

# Registry Data, TMLE, and `ltmle`

Zeyi Wang  
wangzeyi@berkeley.edu

Division of Biostatistics, UC Berkeley School of Public Health

December 15, 2022

# Outline

- 1 Prepare registry data for causal inference
- 2 TMLE and `ltmle`
- 3 Have a successful `ltmle` implementation
- 4 Extensions

## Prepare registry data for causal inference

# Introduction: longitudinal data

Assume a Structural Causal Model (SCM) for

$O = (L_0, A_1, L_1, \dots, A_K, L_K)$ :

$$L_0 = f_{L_0}(U_{L_0})$$

$$A_1 = f_{A_1}(L_0, U_{A_1})$$

$$L_1 = f_{L_1}(L_0, A_1, U_{L_1})$$

...

$$L_K = f_{L_K}(L_0, A_1, L_1, \dots, A_K, U_{L_K}).$$

Suppose we care about an intervention over variables as

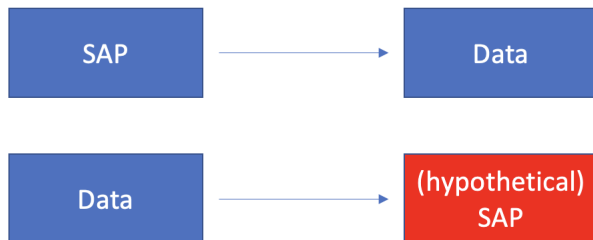
$A_1 = a_1, \dots, A_K = a_K$ , which defines counterfactuals

$$L_k(\bar{a}) = f_{L_k}(L_0, a_1, L_1, \dots, a_k, U_{L_k}), \text{ for } k = 1, \dots, K$$

and satisfies certain identification assumptions, then we can identify and analyze targets such as  $E[L_K(\bar{a})]$  with `ltmle`: ...

# What is special about registry data?

- Data is collected continuously;
  - different implications of time discretization (e.g. with `1tm1e`).
- Your Statistical Analysis Plan (SAP) is written *after* data is collected;
  - the trial in SAP is only hypothetical; no extra data can be collected if it is not already existing in the registry.



## Example: definition of time zero

**Time zero** is the time that a subject enters the study ( $t = 0$ ).

### Example

Suppose we have a database with treated and untreated patients. Both groups were followed up for  $x$  years.

## Example: definition of time zero

**Time zero** is the time that a subject enters the study ( $t = 0$ ).

### Example

Suppose we have a database with treated and untreated patients. Both groups were followed up for  $x$  years.

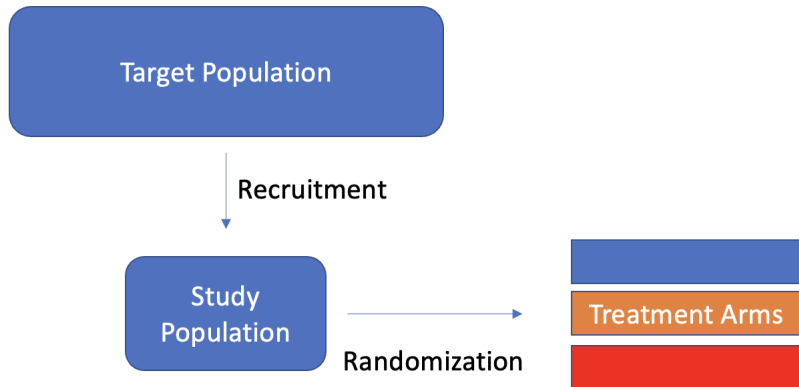
Possible issues:

- Some of the control group might never have a chance of receiving this treatment (not having the same disease; healthy controls; geographic or socioeconomic restrictions).
- Some may develop symptoms and conditions months after entering the database.
- One can end up with completely uncomparable arms with no proper confounding adjustment.
- Bias can be in either direction.

# Principle: define target population first

Even for a (fake) recruitment of a hypothetical trial, target population needs to be specified.

$t = 0$  should be the date of (fake) recruitment and treatment allocation as specified in the hypothetical SAP.





