

Matching Points: Supplementing Instruments with Covariates in Triangular Models

Junlong Feng ¹

¹Hong Kong University of Science and Technology

Motivation & Objective

A Common Problem in Practice. In economic applications, it is common to have a discrete endogenous variable (D) and an instrument (Z) taking on **fewer** values:

- Return to education. D : Multiple levels of education. Z : Whether lived near a college or not (binary).
- Program evaluation. D : Multiple training programs. Z : A lottery granting access to a certain program (binary).

No Sufficient Variation in IV, No Identification. Let Y , \mathbf{X} and U be the outcome variable, covariates, and the unobservable. The model $Y = g(D, \mathbf{X}, U)$ is under-identified by the IV when the cardinality of the support of D is greater than that of Z ($|S(D)| > |S(Z)|$). **For example** (Newey and Powell, 2003):

- $S(D) = \{1, 2, 3\}$ and $S(Z) = \{0, 1\}$.
- Suppose $g(D, \mathbf{X}, U) = \sum_{d \in S(D)} \mathbb{1}(D = d)m_d^*(\mathbf{X}) + U$.
- For some covariates value $\mathbf{X} = \mathbf{x}_0$, if $\mathbb{E}_{U|\mathbf{X}Z}(\mathbf{x}_0, Z) = 0$, then for every $z \in S(Z)$,

$$\sum_{d \in S(D)} m_d^*(\mathbf{x}_0)p_d(\mathbf{x}_0, z) = \mathbb{E}_{Y|\mathbf{X}Z}(\mathbf{x}_0, z) \quad (1)$$

where $p_d(\mathbf{x}, z) \equiv \mathbb{P}(D = d|\mathbf{X} = \mathbf{x}, Z = z)$ is the generalized propensity score.

- TWO linear equations** ($z = 0, 1$) for **THREE unknowns** ($m_1^*(\mathbf{x}_0), m_2^*(\mathbf{x}_0), m_3^*(\mathbf{x}_0)$): **Underidentified.**

Objective: Achieve point identification by supplementing instruments with covariates.

Model

A triangular model with $S(D) = \{1, 2, 3\}$ and $S(Z) = \{0, 1\}$:

$$\text{Outcome eq.: } Y = \sum_d \mathbb{1}(D = d)m_d^*(\mathbf{X}) + U \quad (\text{SP})$$

$$\text{OR: } Y = \sum_d \mathbb{1}(D = d)g_d^*(\mathbf{X}, U) \quad (\text{NSP})$$

$$\text{Selection eq.: } D = d \text{ iff } h_d(\mathbf{X}, Z, \mathbf{V}) = 1, \sum_d h_d = 1 \quad (\text{SL})$$

- The outcome heterogeneity U is a scalar.
 - D -dependent outcome heterogeneity is allowed. See the paper for more details.
- The selection heterogeneity can be vector-valued.
- No notion of monotonicity is required for selection.

Main Idea

Take Model-SP as an example.

- Consider a different value of covariates \mathbf{x}_m such that equation (1) also holds at \mathbf{x}_m .
- Four equations but six unknowns.
- Identification is possible if $\Delta_d(\mathbf{x}_0, \mathbf{x}_m) \equiv m_d^*(\mathbf{x}_m) - m_d^*(\mathbf{x}_0)$ is known for all $d \in S(D)$: One can substitute $m_d^*(\mathbf{x}_m) = m_d^*(\mathbf{x}_0) + \Delta_d(\mathbf{x}_0, \mathbf{x}_m)$ into the equations for $m_d^*(\mathbf{x}_m)$, and then:
 - Only **three** unknowns.

Such special covariate values \mathbf{x}_m are **matching points** of \mathbf{x}_0 . The paper shows how to find them for given \mathbf{x}_0 s and how identification is restored using them.

