

Covariate selection criteria for controlling confounding bias in a causal study

Seethad Thepepomma^a · Ji-Hyun Kim^{a,1}

^aDepartment of Statistics and Actuarial Science, Soongsil University

(Received April 27, 2016; Revised June 4, 2016; Accepted June 10, 2016)

Abstract

It is important to control confounding bias when estimating the causal effect of treatment in an observational study. We illustrated that the covariate selection in the causal inference is different from the variable selection in the ANCOVA model. We then investigated the three criteria of covariate selection for controlling confounding bias, which can be used when we have inadequate information to draw a complete causal graph. VanderWeele and Shpitser (2011) proposed one of them and claimed it was better than the other two. We show by example that their criterion also has limitations and some disadvantages. There is no clear winner; however, their criterion is better (if some correction is made on its condition) than the other two because it can remove the confounding bias.

Keywords: causal graph, back-door criterion, strongly ignorable, confounding

1. 연구목적

확률화(randomization)가 이루어진 실험 자료가 아닌 관측 자료를 이용해 인과연구(causal study)를 하고자 할 때 치료효과(treatment effect)를 제대로 추정하기 위해 조건화해야 할 변수들을 선택하는 문제는 매우 중요하다. LaLonde (1986)의 예에서와 같이 직업훈련 프로그램 이수 여부가 소득 향상에 미치는 효과를 추정하고자 할 때 프로그램 이수 여부를 피실험자가 선택(self-selection)하므로 프로그램을 이수한 집단과 하지 않은 집단의 특성, 예를 들어 교육수준이나 나이 등의 특성이 다르게 되는데, 이러한 특성의 차이에 대한 고려 없이 단순히 프로그램을 이수한 집단과 이수하지 않은 집단의 평균 소득의 차이만으로 교육의 효과를 추정하면 편향된 추정량을 얻게 된다.

직업훈련 프로그램의 이수 여부와 같이 그 인과효과(causal effect)를 알고 싶은 변수를 처리변수라고 하고, 소득과 같이 비교하고 싶은 크기를 나타내는 변수를 반응변수라고 하자. 이 두 변수 이외에 관측되는 변수들을 공변량(covariates 또는 concomitant variables)이라고 할 때, 치료효과와 비편향 추정량을 얻기 위해 조건화해야 하는 공변량들을 선택하는 방법 또는 기준이 중요하다. 반응변수에 관한 모형에서 변수를 선택할 때는 모형에 의한 예측값이 관측값에 얼마나 가까운가를 보는 것이고, 인과연구에서 변수를 선택할 때는 어떻게 하면 치료효과를 편향되지 않게 추정할 수 있을지를 보는 것이기 때문에 변수를 선택하는 목적이 다르다. 따라서 회귀분석 또는 공분산분석(analysis of covariance; ANCOVA)

¹Corresponding author: Department of Statistics and Actuarial Science, Soongsil University, Sangdo-Ro 369, Dongjak-Gu, Seoul 06978, Korea. E-mail: jxk61@ssu.ac.kr

모형에서 쓰이는 변수선택기준은 인과연구에서 처리효과를 제대로 추정하기 위한 공변량선택기준으로 적합하지 않다. 3절에서 이 점을 예증하였다.

관측연구에서의 인과추론(또는 관측 자료를 이용한 인과연구)에 대한 체계적인 접근법으로 크게 두 가지를 들 수 있는데, Rubin (1974, 1990)의 가상결과(counterfactual 또는 possible outcome) 접근법과 Pearl (1993, 1995, 2009)의 인과 그래프(causal graph)를 이용한 접근법이 그것이다. 이 중에서 Pearl (1993)은 처리효과의 비편향 추정량을 얻기 위한 조건인 ‘뒷문기준(back-door criterion)’이라는 정리를 제시하였는데, 이 정리에 따라 공변량을 선택하면 ‘중첩(confounding)’을 피할 수 있다. 하지만 이 정리는 처리변수와 반응변수, 그리고 모든 공변량들의 인과관계를 파악하여 그래프로 표현할 수 있다는 가정을 하고 있다. 실제 상황에서는 이런 가정이 비현실적일 수 있어 보다 실용적으로 적용할 수 있는 기준이 필요하다. 두 가지 통용되는 선택기준이 있었는데 VanderWeele와 Shpitser (2011)가 기존의 대표적인 두 기준에 결함이 있음을 예를 들어 보이고 새로운 기준을 제안하였다. 그리고 그들의 공변량선택기준이 일정한 전제조건이 만족될 때 제대로 작동한다는 사실을 증명하였다.