## Example: time zero

### Example

Consider all type 2 diabetes patients (defined by previous usage of first-line drugs such as metformin) in a healthcare record database who started second-line drugs (different brands/types) between 2010-2019.

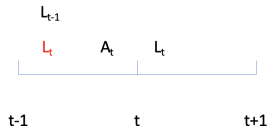
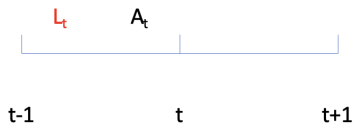
- Suitable for analysis and comparisons of second-line drugs on the market between 2010-2019 without the impact of COVID.
- $t = 0$  is the date of the earliest second-line drug subscription record.
- Randomness in treatment types might be explained by a list of confounders (pre-conditions, education, income, etc.).

# Example: time discretization

## Example (Unintended post-randomization confounders. )

Data is collected at the first and the last week of a month.

An online tutorial uses ordering  $L_0, A_0, L_1, A_1, \dots, L_K, A_K$ .

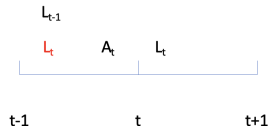
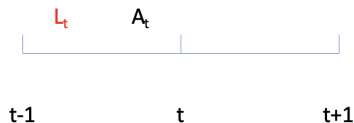


# Example: time discretization

## Example (Unintended post-randomization confounders. )

Data is collected at the first and the last week of a month.

An online tutorial uses ordering  $L_0, A_0, L_1, A_1, \dots, L_K, A_K$ .



## Possible issues

- "Confounding adjustment"  $A_t | L_t$  might be post-randomization.
- Actual effects may be wiped out like a Simpson's paradox.
- **Cautious:** Do NOT blindly discretize variable into bins and enforce preferred ordering.

# Principle: be careful about variable ordering with discretized time

Need to decide between:

- a conservative adjustment possibly biased due to unmeasured confounding;
- a careless adjustment with invalid causal interpretation.

Considerations

- Only adjust for confounders that are guaranteed to be pre-randomization.
- For longitudinal treatment, randomization can happen at multiple time points, and this applies to each time point separately.
- Evaluate the chance that the final confounder list fully explains randomness in treatment allocation.

# Different types of discretized variables

- Status (treatment/censoring/confounder) at the beginning/end of an interval;
- Summary of measurements collected across an interval;
- Incidence count across an interval;
- Indicator that an event has happens within an interval.

For treatment node  $A_t$ , one cannot adjust for any "blocks" overlapping with  $A_t$  definition or any "status" variables with time stamps prior to the end of  $A_t$  definition.

# Width of time bin

- No need to make a time grid finer than how frequent variables can change.
- No need to make the finest time grid either — can cause trouble.
- Continuous time TMLE package is under development (stay tuned).

## Example: an unsupported trial

The whole SAP is hypothetical — the proposed hypothetical trial might not have enough support in the data.

### Example (Off-label use. )

An approved anti-diabetic drug was unofficially used for preventing dementia. Patients were in doubt of the effect; none of them regularly prescribed the drug throughout the whole follow-up period; majority of the users only received one prescription.

- An "always-on" treatment arm would have almost no support in the data.
- Outcome-blind simulation is useful for detecting problems at an early stage and creating a chance to update (reversible) SAP.
- Intention-to-treat is not useful either. (Treatment holiday is discussed in later slides.)

# Principle: adjust SAP according to data

With hypothetical trials:

- Not all interesting questions can be answered by your data;
- SAP is reversible;
- It is allowed to have unpractical hypothetical trials, and they can be interpretable and supported;
- Interpretability is verified conceptually, while data support can be verified with outcome-blind simulation.



## Example: treatment holiday

Treatment holidays are the periods of time that a patient stops receiving a treatment (medicine, chemotherapy, etc.).

- Treatment holidays in practice can be intentional, considering quality of life, tolerance, toxicity etc.
- It is less often to have unscheduled treatment holidays as part of an actual randomized trial.
- But adding treatment holidays to hypothetical trials can gain in data support while preserving interpretability.