Identification of Model-SP

Definition. Matching point \mathbf{x}_m of \mathbf{x}_0 : for $z, z' \in S(Z)$ and for all $d \in S(D)$:

- $h_d(\mathbf{x}_0, z, \mathbf{V}) = h_d(\mathbf{x}_m, z', \mathbf{V})$ a.s.
- $\mathbb{E}_{U|\mathbf{V}\mathbf{X}Z}(\mathbf{V}, \mathbf{x}_m, Z) = \mathbb{E}_{U|\mathbf{V}\mathbf{X}Z}(\mathbf{V}, \mathbf{x}_0, Z)$ a.s.
- $(\mathbf{V}|\mathbf{x}_m, Z)$ and $(\mathbf{V}|\mathbf{x}_0, Z)$ are identically distributed.

Main message of the definition: Selection patterns are the same at the covariates-IV combinations in a matching pair. Covariates need NOT to be exogenous, but equal dependence of (U, \mathbf{V}) on \mathbf{X} at \mathbf{x}_0 and \mathbf{x}_m is required.

Key implication. Under exogeneity of Z , the definition of a matching point implies that:

$$\begin{aligned} \mathbb{E}_{U|\mathbf{D}\mathbf{X}Z}(d, \mathbf{x}_0, z) &= \mathbb{E}_{U|\mathbf{V}}(h_d(\mathbf{x}_0, z, \mathbf{V}) = 1) = \mathbb{E}_{U|\mathbf{V}}(h_d(\mathbf{x}_m, z', \mathbf{V}) = 1) = \mathbb{E}_{U|\mathbf{D}\mathbf{X}Z}(d, \mathbf{x}_m, z') \\ &\Downarrow \\ \mathbb{E}_{Y|\mathbf{D}\mathbf{X}Z}(d, \mathbf{x}_0, z) - m_d^*(\mathbf{x}_0) &= \mathbb{E}_{Y|\mathbf{D}\mathbf{X}Z}(d, \mathbf{x}_m, z') - m_d^*(\mathbf{x}_m) \end{aligned}$$

Identification. Substituting this into the equations in the form of (1) for \mathbf{x}_m yields the following equation system for $\mathbf{m}^*(\mathbf{x}_0) \equiv (m_1^*(\mathbf{x}_0), m_2^*(\mathbf{x}_0), m_3^*(\mathbf{x}_0))$:

$$\begin{pmatrix} p_1(\mathbf{x}_0, z) & p_2(\mathbf{x}_0, z) & p_3(\mathbf{x}_0, z) \\ p_1(\mathbf{x}_0, z') & p_2(\mathbf{x}_0, z') & p_3(\mathbf{x}_0, z') \\ p_1(\mathbf{x}_m, z) & p_2(\mathbf{x}_m, z) & p_3(\mathbf{x}_m, z) \end{pmatrix} \cdot \mathbf{m}^*(\mathbf{x}_0) = \begin{pmatrix} \mathbb{E}_{Y|\mathbf{X}Z}(\mathbf{x}_0, z) \\ \mathbb{E}_{Y|\mathbf{X}Z}(\mathbf{x}_0, z') \\ \mathbb{E}_{Y|\mathbf{X}Z}(\mathbf{x}_m, z) \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ \sum_{d=1}^3 [\mathbb{E}_{Y|\mathbf{D}\mathbf{X}Z}(d, \mathbf{x}_0, z) - \mathbb{E}_{Y|\mathbf{D}\mathbf{X}Z}(d, \mathbf{x}_m, z')] p_d(\mathbf{x}_m, z) \end{pmatrix}$$

- The combination (\mathbf{x}_m, z) is treated as if it was a third instrumental value.
- The term in red offsets the change caused by conditioning on a different value of the non-excluded covariates.
- Identification is achieved as long as the generalized propensity score matrix is full-rank.
- Over-identification is possible because multiple matching points may exist.

Finding the Matching Point

Two different approaches depending on how much we know about Model-SL.

If the selection model is known/specified. Find the matching point by solving $h_d(\mathbf{x}_0, z, \mathbf{v}) = h_d(\mathbf{x}_m, z', \mathbf{v})$.