하지만 이 새로운 기준에도 결함과 제한이 있음을 이 연구를 통해 보이고자 한다. 그리고 기준에 통용되고 있는 선택기준들도 나름의 장점을 가진다는 것도 보이고자 한다. 관측연구의 인과추론에서 공변량선택 방법이 중요한 주제인 만큼 실용적으로 적용할 수 있는 공변량선택기준의 한계와 적용 시 주의할 점을 정확하게 인식할 필요가 있다. 이 연구를 통해 그런 점들을 지적하고자 한다.

논문의 구성은 다음과 같다. 2절에서 관측연구의 인과추론을 위한 체계적인 접근법 중의 하나인 인과 그래프에 대해 간략히 소개하였다. 그리고 처리효과를 제대로 추정하기 위해 어떤 공변량을 선택해야 하는가를 알려주는 정리인 Pearl (1993)의 ‘뒷문기준정리’에 대해서도 설명하였다. 3절에서 인과연구를 위한 공분산분석모형의 공변량선택기준이 일반적인 회귀모형의 변수선택기준과 다르다는 것을 모의실험을 통해 보였다. 4절에서 관측연구의 인과추론에 쓰이고 있는 실용적인 공변량선택기준 세 가지에 대해 정리하였다. 5절에서 VanderWeele와 Shpitser (2011)의 기준에도 한계가 있음을 지적하고 다른 기준이 더 나은 성능을 보일 수 있는 예를 보였다. 6절에서 전체적인 내용을 요약하고 결론을 서술하였다.

2. 인과 그래프(causal graph)와 뒷문기준(back-door criterion)

관측연구의 인과추론에서 처리효과에 대한 비편향추정량을 얻을 수 있는 조건을 서술한 Pearl (1993)의 정리는 인과 그래프를 그릴 수 있다는 것을 가정한다. 인과 그래프는 처리변수와 반응변수, 그리고 공변량들 사이의 인과관계를 나타낸 그래프인데, 중요한 용어들을 Greenland 등 (1999)을 참조하여 정리하였다.

한 변수 B 가 다른 변수 C 의 ‘직접적인 원인(direct cause)’일 때 $B \rightarrow C$ 또는 $C \leftarrow B$ 로 나타낸다. 변수 A 로부터 Y 까지의 ‘경로(path)’는 $A \rightarrow B \leftarrow C \rightarrow Y$ 와 같이 두 변수 사이에 있는 변수들 사이의 인과관계를 모두 나열한 것이다. ‘방향경로(directed path)’는 $A \rightarrow B \rightarrow C \rightarrow Y$ 와 같이 한쪽 방향으로만 향해 있는 경로를 의미하는데, 시작 위치에 있는 변수를 끝 위치에 있는 변수의 ‘선조(ancestor)’ 또는 ‘원인(cause)’이라고 하며, 반대로 끝 위치에 있는 변수를 시작 위치에 있는 변수의 ‘자손(descendant)’이라고 한다. $A \rightarrow B \rightarrow C$ 의 경우 A 와 B 가 C 의 선조이자 원인이며, C 를 A 와 B 의 자손이라고 한다. A 에서 Y 까지의 경로 중에서 $A \leftarrow B \rightarrow C \rightarrow Y$ 와 같이 A 로 들어오는 화살표가 있는 경로를 ‘뒷문경로(back-door path)’라고 부른다. 한편, 경로 $A \leftarrow B \rightarrow C \leftarrow D \rightarrow Y$ 의 변수 C 와 같이 변수의 양쪽에 모두 들어오는 화살표만 있는 경우에 이 변수를 ‘충돌변수(collider)’라고 부른다. 충돌변수가 적어도 하나 있는 경로를 ‘차단경로(blocked path)’라고 부르고 그렇지 않은 경로를 ‘비차단경로(unblocked path)’라고 부른다. 처리변수 A 에서 반응변수 Y 에 이르는 뒷문경로에 충돌변

수가 있으면 이 경로에서 발생하는 효과가 저절로 차단되므로 그런 이름이 붙여졌다 (만약 충돌변수를 조건화하면 오히려 뒷문경로가 열리게 되어 중첩이 일어나게 되므로 주의해야 한다). 비차단경로를 그대로 두면 처리효과를 추정하는 데에 편향이 생기므로 차단해야 하는데, 경로 상에 있는 적어도 한 개의 변수를 조건화하여 값을 고정시킴으로써 차단이 가능하다 (Pearl, 2009, 11.1.2절).

앞에서 든 모든 예와 같이, 인과 그래프의 모든 경로 상의 이웃하는 두 변수가 $B \rightarrow C$ 또는 $B \leftarrow C$ 와 같이 어느 한쪽으로만 방향이 있는 화살표로 연결이 되어 있고(directed edge), 경로 상의 두 변수가 서로 상대 변수의 원인이 되는 경우가 없을 때(acyclic), 이 인과 그래프를 directed acyclic graphs(DAG)라고 부른다. 뒷문기준을 서술한 정리는 (앞으로 이 정리를 ‘뒷문기준정리’라고 부르기로 한다) DAG에 적용되는 정리이다. 이 정리는 처리효과를 ‘중첩(confounding)’ 없이 제대로 추정할 수 있는 조건을 알려주는 정리이다. 이 정리를 서술하기 위해 필요한 개념들을 먼저 살펴보자.

