

R Lab 2 - Identifiability & the Simple Substitution Estimator

Laura B. Balzer

Biostat683 - Intro. to Causal Inference

Goals:

1. Review the steps 1-5 of the Roadmap: (1) specify the causal model, (2) specify the causal question, (3) specify the observed data and its link to the causal model, (4) assess identifiability, and (5) specify a statistical estimand and statistical model.
2. Introduce and implement the simple substitution estimator based on the G-Computation formula.
3. Use simulations to evaluate the properties of estimators.

Next lab:

We will implement the inverse probability of treatment weighted (IPTW) estimator and explore the impact of positivity violations on estimator performance.

Reminder:

This is not an R class. However, software is an important bridge between the statistical concepts and implementation.

1 Background Story

“[The Hunger Games] is written in the voice of sixteen-year-old Katniss Everdeen, who lives in a post-apocalyptic world in the country of Panem where the countries of North America once existed. The Capitol, a highly advanced metropolis, holds hegemony over the rest of the nation. The Hunger Games are an annual event in which one boy and one girl aged 12 to 18 from each of the 12 districts surrounding the Capitol are selected by lottery [as ‘tributes’] to compete in a televised battle in which only one person can survive.” - Source: Wikipedia “The Hunger Games”

Some of the tributes have trained extensively for this tournament. The life experiences of other tributes have resulted in certain abilities/advantages (e.g., strength, tree climbing, marksmanship). Prior to the tournament, a committee of judges assigns a score to each the tribute indicating their probability of winning. Once the tournament starts, forming alliances and sponsorship can aid in survival. A lone victor returns to their district and is showered with wealth and other resources.

Suppose we are interested in the effect of forming an alliance on the probability of surviving through the first 24 hours. We have randomly sampled one tribute from each year of the games. Let $W1$ denote the tribute’s sex with $W1 = 1$ being male and $W1 = 0$ female. Let $W2$ denote the score from the judges. Let A be an indicator that an alliance is formed, and Y be an indicator of survival through the first 24 hours. Finally, let $W3$ be an indicator of whether the tribute receives aid from sponsors during the tournament. Our goal is to evaluate the effect of forming an alliance on the probability of surviving through the first 24 hours.

This study can be translated into the following directed acyclic graph (DAG) shown in Figure 1.

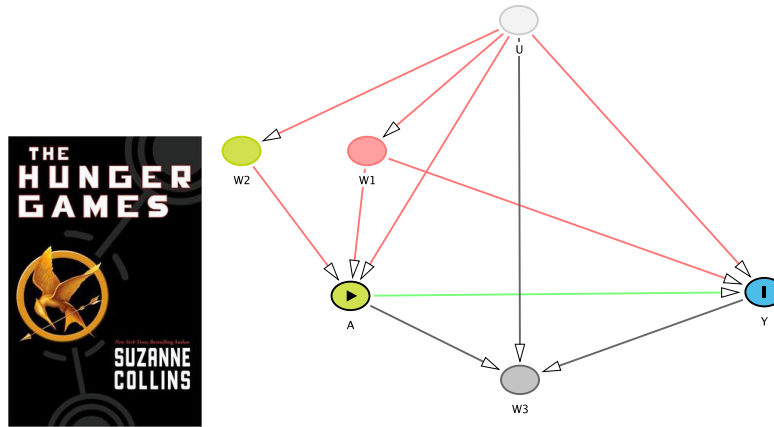


Figure 1: Directed Acyclic Graph for the Hunger Games study.

1. Translate the DAG into the corresponding structural causal model \mathcal{M}^* .
2. Are there any exclusion restrictions?
3. Are there any restrictions on the distribution of the background variables \mathbb{P}_U ? In other words, are there any independence assumptions?
4. Specify the causal question and parameter.
5. Suppose the observed data consist of n independent, identically distributed (i.i.d.) draws of the random variable $O = (W1, W2, A, Y, W3) \sim \mathbb{P}_0$. Specify the link between the SCM and the observed data. Does the SCM place any restrictions on the statistical model \mathcal{M} ?

