

Nonparametric Estimation of IV Model

Jiaqi Yin

May 22, 2018

1 Estimations

We consider binary instrumental variable model Fig.(1) , where we have $A \perp U$, $C \perp A \mid (B, U)$ and A is not independent of B . We have $A, B, C \in \{0, 1\}$. Let $P_{cb \cdot a} = \Pr(C = c, B = b \mid A = a)$. Sufficient and necessary conditions for IV model are

$$\begin{cases} P_{00 \cdot 0} + P_{10 \cdot 1} \leq 1 \\ P_{01 \cdot 0} + P_{11 \cdot 1} \leq 1 \\ P_{10 \cdot 0} + P_{00 \cdot 1} \leq 1 \\ P_{11 \cdot 0} + P_{01 \cdot 1} \leq 1 \end{cases}. \quad (1)$$

Ineq.(1) are also called IV-inequalities.

Suppose our sample size is n , and each individual is i.i.d. Let $n_{cba} = \sum_{i=1}^n \mathbb{I}\{A = a, B = b, C = c\}$. We need to estimate $P_{cb \cdot a}$. The log-likelihood is

$$l(P_{cb \cdot a}) = \sum_{a,b,c} n_{cba} \log P_{cb \cdot a}.$$

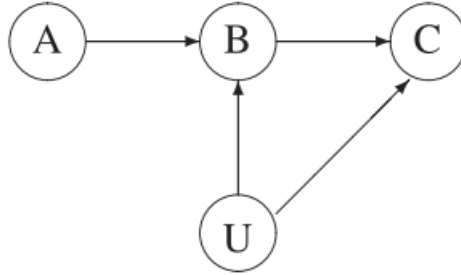


Figure 1: Directed acyclic graph which represents the instrumental variable model.

Of course, we assume $P_{cb \cdot a} > 0$. An optimization question is raised as follows,

$$\left\{ \begin{array}{ll} \min & -\sum_{a,b,c} n_{cba} \log P_{cb \cdot a} \\ \text{subject to} & P_{00 \cdot 0} + P_{10 \cdot 1} \leq 1 \\ & P_{01 \cdot 0} + P_{11 \cdot 1} \leq 1 \\ & P_{10 \cdot 0} + P_{00 \cdot 1} \leq 1 \\ & P_{11 \cdot 0} + P_{01 \cdot 1} \leq 1 \\ & P_{00 \cdot 0} + P_{01 \cdot 0} + P_{10 \cdot 0} + P_{11 \cdot 0} = 1 \\ & P_{00 \cdot 1} + P_{01 \cdot 1} + P_{10 \cdot 1} + P_{11 \cdot 1} = 1 \end{array} \right. ,$$

where objective function and feasible set are convex. The Lagrange function is

$$\begin{aligned} L(P_{abc}, \lambda, \nu) = & -\sum_{a,b,c} n_{cba} \log P_{cb \cdot a} + \lambda_1 (P_{00 \cdot 0} + P_{10 \cdot 1} - 1) + \lambda_2 (P_{01 \cdot 0} + P_{11 \cdot 1} - 1) \\ & + \lambda_3 (P_{10 \cdot 0} + P_{00 \cdot 1} - 1) + \lambda_4 (P_{11 \cdot 0} + P_{01 \cdot 1} - 1) \\ & + \nu_1 (P_{00 \cdot 0} + P_{01 \cdot 0} + P_{10 \cdot 0} + P_{11 \cdot 0} - 1) + \nu_2 (P_{00 \cdot 1} + P_{01 \cdot 1} + P_{10 \cdot 1} + P_{11 \cdot 1} - 1). \end{aligned}$$

