

Application: Causal Mediation Analysis for Stochastic Interventions

Iván Díaz and Nima Hejazi

2018-12-26

Background

We are interested in assessing the Natural Direct Effect (NDE) and the Natural Indirect Effect (NIE), based on the decomposition of the population mediated intervention mean given in Díaz and Hejazi (n.d.).

To proceed, we'll use as our running example a simple data set from an observational study of the relationship between BMI and kids behavior, distributed as part of the `mma` R package on CRAN. First, let's load this data set and take a quick look at it

```
data(weight_behavior)
dim(weight_behavior)

## [1] 691 15

head(weight_behavior)

##      bmi  age sex  race numpeople car gotosch snack tvhours cmpthours
## 1 18.20665 12.2  F OTHER         5   3      2     1      4         0
## 2 22.78401 12.8  M OTHER         4   3      2     1      4         2
## 3 19.60725 12.6  F OTHER         4   2      4     2     NA        NA
## 4 25.56754 12.1  M OTHER         2   3      2     1      0         2
## 5 15.07408 12.3  M OTHER         4   1      2     1      2         1
## 6 22.98338 11.8  M OTHER         4   1      1     1      4         3
##  cellhours sports exercises sweat  overweigh
## 1         0      2         2      1         0
## 2         0      1         8      2         0
## 3        NA <NA>         4      2         0
## 4         0      2         9      1         1
## 5         3      1        12      1         0
## 6         2      1         1      1         0
```

The documentation for the data set describes it as a “database obtained from the Louisiana State University Health Sciences Center, New Orleans, by Dr. Richard Scribner. He explored the relationship between BMI and kids behavior through a survey at children, teachers and parents in Grenada in 2014. This data set includes 691 observations and 15 variables.”

Unfortunately, the data set contains a few observations with missing values. As these are unrelated to the object of our analysis, we'll simply remove these for the time being. Note that in a real data analysis, we might consider strategies to fully make of the observed data, perhaps by imputing missing values. For now, we simply remove the incomplete observations, resulting in a data set with fewer observations but much the same structure as the original:

```
## [1] 567 15

##      bmi  age sex  race numpeople car gotosch snack tvhours cmpthours
## 1 18.20665 12.2  F OTHER         5   3      2     1      4         0
## 2 22.78401 12.8  M OTHER         4   3      2     1      4         2
## 4 25.56754 12.1  M OTHER         2   3      2     1      0         2
```

## 5	15.07408	12.3	M	OTHER	4	1	2	1	2	1
## 6	22.98338	11.8	M	OTHER	4	1	1	1	4	3
## 8	19.15658	12.1	F	OTHER	3	3	2	1	0	0
##	cellhours	sports	exercises	sweat	overweigh					
## 1	0	1		2	1	0				
## 2	0	0		8	2	0				
## 4	0	1		9	1	1				
## 5	3	0		12	1	0				
## 6	2	0		1	1	0				
## 8	1	0		1	3	0				

For the analysis of this observational data set, we focus on the effect of participating in a sports team (`sports`) on the BMI of children (`bmi`), taking several related covariates as mediators (`snack`, `exercises`, `overweigh`) and all other collected covariates as potential confounders. Considering an NPSEM, we separate the observed variables from the data set into their corresponding nodes as follows

```
Y <- weight_behavior_complete$bmi
A <- weight_behavior_complete$sports
Z <- weight_behavior_complete %>%
  select(snack, exercises, overweigh)
W <- weight_behavior_complete %>%
  select(age, sex, race, numpeople, car, gotosch, tvhours, cmpthours,
         cellhours, sweat)
```

Finally, in our analysis, we consider an incremental propensity score intervention (IPSI), as first proposed by Kennedy (2018), wherein the *odds of participating in a sports team* is modulated by some fixed amount ($0 \leq \delta \leq 1$) for each individual. Such an intervention may be interpreted as the effect of a school program that motivates children to participate in sports teams. To exemplify our approach, we postulate a motivational intervention that increases the odds of participating in a sports team by 25% for each individual:

```
delta_shift_ipsi <- 0.25
```