Solution:

1. Endogenous variables: $X = (W1, W2, A, Y, W3)$
 Background variables: $U = (U_{W1}, U_{W2}, U_A, U_Y, U_{W3}) \sim \mathbb{P}_U$
 Structural equations F :

$$\begin{aligned}
 W1 &= f_{W1}(U_{W1}) \\
 W2 &= f_{W2}(U_{W2}) \\
 A &= f_A(W1, W2, U_A) \\
 Y &= f_Y(W1, A, U_Y) \\
 W3 &= f_{W3}(A, Y, U_{W3})
 \end{aligned}$$

2. In our recursive (i.e., time-ordered) SCM, we have made lots of exclusion restrictions. Sex $W1$ does not affect the judge's score $W2$. The outcome Y is not affected by the score $W2$. Whether the tribute receives aid from sponsors is not a function of sex $W1$ or the score $W2$.
3. There are no independence assumptions. (One may reasonably argue that the unmeasured factors determining biological sex U_{W1} are independent of the others, but we have not made that independence assumption here.)
4. The target causal parameter is the difference in the counterfactual probability of survival through the first 24 hours, if all tributes formed alliances, and the counterfactual probability of survival, if all tributes did not form alliances:

$$\Psi^*(\mathbb{P}^*) = \mathbb{P}^*(Y_1 = 1) - \mathbb{P}^*(Y_0 = 1) = \mathbb{E}^*(Y_1) - \mathbb{E}^*(Y_0)$$

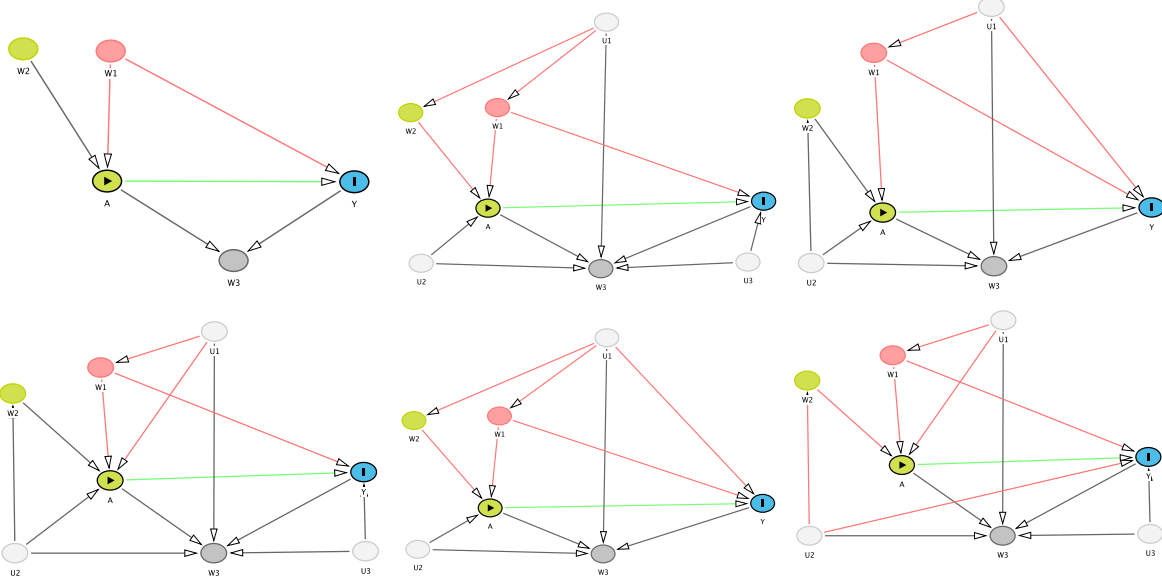


Figure 2: Evaluating the back-door criterion: In the first 4 DAGs, $W1$ alone would satisfy the back-door criterion. In the last 2 DAGs, $W1$ and $W2$ are needed to satisfy the back-door criterion. The needed independence assumptions should be carefully discussed and considered with the help of subject matter experts.

where Y_1 denotes the counterfactual survival under an alliance $A = 1$ and Y_0 denotes the counterfactual survival under no alliance $A = 0$.

5. We assume the observed data $O = (W1, W2, A, Y, W3)$ were generated by sampling n i.i.d. times from a data generating system compatible with \mathcal{M}^* . This provides a link between the causal model \mathcal{M}^* and the observed data O . The distribution of the background variables U and the structural equations F identify the distribution of the endogenous variables X and thus the distribution of the observed data O . We have not placed any restrictions on the statistical model \mathcal{M} , which is thereby non-parametric.

2 Assess identifiability of our causal parameter $\Psi^*(\mathbb{P}^*)$ & Commit to a statistical estimand $\Psi(\mathbb{P}_0)$.