Because Slater's condition is satisfied, strong duality holds. Let $p_{cb \cdot a}^*$ and (λ^*, ν^*) be primal and dual optimal points. We further have KKT conditions,

$$\left\{ \begin{array}{l} p_{00 \cdot 0}^* + p_{10 \cdot 1}^* - 1 \leq 0 \\ p_{01 \cdot 0}^* + p_{11 \cdot 1}^* - 1 \leq 0 \\ p_{10 \cdot 0}^* + p_{00 \cdot 1}^* - 1 \leq 0 \\ p_{11 \cdot 0}^* + p_{01 \cdot 1}^* - 1 \leq 0 \\ p_{00 \cdot 0}^* + p_{01 \cdot 0}^* + p_{10 \cdot 0}^* + p_{11 \cdot 0}^* - 1 = 0 \\ p_{00 \cdot 1}^* + p_{01 \cdot 1}^* + p_{10 \cdot 1}^* + p_{11 \cdot 1}^* - 1 = 0 \\ \lambda_i^* \geq 0, \quad i = 1, 2, 3, 4 \\ \lambda_1^* (p_{00 \cdot 0}^* + p_{10 \cdot 1}^* - 1) = 0 \\ \lambda_2^* (p_{01 \cdot 0}^* + p_{11 \cdot 1}^* - 1) = 0 \\ \lambda_3^* (p_{10 \cdot 0}^* + p_{00 \cdot 1}^* - 1) = 0 \\ \lambda_4^* (p_{11 \cdot 0}^* + p_{01 \cdot 1}^* - 1) = 0 \\ -\frac{n_{000}}{p_{00 \cdot 0}^*} + \lambda_1^* + \nu_1^* = 0; \quad -\frac{n_{101}}{p_{10 \cdot 1}^*} + \lambda_1^* + \nu_2^* = 0; \\ -\frac{n_{010}}{p_{01 \cdot 0}^*} + \lambda_2^* + \nu_1^* = 0; \quad -\frac{n_{111}}{p_{11 \cdot 1}^*} + \lambda_2^* + \nu_2^* = 0; \\ -\frac{n_{100}}{p_{10 \cdot 0}^*} + \lambda_3^* + \nu_1^* = 0; \quad -\frac{n_{001}}{p_{00 \cdot 1}^*} + \lambda_3^* + \nu_2^* = 0; \\ -\frac{n_{110}}{p_{11 \cdot 0}^*} + \lambda_4^* + \nu_1^* = 0; \quad -\frac{n_{011}}{p_{01 \cdot 1}^*} + \lambda_4^* + \nu_2^* = 0. \end{array} \right.$$

We start by noting that λ^* act as slack variables in the last four equations, so it can be eliminated, leaving

$$\left\{ \begin{array}{l} p_{00.0}^* + p_{10.1}^* - 1 \leq 0 \\ p_{01.0}^* + p_{11.1}^* - 1 \leq 0 \\ p_{10.0}^* + p_{00.1}^* - 1 \leq 0 \\ p_{11.0}^* + p_{01.1}^* - 1 \leq 0 \\ p_{00.0}^* + p_{01.0}^* + p_{10.0}^* + p_{11.0}^* - 1 = 0 \\ p_{00.1}^* + p_{01.1}^* + p_{10.1}^* + p_{11.1}^* - 1 = 0 \\ \left(\frac{n_{000}}{p_{00.0}^*} - \nu_1^* \right) (p_{00.0}^* + p_{10.1}^* - 1) = 0; \quad \left(\frac{n_{101}}{p_{10.1}^*} - \nu_2^* \right) (p_{00.0}^* + p_{10.1}^* - 1) = 0 \\ \left(\frac{n_{010}}{p_{01.0}^*} - \nu_1^* \right) (p_{01.0}^* + p_{11.1}^* - 1) = 0; \quad \left(\frac{n_{111}}{p_{11.1}^*} - \nu_2^* \right) (p_{01.0}^* + p_{11.1}^* - 1) = 0 \\ \left(\frac{n_{100}}{p_{10.0}^*} - \nu_1^* \right) (p_{10.0}^* + p_{00.1}^* - 1) = 0; \quad \left(\frac{n_{001}}{p_{00.1}^*} - \nu_2^* \right) (p_{10.0}^* + p_{00.1}^* - 1) = 0 \\ \left(\frac{n_{110}}{p_{11.0}^*} - \nu_1^* \right) (p_{11.0}^* + p_{01.1}^* - 1) = 0; \quad \left(\frac{n_{011}}{p_{01.1}^*} - \nu_2^* \right) (p_{11.0}^* + p_{01.1}^* - 1) = 0 \\ \nu_1^* \leq \min \left\{ \frac{n_{000}}{p_{00.0}^*}, \frac{n_{010}}{p_{01.0}^*}, \frac{n_{100}}{p_{10.0}^*}, \frac{n_{110}}{p_{11.0}^*} \right\} \\ \nu_2^* \leq \min \left\{ \frac{n_{101}}{p_{10.1}^*}, \frac{n_{111}}{p_{11.1}^*}, \frac{n_{001}}{p_{00.1}^*}, \frac{n_{011}}{p_{01.1}^*} \right\} \end{array} \right. .$$