관측연구에서의 인과추론을 체계화한 Rubin (1974, 1990)의 가상결과 접근법을 이용하여 처리효과와 중첩, 그리고 중첩을 제거할 수 있는 경우를 의미하는 ‘강한 무시가능성(strong ignorability)’을 정의하기로 한다. 동일한 관측대상에게 $A = a$ 가 할당됐을 때 Y 의 가상결과값을 나타내는 변수를 Y_a 라고 할 때 처리효과를

$$\theta = E(Y_1) - E(Y_0)$$

라고 정의할 수 있다. 처리변수가 이항형인 관측연구에서 $A = 0$ 인 관측대상으로부터는 Y_0 를, $A = 1$ 인 관측대상으로부터는 Y_1 만을 관측하게 된다. 이 때 $E(Y_1 | A = 1) - E(Y_0 | A = 0)$ 은 θ 와 다르다. 왜냐하면 $A = 1$ 인 관측대상들과 $A = 0$ 인 관측대상들 사이에 A 의 값 이외의 다른 공변량들의 값의 체계적인 차이가 있을 수 있기 때문이다. 이 때 처리효과, 즉 A 의 Y 에 대한 효과가 ‘중첩’되었다고 한다. 만약 조건부로 주어졌을 때 A 와 Y_a 가 독립이 되는 공변량 (또는 다변량 공변량) S 를 찾을 수 있다면

$$\begin{aligned} & \sum_S [E(Y_1 | A = 1, S = s) - E(Y_0 | A = 0, S = s)]P(S = s) \\ &= \sum_S [E(Y_1 | S = s) - E(Y_0 | S = s)]P(S = s) \\ &= E(Y_1) - E(Y_0) = \theta \end{aligned}$$

이 된다. 위 식은 $S = s$ 이면서 $A = 1$ 인 관측대상의 Y 의 관측값과 $S = s$ 이면서 $A = 0$ 인 Y 의 관측값들을 이용하여 $E(Y_1 | A = 1, S = s) - E(Y_0 | A = 0, S = s)$ 를 추정한 다음, S 에 관해 평균을 추정하면, 즉,

$$\sum_S [E(Y_1 | A = 1, S = s) - E(Y_0 | A = 0, S = s)]P(S = s)$$

를 추정하면 처리효과 θ 를 편향 없이 추정할 수 있게 됨을 의미한다. 이러한 S 가 존재할 때 처리할 당(treatment assignment)은 ‘강하게 무시가능(strongly ignorable)’하다고 정의하고, $A \perp\!\!\!\perp Y_a | S$ 로 표기한다 (엄격하게 정의하려면 $0 < P(A = 1 | S) < 1$ 라는 조건을 추가해야 한다).

정리 2.1 (뒷문기준정리, Pearl (1993)) 공변량들의 집합 S 가 ‘뒷문기준’이라고 부르는 다음 두 조건을 만족할 때 $A \perp\!\!\!\perp Y_a | S$ 가 성립한다.

- (i) S 에 속하는 공변량은 A 의 자손(descendant)이 아니어야 한다.
- (ii) S 는 A 부터 Y 까지의 뒷문경로를 모두 차단한다.

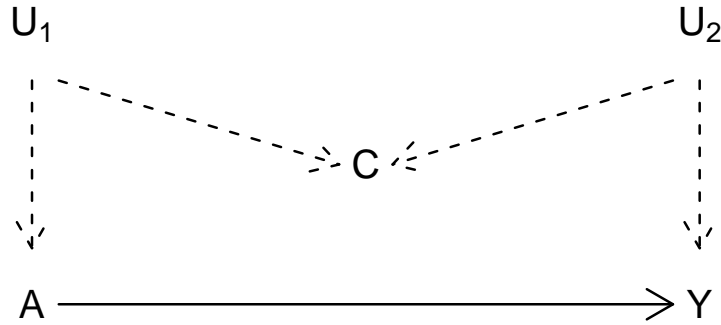


Figure 3.1. Causal graph showing that the variable selection criterion in regression analysis and the covariate selection criterion in causal study may produce different models. The variable C is a collider, and variables U_1 and U_2 are unobserved.

변수들 사이의 인과관계를 모두 파악해서 DAG로 표현할 수 있다면 뒷문기준정리를 이용해서 조건화해야 할 공변량을 선택할 수 있지만, 그렇지 못할 경우 공변량을 어떤 기준으로 선택하면 좋은가에 대한 실용적인 기준이 필요하다. 어떤 실용적인 기준들이 있으며 어떤 장단점들을 갖고 있는가를 알아보기 전에 회귀모형 또는 공분산분석 모형에서의 변수선택과 인과연구에서의 공변량선택이 어떻게 다른가를 먼저 살펴보기로 한다.

