [1] 0 0 0 0 1 1
```

```
mean(qA)
```

```
## [1] 0.2501
```

```
mean(a$qsmk)
```

```
## [1] 0.2452575
```

SIMULATE OUTCOMES UNDER BASELINE CONDITIONS

Taking new simulated exposure qA status, we simulate the outcome under observed conditions:

```
pY <- predict(model_Y, newdata=data.frame(MC, qsmk=qA), type="response")
```

```
mean(pY)
```

```
## [1] 0.6733915
```

```
mean(a$delta)
```

```
## [1] 0.6720867
```

Note: While this is neither necessary nor sufficient to demonstrate a valid model, deviations from the observed outcome can indicate modelling or data errors.

EFFECT OF SMOKING CESSATION

We can now estimate the effect of smoking on weight gain risk by:

1. Setting $q_{smk} = 1$, predict the probability of weight gain, taking the mean
2. Repeat with q_{smk} set to 0
3. Taking the difference (and ratio)

```
pY_1 <- predict(model_Y, newdata=data.frame(MC, qsmk=1), type="response")
mY_1 <- mean(pY_1)

pY_0 <- predict(model_Y, newdata=data.frame(MC, qsmk=0), type="response")
mY_0 <- mean(pY_0)

(RD <- round((mY_1 - mY_0)*100, 2))
```

```
## [1] 11.39
```

```
(RR <- round(mY_1 / mY_0, 2))
```

```
## [1] 1.18
```

INTERPRETATION

Can we interpret $RD = 11.39$ more cases of weight gain per 100 persons (or $RR = 1.18$ times risk) causally?

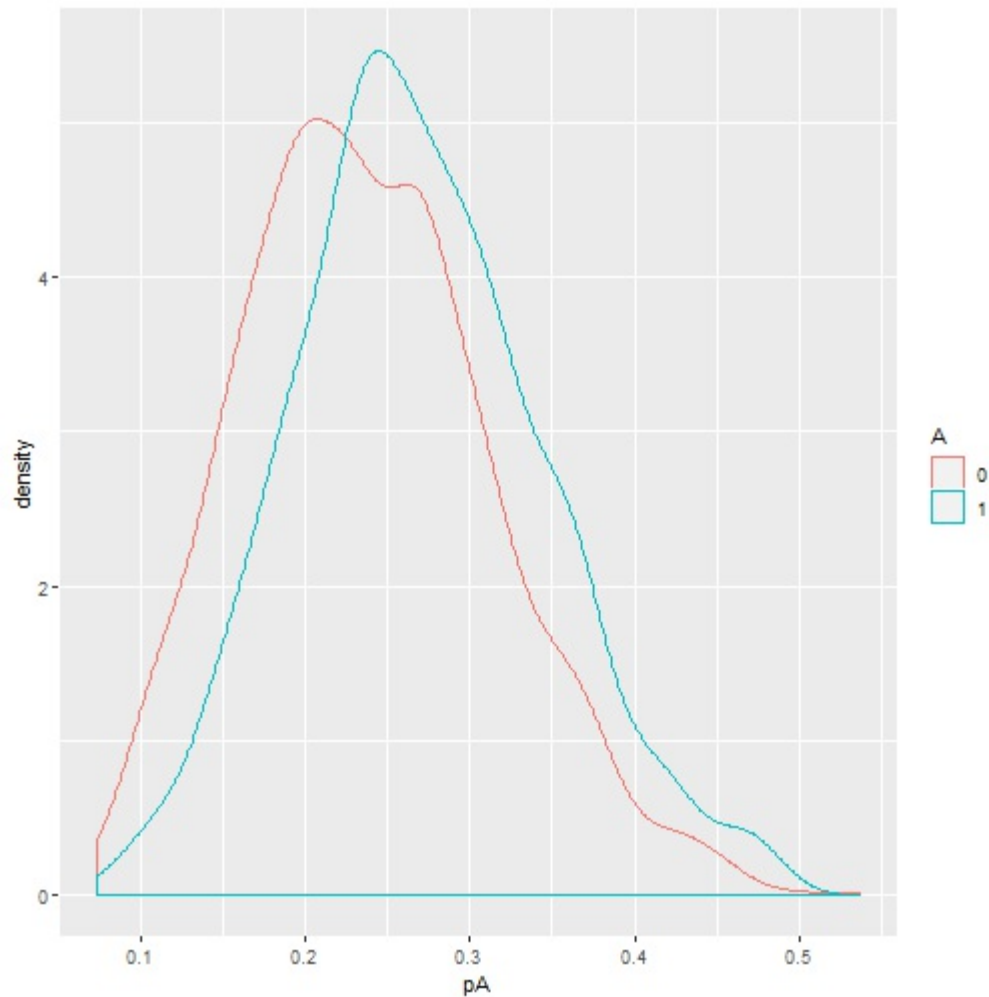
Evaluate causal assumptions:

- *Consistency*: Will different ways to quit smoking lead to different effects?
- Correct model specification and *exchangeability*: Did we include all the possible confounders? Residual confounding? Interaction?
- *Interference* (NEW): If a given person quit smoking, will it affect the outcome of another person in the study?
- *Positivity*: Are there are exposed and unexposed individuals in each confounder level?

POSITIVITY CHECK