If $p_{00.0}^* + p_{10.1}^* - 1 < 0$; $p_{01.0}^* + p_{11.1}^* - 1 < 0$; $p_{10.0}^* + p_{00.1}^* - 1 < 0$; $p_{11.0}^* + p_{01.1}^* - 1 < 0$, we will have

$$\begin{aligned} \nu_1^* &= \frac{n_{000}}{p_{00.0}^*} = \frac{n_{010}}{p_{01.0}^*} = \frac{n_{100}}{p_{10.0}^*} = \frac{n_{110}}{p_{11.0}^*}, \\ \nu_2^* &= \frac{n_{101}}{p_{10.1}^*} = \frac{n_{111}}{p_{11.1}^*} = \frac{n_{001}}{p_{00.1}^*} = \frac{n_{011}}{p_{01.1}^*}. \end{aligned}$$

Further, $\nu_1^* = n_0$, $\nu_2^* = n_1$, and $p_{cb.a}^* = \frac{n_{cba}}{n_a}$.

However, once we consider sampling variability, we could have the case such that $\frac{n_{cba}}{n_a} + \frac{n_{(1-c)b(1-a)}}{n_{1-a}} > 1$, which violates IV Inequalities if we estimate $p_{cb.a}^* = \frac{n_{cba}}{n_a}$. When we have such case, let $p_{c'b'.a'}^* + p_{(1-c')b'.(1-a')}^* = 1$ for some a', b', c' . Consider nonzero $p_{cb.a}^*$ and $\sum_{bc} p_{cb.a}^* = 1$, we could not have $p_{(1-c')b'.a'}^* + p_{c'b'.(1-a')}^* = 1$, $p_{(1-c')(1-b').a'}^* + p_{c'(1-b').(1-a')}^* = 1$ or $p_{c'(1-b').a'}^* + p_{(1-c')(1-b').(1-a')}^* = 1$ (once any one of it holds, we will have zero probability). WOLG, let $c' = 0, b' = 0, a' = 0$, and $\frac{n_{000}}{n_0} + \frac{n_{101}}{n_1} > 1$. Let $p_{00.0}^* + p_{10.1}^* = 1$, $p_{01.0}^* + p_{11.1}^* < 1$, $p_{10.0}^* + p_{00.1}^* < 1$, and $p_{11.0}^* + p_{01.1}^* < 1$.

Further,

$$\begin{aligned} \frac{n_{010}}{p_{01.0}^*} &= \frac{n_{100}}{p_{10.0}^*} = \frac{n_{110}}{p_{11.0}^*} = \nu_1^* < \frac{n_{000}}{p_{00.0}^*} \\ \frac{n_{111}}{p_{11.1}^*} &= \frac{n_{001}}{p_{00.1}^*} = \frac{n_{011}}{p_{01.1}^*} = \nu_2^* < \frac{n_{101}}{p_{10.1}^*}. \end{aligned} \tag{2}$$