The causal risk difference $\Psi^*(\mathbb{P}^*)$ is not identified under our causal model \mathcal{M}^* . A sufficient, but *not* minimal, identifiability assumption is that all of the unmeasured factors are independent. If the all unmeasured factors were conveniently independent, the back-door criterion would hold conditional on $W1$. Equivalently, the counterfactual outcome Y_a would be conditionally independent of the treatment A , given $W1$; this is the randomization assumption and also called “conditional exchangeability”. Other possible independence assumptions and the corresponding sufficient sets are given in Figure 2.

For the statistical estimand $\Psi(\mathbb{P}_0)$ to be well-defined, we need additional condition of data support, known as the positivity assumption. There must be a positive probability of each exposure condition within each possible values of the adjustment variable $W1$:

$$\begin{aligned} \mathbb{P}_0(A = 1|W1 = 1) &> 0 & \mathbb{P}_0(A = 1|W1 = 0) &> 0 \\ \mathbb{P}_0(A = 0|W1 = 1) &> 0 & \mathbb{P}_0(A = 0|W1 = 0) &> 0 \end{aligned}$$

For this specific example, we need a positive probability of forming and not forming an alliance for both men and women.

Despite lack of identifiability (see the Figure 1), we can still “commit” to an interesting statistical estimand inspired by our scientific/causal question:

$$\begin{aligned}\Psi(\mathbb{P}_0) &= \mathbb{E}_0[\mathbb{E}_0(Y|A=1, W1) - \mathbb{E}_0(Y|A=0, W1)] \\ &= \sum_{w1} [\mathbb{E}_0(Y|A=1, W1=w1) - \mathbb{E}_0(Y|A=0, W1=w1)] \mathbb{P}_0(W1=w1)\end{aligned}$$

Formally, the parameter Ψ is a mapping from the statistical model \mathcal{M} to the parameter space $\Psi : \mathcal{M} \rightarrow \mathbb{R}$. In other words, Ψ is a function with input as a distribution in \mathcal{M} and output a value in the parameter space (e.g., a number). We have not made any new assumptions during identifiability; therefore, our statistical model \mathcal{M} remains non-parametric.

Note: Alternatively, we could have specified the following as our statistical estimand

$$\Psi^{alt}(\mathbb{P}_0) = \mathbb{E}_0 \left[\mathbb{E}_0(Y|A=1, W1, W2) - \mathbb{E}_0(Y|A=0, W1, W2) \right]$$

which would equal the causal effect of interest if the independence assumptions in the last 2 DAGs in Figure 2 conveniently held. However, adjusting for both $(W1, W2)$ would require a different and stronger positivity assumption. For simplicity, we will focus on the previous statistical estimand $\Psi(\mathbb{P}_0)$.

3 A specific data generating process

The above SCM is compatible with many possible data generating processes. Recall \mathcal{M}^* is a causal model for the set of possible distributions $\mathbb{P}_{U,X}$ for (U, X) . Now, consider the a specific data generating process, where each of the exogenous nodes U_{X_i} is drawn independently from the following distributions:

$$\begin{aligned}U_{W1} &\sim Uniform(0, 1) \\ U_{W2} &\sim Normal(\mu = 1, \sigma^2 = 2^2) \\ U_A &\sim Uniform(0, 1) \\ U_Y &\sim Uniform(0, 1) \\ U_{W3} &\sim Uniform(0, 1)\end{aligned}$$

Given the U s, the endogenous variables are deterministically generated as:

$$\begin{aligned}W1 &= \mathbb{I}[U_{W1} < 0.45] \\ W2 &= 0.75 * U_{W2} \\ A &= \mathbb{I}[U_A < \text{logit}^{-1}(-1 + 2.6 * W1 + 0.9 * W2)] \\ Y &= \mathbb{I}[U_Y < \text{logit}^{-1}(-2 + A + 0.7 * W1)] \\ W3 &= \mathbb{I}[U_{W3} < \text{logit}^{-1}(-1 + 1.3 * A + 2.9 * Y)]\end{aligned}$$

The logit^{-1} function is the inverse of the logistic function and given by the `plogis` function in R:

$$\text{logit}(x) = \log\left(\frac{x}{1-x}\right) \quad \text{and} \quad \text{logit}^{-1}(x) = \frac{1}{1 + e^{-x}}$$