To easily incorporate ensemble machine learning into the estimation procedure, we rely on the facilities provided in the `sl3` R package (Coyle et al. 2018). For a complete guide on using the `sl3` R package, consider consulting <https://tlverse.org/sl3>, or <https://tlverse.org> (and <https://github.com/tlverse>) for the `tlverse` ecosystem, of which `sl3` is a major part. We construct an ensemble learner using a handful of popular machine learning algorithms below

```
# SL learners used for continuous data (the nuisance parameter M)
hal_contin_lrnr <- Lrnr_hal9001$new(n_folds = 5, fit_type = "glmnet",
                                family = "gaussian",
                                lambda.min.ratio = 0.00001,
                                type.measure = "deviance")
xgboost_lrnr_50_contin <- Lrnr_xgboost$new(nrounds = 50)
xgboost_lrnr_100_contin <- Lrnr_xgboost$new(nrounds = 100)
enet_contin_lrnr <- Lrnr_glmnet$new(alpha = 0.5, family = "gaussian")
ridge_contin_lrnr <- Lrnr_glmnet$new(alpha = 0, family = "gaussian")
lasso_contin_lrnr <- Lrnr_glmnet$new(alpha = 1, family = "gaussian")
fglm_contin_lrnr <- Lrnr_glm_fast$new(family = gaussian())
contin_lrnr_lib <- Stack$new(enet_contin_lrnr, ridge_contin_lrnr,
                           lasso_contin_lrnr, fglm_contin_lrnr,
                           xgboost_lrnr_50_contin, xgboost_lrnr_100_contin,
                           hal_contin_lrnr)
sl_contin_lrnr <- Lrnr_sl$new(learners = contin_lrnr_lib,
                             metalearner = Lrnr_nnls$new())
```

```

# SL learners used for binary data (nuisance parameters G and E in this case)
hal_binary_lrnr <- Lrnr_hal9001$new(n_folds = 5, fit_type = "glmnet",
                                  family = "binomial",
                                  lambda.min.ratio = 0.00001,
                                  type.measure = "deviance")
xgboost_lrnr_50_binary <- Lrnr_xgboost$new(nrounds = 50,
                                           objective = "reg:logistic")
xgboost_lrnr_100_binary <- Lrnr_xgboost$new(nrounds = 100,
                                           objective = "reg:logistic")
enet_binary_lrnr <- Lrnr_glmnet$new(alpha = 0.5, family = "binomial")
ridge_binary_lrnr <- Lrnr_glmnet$new(alpha = 0, family = "binomial")
lasso_binary_lrnr <- Lrnr_glmnet$new(alpha = 1, family = "binomial")
fglm_binary_lrnr <- Lrnr_glm_fast$new(family = binomial())
binary_lrnr_lib <- Stack$new(enet_binary_lrnr, ridge_binary_lrnr,
                             lasso_binary_lrnr, fglm_binary_lrnr,
                             xgboost_lrnr_50_binary, xgboost_lrnr_100_binary,
                             hal_binary_lrnr)
sl_binary_lrnr <- Lrnr_sl$new(learners = binary_lrnr_lib,
                             metalearner = Lrnr_nnls$new())

```

Decomposing the population intervention mean

We may decompose the PIE in terms of a *direct effect (DE)* and an *indirect effect (IE)*:

$$\psi(\delta) = \overbrace{\mathbb{E}\{Y(g, q) - Y(g_\delta, q)\}}^{\text{DE}} + \overbrace{\mathbb{E}\{Y(g_\delta, q) - Y(g_\delta, q_\delta)\}}^{\text{IE}}$$

This decomposition of the PIE as the sum of direct and indirect effects has an interpretation analogous to the corresponding standard decomposition of the average treatment effect. In the sequel, we will compute each of the components of the direct and indirect effects above using appropriate estimators as follows

- For $\mathbb{E}\{Y(g, q)\}$, the sample mean $\frac{1}{n} \sum_{i=1}^n Y_i$ is sufficient;
- for $\mathbb{E}\{Y(g_\delta, q)\}$, a one-step efficient estimator for the pure effect (of altering the exposure mechanism but not the mediation mechanism), as proposed in Díaz and Hejazi (n.d.); and,
- for $\mathbb{E}\{Y(g_\delta, q_\delta)\}$, a one-step efficient estimator for the joint effect (of altering both the exposure and mediation mechanisms), as proposed in Kennedy (2018) and implemented in the `npcausal` R package.

Estimating the joint (mediated) effect

To estimate $\psi_0(\delta) = \mathbb{E}\{Y(g_\delta, q_\delta)\}$, the effect of altering both the exposure and mediation mechanisms, we rely on a one-step nonparametric-efficient estimator, denoted $\hat{\psi}(\delta)$, first proposed by Kennedy (2018). In the case of an incremental propensity score intervention, the estimator $\hat{\psi}(\delta)$ is implemented in the open source and freely available `npcausal` R package: <https://github.com/ehkennedy/npcausal>

```

# let's compute the parameter where both A and Z are shifted
psi_ipsi_fit <- ipsi(y = Y, a = A, x.trt = W, x.out = W,
                   time = rep(1, length(Y)), id = seq_along(Y),
                   delta.seq = delta_shift_ipsi, nsplits = 10)