The constraints are updated as

$$\begin{cases} p_{00\cdot0}^* + p_{10\cdot1}^* = 1 \\ p_{01\cdot0}^* + p_{11\cdot1}^* + p_{10\cdot0}^* + p_{00\cdot1}^* + p_{11\cdot0}^* + p_{01\cdot1}^* = 1 \\ \frac{n_{010}}{p_{01\cdot0}^*} = \frac{n_{100}}{p_{10\cdot0}^*} = \frac{n_{110}}{p_{11\cdot0}^*} = \nu_1^* < \frac{n_{000}}{p_{00\cdot0}^*} \\ \frac{n_{111}}{p_{11\cdot1}^*} = \frac{n_{001}}{p_{00\cdot1}^*} = \frac{n_{011}}{p_{01\cdot1}^*} = \nu_2^* < \frac{n_{101}}{p_{10\cdot1}^*} \\ p_{01\cdot0}^* + p_{10\cdot0}^* + p_{11\cdot0}^* < 1 \\ p_{01\cdot1}^* + p_{10\cdot1}^* + p_{11\cdot1}^* < 1 \end{cases} \quad (3)$$

From Eq.(3)-1,2, we have $p_{00\cdot0}^* = 1 - p_{10\cdot1}^* = p_{00\cdot1}^* + p_{11\cdot0}^* + p_{01\cdot1}^*$, and $p_{10\cdot1}^* = 1 - p_{00\cdot0}^* = p_{01\cdot0}^* + p_{11\cdot1}^* + p_{10\cdot0}^*$. Plug them into (3)-3,4, we further have

$$\frac{n_{010}}{p_{01\cdot0}^*} = \frac{n_{100}}{p_{10\cdot0}^*} = \frac{n_{110}}{p_{11\cdot0}^*} = \frac{n_0 - n_{000}}{1 - p_{00\cdot0}^*} = \frac{n_0 - n_{000}}{p_{10\cdot1}^*}. \quad (4)$$

$$\frac{n_{111}}{p_{11\cdot1}^*} = \frac{n_{001}}{p_{00\cdot1}^*} = \frac{n_{011}}{p_{01\cdot1}^*} = \frac{n_1 - n_{101}}{1 - p_{10\cdot1}^*} = \frac{n_1 - n_{101}}{p_{00\cdot0}^*}. \quad (5)$$

From above equations, it yields

$$p_{01\cdot0}^* = \frac{n_{010}}{n_0 - n_{000}} p_{10\cdot1}^*, \quad p_{10\cdot0}^* = \frac{n_{100}}{n_0 - n_{000}} p_{10\cdot1}^*, \quad p_{11\cdot0}^* = \frac{n_{110}}{n_0 - n_{000}} p_{10\cdot1}^*; \quad (6)$$

$$p_{11\cdot1}^* = \frac{n_{111}}{n_1 - n_{101}} p_{00\cdot0}^*, \quad p_{00\cdot1}^* = \frac{n_{001}}{n_1 - n_{101}} p_{00\cdot0}^*, \quad p_{01\cdot1}^* = \frac{n_{011}}{n_1 - n_{101}} p_{00\cdot0}^*. \quad (7)$$

From Ineq.(3)-3,4,5,6, we further have

$$\begin{cases} \frac{n_0 - n_{000}}{1 - p_{00\cdot0}^*} < \frac{n_{000}}{p_{00\cdot0}^*} \\ \frac{n_1 - n_{101}}{1 - p_{10\cdot1}^*} < \frac{n_{101}}{p_{10\cdot1}^*} \\ \left(\frac{n_{010}}{n_0 - n_{000}} + \frac{n_{100}}{n_0 - n_{000}} + \frac{n_{110}}{n_0 - n_{000}} \right) p_{10\cdot1}^* < 1 \\ \left(\frac{n_{111}}{n_1 - n_{101}} + \frac{n_{001}}{n_1 - n_{101}} + \frac{n_{011}}{n_1 - n_{101}} \right) p_{00\cdot0}^* < 1 \end{cases} \quad .$$

Finally, we have

$$\max \left\{ 0, 1 - \frac{n_{101}}{n_1} \right\} < p_{00\cdot0}^* < \min \left\{ \frac{n_{000}}{n_0}, 1 \right\} \quad (8)$$

Next, we need to maximize $l(p_{cb\cdot a}^*) = \sum_{a,b,c} n_{cba} \log p_{cb\cdot a}^*$, and we further plug Eq.(6)(7) into the log-likelihood