```
prop <- model_A$fitted.values
propD <- data.frame(
  A=as.factor(a$qsmk),pA=prop)

ggplot(propD, aes(x=pA,color=A)) +
  geom_density()
```



Propensity score overlap:

Because of the number of confounding variables, we cannot use 2×2 tables.

Instead we can examine predicted probability of exposure.

Reasonable overlap between the two groups and likely no rare confounder combinations.

CI BY BOOTSTRAP

Resample (with replacement) original data re-fitting outcome model; MC to compute contrast.

Repeat x 100 saving contrasts.

```
res <- NULL
for(i in 1:100){
  index <- sample(1:nrow(a),nrow(a),replace=T)
  boot_dat <- a[index,]
  model_Y <- glm(delta ~ qsmk+sex+age+race+income+marital+school+active+hf+hbpmed+asthma+bronch+
                 smokeintensity+exercise+diabetes+hbp_71,data=boot_dat,family="logit")

  index <- sample(1:nrow(a),size=1e4,replace=T)
  MC <- boot_dat[index,]
  MC$qsmk<-NULL

  mY_1 <- mean(predict(model_Y,newdata=data.frame(MC,qsmk=1),type="response"))
  mY_0 <- mean(predict(model_Y,newdata=data.frame(MC,qsmk=0),type="response"))

  RD <- (mY_1 - mY_0)*100
  logRR <- log(mY_1 / mY_0)

  res <- rbind(res,cbind(RD,logRR))
}
```


CI BY BOOTSTRAP

Take take SD across bootstraps as SE; calculate Wald confidence limits:

```
head(res)
```

```
##           RD      logRR
## [1,] 11.348235 0.16168739
## [2,]  9.082945 0.13487830
## [3,]  5.734421 0.08472272
## [4,] 12.565192 0.17826077
## [5,] 10.325501 0.14892598
## [6,] 11.409032 0.16151644
```