```

```
##
```

```
|
|
|
```

```
| 0%
```

===	4%
=====	9%
=====	13%
=====	17%
=====	22%
=====	26%
=====	30%
=====	35%
=====	39%
=====	43%
=====	48%
=====	52%
=====	57%
=====	61%
=====	65%
=====	70%
=====	74%
=====	78%
=====	83%
=====	87%
=====	91%
=====	96%
=====	100%

```
# print return object and extract point estimate
psi_ipsi <- psi_ipsi_fit$res$est
psi_ipsi
```

```
## [1] 18.88448
```

Estimating the pure non-mediated effect

As given in Díaz and Hejazi (n.d.), the statistical functional identifying the pure effect $\mathbb{E}\{Y(g_\delta, q)\}$, which corresponds to altering the exposure mechanism while keeping the mediation mechanism (and its reliance on the exposure) fixed, is

$$\theta_0(\delta) = \int m_0(a, z, w) g_{0,\delta}(a | w) p_0(z, w) d\nu(a, z, w),$$

for which a one-step nonparametric-efficient estimator is available. The corresponding *efficient influence function* (EIF) with respect to the nonparametric model \mathcal{M} is $D_{\eta,\delta}(o) = D_{\eta,\delta}^Y(o) + D_{\eta,\delta}^A(o) + D_{\eta,\delta}^{Z,W}(o) - \theta(\delta)$. The one-step estimator may be computed using the EIF estimating equation, making use of cross-fitting to circumvent any need for entropy conditions. The resultant estimator is

$$\hat{\theta}(\delta) = \frac{1}{n} \sum_{i=1}^n D_{\hat{\eta}_{j(i)},\delta}(O_i) = \frac{1}{n} \sum_{i=1}^n \left\{ D_{\hat{\eta}_{j(i)},\delta}^Y(O_i) + D_{\hat{\eta}_{j(i)},\delta}^A(O_i) + D_{\hat{\eta}_{j(i)},\delta}^{Z,W}(O_i) \right\},$$

which is implemented in the `medshift` R package (Hejazi and Díaz, n.d.). We make use of that implementation to estimate $\mathbb{E}\{Y(g_\delta, q)\}$ via its one-step estimator $\hat{\theta}(\delta)$ below

```
# let's compute the parameter where A (but not Z) are shifted
theta_eff <- medshift(W = W, A = A, Z = Z, Y = Y,
  delta = delta_shift_ipsi,
  g = sl_binary_lnr,
  e = sl_binary_lnr,
  m = sl_contin_lnr,
  phi = fglm_contin_lnr,
  estimator = "onestep")

theta_eff

## [1] 19.03938
```

Estimating the Natural Direct Effect

Recall that, based on the decomposition outlined previously, the *natural direct effect* (NDE) may be denoted $\beta_{\text{NDE}}(\delta) = \mathbb{E}Y - \theta_0(\delta)$. Thus, an estimator of the NDE, $\hat{\beta}_{\text{NDE}}(\delta)$ may be expressed as a composition of estimators of its constituent parameters:

$$\hat{\beta}_{\text{NDE}}(\delta) = \frac{1}{n} \sum_{i=1}^n Y_i - \hat{\theta}(\delta).$$

Based on the above, we may construct an estimator of the NDE using the quantities already computed

```
# natural direct effect = EY - estimated quantity
nde_est <- mean(Y) - theta_eff
nde_est

## [1] 0.08772927
```

As given above, we have for our estimate of the natural direct effect $\hat{\beta}_{\text{NDE}}(\delta) = 0.088$.

Estimating the Natural Indirect Effect

Similarly to the direct effect, the *natural indirect effect* (NIE) may be denoted $\beta_{\text{NIE}}(\delta) = \theta_0(\delta) - \psi_0(\delta)$, and an estimator of the NIE, $\hat{\beta}_{\text{NIE}}(\delta)$, may be expressed

$$\hat{\beta}_{\text{NIE}}(\delta) = \hat{\theta}(\delta) - \hat{\psi}(\delta),$$

which may be constructed from the quantities already computed

```
# natural indirect effect = estimated quantity - Edward's estimate
nie_est <- theta_eff - psi_ipsi
nie_est
```

```
## [1] 0.1549016
```

Thus, we have for our estimate of the natural indirect effect $\hat{\beta}_{\text{NIE}}(\delta) = 0.155$.

References

- Coyle, Jeremy R, Nima S Hejazi, Ivana Malenica, and Oleg Sofrygin. 2018. *sl3: Pipelines for Machine Learning and Super Learning in R*. <https://github.com/tlverse/sl3>.
- Díaz, Iván, and Nima S Hejazi. n.d. “Causal Mediation Analysis for Stochastic Interventions.”
- Hejazi, Nima S, and Iván Díaz. n.d. *medshift: Causal Mediation Analysis for Stochastic Interventions in R*. <https://github.com/nhejazi/medshift>.
- Kennedy, Edward H. 2018. “Nonparametric Causal Effects Based on Incremental Propensity Score Interventions.” *Journal of the American Statistical Association*. Taylor & Francis.