Example. Ordered choice with linear cutoffs. $h_1(\mathbf{X}, Z, V) = \mathbb{1}(V < \kappa_1 + Z\alpha + \mathbf{X}'\beta)$ and $h_3(\mathbf{X}, Z, V) = \mathbb{1}(V \geq \kappa_2 + Z\alpha + \mathbf{X}'\beta)$ with $\kappa_1 < \kappa_2$. Then matching points can be obtained by solving $z\alpha + \mathbf{x}'_0\beta = z'\alpha + \mathbf{x}'_m\beta$.

If the selection model is unknown/unspecified. Obtain the matching points by matching the generalized propensity scores under the following assumption:

$$p_d(\mathbf{x}_0, z) = p_d(\mathbf{x}_m, z') \forall d \implies h_d(\mathbf{x}_0, z, \mathbf{V}) = h_d(\mathbf{x}_m, z', \mathbf{V}) \text{ a.s. } \forall d$$

The assumption holds in many widely used selection models like parametric/nonparametric ordered choice or discrete choice models.

Remark. These two approaches apply to both Model-SP and Model-NSP.

Identification of Model-NSP

- $U \sim \text{Unif}[0, 1]$.
- $g_d^*(\mathbf{X}, \cdot) : [0, 1] \mapsto S(Y|d, \mathbf{X})$ strictly increasing for all realizations of \mathbf{X} for all d .
- Matching point (strengthened): $(U, \mathbf{V}|\mathbf{x}_0, Z)$ and $(U, \mathbf{V}|\mathbf{x}_m, Z)$ are identically distributed.

Key implication. For the conditional CDFs $F_{U|\mathbf{D}\mathbf{X}Z}$ and $F_{Y|\mathbf{D}\mathbf{X}Z}$ for all $u \in [0, 1]$ and $d \in S(D)$:

$$\begin{aligned} F_{U|\mathbf{D}\mathbf{X}Z}(u|d, \mathbf{x}_0, z) &= F_{U|\mathbf{V}}(u|h_d(\mathbf{x}_0, z, \mathbf{V}) = 1) = F_{U|\mathbf{V}}(u|h_d(\mathbf{x}_m, z', \mathbf{V}) = 1) = F_{U|\mathbf{D}\mathbf{X}Z}(u|d, \mathbf{x}_m, z') \\ &\Downarrow \\ F_{Y|\mathbf{D}\mathbf{X}Z}(g_d^*(\mathbf{x}_0, u)|d, \mathbf{x}_0, z) &= F_{Y|\mathbf{D}\mathbf{X}Z}(g_d^*(\mathbf{x}_m, u)|d, \mathbf{x}_m, z') \end{aligned}$$

Inverting the CDF yields the relationship between $g_d^*(\mathbf{x}_0, u)$ and $g_d^*(\mathbf{x}_m, u)$ for all u and d . Then for each u , obtain three moment equations $\Psi(\mathbf{g}^*(\mathbf{x}_0, u)) = \mathbf{u}$.

- $\Psi(\cdot) : \mathbb{R}^3 \mapsto \mathbb{R}^3$ is continuous everywhere and strictly increasing on $\Pi_d S(Y|d, \mathbf{x}_0)$.
- A typical element in $\Psi(\mathbf{g}^*(\mathbf{x}_0, u))$ is

$$\sum_{d=1}^3 p_d(\mathbf{x}, z) \cdot F_{Y|\mathbf{D}\mathbf{X}Z}(\varphi_d(g_d^*(\mathbf{x}_0, u); \mathbf{x})|d, \mathbf{x}, z)$$

- $(\mathbf{x}, z) = (\mathbf{x}_0, z), (\mathbf{x}_0, z'), (\mathbf{x}_m, z)$.
- $\varphi_d(y; \mathbf{x}) = F_{Y|\mathbf{D}\mathbf{X}Z}^{-1}(F_{Y|\mathbf{D}\mathbf{X}Z}(y|d, \mathbf{x}_0, z)|d, \mathbf{x}, z_1)$ where (\mathbf{x}_0, z) and (\mathbf{x}, z_1) are a matching pair.