A	B	C	Count
0	0	0	150
0	1	0	50
0	0	1	100
0	1	1	200
1	0	0	50
1	1	0	325
1	0	1	100
1	1	1	25

Table 1: fake data set

$$\begin{aligned}
l(p_{00.0}^*) &= \sum_{a,b,c} n_{cba} \log p_{cb.a}^* \\
&= n_{101} \log p_{10.1}^* + n_{010} \log \left(\frac{n_{010}}{n_0 - n_{000}} p_{10.1}^* \right) + n_{100} \log \left(\frac{n_{100}}{n_0 - n_{000}} p_{10.1}^* \right) + n_{110} \log \left(\frac{n_{110}}{n_0 - n_{000}} p_{10.1}^* \right) \\
&\quad + n_{000} \log p_{00.0}^* + n_{111} \log \left(\frac{n_{111}}{n_1 - n_{101}} p_{00.0}^* \right) + n_{001} \log \left(\frac{n_{001}}{n_1 - n_{101}} p_{00.0}^* \right) + n_{011} \log \left(\frac{n_{011}}{n_1 - n_{101}} p_{00.0}^* \right) \\
&= (n_{101} + n_{010} + n_{100} + n_{110}) \log (1 - p_{00.0}^*) + n_{010} \log \left(\frac{n_{010}}{n_0 - n_{000}} \right) + n_{100} \log \left(\frac{n_{100}}{n_0 - n_{000}} \right) + n_{110} \log \left(\frac{n_{110}}{n_0 - n_{000}} \right) \\
&\quad + (n_{000} + n_{111} + n_{001} + n_{011}) \log p_{00.0}^* + n_{111} \log \left(\frac{n_{111}}{n_1 - n_{101}} \right) + n_{001} \log \left(\frac{n_{001}}{n_1 - n_{101}} \right) + n_{011} \log \left(\frac{n_{011}}{n_1 - n_{101}} \right)
\end{aligned}$$

The derivative of $l(p_{00.0}^*)$ is

$$l'(p_{00.0}^*) = \frac{n_{000} + n_{111} + n_{001} + n_{011}}{p_{00.0}^*} - \frac{n_{101} + n_{010} + n_{100} + n_{110}}{1 - p_{00.0}^*}.$$

Let $l'(p_{00.0}^*) = 0$, we have maximizer $\tilde{p}_{00.0} = \frac{n_{000} + n_{111} + n_{001} + n_{011}}{n}$. When $p_{00.0}^* < \tilde{p}_{00.0}$, $l(p_{00.0}^*)$ increases; when $p_{00.0}^* > \tilde{p}_{00.0}$, $l(p_{00.0}^*)$ decreases.

Denote the set of $p_{00.0}^*$ in (8) as \mathcal{C} . Follow the following steps to have $p_{00.0}^*$. First, according to the data, have the constraint set \mathcal{C} of $p_{00.0}^*$; second, $p_{00.0}^* = \arg \min_{p \in \mathcal{C}} |\tilde{p}_{00.0} - p|$.

2 Simulation

Suppose we have dataset (fake).

First, we calculate the empirical distribution

$$\begin{aligned}
\pi_{emp} &= (\hat{p}_{00.0}, \hat{p}_{10.0}, \hat{p}_{01.0}, \hat{p}_{11.0}, \hat{p}_{00.1}, \hat{p}_{10.1}, \hat{p}_{01.1}, \hat{p}_{11.1}) \\
&= \left(\frac{150}{500}, \frac{50}{500}, \frac{100}{500}, \frac{200}{500}, \frac{50}{500}, \frac{325}{500}, \frac{100}{500}, \frac{25}{500} \right) \\
&= (0.3, 0.1, 0.2, 0.4, 0.1, 0.65, 0.2, 0.05).
\end{aligned}$$