3. 공분산분석 모형의 변수선택기준과 인과연구의 공변량선택기준

관측 자료로부터 인과추론을 하고자 할 때 가장 대중적으로 쓰이는 회귀모형인 공분산분석 모형에서 공변량을 선택하는 다양한 방법들이 있다. 하지만 공분산분석 모형은 반응변수에 관한 모형으로서 모형이 반응변수를 얼마나 효과적으로 설명하는가 하는 기준에서 공변량을 선택한다. 반면에 인과추론이 목적인 경우 어떻게 처리효과를 증첩되지 않게 추정할 수 있는가 하는 기준에서 공변량을 선택하므로 두 선택기준의 목적이 다르다. 따라서 인과추론을 하고자 하는 경우 공분산분석 모형의 일반적인 변수선택기준으로 공변량을 선택하면 증첩을 제거하지 못하게 되는 경우가 발생하는데, 이 사실을 예로써 입증하고자 한다.

Figure 3.1의 경우 공변량 C 는 C 의 원인이 되는 서로 다른 두 변수가 존재하는 충돌변수이다. 충돌변수가 존재하는 경로는 차단할 필요가 없으며, 충돌변수를 조건화하면 A 의 효과를 추정함에 있어 오히려 편향을 초래하게 됨이 알려져 있다. 하지만 공변량 C 는 처리변수 A 와 반응변수 Y 둘 다와 관련성(association)이 있으므로 공분산분석 모형에서 통상적인 변수선택을 하면 A 와 함께 모형에 포함되어야 할 공변량으로 잘못 선택된다. 모의실험으로 이를 확인해보자.

Figure 3.1에 대응하는 자료생성 모형으로 다음과 같은 모형을 가정하였다.

$$\begin{aligned} U_1, U_2 &\sim N(0, 1), \\ C &= 2U_1 + 2U_2 + \epsilon_1, \quad \epsilon_1 \sim N(0, 1), \\ A &= I\left(\left(1 + e^{-0.25 U_1}\right)^{-1} > \epsilon_2\right), \quad \epsilon_2 \sim U(0, 1), \\ Y &= 1 + 2A + 2U_2 + \epsilon_3, \quad \epsilon_3 \sim N(0, 2^2). \end{aligned}$$

위 식에서 I 는 지시함수(indicator function)이다. 위 모형에 의해 생성된 자료에 공변량 C 가 포함된

Table 3.1. Results of 500 simulations of fitting two models to data generated by the model of Figure 3.1 with a collider

| Fitted model | Mean and std. error of $\hat{\beta}_1$ | Mean and std. error of $\hat{\beta}_2$ |
|---|--|--|
| $E(Y A, C) = \beta_0 + \beta_1 A + \beta_2 C$ | 1.783 (0.005) | 24.084 (0.049) |
| $E(Y A) = \beta_0 + \beta_1 A$ | 2.008 (0.006) | |

The true value of β_1 is 2.

모형과 포함되지 않은 두 모형

$$E(Y|A, C) = \beta_0 + \beta_1 A + \beta_2 C, \quad (3.1)$$

$$E(Y|A) = \beta_0 + \beta_1 A \quad (3.2)$$

을 적합시켜 결과를 비교하였다. 추정량의 표집오차가 미치는 효과를 작게 하기 위해 표본크기를 2,000으로 크게 하였으며 평균적인 결과를 비교하기 위해 실험을 500번 반복하였다. 뒷문기준정리에 의하면 공변량 C 는 처리효과를 제대로 추정하기 위해서 조건화시키지 말아야 할 변수이다. 이 공변량을 잘못 포함시킨 모형 (3.1)의 결과를 보면 (Table 3.1) 처리효과의 추정값 500개의 평균이 1.783, 표준오차가 0.005로서 참값 2에 대한 편향 추정값을 얻게 됨을 알 수 있다. 그리고 공변량 C 의 계수 추정값 500개의 평균이 24.084, 표준오차가 0.049인데, 이는 공분산분석에서 어떤 변수선택기준을 쓰더라도 공변량 C 가 모형에 포함됨을 의미한다. 즉, 공분산분석 모형의 변수선택기준을 쓰면 공변량 C 는 모형에 포함되고 따라서 처리효과를 편향 추정하게 됨을 알 수 있다. 한편 ‘뒷문기준정리’나 뒤에 설명할 인과연구의 변수선택기준에 의하면 C 는 선택하지 말아야 할 변수가 되고, 따라서 Y 와의 높은 관련성에도 불구하고 C 를 포함시키지 않은 모형을 적용시키면 A 의 효과는 제대로 추정됨을 Table 3.1에서 알 수 있다.

