# final project

## Fanding Zhou

## 4/29/2022

## Step 1: State scientific question.

(a) The causal question

What is the effect of a large-scale primary care redesign, the Comprehensive Primary Care Plus (CPC+) Initiative, on the monthly Medicare patient expenditures?

(b) The target population

Patients in different primary care practices in 18 regions.

## Step 2: Specify a structural causal model (SCM).

(a) The endogenous variables V are variables that are meaningful for the scientific question and affected by other variables in this model. For this study, we define the following endogenous variables:

- C = {V1, V2, V3, V4, V5}

- V1: age, continuous variable supported on [0,100]

- V2: income groups, ordinal variable encoded into 0~14 integers.

- V3: sex, 0 for female, 1 for male, approximately equally distributed.

- V4: mean standardized Hierarchical Condition Category (HCC) score, approximately standard normal distribution.

- V5: race categories, white/black/all other, 70:20:10.

- T = Initiative of Comprehensive Primary Care Plus (CPC+)

- Y = Monthly Medicare patient expenditures

Here $V = (C, T, Y)$

(b) The exogenous U include all the unmeasured or unknown factors not included in V that impact the values that the V variables take.

The exogenous nodes are $U = (U_T, U_C, U_Y)$.

We would need to place some independence assumptions on the distribution of unmeasured factors $P_U$. Specifically, we first need $U_T \perp\!\!\!\perp U_Y$ and then we also need one of the following assumption: $U_C \perp\!\!\!\perp U_T$ or $U_C \perp\!\!\!\perp U_Y$.

Here we let $U_1 = \{U_C, U_T\}$ and $U_2 = \{U_C, U_Y\}$. Under these two sets of assumptions, we can plot DAGs as follows. From the DAGs, we notice that when conditioning on C, the backdoor path is blocked and the target parameter is therefore identifiable.

(c) This would suggest the following structural equations F:

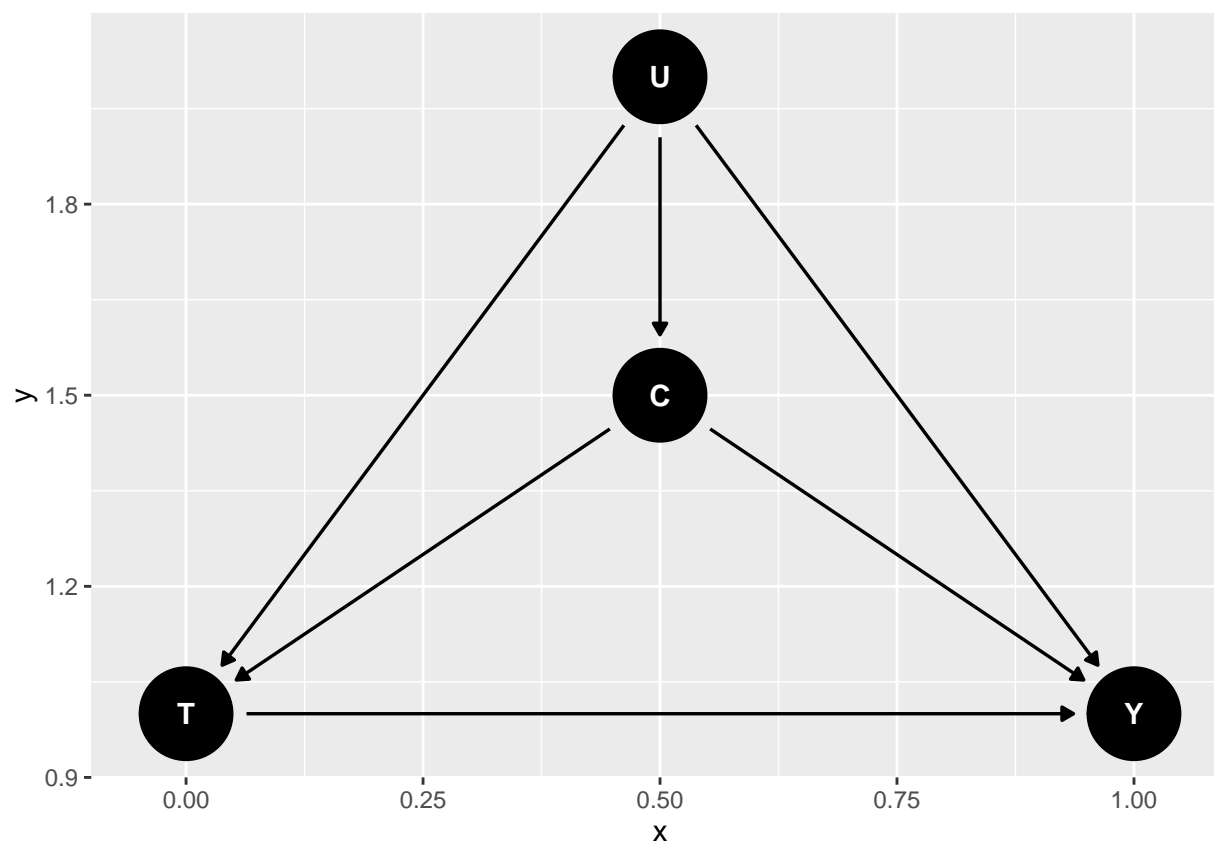$$C = f_C(U_C)$$
$$T = f_T(C, U_T)$$
$$Y = f_Y(C, T, U_Y)$$