### Example

Regime option 1: always taking a drug (despite it being “off-label” use of an approved drug, or with other side effects).

Regime option 2: Stay on treatment more than 50% of the time with  $< 1$  month treatment holidays.

## Example: length of treatment

- Intention-to-treat might be a meaningless group for example if patients are in doubt are dropping off quickly while it is only effective with cumulative usage.
- Full follow-up might lack data support.

### Example

Regime option 1: Stay on treatment more than 80% of the time with  $< 1$  month treatment holidays for 1 year.

Regime option 2: Stay on treatment more than 80% of the time with  $< 1$  month treatment holidays for 2 year.

...

## Example: mediation

- Direct effect: treated vs control, while holding a specific mediator distribution;
- Indirect effect: without changing treatment condition, enforcing mediators to be *like* their treatment group vs control group counterfactuals.

Mediation is a combination of intervention over treatment variables and mediator variables.

The intervention over mediators is almost always hypothetical/unrealistic unless for few instances such as separable effects (effect of smoking and effect of nicotine).

## Example: censoring

Data censoring/missing is common even in the analysis of actual trials in the real world (end of study, drop-outs, etc.).

Analyzing *treatment effects among uncensored subjects* is essentially defining a hypothetical trial that combine the actual trial with an additional intervention on censoring indicators.

Example (Average treatment effect among the uncensored)

$$E[Y(\text{treated, uncensored}) - Y(\text{untreated, uncensored})]$$

Missing value indicator vs multiple imputation.

- Can create a bivariate version  $\tilde{X} = (X, \delta_X)$  of a variable  $X$ , where  $\delta_X$  is an indicator of whether  $X$  is observed, and  $\tilde{X} = (\text{median}(X), 0)$  when  $\delta_X = 0$ .
- Observed data information is fully preserved for training of learners.
- Multiple imputation is expensive and not needed for follow-up analysis to achieve double robustness.

# Different types of regimes/interventions

Intervention (of an SCM) is a manipulation of variable status, which defines counterfactual variables.

- Counterfactuals defined with SCM implicitly satisfy the *consistency* assumption (as part of SUTVA; a theorem in SCM, an assumption for potential outcomes).

## Example (Intervention in SCM)

$O = (W, A, Y)$ .

$$W = f_W(U_W)$$

$$A = f_A(W, U_A)$$

$$Y = f_Y(W, A, U_Y)$$

$$W = f_W(U_W)$$

$$A = a$$

$$Y(a) = f_Y(W, a, U_Y)$$

$Y(a)$  is the counterfactual under intervention  $A = a$ .  $Y(a) = Y$  if  $A = a$ .

# Dynamic regimes

Instead of a specific value  $a$ , one can intervene on SCM with a deterministic function

$d$  : history up till treatment  $\rightarrow$  next treatment value; e.g.  $A = d(W)$ .

- $d(l_{t-1}) = 1$  if  $l_{t-1} > (120, 80)$  mmHg.
- $d_0 = \arg \max_d E[Y(d)]$  (optimal treatment regimes).
- $d_t(a_{t-1}) = \begin{cases} a_t, & a_{t-1} = 1 \\ 1, & a_{t-1} = 0. \end{cases}$  (max 1 consecutive treatment holiday)
- Dynamic regimes are used in both real-world and hypothetical trials.

# Static vs random interventions

Instead of assigning a static value such as  $A = a$  or  $A = d(W)$ , one could intervene on SCM by a random draw such as  $A \sim g^*(W)$ .

- Currently random/stochastic intervention is not implemented in `ltmle` but there exists packages that handle specific problems involving random interventions.
- Natural direct/indirect effects can be defined with random interventions whereas controlled direct effects are defined with static intervention.
- Treatment holiday (natural treatment value) can be used as either static or random intervention.
- Caution: random allocation of a static intervention is NOT random intervention.
- Dynamic regimes are static (assigning  $A = d(\cdot)$ ) and can be implemented with `ltmle`.