4. 인과추론을 위한 세 가지 공변량선택기준

Pearl (1993)의 뒷문기준정리는 증첩을 제거하기 위해 조건화해야 할 변수를 알려주지만, 인과 그래프, 즉, 처리변수와 반응변수 그리고 공변량들 사이의 관계를 모두 파악하고 있다는 것을 전제로 한다. 모든 변수들 사이의 인과관계를 알지 못할 때에도 쓸 수 있는 보다 실용적인 기준이 필요한데, VanderWeele과 Shpitser (2011)는 기존의 두 기준의 결함을 지적하고 이를 극복하는 새로운 기준을 제안하였다. 기준들의 장단점을 비교하기 전에 먼저 이 세 가지 기준을 정리해 보았다.

그 동안 써 왔던 두 기준이 있는데, 먼저 ‘처리전기준(pretreatment criterion)’이라고 부르는 기준은 처리변수의 값이 정해지기 전에 측정된 공변량들을 모두 선택해야 한다는 기준이다 (Rubin, 2009). 각 공변량이 처리변수나 반응변수와 어떤 관계에 있던 상관하지 않고 측정시점이 처리변수의 측정시점보다 앞서기만 하면 모두 조건화해야 한다는 기준이다. 다음으로 ‘공통원인기준(common cause criterion)’이 있는데, 처리변수의 원인인 동시에 반응변수의 원인인 공변량들을 모두 조건화해야 한다는 기준이다. 이 기준은 처리변수나 반응변수 중 어느 하나에만 영향을 미치는 변수는 조건화할 필요가 없다고 주장하는 기준이다.

VanderWeele과 Shpitser (2011)는 앞의 두 기준을 적용했을 때 모두 편향된 처리효과를 얻게 되는 예를 인과 그래프로 표현하고 모의실험을 통해 입증하였다. 그리고 ‘분리성기준(disjunctive cause criterion)’이라고 이름붙인 새로운 기준을 제안하였는데 이 새로운 기준은 ‘처리 전 측정된 공변량

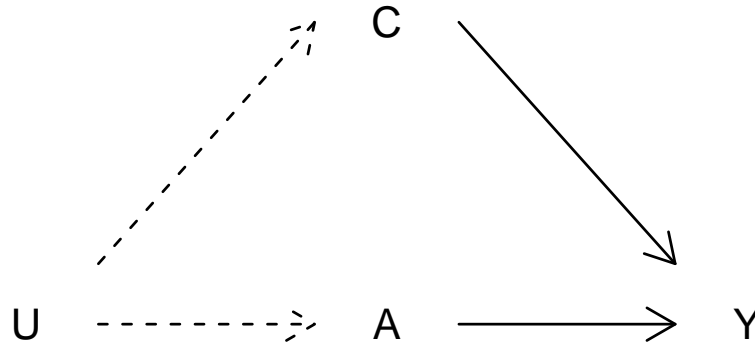


Figure 5.1. Causal graph showing that the disjunctive cause criterion may not be appropriate. The covariate C is assumed to be measured after the treatment is assigned. The variable U is unobserved.

들(pretreatment covariates)’ 중에서 처리변수와 반응변수 둘 중 하나의 원인이거나 또는 둘 모두의 원인인 공변량들을 전부 선택해야 한다는 것이다. 이 기준에 의해 조건화해야 할 공변량들을 선택하면 비편향 추정량을 얻게 된다는 것을 다음과 같은 정리로 서술하고 증명하였다.

정리 4.1 (VanderWeele과 Shpitser, 2011) 처리 할당이 이루어지기 전에 측정되는 모든 공변량들의 집합을 T 라고 하자. 조건화를 통해 중첩을 제거할 수 있는 공변량들의 집합 W 가 존재하며 $W \subseteq T$ 라고 하자. 그러면 분리성기준으로 변수들을 선택해 조건화하면 중첩을 제거할 수 있다.

위 정리에서 분리성기준이 정당한 기준이 될 수 있기 위해서는 처리 전 측정된 공변량들만으로 중첩을 제거할 수 있다는 가정이 필요하다. 이 가정이 만족될 때 그 공변량들 중에서 분리성기준을 만족하는 공변량을 선택하면 된다는 것인데, 이 가정은 아주 강한 가정으로서 이 가정이 만족되지 않을 때 분리성기준에 의한 공변량들이 중첩을 제거한다는 보장이 없다. 또한 이 가정이 만족된다고 하더라도 분리성기준에 의해 선택된 공변량들 중에 불필요한 게 있을 수도 있다. 이러한 한계를 구체적인 예를 들어 보자 한다.