We can evaluate the statistical parameter $\Psi(\mathbb{P}_0)$ in closed form:

$$\begin{aligned}\Psi(\mathbb{P}_0) &= \mathbb{E}_0[\mathbb{E}_0(Y|A=1, W1) - \mathbb{E}_0(Y|A=0, W1)] \\ &= \sum_{w1} [\mathbb{E}_0(Y|A=1, W1=w1) - \mathbb{E}_0(Y|A=0, W1=w1)] \mathbb{P}_0(W1=w1)\end{aligned}$$

in closed form. In this particular data generating system (one of many compatible with the SCM), $W1$ (sex) is a Bernoulli random variable with mean 0.45:

$$\mathbb{P}_0(W1 = 1) = \mathbb{E}_0[W1] = 0.45$$

For a given tribute, random error U_{W1} determines whether $W1$ is 1 (male) or 0 (female). Likewise, the binary outcome Y (survival or not) is a Bernoulli random variable with mean given by the logit^{-1} of a function of A and $W1$. Random error U_Y determines whether Y is 1 (survival) or 0 (death). In other words, we know the conditional mean of Y , given A and $W1$:

$$\mathbb{P}_0(Y = 1|A, W) = \mathbb{E}_0(Y|A, W) = \text{logit}^{-1}(-2 + A + 0.7W1)$$

Plugging these functions into the G-Computation formula and evaluating $\Psi(\mathbb{P}_0)$ in closed form, we have:

$$\begin{aligned}\Psi(\mathbb{P}_0) &= \sum_{w1} [\mathbb{E}_0(Y|A = 1, W1 = w1) - \mathbb{E}_0(Y|A = 0, W1 = w1)] P(W1 = w1) \\ &= [\text{logit}^{-1}(-2 + 1 + 0.7 * 1) - \text{logit}^{-1}(-2 + 0 + 0.7 * 1)] 0.45 \\ &\quad + [\text{logit}^{-1}(-2 + 1 + 0.7 * 0) - \text{logit}^{-1}(-2 + 0 + 0.7 * 0)] (1 - 0.45) \\ &= 0.1775\end{aligned}$$

```
> # in R the logit^{-1} function is equal to plogis
> Psi.P0<- (plogis(-2+1+0.7*1) - plogis(-2+0+0.7*1) )*0.45 +
+ (plogis(-2+1+0.7*0) - plogis(-2+0+0.7*0))* 0.55
> Psi.P0
```

```
[1] 0.1774828
```

We can interpret $\Psi(\mathbb{P}_0)$ as the difference in the sex-specific probability of survival with and without an alliance, averaged with respect to the distribution of sex, is **0.1775**. Since the wished-for identifiability assumptions (randomization + positivity) did not hold in our original causal model \mathcal{M}^* , we cannot interpret this parameter causally.

4 Translate this data generating process into simulations

1. **Write a function to generate n i.i.d. observations of random variable $O = (W1, W2, A, Y, W3) \sim \mathbb{P}_0$.** As input to this function, use the sample size n . Within this function, simulate the background factors U and evaluate the structural equations F . Recall the logit^{-1} function in R is `plogis`. Give as output of this function, a data frame (`data.frame`) to hold the observed data.
2. **Set the seed to 252 and generate $n = 5000$ observations using your function. Call the output `Obs`. Use the head and summary functions to examine the output.** The rows are the n repetitions of the data generating process and the columns are the random variables. In other words, the rows are the n participants and the columns are their characteristics.

Solution:

```
> # 1. function to generate the observed data
> generate.data <- function(n){
+   U.W1<- runif(n, min=0, max=1)
+   U.W2<- rnorm(n, mean=1, sd=2)
```

```

+   U.A <- runif(n, min=0, max=1)
+   U.Y <- runif(n, min=0, max=1)
+   U.W3<- runif(n, min=0, max=1)
+   #
+   W1<- as.numeric( U.W1 < 0.45)
+   W2<- 0.75*U.W2
+   A <- as.numeric( U.A < plogis(-1+2.6*W1+0.9*W2))
+   Y <- as.numeric( U.Y < plogis(-2+A+0.7*W1))
+   W3<- as.numeric( U.W3 < plogis(-1 + 1.3*A + 2.9*Y))
+   data.frame(cbind(W1, W2, A, Y, W3))
+ }

> # set the seed & draw the data
> set.seed(252)
> n <- 5000
> Obs <- generate.data(n)
> head(Obs)

  W1      W2 A Y W3
1  0  0.9012694 1 0  1
2  0 -0.2273375 0 0  1
3  1  0.5045215 1 1  1
4  0  1.2643669 1 0  1
5  0  5.3133210 1 1  1
6  1  1.2645811 1 0  0

> summary(Obs)

      W1      W2      A      Y
Min.   :0.0000 Min.   :-4.0200 Min.   :0.0000 Min.   :0.0000
1st Qu.:0.0000 1st Qu.: -0.2676 1st Qu.:0.0000 1st Qu.:0.0000
Median :0.0000 Median :  0.7765 Median :1.0000 Median :0.0000
Mean   :0.4362 Mean   :  0.7356 Mean   :0.6168 Mean   :0.2712
3rd Qu.:1.0000 3rd Qu.:  1.7261 3rd Qu.:1.0000 3rd Qu.:1.0000
Max.   :1.0000 Max.   :  7.5540 Max.   :1.0000 Max.   :1.0000

      W3
Min.   :0.0000
1st Qu.:0.0000
Median :1.0000
Mean   :0.5692
3rd Qu.:1.0000
Max.   :1.0000

>
> #-----