Three moments conditions for three unknowns at each u : Under-identification solved!

A New Global Uniqueness Theorem. These moment conditions nest Chernozhukov and Hansen (2005). We prove a new global uniqueness theorem under weaker conditions.

Theorem. If Y is bounded at least from one side and the Jacobian of Ψ is full rank at $\mathbf{g}^*(\mathbf{x}_0, u)$ for all u , then $\mathbf{g}^*(\mathbf{x}_0, \cdot)$ is the unique solution (path) to $\Psi(\mathbf{y}(u)) = \mathbf{u}$ for all u among all increasing functions on $[0, 1]$.

- Exploited monotonicity and continuity of $g_d^*(\mathbf{X}, \cdot)$ and Ψ .
- The condition is minimal: Only full-rankness is required without which local identification at some u is even not guaranteed.

Idea of Proof. Below is the key idea to show uniqueness among increasing functions whose range is $\Pi_d S(Y|d, \mathbf{x}_0)$. The general case is shown in the paper.

By boundedness of Y , all candidate solution paths start or end at the same point, so they must intersect with the true solution paths at least for once. There are then two cases.

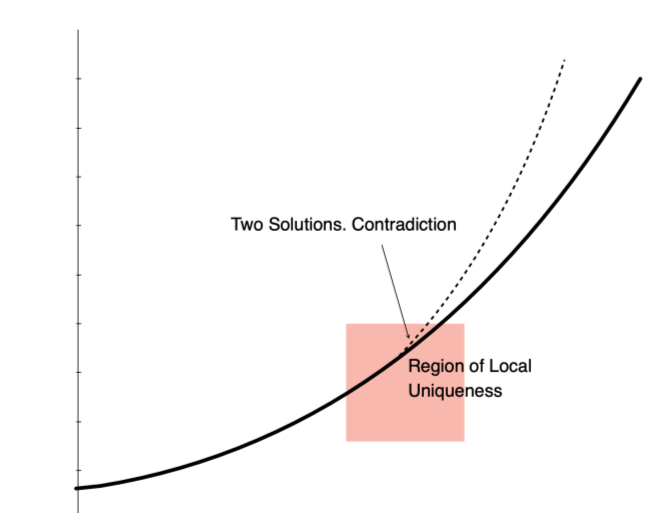


Figure 1. Case 1.

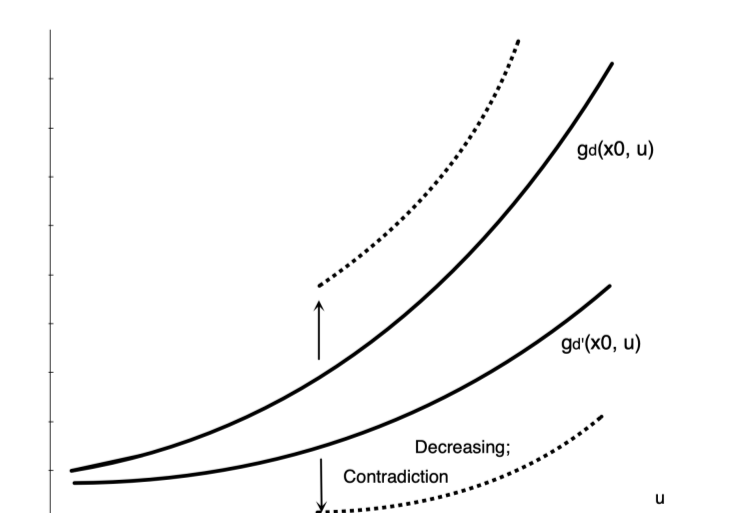


Figure 2. Case 2.

- Case 1: The difference between two solution paths is continuous. Impossible because they must enter the local-uniqueness region by continuity.
- Case 2: The difference is discontinuous. Then the alternative solution path must jump up for some d . But by monotonicity and continuity of Ψ , there must exist a d' for which the alternative solution path jumps down to make the equation still hold, violating monotonicity.