5. 분리성기준의 한계

VanderWeele과 Shpitser (2011)는 처리전기준과 공통원인기준을 적용하면 편향추정량을 얻게 되는 예를 보였다. 그리고 이 예에서 그들이 제안한 분리성기준을 적용하면 비편향추정량을 얻게 되므로 분리성기준이 우월한 기준이라고 주장하였다. 하지만 이 기준에도 한계가 있는데, 예를 들어 살펴보자. Figure 5.1에서 공변량 U 는 관측되지 않는 변수이고 C 는 처리변수 A 의 값이 할당된 후에 관측되는 공변량이라고 하자 (A 를 알코올 섭취량, Y 를 심장질환, U 를 유전적 요인, C 를 운동 습관이라고 할 때, 코호트 연구의 경우 꾸준한 운동 습관에 관한 측정은 알코올 섭취량을 측정된 후에 이루어질 수도 있다). 변수들과의 인과관계가 Figure 5.1과 같을 때, 분리성기준에 의해 공변량을 선택하게 되면 C 는 처리 전(pretreatment) 공변량이 아니므로 선택 대상에서 빠지게 된다 (이 사실은 처리전기준이나 공통원인기준도 마찬가지이다). 하지만 C 를 조건화하지 않으면 뒷문경로 $A \leftarrow U \rightarrow C \rightarrow Y$ 가 차단되지 않으므로 편향이 발생하게 된다. 뒷문 경로를 차단하지 않으면 편향이 발생한다는 사실은 Pearl (1993)의 뒷문기준정리에 의해 이미 알려져 있지만 수치적 증거를 제시하기 위해 모의실험을 실시하였다.

Table 5.1. Results of 500 simulations of fitting two models to data generated by model (5.1) of Figure 5.1 with a back-door path

| Fitted model | Mean and standard error of $\hat{\beta}_1$ |
|--|--|
| Model by disjunctive cause criterion $E(Y A) = \beta_0 + \beta_1 A$ | 7.848 (0.007) |
| Model which correctly blocks the back-door path $E(Y A, C) = \beta_0 + \beta_1 A + \beta_2 C$ | 3.000 (0.002) |

The true value of β_1 is 3.

모의실험에서 자료생성을 위한 모형으로 VanderWeele과 Shpitser (2011)에 있는 모형을 변형하였다.

$$\begin{aligned}
 U &\sim N(0, 1), \\
 C &= 1 + 2U + \epsilon_1, \quad \epsilon_1 \sim N(0, 0.5^2), \\
 A &= I\left(\left(1 + e^{-2U}\right)^{-1} > \epsilon_2\right), \quad \epsilon_2 \sim U(0, 1), \\
 Y &= 3 + 3A + 2C + \epsilon_3, \quad \epsilon_3 \sim N(0, 1).
 \end{aligned} \tag{5.1}$$

표본크기는 표집오차의 영향을 줄이기 위해 2,000으로 하였다. 자료생성 모형을 모르는 상태에서 처리 효과를 추정하기 위해 분리성기준을 적용하여 공변량을 선택하면 공변량 C 를 선택하지 않게 된다. 따라서 처리변수 A 의 효과를 추정하기 위한 모형은

$$E(Y|A) = \beta_0 + \beta_1 A \tag{5.2}$$

가 된다. 하지만 뒷문경로를 차단해야 처리효과와 비편향추정량을 얻을 수 있으므로 올바른 모형은

$$E(Y|A, C) = \beta_0 + \beta_1 A + \beta_2 C \tag{5.3}$$

이다. 모형 (5.2)와 (5.3)를 모형 (5.1)에 의해 생성된 자료에 각각 적합시켜 β_1 을 추정하였으며, 이 작업을 500번 반복해서 실험하였다. 처리효과와 실제 크기는 3이다. 분리성기준에 의한 모형 (5.2)와 제대로 된 모형 (5.3)를 적용하여 얻은 처리효과 β_1 의 추정값의 평균은 각각 7.848과 3.000이었고, 표준오차는 각각 0.007, 0.002였다 (Table 5.1). 따라서 Figure 5.1의 인과 그래프에서 분리성기준은 3이라는 처리효과를 7.85로 편향되게 추정하게 한다는 사실을 확인할 수 있었다. Figure 5.1과 같은 분리성기준의 한계는 ‘처리 전(pretreatment)’ 공변량이라는 분리성기준의 전제조건을 처리 변수의 ‘자손이 아닌(nondescendant)’ 공변량(또는 처리 변수에 ‘영향을 받지 않는(unaaffected by)’ 공변량)으로 바꾸면 해결되지만 전제조건을 만족 여부에 대한 판단이 조금 더 어려워질 수 있다.