```
res_sd <- apply(res,2,sd)
```

```
lclRD <- RD - 1.96*res_sd[1]
```

```
uclRD <- RD + 1.96*res_sd[1]
```

```
lclRR <- exp(log(RR) - 1.96*res_sd[2])
```

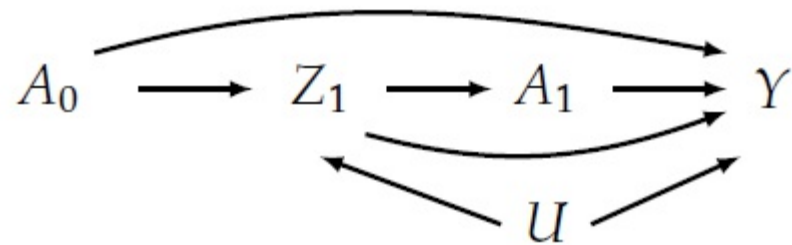
```
uclRR <- exp(log(RR) + 1.96*res_sd[2])
```

This bootstrap estimator yields 95% CIs of [2.51, 11.37] for the risk difference, and [1.11, 1.25] for the risk ratio.

CASE STUDY 2

ART treatment on CD4 count

ART treatment on CD4 count



Working DAG

Where:

A_0 = ART treatment at $t = 0$

Z_1 = HIV viral load >200 copies/mL just before $t = 1$

A_1 = ART treatment at $t = 1$

Y = CD4 count (cells / mm^3)

U = unmeasured health status

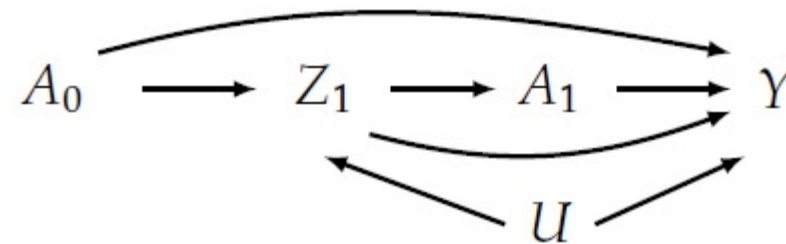
Observed data on $N = 100,000$ patients:

A_0	Z_1	A_1	Y	N
0	0	0	87.29	20,927
0	0	1	112.11	9,378
0	1	0	119.65	6,065
0	1	1	144.84	13,630
1	0	0	105.28	13,478
1	0	1	130.18	6,079
1	1	0	137.72	9,390
1	1	1	162.83	21,053

(Hypothetical cohort study)

Question and statistical model

What is the effect of $\bar{A}_1 = \{1, 1\}$ on CD4 count (Y)?



Variable	Model
Y	$E(Y \mid A_1, Z_1, A_0) = \alpha_0 + \alpha_1 A_1 + \alpha_2 Z_1 + \alpha_3 A_0$
A_1	$P(A_1 \mid Z_1) = \text{expit}(\beta_0 + \beta_1 Z_1)$
Z_1	$P(Z_1 \mid A_0) = \text{expit}(\gamma_0 + \gamma_1 A_0)$
A_0	$P(A_0) = \text{expit}(\theta_0)$

Statistical model to g-formula

$$E(Y) = \sum_{A_1} \sum_{Z_1} \sum_{A_0} E(Y \mid A_1, Z_1, A_0) Pr(A_1 \mid Z_1) P(Z_1 \mid A_0) Pr(A_0)$$

Under causal assumptions, we can:

- replace all A_0 and A_1 with a_0 and a_1 (respectively),
- remove models for A_0 and A_1 (because they are assigned),
- and equate observed conditional outcome with their PO

$$E(Y^{a_0, a_1}) = \sum_{Z_1} E(Y \mid A_1 = a_1, Z_1, A_0 = a_0) Pr(Z_1 \mid A_0 = a_0)$$

Then the objective is to estimate the conditional means and probabilities.

INPUT DATA

```
a0<-c(0,0,0,0,1,1,1,1)
z1<-c(0,0,1,1,0,0,1,1)
a1<-c(0,1,0,1,0,1,0,1)
y<-c(87.29,112.11,119.65,144.84,105.28,130.18,137.72,162.83)
N<-c(20927,9378,6065,13630,13478,6079,9390,21053)
#N<-c(209271,93779,60654,136293,134781,60789,93903,210530)
D<-NULL
for(i in 1:8){
  d<-data.frame(cbind(rep(a0[i],N[i]),rep(z1[i],N[i]),
                      rep(a1[i],N[i]),rep(y[i],N[i]))))
  D<-rbind(D,d)
}
names(D)<-c("a0","z1","a1","y")
```

Sample size: 100000

Initially treated: 0.5

Mean CD4 count: 125.0951071

We fit each of the statistical models implied by our DAG:

```
mA0<-glm(a0~1,data=D,family=binomial("logit"))
mZ1<-glm(z1~a0,data=D,family=binomial("logit"))
mA1<-glm(a1~z1,data=D,family=binomial("logit"))
mY<-glm(y~a1+z1+a0,data=D,family=gaussian("identity"))
```

Predict each variable in the correct order:

```
# A0 (50% as observed)
pA0<-predict(mA0,type="response")
# Z1 <- A0
pZ1<-predict(mZ1,newdata=data.frame(a0=pA0),type="response")
# A1 <- Z1
pA1<-predict(mA1,newdata=data.frame(z1=pZ1),type="response")
# Y <- A0, Z1, and A1
pY<-predict(mY,newdata=data.frame(a0=pA0,z1=pZ1,a1=pA1),type="response")
```

Predicted initially treated: 0.5

Simulate mean CD4 count under "natural course": 125.102381

Sample mean CD4 count: 125.0951071

ESTIMATE TARGET CONTRAST

Setting $A\{1,1\}$:

```
pZ_1<-predict(mZ1,newdata=data.frame(a0=1),type="response")  
pY_1<-predict(mY,newdata=data.frame(a0=1,z1=pZ_1,a1=1),type="response")  
mY_1<-mean(pY_1)
```

Mean CD4 count until always treat: 150.0640176

Setting $A\{0,0\}$:

```
pZ_0<-predict(mZ1,newdata=data.frame(a0=0),type="response")  
pY_0<-predict(mY,newdata=data.frame(a0=0,z1=pZ_0,a1=0),type="response")  
mY_0<-mean(pY_0)
```

Mean CD4 count until always treat: 100.0562062

Mean difference = 50 cells/mL (+25 cell/mL per time-point).

(Causal if assumptions hold to allow interpretation of g-formula.)

How about simple regression?

Effect of $A_1 | Z_1$ (true effect = +25 cells/ML):

```
round(coef(glm(y~a1+z1,data=D,family=gaussian("identity"))),1)
```

```
## (Intercept)          a1          z1  
##          94.3          25.0          36.4
```

Effect of A_0 (true effect = +25 cells/ML):

```
round(coef(glm(y~a0,data=D,family=gaussian("identity"))),1)
```

```
## (Intercept)          a0  
##          111.6          27.1
```

Overestimates the true effect by 2.1 cells/mL, because it includes both the direct effect ($A_0 \rightarrow Y$) **plus** indirect effect ($A_0 \rightarrow A_1 \rightarrow Y$).

Summing to get the total effect of $A\{1,1\}$ would give incorrect estimate!

CASE STUDY 3

Aspirin on live birth (EAGeR Trial)

Estimating per-protocol effects

ITT estimates effect of treatment randomization:

- may not be good estimate of drug efficacy with non-compliance

Standard regression approaches to per-protocol effect biased:

- non-compliance caused by treatment, related to POs
- stratification by compliance disrupts randomization
- may be temporary non-compliance

EAGeR Study

- *Effect of Aspirin in Gestation and Reproduction (EAGeR)*
- Multicenter, block randomized, double-blind RCT
- ASA versus placebo on live birth
- Women (18-40 y/o) with 1-2 pregnancy loss
- 81 mg ASA + folic acid / day (N = 614)
- folic acid only (N = 614)

Other parameters

- Allowed up to 6 cycles to conceive
- Treatment up to 36 weeks gestation
- Baseline: age, BMI, income, race, education, marital status, employment, study site
- Time-varying: compliance, vaginal bleeding, GI symptoms, TTP
- Outcomes: live birth, pregnancy loss, loss-to-follow-up

Question of interest:

What is the effect of continual treatment by ASA on live birth among target population if compliance was maintained through end of follow-up (by birth, pregnancy loss, or withdrawal)?

READ DATA AND DESCRIBE