```

It is worth re-iterating that R generates the exogenous input U by calling a pseudorandom number generator to simulate a $Uniform(0, 1)$ variable and then transforming it to correspond with a draw from the specified distribution. We went through an equivalent process above when we drew U_{W1} from a Uniform distribution and then using an indicator function to generate the Bernoulli random variable $W1$ with probability $p = 0.45$. In some cases, we can simplify the R code and directly simulate the observed data as follows:

```

> W1 <- rbinom(n, size=1, prob=0.45)
> W2 <- 0.75*rnorm(n, mean=1, sd=2)
> A <- rbinom(n, size=1, prob=plogis(-1+2.6*W1+0.9*W2))
> Y <- rbinom(n, size=1, prob=plogis(-2+A+0.7*W1))
> W3 <- rbinom(n, size=1, prob=plogis(-1+1.3*A+2.9*Y))

```

Note: this is NOT recommended, because we are changing the U s with each draw. As a result, we will not have consistency: $Y_a \neq Y$ when $A = a$.

5 Simple substitution estimator based on the G-Computation formula (a.k.a., parametric G-comp)

In reality, we usually do not know the true distribution of the observed data \mathbb{P}_0 . Instead, we only have a sample of n i.i.d. observations of O from \mathbb{P}_0 . An intuitive estimator of the statistical estimand $\Psi(\mathbb{P}_0)$ is the simple substitution estimator based on the G-Computation formula. Briefly, the algorithm estimates the relevant parts of the observed data distribution \mathbb{P}_0 and plugs them into the parameter mapping Ψ :

1. Estimate the conditional mean $\mathbb{E}_0(Y|A, W)$ using the observed data as input.
2. Estimate the marginal distribution of baseline covariates $\mathbb{P}_0(W)$ using the observed data as input.
3. Substitute these estimates into the target parameter mapping:

$$\hat{\Psi}(\mathbb{P}_n) = \sum_w [\hat{\mathbb{E}}(Y|A=1, W=w) - \hat{\mathbb{E}}(Y|A=0, W=w)] \hat{\mathbb{P}}(W=w)$$

where \mathbb{P}_n denotes the empirical distribution, which puts weight $1/n$ on each copy O_i , $i = 1, \dots, n$, and where W denotes our adjustment set.

We will always use the sample proportion to estimate the covariate distribution $\mathbb{P}_0(W)$ and therefore can express our simple substitution estimator as

$$\begin{aligned} \hat{\Psi}(\mathbb{P}_n) &= \sum_w [\hat{\mathbb{E}}(Y|A=1, W=w) - \hat{\mathbb{E}}(Y|A=0, W=w)] \times \frac{1}{n} \sum_{i=1}^n \mathbb{I}(W_i = w) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_w [\hat{\mathbb{E}}(Y|A=1, W=w) - \hat{\mathbb{E}}(Y|A=0, W=w)] \times \mathbb{I}(W_i = w) \\ &= \frac{1}{n} \sum_{i=1}^n [\hat{\mathbb{E}}(Y|A=1, W_i) - \hat{\mathbb{E}}(Y|A=0, W_i)] \end{aligned}$$

Formally, an estimator $\hat{\Psi}$ is a mapping from the set of possible empirical distributions \mathbb{P}_n to the parameter space (\mathbb{R}) . In other words, $\hat{\Psi}$ is a function with input as the observed data (a realization of \mathbb{P}_n) and output a value in the parameter space (e.g., a number). The estimator should respect the statistical model \mathcal{M} , which is non-parametric. In other words, we should not make any unfounded assumptions about the observed data distribution \mathbb{P}_0 .