분리성기준의 한계를 알려주는 또 다른 예를 들어보자. Figure 5.2에서 공변량 C_1 은 A 의 원인이므로 분리성기준을 적용했을 때 C_2 와 함께 조건화해야 할 변수로 선택된다 (반면에 공통원인기준을 적용하면 C_2 만 선택된다). C_1 은 처리변수를 통해서만 반응변수에 영향을 미치는 변수로서 ‘도구변수(instrumental variable)’라고 부르기도 한다. C_1 은 뒷문경로 상에 있지 않으므로 조건화할 필요가 없는 변수인데 조건화한다고 해서 편향을 가져오지는 않는다. 하지만 반응변수와 무관한 변수를 조건화하면 처리효과와 추정량의 표준오차를 크게 하는 경향이 있다고 알려져 있는데, 이를 모의실험을 통해 확인하고자 한다. 자료 생성 모형은 다음과 같다.

$$C_1, C_2 \sim N(0, 1),$$

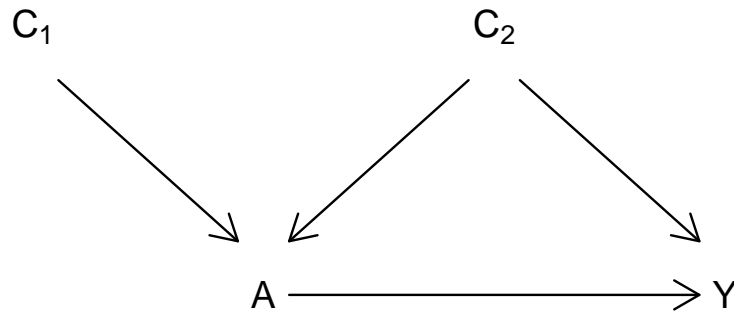


Figure 5.2. Causal graph showing the common cause criterion may be better than the disjunctive cause criterion.

Table 5.2. Results of 500 simulations of fitting two models to data generated by model (5.4) of Figure 5.2 with an instrumental variable

| Fitted model | Mean and std. err. of $\hat{\beta}_A$ | Mean and std. err. of $SE(\hat{\beta}_A)$ |
|--------------------------------------|---------------------------------------|---|
| Model by disjunctive cause criterion | 2.012 (0.017) | 0.397 (0.001) |
| Model by common cause criterion | 2.015 (0.013) | 0.291 (0.001) |

Both models provide unbiased estimates, but the standard errors are different.

$$A = I\left(\left(1 + e^{-(3C_1 + C_2)}\right)^{-1} > \epsilon_1\right), \quad \epsilon_1 \sim U(0, 1), \quad (5.4)$$

$$Y = 1 + 2A + 2C_2 + \epsilon_2, \quad \epsilon_2 \sim N(0, 2^2).$$

이 경우 분리성기준에 의해 선택된 모형

$$E(Y|A, C_1, C_2) = \beta_0 + \beta_A A + \beta_{C_1} C_1 + \beta_{C_2} C_2$$

과 공통원인기준에 의해 선택된 모형

$$E(Y|A, C_2) = \beta_0 + \beta_A A + \beta_{C_2} C_2$$

을 생성된 자료에 각각 적합시켰다. 이 모의실험의 목적은 치료효과 추정량의 표준오차의 크기 비교에 있으므로 표본크기를 200으로 줄였다. 그리고 500번의 반복실험에서 매번 추정량의 표준오차를 구하였다. 표준오차의 크기를 비교하기 위해 500개의 표준오차의 평균과 그 평균의 표준오차를 구하여 Table 5.2에 정리하였다. 분리성기준에 의해 선택된 불필요한 공변량 C_1 을 조건화하면 치료효과의 추정량의 표준오차가 평균적으로 1.36배($0.397/0.291 \doteq 1.36$) 커져서 공통원인기준을 적용했을 때보다 추정의 정밀도가 떨어지게 됨을 알 수 있다.

6. 결론 및 요약

관측연구에서 얻어진 자료로 치료효과에 대한 인과추론을 해야 할 때 증첩을 제거하여 비편향추정량을 얻는 것이 매우 중요하다. 따라서 어떤 공변량들을 조건화해야 하는지를 알려주는 공변량선택기준이 중요한데, 먼저 이 공변량선택의 문제는 회귀모형의 변수선택의 문제와 다르다는 것을 예를 통해 보였다. 만약 변수들 사이의 인과관계를 모두 파악할 수 있다면 ‘뒷문기준정리’에 의해 모든 뒷문경로가 차단되

도록 변수를 선택하면 된다. 이 연구에서 논의한 세 가지 기준은 변수들 사이의 인과관계를 모두 파악하지는 못하고 공변량과 처리변수, 그리고 공변량과 반응변수와의 인과관계만 알고 있는 경우에 적용할 수 있는 실용적 기준이다. 이 실용적 기준들의 한계와 적용 시 주의할 점을 지적하여 관측연구의 인과추론을 위해 공변량을 선택해야 하는 문제에 직면한 연구자에게 도움을 주고자 하였다.