```
aspirin<-read.table("aspirin2.txt",header=T,sep="\t")
```

```
aspirin %>% names()
```

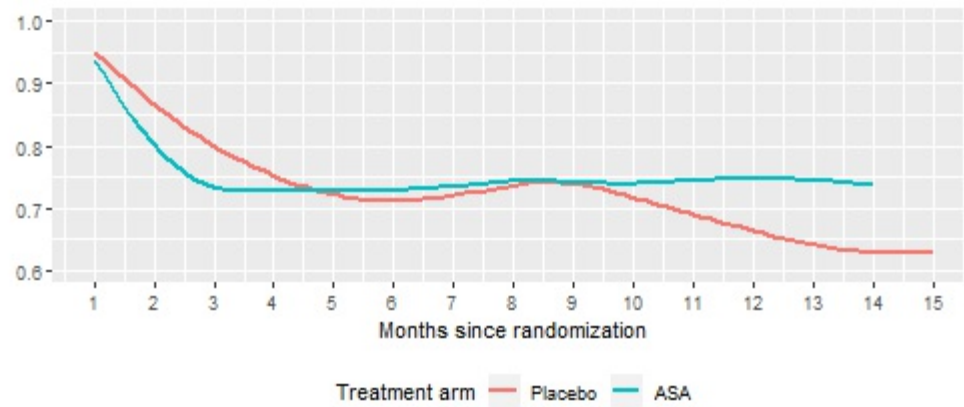
```
## [1] "id"          "study_month"  "eligibility"  "age"
## [5] "income"      "education"    "white"        "marital"
## [9] "employed"    "BMI"          "compliance"   "treatment"
## [13] "bleeding"    "gastro"       "conceived"    "efuwp"
## [17] "pregnancy_loss" "live_birth"   "last"         "site"
```

```
aspirin %>% filter(last == 1) %>% group_by(treatment) %>%
  summarize(Y_Birth = mean(live_birth), D_Loss = mean(pregnancy_loss),
            S_Censor = mean(efuwp))
```

```
## # A tibble: 2 x 4
##   treatment Y_Birth D_Loss S_Censor
##   <int>    <dbl> <dbl>    <dbl>
## 1         0  0.485  0.194  0.303
## 2         1  0.549  0.192  0.259
```

Compliance?

```
aspirin %>% mutate(tx = factor(treatment, levels = c(0,1), labels = c("Placebo", "ASA"))) %>%  
  ggplot() + geom_smooth(aes(x = study_month, y = compliance, color = tx), se = F) +  
  scale_y_continuous(limits = c(0,1)) + scale_x_continuous(breaks = c(1:15)) +  
  labs(y = "", x = "Months since randomization", color = "Treatment arm") + theme(legend.position = "bottom") + coo
```



FIT MODELS BASED ON ORDERING

Order	Variable	Description
8	Y	Live birth
7	D	pregnancy loss
6	S	No pregnancy
5	C	Withdrawal
4	Z	Conception
3	X	Compliance
2	N	GI Symptoms
1	B	Bleeding

FIT MODELS BASED ON ORDERING

Order	Variable	Description
8	Y	Live birth
7	D	pregnancy loss
6	S	No pregnancy
5	C	Withdrawal
4	Z	Conception
3	X	Compliance
2	N	GI Symptoms
1	B	Bleeding

e.g. Compliance X at each month k :

