

# Estimation of probability associated with collective counterfactual outcomes

*Youjin Lee*

## Approximation of Directed Acyclic Graph (DAG)

In social network, treatments may spill over from the treated individual to his or her social contacts and outcomes may be contagious. Therefore, causal inference using observational data from a single social network requires observing longitudinal data on treatments and outcomes as they evolve in real-time, so that each spillover or contagious event appears in the data. In these settings, directed acyclic data (DAG) has been used to represent interference and contagion, but DAGs can be quickly be cumbersome. Especially in most cases it is impossible to collect the kind of real-time data required. If it is, the resulting model will generally high-dimensional and often too big to fit the available data.

Using `simcausal` R package, we can define DAG object where outcome variables of  $(Y_1^t, Y_2^t, Y_3^t)$  from  $t = 0, 1, 2, \dots, 20$ . Figure~1 illustrates data generating mechanism or causal relationship that a DAG object of `LDAG` stands for.

```
library(simcausal)
D <- DAG.empty()
D <- D +
  node("A1", t = 0, distr = "rbern", prob = 0.3) +
  node("A2", t = 0, distr = "rbern", prob = 0.5) +
  node("A3", t = 0, distr = "rbern", prob = 0.7) +
  node("Y1", t = 0, distr = "rbern", prob = plogis(-1 + A1[0])) +
  node("Y2", t = 0, distr = "rbern", prob = plogis(0 - 0.5*A2[0])) +
  node("Y3", t = 0, distr = "rbern", prob = plogis(1 - 0.5*A3[0]))
t.end <- 20
D <- D +
  node("Y1", t = 1:t.end, distr = "rbern",
    prob = plogis(3*Y1[t-1] - 1.5 + A1[0] + Y2[t-1])) +
  node("Y2", t = 1:t.end, distr = "rbern",
    prob = plogis(3*Y2[t-1] - 2 - 0.5*A2[0] + Y1[t-1] - Y3[t-1])) +
  node("Y3", t = 1:t.end, distr = "rbern",
    prob = plogis(3*Y3[t-1] - 1 - 0.5*A3[0] - Y2[t-1]))
LDAG <- set.DAG(D, verbose = FALSE)
```

Note that for causal inference, DAGs are only useful when if we observe all of the  $Y_i^t$  variables. Instead chain graph includes directed edges and undirected edges, with no directed cycles. If  $Y_i^t$  is approximated equal to  $Y_i^{t-1}$ , the conditional independence implied by the chain graph holds approximately in the DAG model. In this case, conditional (in)dependence structure in DAG from Figure~1 (approximately) holds for those of chain graph in Figure~2.

```
par(mar = c(4, 4, 1, 1))
plotDAG(LDAG, tmax = 10, xjitter = 0.3, yjitter = 0.1,
  edge_attrs = list(width = 0.2, arrow.width = 0.5, arrow.size = 0.3),
  vertex_attrs = list(size = 10, label.cex = 0.5, label.color = "dodgerblue"),
  verbose = FALSE)

library(netchain)
par(mar = c(4, 4, 1, 1))
```

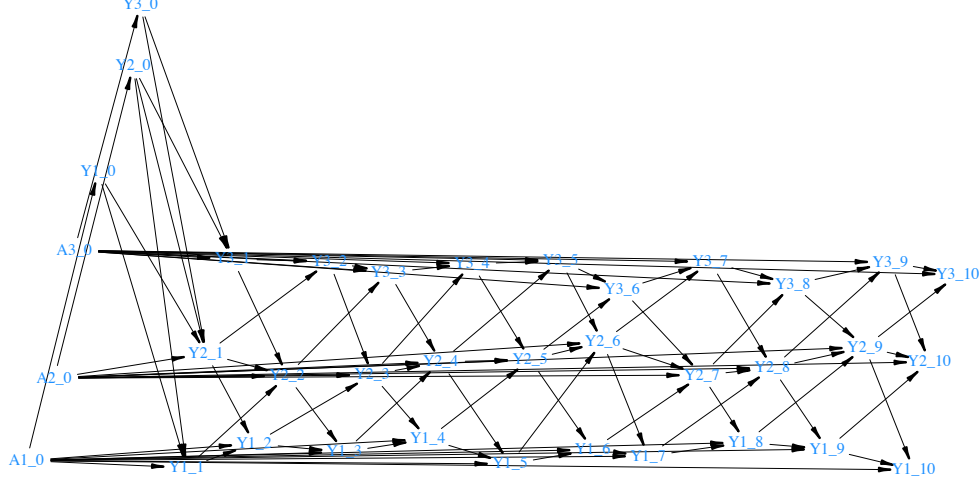


Figure 1: plotting an example of evolving outcome variables ( $Y_1, Y_2, Y_3$ ) and treatment variables ( $A_1, A_2, A_3$ ) following directed cyclic graph of 1DAG.

```
plotCG(1DAG)
```

## Input observations

Chain graphs can be used to define hybrid graphical models combining features of both log-linear models on undirected graphs and DAG models used to represent causal relationships. Therefore, chain graphs provide justification of parsimonious conditional log-linear models to identify causal effects on the collective outcomes.