## TMLE and ltmle

# Longitudinal data, revisited

$O = (L_0, A_1, L_1, A_2, L_2, \dots, A_K, L_K) \sim P_0 \in \mathcal{M}$ .  $O_1, \dots, O_n$  IID.

## Example (Static intervention with longitudinal data)

$$L_0 = f_{L_0}(U_{L_0})$$

$$A_k = a_k$$

$$L_k(\bar{a}) = f_{L_k}(L_0, a_1, L_1, \dots, a_k, U_{L_k}).$$

## Example (Dynamic intervention with longitudinal data)

$$L_0 = f_{L_0}(U_{L_0})$$

$$A_k(\bar{d}) = d_k(Pa(A_k))$$

$$L_k(\bar{d}) = f_{L_k}(L_0, A_1(\bar{d}), L_1(\bar{d}), \dots, A_k(\bar{d}), U_{L_k}).$$

# Identification assumptions

(Consistency assumption is implicit in SCM.)

- 1 Sequential ignorability:  $L_s(\bar{a}) \perp A_t | L_0, A_1, \dots, L_{k-1}$  for  $s \geq t$ .
- 2 Positivity:  $P_0(A_t = a_t | L_0 = l_0, a_1, \dots, L_{k-1} = l_{k-1}) > 0$  for all  $P_0(L_0 = l_0, A_1 = a_1, \dots, L_{k-1} = l_{k-1}) > 0$  at the true distribution  $P_0$ .

Note: both are important. The first is verified conceptually. The second needs to be verified on finite-sample estimators  $\hat{P}(A_t | \dots)$  for  $P_0(A_t | \dots)$ .

# Target parameter

G-computation formula for an outcome  $Y \in L_K$  under aforementioned assumptions:

$$E[Y(\bar{a})] = \sum_{l_1, \dots, l_K} y P(L_0 = l_0) \prod_{t=1}^K P(L_t = l_t | l_0, a_1, \dots, a_t)$$

Define  $\Psi : \mathcal{M} \rightarrow \mathbb{R}$  as RHS above.

# Sequential regression (iterated conditional expectation)

Define

$$\begin{aligned}Q_{K+1}(\bar{a}) &= Y \in L_K \\Q_K(\bar{a})(\bar{L}_{K-1}) &= E[Q_{K+1}(\bar{a}) | L_0, a_1, L_1, \dots, a_K] \\Q_{K-1}(\bar{a})(\bar{L}_{K-2}) &= E[Q_K(\bar{a}) | L_0, a_1, L_1, \dots, a_{K-1}] \\&\dots \\Q_1(\bar{a})(L_0) &= E[Q_2(\bar{a}) | L_0, A_1 = a_1] \\ \Psi(P) &= E[Q_1(L_0)]\end{aligned}$$

- **Note:** the regression at  $K - 1 \dots, 1$  is continuous, whereas at regression  $K$  the outcome may be binary.
- The  $k$ -th regression model is always against the history before  $L_k$ .
- These are the  $Q$  regression models that need to be specified in `ltmle`.

## $Q$ and $g$ components

- Propensity score functions  $g_t(a_t|l_0, a_1, l_1, \dots, l_{k-1}) = P(A_t = a_t|L_0 = l_0, A_1 = a_1, \dots, L_{k-1} = l_{k-1})$  are the  $g$  components in `ltmle`.
- $\Psi(P) = \Psi(Q)$  is a function of the  $Q$  components.
- If  $Q$  is the truth, might as well just do plug-in  
 $E_{P_n}[Q_1(L_0)] = \frac{1}{n} \sum_{i=1}^n Q_1(L_{0,i})$ .
- In randomized control trials,  $g$  components are correct, and can be used to improve  $Q$  estimation.
- In observational studies, doubly robust estimators improve the chance of getting asymptotically unbiased estimation.
- TMLE for  $\Psi(Q)$  is a locally efficient doubly robust estimator.

# Targeted Maximum Likelihood Estimation (TMLE)

TMLE is a two-stage estimator that maps an initial estimator  $Q$  into an updated final estimation  $Q^*$ .