```
fitX<-glm(X~V1+V2+V3+V4+V5+V6+V7+
ns(V8,df=3)+V9+V10+V11+V12+V13+
ns(V14,df=3)+X1+X11+B+B1+B11+N+N1+
N11+Z+Z1+Z11+ns(jj,df=3),
family=binomial,data=boot,subset=R==k)
```

Modelled by natural cubic splines:

$V8$ = maternal age @ baseline

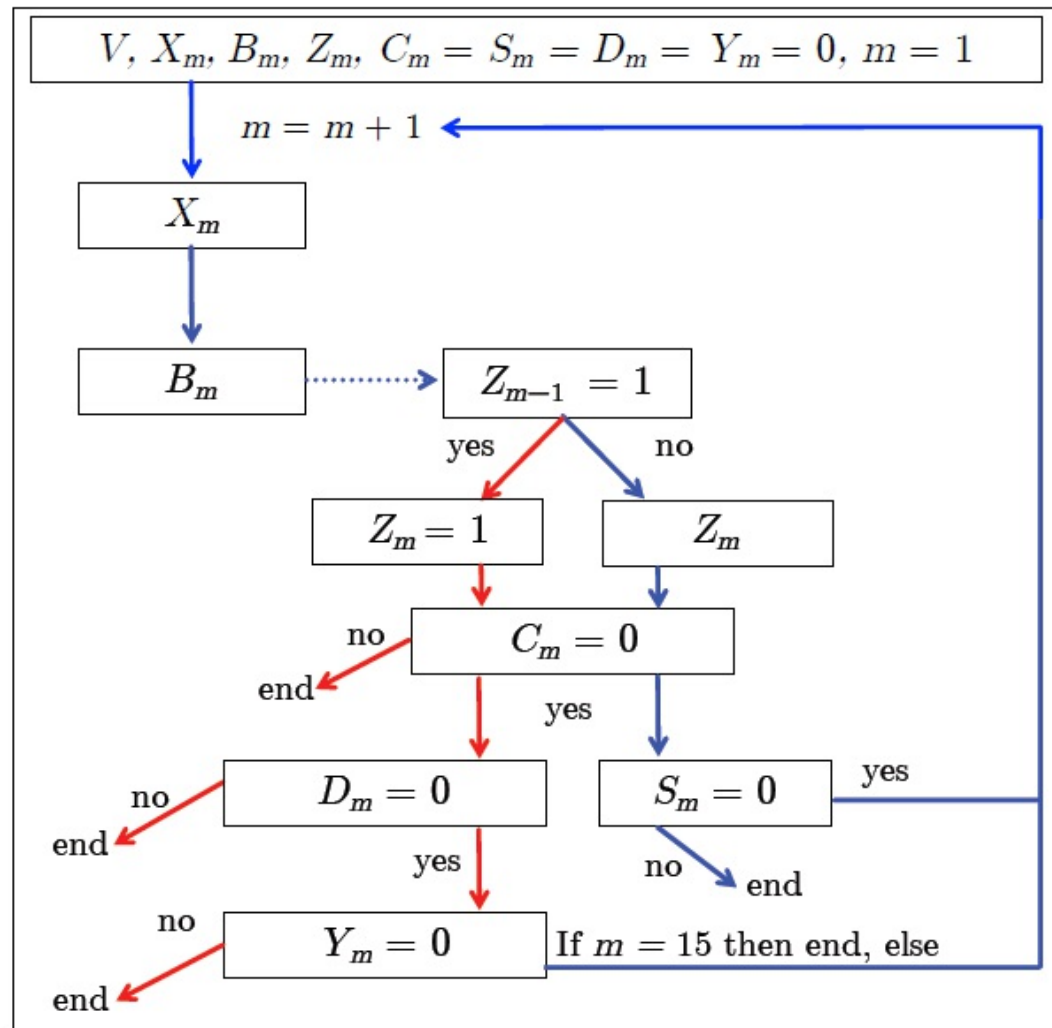
$V14$ = maternal BMI @ baseline

jj = months since enrollment (mean-centered)

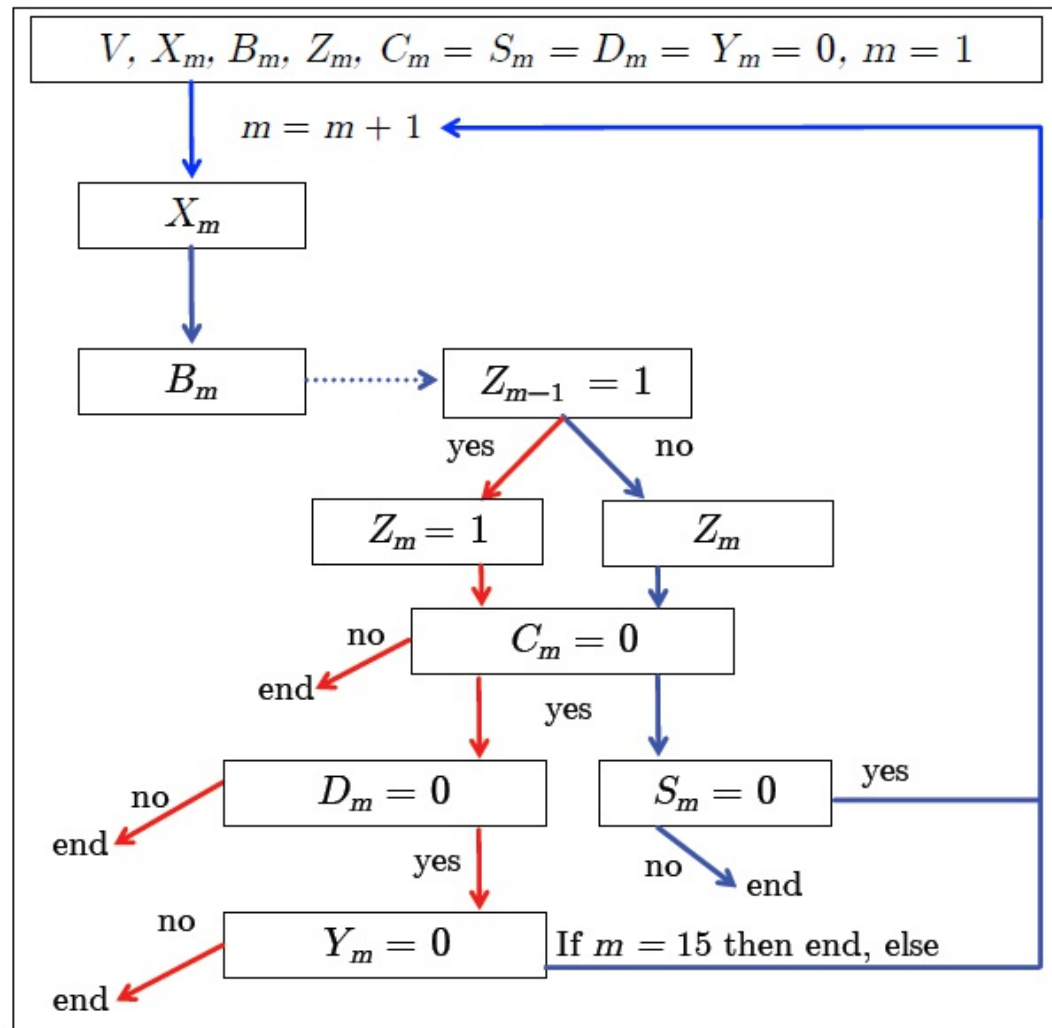
Lag - two previous values: $X1$, $X11$, $B1$, $B11$, $N1$, $N11$, $Z1$, $Z11$

(Note: Can't be influenced by later events.)

Resample and iterate based on flowchart



Resample and iterate based on flowchart



Simulation decisions based on trial characteristics:

- set live birth (Y) to 0 for all months < 8
- set pregnancy loss (D) to 0 for months < 2
- censor (S) at 6 months if no conception
- follow to 15 months

RESULTS

MC resampled $500 \times N$ ($500 \times 1228 = 614,000$)

Simulate "natural course": g-formula: **0.5381526** (observed: **0.517101**)