To estimate the parameters in conditional log-linear model, we require enough independent replicates of observations. Observations must take a form of a  $m \times n$  matrix for binary  $Y$  and  $A$  respectively, where  $m$  is the number of units and  $n$  is the number of independent observations. Depending on the contexts of your scenario, we can enter no confounders (Defaults to `listC = NULL`), a  $m \times n$  matrix for `listC`, or a list of  $m \times n$  matrices for `listC` for which of each matrix stands for observations for one confounding factors (we later call them  $C_1, C_2, \dots, C_q$ ).

We provide a simple function `simGibbs` to generate binary  $(Y, A, C)$  from chain graph model under simple scenario, which requires additional structural information (`weight.matrix`, `treat.matrix`, and `cov.matrix`). Consider the following conditional log-linear model following chain graph model defined up-to two-way interaction effects:

$$p(Y = (y_1, y_2, \dots, y_m) \mid A = a, C = c) = \frac{1}{Z} \exp \left\{ \sum_{i,j=1}^m w_{ij} y_i y_j + \sum_{i,j=1}^m k_{ij} a_i y_j + \sum_{i,j=1}^m h_{ij} c_i y_j \right\}$$

A `weight.matrixij` ( $w_{ij}$ ) indicates two-way interaction effects between  $Y_i$  and  $Y_j$  when  $i \neq j$  and main effect of each outcomes when  $i = j$ . Note that by the construction of chain graph model, `weight.matrix` should be symmetric. A `treat.matrixij` ( $t_{ij}$ ) indicates a direct causal effect of  $A_i$  on  $Y_j$ ; a `cov.matrixij` ( $h_{ij}$ ) indicates a direct causal effect of  $C_i$  on  $Y_j$ . By the definition, `treat.matrix` and `cov.matrix` are directional.

The following codes generate  $n = 1000$  (`n.gibbs` × `n.sample` = 500) observations of  $m = 3$  (`n.unit` = 3) units assuming chain graph of  $Y_1 - Y_2 - Y_3$  and direct effect from  $A_i \rightarrow Y_i$  and  $C_i \rightarrow Y_i$  for  $i = 1, 2, 3$ .

```
library(netchain)
weight.matrix = matrix(c(0.5, 1, 0, 1, 0.3, 0.5, 0, 0.5, -0.5), 3, 3)
```

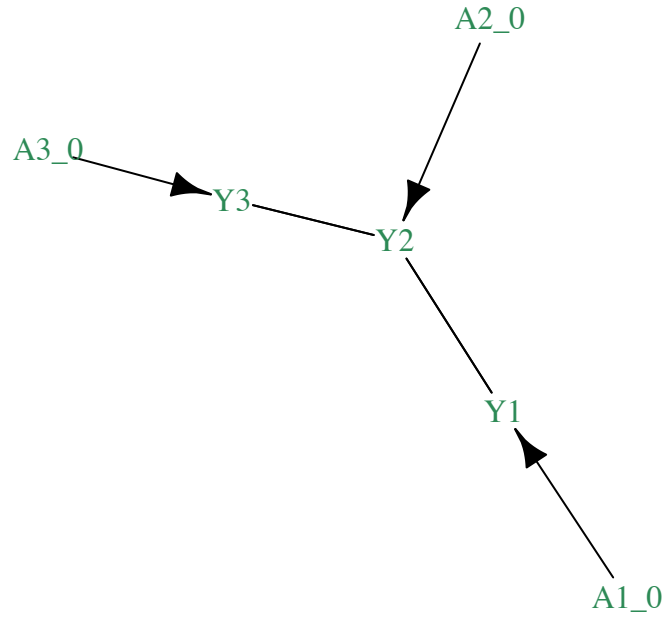


Figure 2: Plotting chain graph component from 1DAG.

```

simobs = simGibbs(n.unit = 3, n.gibbs = 10, n.sample = 10,
                  weight.matrix,
                  treat.matrix = 0.5*diag(3), cov.matrix= (-0.3)*diag(3) )

inputY = simobs$inputY
inputA = simobs$inputA
inputC = simobs$inputC
head(inputY)
#>      [,1] [,2] [,3]
#> [1,]    1    1    1
#> [2,]    1    1    0
#> [3,]    1    1    0
#> [4,]    1    1    1
#> [5,]    1    1    1
#> [6,]    1    1    1
head(inputA)
#>      [,1] [,2] [,3]
#> [1,]    0    1    0
#> [2,]    0    1    0
#> [3,]    0    1    0
#> [4,]    0    1    0
#> [5,]    0    1    0
#> [6,]    0    1    0
head(inputC)
#>      [,1] [,2] [,3]
#> [1,]    0    1    1
#> [2,]    0    1    1
#> [3,]    0    1    1
#> [4,]    0    1    1
#> [5,]    0    1    1
#> [6,]    0    1    1