(d) Exclusion restrictions.

An exclusion restriction is when a variable is not directly affected by another variable that precedes it in the structural causal model. For now we do not have a restriction regarding the endogenous variables.

(f) Causal graph.

```
dag_m <- dagitty('dag {
    T [pos="0,1"]
    Y [pos="1,1"]
    U [pos="0.5,2"]
    C [pos="0.5,1.5"]
    U->T
    U->C
    U->Y
    C->T
    C->Y
    T->Y
    }')
ggdag(dag_m, layout = "circle")
```



## Step 3. Specify the target parameter of the observed data distribution

(a) The counterfactuals of interest are Y(1) and Y(0) and they can be defined as

- Y(1): Monthly Medicare patient expenditures if all patients had joined CPC+ for 2 years

- Y(0): Monthly Medicare patient expenditures if all patients hadn't joined CPC+ for 2 years.

(b) Our target parameter is the average causal effect:

$$\Psi_{U,V}(P_0) = \mathbb{E}_{U,V}[Y(1) - Y(0)]$$

This means the target causal parameter is the difference in the monthly Medicare patient expenditures for patients enrolled in CPC+ for 2 years and not enrolled in CPC+.

## Step 4. Specify the observed data.

Suppose the data is from i.i.d sample from n = 329250 randomly sampled patients.

- $T_i$: whether patient i has enrolled in CPC+ for 2 years
- $C_i$: baseline confounders collected for patient i
- $Y_i$: Monthly Medicare patient expenditures for patient i

Further assume that the collected data are independent and identically distributed.

## Step 5. Identify the targeted causal effect with the observed data.

(a) Conditional independent assumption

$Y_i(0), Y_i(1) \per\!\!\!\perp T_i | C_i, i = 1, ..., n,$

(b) write our causal parameter with the observed data by G-computation:

$$\Psi(P_0) = \mathbb{E}[\mathbb{E}[Y|T = 1, X]] - \mathbb{E}[\mathbb{E}[Y|T = 1, X]]$$

Under our causal model and assumptions, average treatment effect equals to the observed difference in mean outcome within confounder strata, standardized to distribution of confounders.

(c) Relevant positivity assumption

## Step 6. Estimation and Statistical Inference.

```
set.seed(252)
ObsData = read.table("filtered_patient_1.csv")
```

result:

```
ltmle.SL = readRDS("result1.rds")
summary(ltmle.SL)
```

```
## Estimator:  tmle
## Call:
## ltmle(data = ObsData, Anodes = "Z", Ynodes = "Y", abar = list(1,
##     0), SL.library = SL.library)
##
## Treatment Estimate:
##    Parameter Estimate:  1317.1
##     Estimated Std Err:  6.6059
##               p-value:  <2e-16
##     95% Conf Interval: (1304.2, 1330.1)
##
## Control Estimate:
##    Parameter Estimate:  1262.7
```

```
##      Estimated Std Err:  5.9451
##               p-value:  <2e-16
##      95% Conf Interval: (1251.1, 1274.4)
##
## Additive Treatment Effect:
##    Parameter Estimate:  54.378
##      Estimated Std Err:  8.7306
##               p-value:  4.7114e-10
##      95% Conf Interval: (37.266, 71.49)
```