```
mean(bn$Y,na.rm=T); mean(a2[a2$last==1,]$Y)
```

Per-protocol estimate, ASA: g-formula: **0.511022** (observed, ITT: **0.5488599**)

```
mean(b1$Y,na.rm=T); mean(a2[a2$last==1&R==1,]$Y)
```

Per-protocol estimate, placebo: g-formula: **0.448** (observed, ITT: **0.485342**)

```
mean(b0$Y,na.rm=T); mean(a2[a2$last==1&R==0,]$Y)
```

Per-protocol effect estimate: +6.3% live birth

(to run: *pgf_annotated.R*)

Review: Estimation Steps

- Draw a DAG representing causal structure
- Model temporal / structural relationships (Q-model) implied by DAG
- Sample baseline values simulate POs using model
- Repeat setting exposures to levels reflecting desired contrast
- Bootstrap statistical model to obtain SE

Other considerations, methods, Q&A

Other estimation considerations

Other estimation considerations

1. *Missing data*

- Single imputation within each bootstrap

Other estimation considerations

1. *Missing data*

- Single imputation within each bootstrap

2. *No interference*

- Defining PO based on "distance" to other units

Other estimation considerations

1. *Missing data*

- Single imputation within each bootstrap

2. *No interference*

- Defining PO based on "distance" to other units

3. *Test sensitivity to causal ordering*

- Re-estimate with different Q-models

Other estimation considerations

1. *Missing data*
 - Single imputation within each bootstrap
2. *No interference*
 - Defining PO based on "distance" to other units
3. *Test sensitivity to causal ordering*
 - Re-estimate with different Q-models
4. *Test sensitivity to simulation error*
 - Try out larger MC samples

Other estimation considerations

1. *Missing data*
 - Single imputation within each bootstrap
2. *No interference*
 - Defining PO based on "distance" to other units
3. *Test sensitivity to causal ordering*
 - Re-estimate with different Q-models
4. *Test sensitivity to simulation error*
 - Try out larger MC samples
5. *Natural course*
 - Visualize means, medians, distributions

Related approaches

Inverse-Probability Weighted Marginal Structural Models

- standardization by propensity score weighting
- *e.g.* counterfactual means for a pseudo-population
- (rather than standardized to observed covariate distributions)
- generate IPW for each treatment time
- slightly easier to implement

Related approaches

Inverse-Probability Weighted Marginal Structural Models

- standardization by propensity score weighting
- *e.g.* counterfactual means for a pseudo-population
- (rather than standardized to observed covariate distributions)
- generate IPW for each treatment time
- slightly easier to implement

Targeted Maximum Likelihood-Based Estimation (TMLE)

- a g-computation approach that adds exposure model to "target" estimate
- usually implements cross-validated library of algorithms (SuperLearner) to estimate the Q-model:
 - lasso, regularized GLM,
 - K nearest neighbors,
 - support vector machine,
 - random forests, *etc.*

Some Helpful Resources

- Daniel RM, De Stavola BL, Cousens SN. gformula: Estimating causal effects in the presence of time-varying confounding or mediation using the g-computation formula. *The Stata Journal*. 2011; 11(4): 479-517.
- Hernán MA, Robins JM (2019). *Causal Inference*. Boca Raton: Chapman & Hall/CRC, forthcoming.
 - Full working draft accessible at [Miguel's site](#)
- Naimi AI, Cole SR, Kennedy EH. An Introduction to G Methods. *Int J Epi*. 2016.
- Datasets & examples from Ashley Naimi's SER shortcourses.
 - Accessible at [Ashley's github repo](#)

Parting thoughts

- Clear thinking about desired estimands are essential
- Draw the DAG that best represents prior knowledge
- Natural intervals may or may not be defined by data
- Careful consideration of causal assumptions is critical
- Dynamic treatment regimes (Friday!)
- Interpret carefully!
- Any questions / comments: [Email me!](#)

Final Word



Judea Pearl @yudapearl · 20 Nov 2018

Contrary to expectations, the definition of "causal modeling" is fairly easy to articulate. To me, "causal model" is a set of assumptions about the data generating process, which cannot be expressed as properties of the joint distribution of observed variables. [#Bookofwhy](#)



5



26



100



ADDITIONAL SLIDES

Hernán, Hsu, Healy (2018)

Description - Quantitative summary of features of the world

Prediction - Map input features onto output features

Causal Inference - Predict different worlds given change of certain features (counterfactual)

Data science is science's second chance to get causal inference right: A classification of data science tasks.
(arXiv:1804.10846)