```

## Input edge information

With  $(\text{outcomes}, \text{treatment}, \text{confounders}) = (\text{inputY}, \text{inputA}, \text{listC})$ , we ultimately aim to calculate probability associated with counterfactual collective outcomes given a certain unit-specific treatment of length  $m$ ,  $\mathbf{a}$  (`treatment`):

$$P(\mathbf{Y}(\mathbf{a}) = \mathbf{y}).$$

Depending on the form of target outcomes  $\mathbf{y}$ , `targetoutcome` vary: in the above example of  $m = 3$  units, `targetoutcome = c(1,0,0)` indicates we want to calculate  $P(\mathbf{Y}(\mathbf{a}) = c(1, 0, 0))$ . If we want to calculate the probability with multiple outcomes, e.g. homogeneous outcomes of  $(0,0,0)$  or  $(1,1,1)$ , you can input a matrix `targetoutcome = rbind(c(0,0,0), (1,1,1))`. If you want to calculate the probability when only one unit has one, instead of listing all the possible outcomes `targetoutcome = rbind(c(1,0,0), c(0,1,0), c(0,0,1))`, you can specify the number of one (or maximum value of `inputY`): `targetoutcome = 1`.

Most tricky and important part is edge information (`R.matrix`, `E.matrix`, and `edgeinfo`), which determines the factors put in the conditional log-linear model.

A `R.matrix` is  $m \times m$  relational symmetric matrix where `R.matrixij = 1` indicates  $Y_i$  and  $Y_j$  are adjacent, meaning the existence of interactions of feedback between the two. A `E.matrix` a  $m \times m$  matrix where `E.matrixij` indicates  $A_i$  has a direct causal effect on  $Y_j$ ; defaults to diagonal matrix, which indicates no interference. On the other hand, `edgeinfo` a list of matrix specifying additional directed edges (from confounders (specified in `listC`) or treatment (`inputA`) to the outcomes (`inputY`)) information. Assume that `edgeinfo = list(mat1, mat2)` where both `mat1` and `mat2` are matrices. For each matrix the first column specifies the types of variables that must involve outcome variable, "Y" while "A" indicates treatment variable and "C" indicates confounder. If you input multiple confounders in `listC` as a list, "C1" indicates the first confounders, "C2" indicates the second, and so on. For each matrix the second column specifies an index for unit corresponding to the variable in the first column. First example,

```
mat1 = rbind(c("Y", 1), c("A", 1), c("A", 2))
```

implies that  $A_1$  and  $A_2$  have a direct causal effect on  $Y_1$  and these two have a interaction effect; this will lead to three-way interaction term of  $Y_1$ ,  $A_1$ , and  $A_2$  in the log-linear model. (In this case a direct effect of  $A_1$  on  $Y_1$  and  $A_2$  on  $Y_1$  should be also specified.)

```
mat2 = rbind(c("Y", 1), c("C1", 1), c("C2", 1))
```

When you have multiple confounders and these confounders interactively have an effect on the outcomes, you can put edge information as `mat2`, which implies a direct effect of  $C_{11}$  and  $C_{21}$  on  $Y$  but these effects are dependent each other. (In this case, you should include `rbind(c("Y", 1), c("C1",1))` and `rbind(c("Y", 1), c("C2", 1) )` in `edgeinfo`.)

## Counterfactual outcomes

With those observations and edge information, the next step is to infer parameters in the conditional log-linear models, and with those parameters we generate counterfactual outcomes using Gibbs sampling. Gibbs sampling requires given `treatment` assignment, the number of independent Gibbs samplings (`n.gibbs`), the number of samples per each iteration (`n.sample`), and the number of burn-in samples (`n.burn`) per each iteration.

```
set.seed(1234)
library(netchain)
weight.matrix = matrix(c(0.5, 1, 0, 1, 0.3, 0.5, 0, 0.5, -0.5), 3, 3)
simobs = simGibbs(n.unit = 3, n.gibbs = 10, n.sample = 10,
                  weight.matrix,
                  treat.matrix = 0.5*diag(3), cov.matrix = (-0.3)*diag(3) )
inputY = simobs$inputY
```

```

inputA = simobs$inputA
inputC = simobs$inputC

R.matrix = ifelse(weight.matrix==0, 0, 1)
diag(R.matrix) = 0

#result = chain.causal.multi(targetoutcome = "mean", treatment = c(1,0,0),
#inputY, inputA, listC = inputC, R.matrix = R.matrix, E.matrix = diag(3),
#edgeinfo = list(rbind(c("Y", 1), c("C", 1)), rbind(c("Y", 2), c("C", 2)),
#rbind(c("Y", 3), c("C", 3))), n.obs = 1000, n.burn = 100)

#print(result)