- Plug-in at the final TMLE update  $\Psi(Q^*)$  is asymptotically linear and satisfies local efficiency properties.
- Like MLE would solve score equations, TMLE solves  $E_{P_n} D(Q^*, \hat{g}) \approx 0$  where  $D : \mathcal{M} \rightarrow \mathbb{R}$  is the **efficient influence curve (EIC)** (for static and dynamic regimes, this can be done inside `ltmle`).
- Note that  $\hat{g}$  is not updated (robust variance option in `ltmle`).

# Impact of using super learners in initial estimation $\hat{Q}, \hat{g}$

Assuming  $P_n D(Q^*, \hat{g}) = 0$ , we have exact expansion

$$\Psi(Q^*) - \Psi(Q_0) = (P_n - P_0)D(Q^*, \hat{g}) + R(Q^*, Q_0, \hat{g}, g_0).$$

- Residual  $R(Q^*, Q_0, \hat{g}, g_0)$  involves products of errors of  $Q^*$  and  $\hat{g}$ . This is the source of double robustness.
- Improving accuracy of  $\hat{Q}$  and  $\hat{g}$  with super learner leads to better controlled residual terms and better chance to be  $o_P(n^{-1/2})$ .
- $(P_n - P_0)D(Q^*, \hat{g}) = (P_n - P_0)D(Q^0, g^0) + (P_n - P_0)(D(Q^*, \hat{g}) - D(Q^0, g^0)).$ 
  - If also  $\hat{g}$  is consistent for  $g^0 = g_0$  or  $\hat{Q}$  is consistent for  $Q^0 = Q_0$ , we have asymptotic linearity, so long as  $D(Q^*, \hat{g})$  is not "too complex" or using CV-TMLE (not implemented in `ltmle` but in `tlverse`);
  - Be cautious for too aggressive SL (simulation; CV-TMLE; discrete super learner; highly adaptive lasso).



# Summary

- 1 Previously discussed static and dynamic interventions (including treatment holidays, censoring, controlled direct effects, etc.) on longitudinal data  $O = (L_0, A_1, L_1, \dots, L_{K-1}, A_K, L_K)$  and binary  $A_1, \dots, A_K$  can be analyzed with `ltmle`.
- 2 Asymptotically unbiased estimation and valid inference for ML/SL algorithms that are otherwise difficult to quantify in uncertainty.
- 3 SL can be used as a tool for objective model selection.
- 4 SL combined with DR properties improves the chance of achieving asymptotically unbiased estimation and valid inference (cautious about learners, or use CV-TMLE).
- 5 In general, improved initial estimation with SL leads to improved testing power and reduced minimum sample sizes for detecting effects.

## TMLE on meta level data (extension; optional)

Consider  $L_0, A_1, L_1, \dots, L_K, Y$ .

Each  $L_t$  can be summarized into lower dimensional  $L_t^r$  as validation set predictions of complex super learners. In practice consider prediction of both  $Q_t$  and  $A_t$ .

Run simpler learners and TMLE on the new data  $L_0^r, A_1, L_1^r, \dots, L_K^r, Y$ .

Have a successful `ltmle` implementation

# Basics: $W, A, Y$

```
rexpit <- function(x) rbinom(n=length(x), size=1, prob=plogis(x))
n <- 10000
W <- rnorm(n)
A <- rexpit(W)
Y <- rexpit(W + A)
data <- data.frame(W, A, Y)
ltmle(data = data,
      Anodes = "A", Lnodes = "W", Ynodes = "Y",
      abar = list(1,0),
      estimate.time = F,
      variance.method = "ic")
```

- Clean up data frame with column names  $W, A, Y$ .
- `Ynodes` needs to be the last variable. (Need to prepare another data frame if to analyze outcome at another time point. )
- `estimate.time` takes a random subsample of size 50; not necessary for small data.
- `variance.method = "ic"` is the standard error estimator by plugging in  $\sqrt{\text{var}_n(D(Q^*, \hat{g}))/n}$  without targeting  $\hat{g}$ .
- See <https://cran.r-project.org/web/packages/ltmle/vignettes/ltmle-intro.html> for more usage examples.