Find that $\hat{p}_{11.0} + \hat{p}_{01.1} = 0.4 + 0.65 = 1.05 > 1$, which violates the IV-inequality. After applying Section 1 result,

we have the restricted estimation of probability distribution as

$$\begin{aligned}\tilde{\pi}_{obs} &= (\tilde{p}_{00,0}, \tilde{p}_{10,0}, \tilde{p}_{01,0}, \tilde{p}_{11,0}, \tilde{p}_{00,1}, \tilde{p}_{10,1}, \tilde{p}_{01,1}, \tilde{p}_{11,1}) \\ &= (0.312, 0.104, 0.208, 0.375, 0.107, 0.625, 0.214, 0.054).\end{aligned}$$

We further generate data according to multinomial distribution $(\tilde{n}_{000}, \tilde{n}_{100}, \tilde{n}_{010}, \tilde{n}_{110}) \sim Multinom(n_0, (\tilde{p}_{00,0}, \tilde{p}_{10,0}, \tilde{p}_{01,0}, \tilde{p}_{11,0}))$, $(\tilde{n}_{001}, \tilde{n}_{101}, \tilde{n}_{011}, \tilde{n}_{111}) \sim Multinom(n_1, (\tilde{p}_{00,1}, \tilde{p}_{10,1}, \tilde{p}_{01,1}, \tilde{p}_{11,1}))$. The null hypothesis

$$H_0 : \text{the observed data from IV model.}$$

The simulation size is 10,000. We define the statistics as

$$\Lambda = 2 [\ln(\text{likelihood for alternative model}) - \ln(\text{likelihood for null model})].$$

The null model here is the restricted model, and alternative model is the empirical model. We denote $\Lambda(\pi', \tilde{\pi}_{obs}) = 2 \left(\sum_{a,b,c} n'_{cba} \log \frac{n'_{cba}}{n_a} - \sum_{a,b,c} n_{cba} \log \tilde{p}_{cb \cdot a} \right)$ in order to emphasize the parameters (distributions). We finally have

$$P(\Lambda(\pi', \tilde{\pi}_{obs}) \geq \Lambda(\pi_{emp}, \tilde{\pi}_{obs})) = 0.1168$$

3 Average Causal Effect

$$\begin{aligned}\text{ACE}(D \rightarrow Y) &\geq \max \left\{ \begin{array}{c} p_{00 \cdot 0} + p_{11 \cdot 1} - 1 \\ p_{00 \cdot 1} + p_{11 \cdot 1} - 1 \\ p_{11 \cdot 0} + p_{00 \cdot 1} - 1 \\ p_{00 \cdot 0} + p_{11 \cdot 0} - 1 \\ 2p_{00 \cdot 0} + p_{11 \cdot 0} + p_{10 \cdot 1} + p_{11 \cdot 1} - 2 \\ p_{00 \cdot 0} + 2p_{11 \cdot 0} + p_{00 \cdot 1} + p_{01 \cdot 1} - 2 \\ p_{10 \cdot 0} + p_{11 \cdot 0} + 2p_{00 \cdot 1} + p_{11 \cdot 1} - 2 \\ p_{00 \cdot 0} + p_{01 \cdot 0} + p_{00 \cdot 1} + 2p_{11 \cdot 1} - 2 \end{array} \right\} \\ \text{ACE}(D \rightarrow Y) &\leq \min \left\{ \begin{array}{c} 1 - p_{10 \cdot 0} + p_{01 \cdot 1} \\ 1 - p_{01 \cdot 0} + p_{10 \cdot 1} \\ 1 - p_{01 \cdot 0} + p_{10 \cdot 0} \\ 1 - p_{01 \cdot 1} + p_{10 \cdot 1} \\ 2 - 2p_{01 \cdot 0} - p_{10 \cdot 0} - p_{10 \cdot 1} - p_{11 \cdot 1} \\ 2 - p_{01 \cdot 0} - 2p_{10 \cdot 0} - p_{00 \cdot 1} - p_{01 \cdot 1} \\ 2 - p_{10 \cdot 0} - p_{11 \cdot 0} - 2p_{01 \cdot 1} - p_{10 \cdot 1} \\ 2 - p_{00 \cdot 0} - p_{01 \cdot 0} - p_{01 \cdot 1} - 2p_{10 \cdot 1} \end{array} \right\}\end{aligned}$$