5.1 Implementation with parametric regression

Consider the following parametric regression to describe the conditional expectation of the outcome $\mathbb{E}_0(Y|A, W1)$:

$$\mathbb{E}_0(Y|A, W1) = \mathbb{P}_0(Y=1|A, W1) = \text{logit}^{-1}(\beta_0 + \beta_1 A + \beta_2 W1 + \beta_3 A \times W1)$$

Since both $W1$ and A are binary, this regression is “saturated”. Fitting its coefficients and predicting outcomes would be equivalent to using the NPMLE (i.e., taking the mean outcome within in each strata of $(W1, A)$.)

1. Use the `glm` function to fit the conditional mean function $\mathbb{E}_0(Y|A, W1)$ with logistic regression. Be sure to specify the arguments `family='binomial'` and `data=Obs`.
Hint: To get interaction terms, try the formula $Y \sim A + W1 + A * W1$.
2. Copy the data set `Obs` into two new data frames `txt` and `control`. Then set `A=1` for all units in `txt` and `A=0` for all units in `control`.
Hint: Columns of a data frame can be accessed with the `$` operator.
3. Use the `predict` function to get the expected outcome for each individual i under the intervention $\hat{\mathbb{E}}(Y|A = 1, W1_i)$. Be sure to specify the arguments `newdata=txt` and the `type='response'`. We use `type='response'` to get back predicted probabilities (instead of log-odds ratios).
4. Use the `predict` function to get the expected outcome for each individual i under the control $\hat{\mathbb{E}}(Y|A = 0, W1_i)$. Be sure to specify the arguments `newdata=control` and the `type='response'`.
5. Evaluate the statistical parameter by substituting the predicted outcomes into the G-Computation formula. As previously discussed, the sample proportion is a non-parametric maximum likelihood estimator of the marginal distribution of $W1$. So we can just take the empirical mean of the difference in the predicted outcomes for each participant:

$$\hat{\Psi}(\mathbb{P}_n) = \frac{1}{n} \sum_{i=1}^n \left[\hat{\mathbb{E}}(Y|A = 1, W1_i) - \hat{\mathbb{E}}(Y|A = 0, W1_i) \right]$$

Solution:

```
> #1. Estimate the conditional mean of Y given the treatment A and W1
> reg.model<- glm(Y ~A + W1+ A*W1, family='binomial', data=Obs)
> reg.model
```

```
Call: glm(formula = Y ~ A + W1 + A * W1, family = "binomial", data = Obs)
```

Coefficients:

(Intercept)	A	W1	A:W1
-1.9857	0.9326	0.6103	0.1258

Degrees of Freedom: 4999 Total (i.e. Null); 4996 Residual

Null Deviance: 5844

Residual Deviance: 5431 AIC: 5439

```
> #2. Copy the original dataset Obs into two new dataframes txt and control
> txt<- control <- Obs
> # set A=1 in the txt dataframe and A=0 in control dataframe
> txt$A <-1
> control$A <- 0
```

```
> # 3 predict the mean outcome for each individual in the sample under the treatment
> predictY.txt<- predict(reg.model, newdata = txt, type='response')
```

```
> # 4. predict the mean outcome for each individual in the sample under the control
> predictY.control<- predict(reg.model, newdata = control, type='response')
> #
> head(cbind(Obs, predictY.txt, predictY.control))
```



```

      W1      W2 A Y W3 predictY.txt predictY.control
1  0  0.9012694 1 0 1    0.2586345    0.1207116
2  0 -0.2273375 0 0 1    0.2586345    0.1207116
3  1  0.5045215 1 1 1    0.4214247    0.2017544
4  0  1.2643669 1 0 1    0.2586345    0.1207116
5  0  5.3133210 1 1 1    0.2586345    0.1207116
6  1  1.2645811 1 0 0    0.4214247    0.2017544

> tail(cbind(Obs, predictY.txt, predictY.control))

      W1      W2 A Y W3 predictY.txt predictY.control
4995  1  2.7779075 1 0 1    0.4214247    0.2017544
4996  0  1.1455181 0 0 0    0.2586345    0.1207116
4997  0  2.0852541 0 0 0    0.2586345    0.1207116
4998  0  0.5176812 1 1 1    0.2586345    0.1207116
4999  0 -1.5372069 0 0 1    0.2586345    0.1207116
5000  1  0.2457420 1 1 1    0.4214247    0.2017544

> # 5. take the mean of the predicted outcomes to average over the distribution of W1
> mean(predictY.txt - predictY.control)

[1] 0.1735812

```

The estimated difference in the sex-specific probability of survival with and without an alliance, averaged with respect to the distribution of sex, was 17.4%.