# Two time points

```
rexpit <- function(x) rbinom(n=length(x), size=1, prob=plogis(x))
n <- 10000
L0 <- rnorm(n)
A1 <- rexpit(L0)
L1 <- rexpit(L0 + A1)
A2 <- rexpit(L0 + L1)
Y <- rexpit(L0 + A1 + A2)
data <- data.frame(L0, A1, L1, A2, Y)
loc_A <- grep("^A", names(data))
loc_L <- grep("^L", names(data))
loc_Y <- grep("^Y", names(data))
ltmle(data = data, Anodes = loc_A, Lnodes = loc_L, Ynodes = loc_Y,
      abar = list(c(1, 1), c(0, 0)),
      estimate.time = F,
      variance.method = "ic")
```

- ltmle automatically picks up the ordering of variables; do not need to name all the treatment nodes in a same way. But after  $L_0$  everything has to be either Anodes, Cnodes, Lnodes, or Ynodes.
- Can use column number for node specification.
- abar needs to be a list of vectors or matrices.

# Survival outcome with right-censoring

```
names_all <- colnames(d)
loc_C <- grep('^censor', names_all)
loc_A <- grep('^A', names_all)
loc_Y <- grep('^event', names_all)
loc_L <- setdiff(seq_along(names_all), c(loc_C, loc_A, loc_Y))
```

```
ate_call <- ltmle(data=d,
  Anodes = loc_A, Cnodes = loc_C, Lnodes= loc_L, Ynodes= loc_Y,
  survivalOutcome = T,
  abar = list(rep(1,length(loc_A)), rep(0,length(loc_A))),
  estimate.time = F,
  variance.method = "ic")
```

- Suppose that censoring nodes start with "censor", treatment nodes start with "A", outcome process start with "event", and all other variables are baseline/time-varying confounders.
- Need to clean data first:
  - After censoring (default value 0 happens, remaining censoring nodes should have value 0, and all other variables should have value NA);
  - After event happens (default 1), remaining Ynodes should be 1, and others should be NA.

# Specify $Q$ and $g$ models

```
ate_call <- ltmle(data=d,  
  Anodes = loc_A, Cnodes = loc_C, Lnodes= loc_L, Ynodes= loc_Y,  
  survivalOutcome = T,  
  abar = list(rep(1,length(loc_A)), rep(0,length(loc_A))),  
  estimate.time = F,  
  variance.method = "ic",  
  SL.library = c("SL.mean", "SL.glmnet", "SL.glm"),  
  gform = vec_gform,  
  Qform = vec_Qform,  
)
```

# Specify Q and g models

```
names_all <- names(data)
loc_AC <- grep("^A_|^censor-", names_all)
names_AC <- names_all[loc_AC]
loc_C <- grep("^censor-", names_all)
vec_gform <- sapply(seq_along(loc_AC), function(u) {
  loc_node <- loc_AC[u] # the current
  if (loc_node > 1) { # interaction between At and Pa(At) \ AC nodes
    loc_AC_node_previous <- loc_AC[loc_AC < loc_node] %>% setdiff(loc_C)
    list_nonAC_past <- lapply(loc_AC_node_previous, function(iii) setdiff(1:iii, loc_AC))
  } else loc_AC_node_previous <- NULL
  name_node <- names_AC[u]
  # main terms
  form <- names_all[(1:(loc_node-1)) %>% setdiff(loc_C)] %>% paste0(collapse = " + ")
  form <- paste0(name_node, " ~ ", form)
  # interaction involving only A_t and its nonAC past
  if (!is.null(loc_AC_node_previous)) {
    for (iii in seq_along(loc_AC_node_previous)) {
      form <- paste0(form,
        " + ",
        paste0(names_all[loc_AC_node_previous[iii]], " : ",
          names_all[list_nonAC_past[[iii]]]) %>% paste0(collapse = " + ")
    )
  }
  return(form)
})
names(vec_gform) <- names_AC
```

- Cnodes do not enter the formulas.
- Write desired regression formula for each learners for each Anodes and Cnodes
- $\hat{Q}$  models (omitted) start with an outcome name Q.kplus1.