```

## Approximate DAG process

```

sample.size <- 1000
simcausaldata <- sim(DAG = lDAG, n = sample.size, rndseed = 1234)
#> simulating observed dataset from the DAG object
dat <- cbind(simcausaldata$A1_0, simcausaldata$A2_0, simcausaldata$A3_0,
             simcausaldata$Y1_20, simcausaldata$Y2_20, simcausaldata$Y3_20)
colnames(dat) <- c("A1", "A2", "A3", "Y1", "Y2", "Y3")
dat <- as.data.frame(dat)

treatments <- c(1, 0, 1)
act_theta <- c(node("A1", t = 0, distr = "rbern", prob = alpha),
              node("A2", t = 0, distr = "rbern", prob = beta),
              node("A3", t = 0, distr = "rbern", prob = gamma))
lDAG <- lDAG + action("all_theta101", nodes = act_theta,
                    alpha = treatments[1], beta = treatments[2], gamma = treatments[3])
Y1_treat101 <- set.targetE(DAG = lDAG, outcome = "Y1", t = t.end,
                          param = "all_theta101")
Y2_treat101 <- set.targetE(DAG = lDAG, outcome = "Y2", t = t.end,
                          param = "all_theta101")
Y3_treat101 <- set.targetE(DAG = lDAG, outcome = "Y3", t = t.end,
                          param = "all_theta101")
results1 <- eval.target(Y1_treat101, n = 10000)$res
#> data not specified, simulating full data
#> no actions specified, sampling full data for ALL actions from the DAG
#> evaluating the target on 10000 simulated samples per action
results2 <- eval.target(Y2_treat101, n = 10000)$res
#> data not specified, simulating full data
#> no actions specified, sampling full data for ALL actions from the DAG
#> evaluating the target on 10000 simulated samples per action
results3 <- eval.target(Y3_treat101, n = 10000)$res
#> data not specified, simulating full data
#> no actions specified, sampling full data for ALL actions from the DAG
#> evaluating the target on 10000 simulated samples per action
print(c(results1, results2, results3))
#> Mean_Y1_20 Mean_Y2_20 Mean_Y3_20
#>      0.9007      0.5376      0.3183

library(gtools)
R.matrix <- matrix(c(0, 1, 0, 1, 0, 1, 0, 1, 0), 3, 3)

```

```

permat <- permutations(2, 3, c(0,1), repeats.allow = T)
marginalY3 <- permat[which(permat[,3] == 1),]
result.chain3 <- chain.causal.multi(targetoutcome = marginalY3,
                                   treatment = treatments,
                                   inputY = cbind(dat$Y1, dat$Y2, dat$Y3),
                                   inputA = cbind(dat$A1, dat$A2, dat$A3),
                                   listC = NULL, R.matrix = R.matrix,
                                   E.matrix = diag(3), edgeinfo = NULL,
                                   n.obs = 500, n.burn = 100)$causalprob
print(c(true = results3, estimate = result.chain3))
#> true.Mean_Y3_20      estimate
#>      0.318300      0.347096

```

## Identifying influential units

A function `causal.influence` takes the same input as `chain.causal.multi` except that a specific treatment is replaced by `Avalues` that In `causal.influence` we evaluate the average of collective outcomes under each treatment assignment  $E[\mathbf{Y}(\mathbf{a}_j)]$  as a measure of influence of unit  $j$ , where  $\mathbf{a}_j$  indicates the sole intervention of unit  $j$ , but `targetoutcome` can be a vector, matrix or an integer as that in `chain.causal.multi`.

```

set.seed(1234)
edgeinfo = list(rbind(c("Y", 1), c("C", 1)), rbind(c("Y", 2), c("C", 2)), rbind(c("Y", 3), c("C", 3)))
#influence = causal.influence(targetoutcome = "mean", Avalues = c(1,0),
#                             inputY, inputA, listC = inputC, R.matrix, #E.matrix = diag(3),
#                             edgeinfo = edgeinfo, n.obs = 100, n.burn = 10)
#print(influence)

```

The above result says that  $E[\mathbf{Y}(\mathbf{a}_3)] \approx 0.71$  achieves the highest influence, indicating that among three, treating unit  $i = 3$  would result highest average potential outcomes.