6 Estimate bias, variance and mean squared error (MSE).

Simulations are useful for evaluating the properties of estimators. We will focus on estimating the bias, variance and mean squared error of the simple substitution estimator. Specifically, for $R = 500$ iterations, we will sample $n = 200$ i.i.d. observations from \mathbb{P}_0 , implement the simple substitution estimator based on the G-Computation formula, and save the resulting estimate ψ_n .

1. **Reset the seed to 252; set R to 500 and n to 200.**
2. **Create a vector estimates of length $R = 500$ to hold the estimated values ψ_n obtained at each iteration.**
Hint: Use the `rep` function to create a vector of missing values NA.
3. **Inside a for loop from 1 to $R = 500$, sample n i.i.d. observations of random variable $O = (W1, W2, A, Y, W3)$; implement the simple substitution estimator using the saturated regression model (adjusting for A , $W1$ and their interaction), and save the resulting estimate ψ_n as an entry in the vector estimates.**

Hint: A simple example of a `for` loop is given below. More information on the syntax can be found with `?for`.

```

> # this code creates an empty vector "temp" of length 10
> # in the for loop, the empty values are replaced by 2*index

```

```

> temp<- rep(NA, 10)
> for(i in 1:10) {
+   temp[i]<- 2*i
+ }
> temp

[1]  2  4  6  8 10 12 14 16 18 20

```

4. **What is the average value of the estimates over $R = 500$ repetitions of the data generating process ?**
5. **Estimate the bias of the estimator.** What is the average deviation of the estimate and the truth $\Psi(\mathbb{P}_0)$? Hint: use the `mean` function.

$$\text{Bias}(\hat{\Psi}(\mathbb{P}_n)) = \mathbb{E}_0(\hat{\Psi}(\mathbb{P}_n) - \Psi(\mathbb{P}_0))$$

6. **Estimate the variance of the estimator.** How much do the estimates vary across samples? Hint: use the `var` function.

$$\text{Variance}(\hat{\Psi}(\mathbb{P}_n)) = \mathbb{E}_0 \left(\left(\hat{\Psi}(\mathbb{P}_n) - \mathbb{E}_0[\hat{\Psi}(\mathbb{P}_n)] \right)^2 \right)$$

7. **Estimate the mean squared error of the estimator.** On average, how far are the estimates from the truth?

$$\begin{aligned} \text{MSE}(\hat{\Psi}(\mathbb{P}_n)) &= \mathbb{E}_0 \left(\left(\hat{\Psi}(\mathbb{P}_n) - \Psi(\mathbb{P}_0) \right)^2 \right) \\ &= \text{Bias}^2 + \text{Variance} \end{aligned}$$

Solution:

```

> # 1. setting the seed, number of iterations, and the sample size
> set.seed(252)
> R <- 500
> n <- 200

> # 2. create a vector for the estimates
> estimates<- rep(NA,R)

> # 3. repeat the data generating experiment and estimation algorithm R times
> for(i in 1:R){
+
+   # 1. simulate the sample of n observations
+   Obs<- generate.data(n)
+
+   #2 Estimate the conditional mean of Y given the treatment A and W1
+   reg.model<- glm(Y ~ A + W1 + A*W1, family='binomial', data=Obs)
+
+   #3. Copy the original dataset O into two new dataframes txt and control.
+   txt<- control <- Obs
+   # set A=1 in the txt dataframe and A=0 in control dataframe
+   txt$A <-1

```

```

+   control$A <- 0
+
+   # 4 predict the outcome for each individual in the sample under the treatment
+   predictY.txt<- predict(reg.model, newdata = txt, type='response')
+
+   # 5 predict the outcome for each individual in the sample under the control
+   predictY.control<- predict(reg.model, newdata = control, type='response')
+
+   # 6. take the mean of the predicted outcomes over the distribution of W1
+   estimates[i]<- mean(predictY.txt - predictY.control)
+ }

> # 3-6 average value, bias, variance, and MSE of the estimator
> meanEst<- mean(estimates)
> meanEst

[1] 0.1778601

> bias<- mean(estimates - Psi.P0)
> bias

[1] 0.000377349

> var<- var(estimates)
> var

[1] 0.005127227

> mse<- mean( (estimates-Psi.P0)^2)
> mse

[1] 0.005117115

> # check that mse=bias^2 + var

```

Over $R = 500$ repetitions, the average value of our estimate is 17.79%. The true value of the statistical estimand was 17.75%. The estimator has very low bias of 0.038% and a variance of 0.005. Indeed, its mean squared error is 0.005 and is dominated by the variance.