VanderWeele과 Shpitser (2011)는 처리전기준과 공통원인기준이 편향된 추정량을 얻게 하는 예를 보이며서 분리성기준이 다른 두 기준보다 더 나은 기준이라고 주장하였다. 분리성기준은 정리 2의 전제조건이 만족되면 다른 두 기준과 달리 비편향추정량을 얻을 수 있게 한다는 점에서 장점을 갖는다. 하지만 이 전제조건은 아주 강한 조건으로서, 관측된 ‘처리 전 변수’들 중에 뒷문을 차단하는 데 필요한 변수가 빠져 있다면 다른 기준들과 마찬가지로 증첩을 제거하지 못하게 됨을 Figure 5.1과 Table 5.1로 보였다. 또한 공통원인기준과 비교했을 때 추정량의 분산이 더 커질 수 있음도 Figure 5.2와 Table 5.2로 보였다. 한편 처리전기준은 변수들 사이의 인과관계에 대한 가정이 필요 없다는 점에서 장점을 가지므로 세 기준 중에서 완전한 승자는 없다. 하지만 관측된 공변량들 중에서 모든 뒷문경로를 차단할 수 있는 해가 있다는 전제조건이 만족되고, ‘처리 전’ 공변량이라는 조건을 처리변수의 ‘자손이 아닌’ 또는 처리변수에 ‘영향을 받지 않는’ 공변량이라는 조건으로 수정한다면 분리성기준은 다른 두 기준과 달리 비편향추정량을 제공한다는 점에서 좀 더 나은 기준이라고 할 수 있다.

지금까지의 논의와 결론은 뒷문경로를 차단하기 위해 필요한 공변량들이 모두 관측된다는 것을 전제하며, 만약 필요한 공변량들이 관측되지 않으면 어떤 기준을 적용하더라도 뒷문경로를 차단하여 증첩편향을 제거할 수는 없다는 인과연구의 한계를 알고 있어야 한다.

References

- Greenland, S., Pearl, J., and Robins, J. M. (1999). Causal diagrams for epidemiologic research, *Epidemiology*, **10**, 37–48.
- LaLonde, R. J. (1986). Evaluating the econometric evaluations of training programs with experimental data, *American Economic Review*, **76**, 604–620.
- Pearl, J. (1993). Comment: graphical models, causality, and intervention, *Statistical Science*, **8**, 266–269.
- Pearl, J. (1995). Causal diagrams for empirical research, *Biometrika*, **82**, 669–688.
- Pearl, J. (2009). *Causality: Models, Reasoning, and Inference*, Cambridge University Press, Cambridge.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies, *Journal of Educational Psychology*, **66**, 688–701.
- Rubin, D. B. (1990). Formal modes of statistical inference for causal effects, *Journal of Statistical Planning and Inference*, **25**, 279–292.
- Rubin, D. B. (2009). Author’s reply (to Pearl’s, Arvid’s and Sjolander’s letters to the editor), *Statistics in Medicine*, **28**, 1420–1423.
- VanderWeele, T. J. and Shpitser, I. (2011). A new criterion for confounder selection, *Biometrics*, **67**, 1406–1413.

인과연구에서 중첩편향을 제거하기 위한 공변량선택기준

Seethad Thepepomma^a · 김지현^{a,1}

^a승실대학교 정보통계보험수리학과

(2016년 4월 27일 접수, 2016년 6월 4일 수정, 2016년 6월 10일 채택)

요약

관측 자료를 이용한 인과연구에서 관심 있는 처리변수의 효과가 다른 공변량의 효과와 중첩되지 않도록 조건화할 공변량을 선택하는 것이 중요하다. 인과연구에서의 공변량선택 문제는 공분산분석 모형에서의 변수선택 문제와 다르다는 것을 예를 들어 설명하였다. 그리고 모든 변수들 사이의 인과관계를 파악하지 않고도 적용할 수 있는 실용적인 공변량선택기준에 대해 살펴보았다. VanderWeele과 Shpitser (2011)가 새로운 기준을 제안하면서 새로운 기준이 다른 두 기준보다 나은 성능을 보인다고 주장하였는데, 이 기준에도 한계와 단점이 있음을 예증하였다. 새로운 기준이 완전한 기준은 아니지만 조건을 조금 수정하면 다른 두 기준과 달리 중첩을 제거할 수 있다는 점에서 좀 더 나은 기준이라고 할 수 있다.

주요용어: 인과 그래프, 뒷문기준, 강한 무시가능성, 중첩

¹교신저자: (06978) 서울시 동작구 상도로 369, 승실대학교 정보통계보험수리학과. E-mail: jxk61@ssu.ac.kr