# Treatment holidays with dynamic regimes

```
abbar <- d[, loc_A]
for (k in 2:K) {
  abbar[,k] = ifelse(abbar[, k-1] == 1, abbar[, k], 1)
}
abbar <- as.matrix(abbar)
ate_call <- ltmle(data=d,
  Anodes = loc_A, Cnodes = loc_C, Lnodes= loc_L, Ynodes= loc_Y,
  survivalOutcome = T,
  abbar = abbar,
  estimate.time = F,
  variance.method = "ic",
  SL.library = c("SL.mean", "SL.glmnet", "SL.glm"),
  gform = vec_gform,
  Qform = vec_Qform,
  deterministic.g.function = det_g_function
)
```

- This example allows up to 1 treatment holiday in the treated arm.
- $\text{abbar}$  is the observed  $A_t(\bar{d})$  values ( $n$  by  $\text{num.Anodes}$ ).

# Treatment holidays with dynamic regimes

If  $A_{t-1} = 1$ , a subject remains in the treatment arm with probability 1. (Note that this deterministic g function asks for return values to be probability of  $A_t = 1$ , not necessarily be the same as the actual  $A_t(\bar{d})$  value.)

```
# use deterministic g function to allow treatment holidays
det_g_function <- function(data, current.node, nodes) {
  # if not A node, skip
  # if there is no previous node, then treat as usual
  Anodes <- nodes$A
  if (!(current.node %in% Anodes)) return(NULL)
  if (!(any(Anodes < current.node))) return(NULL)

  prev.a.node <- max(Anodes[Anodes < current.node])
  is.deterministic <- ifelse(!is.na(data[, current.node]),
    data[, prev.a.node] == 1, F)
  prob1 <- data[, current.node][is.deterministic]
  return(list(is.deterministic=is.deterministic, prob1=prob1))
}
```

# Test runs and robust variance options

When there is finite-sample positivity violation, which makes  $\hat{g}$  extremely close to 0, then the sample variance of  $D(Q^*, \hat{g})$  becomes an unstable estimator for the true standard error. In this case, a targeted estimation of this asymptotic standard error is advised (details omitted).

To implement, keep `variance.method = "tmle"`.

This can take longer time to run, so it is recommended to have test runs on smaller samples first. Despite being more robust against finite-sample positivity violation, robust variance might as well be meaningless (e.g. CI going beyond  $[0, 1]$  for binary outcome) if positivity assumption is truly violated. In those cases, revision of SAP is advised.

# Collecting diagnostics

- `ltmle` object returns `cum.g` which is by default bounded above 0.01.
- `cum.g.unbounded` can also be used to evaluate the finite sample violation of positivity assumptions.

# Extensions

# Outcome-blind simulation

Simulation with fake data generating process is crucial for valid registry causal inference, especially for effective modification of SAP that depends on the data.

- Detect positivity or rare outcome issues;
- Evaluate method validity given those challenges;
- Detect if SAP is lacking data support and if should switch to other target parameters instead.

However, the purpose of such simulation is not for testing the correctness of  $\hat{Q}$ ,  $\hat{g}$  models (with known DGP it can easy to guess which one is more correct).

# Outcome-blind simulations

Typical setup:

- Can use summary statistics of true data in data generating process especially for propensity score models, but avoid using actual outcomes.
- Create tuning hyperparameters for 1) positivity, 2) rareness of outcomes, 3) effect sizes.
- evaluate bias/SD/MSE and confidence interval coverage with repeated iterations.

# Long format vs wide format

Current implementation of `ltmle` relies on wide format data which can be less efficient for survival outcomes with right-censoring (need to have space holders after death and censoring).

Some `tlverse` packages and also `stremr` have options to implement long format implementation of the same algorithm.

One future direction is to have adaptively defined confounder per subject and then have pooled estimation. Stay tuned for future development.



# Thank You

**Thank you!**

**Questions?**

**Happy to discuss more!**

**Email: [wangzeyi@berkeley.edu](mailto:wangzeyi@berkeley.edu)**