7 More practice

Suppose the Capitol (people in charge of the Hunger Games) demand that you estimate the conditional mean outcome, according to following parametric regression model:

$$\mathbb{E}(Y|A, W1, W2, W3) = \text{logit}^{-1}(\beta_0 + \beta_1 A + \beta_2 W1 + \beta_3 W2 + \beta_4 W3)$$

In other words, they believe that conditional probability of survival through the first 24 hours is a linear (on the logit scale) function of the exposure (alliance), all the pre-exposure covariates ($W1, W2$) and a post-exposure covariate ($W3$). This “knowledge” changes our SCM \mathcal{M}^* , because it restricts the set of allowed functions f_Y . This “knowledge” also changes our statistical model \mathcal{M} , because it restricts the allowed conditional distributions for Y given $(A, W1, W2, W3)$.

1. Does the back-door criterion hold conditional on $W1, W2$ and $W3$ (assuming independence of the U s)?
2. For $R = 500$ iterations, repeat the above process of sampling $n = 200$ observations, fitting the conditional mean outcome with a main terms logistic model (adjusting now for $A, W1, W2$ and $W3$), obtaining the predicted values under $A = 1$ and $A = 0$, and substituting the estimates into the target parameter mapping. Don't forget to reset the seed.
3. Compare the bias, variance and mean squared error of the substitution estimators when using the previous saturated model (adjusting for $A, W1$ and their interaction) and the Capitol's parametric model (adjusting for $A, W1, W2, W3$ with main terms) to estimate the conditional mean outcome $\mathbb{E}_0(Y|A, W)$.

Solution:

1. No. The back-door criterion does not hold conditional on $W1, W2$ and $W3$. $W3$ is a collider of the intervention A and the outcome Y . By adjusting for $W3$, we are inducing an association between A and Y . The resulting estimate will not equal the causal risk difference - even if the independence assumptions in Figure 2 held!

```
> set.seed(252)
> estimates.miss <- rep(NA,R)
> for(i in 1:R){
+   # simulate the sample of n=200 observations
+   Obs<- generate.data(n=n)
+
+   #2 Estimate the conditional mean of Y given the treatment A, W1, W2,W3
+   miss.reg.model<- glm(Y ~ A + W1 +W2+W3, family='binomial', data=Obs)
+
+   #3. Copy the original dataset O into two new dataframes txtPop and controlPop.
+   txt<- control <- Obs
+   # set A=1 in the txt dataframe and A=0 in control dataframe
+   txt$A <-1
+   control$A <- 0
+
+   # 4 predict the outcome for each individual in the sample under the treatment
+   predictY.txt<- predict(miss.reg.model, newdata = txt, type='response')
+
+   # 5. predict the outcome for each individual in the sample under the control
+   predictY.control<- predict(miss.reg.model, newdata = control, type='response')
+
+   # 6. take the mean of the predicted outcomes over the distribution of W1
+   estimates.miss[i]<- mean(predictY.txt - predictY.control)
+ }

> # Evaluating the estimator
> meanEst.miss<- mean(estimates.miss)
```

```

> bias.miss<- mean(estimates.miss - Psi.P0)
> var.miss<- var(estimates.miss)
> mse.miss<- mean( (estimates.miss-Psi.P0)^2)

> # 3. compare bias, variance, mse of substitution estimators
> estComparison<- data.frame(rbind( c(meanEst, bias, var, mse),
+   c(meanEst.miss, bias.miss, var.miss, mse.miss) ) )
> rownames(estComparison)<- c('Saturated.W1', 'Capitol.W1W2W3')
> colnames(estComparison)<- c('Mean estimate', 'Bias', 'Var', 'MSE')
> signif(estComparison, 2)

```

	Mean estimate	Bias	Var	MSE
Saturated.W1	0.180	0.00038	0.0051	0.0051
Capitol.W1W2W3	0.052	-0.12000	0.0069	0.0230

When the Capitol's regression is used to estimate the conditional mean outcome given the exposure and covariates, the resulting substitution estimator is biased. This is unsurprising as $W3$ (receiving sponsorship) is a collider of the exposure (alliance or not) and the outcome (survival through the first 24 hours). Indeed, the average estimate using this regression model is 5%. Its absolute bias over $R = 500$ repetitions of the experiment is over 300 times higher than when using the previous, saturated regression.