

Informational Content of Factor Structures in Simultaneous Binary Response Models

Shakeeb Khan
Boston College

Arnaud Maurel
Duke University, NBER and IZA

Yichong Zhang
Singapore Management University

September 2019

Abstract

We study the informational content of factor structures in discrete triangular systems. Factor structures have been employed in a variety of settings in cross sectional and panel data models, and in this paper we formally quantify their identifying power in a bivariate system often employed in the treatment effects literature. Our main findings are that imposing a factor structure yields point identification of parameters of interest, such as the coefficient associated with the endogenous regressor in the outcome equation, under weaker assumptions than usually required in these systems. In particular, we show that an exclusion restriction, requiring an explanatory variable in the outcome equation to be excluded from the treatment equation, is no longer necessary for identification. Under such settings, we propose a rank estimator for both the factor loading and the causal effect parameter that are root- n consistent and asymptotically normal. The estimator's finite sample properties are evaluated through a simulation study. We also establish identification results in models with more general factor structures, that are characterized by nonparametric functional forms and multiple idiosyncratic shocks.

Keywords: Factor Structures, Discrete Choice, Causal Effects.

1 Introduction

Factor models (or structures) see widespread and increasing use in various areas of econometrics. This type of structure has been employed in a variety of settings in cross sectional, panel and time series models, and have proven to be a flexible way to model the behavior of and relationship between unobserved components of econometric models. The baseline idea behind factor models is to assume that the dependence across the unobservables is generated by a low-dimensional set of mutually independent random variables (or factors). The applied and theoretical research in econometrics employing factor structures is extensive. In particular, these models are often used in the treatment

We are thankful to seminar participants at Arizona State University, Emory, Michigan State, Shanghai University of Finance and Economics, University of Arizona, and conference participants at the 2015 SEA meetings for helpful comments. Zhang acknowledges the financial support from Singapore Ministry of Education Tier 2 grant under grant no. MOE2018-T2-2-169 and the Lee Kong Chian fellowship.

effect literature as a way to identify the joint distribution of potential outcomes from the marginals, and then recover the distribution of treatment effects from this joint distribution.¹ Factor models have been used in a number of different contexts in applied microeconomics. Notably, factor models have been used in the context of earnings dynamics (Abowd and Card 1989, Bonhomme and Robin 2010), estimation of returns to schooling and work experiences (Ashworth, Hotz, Maurel, and Ransom 2017), as well as cognitive and non-cognitive skill production technology (Cunha, Heckman, and Schennach 2010), among others. All of these papers, with the notable exception of Cunha, Heckman, and Schennach (2010), rely on linear factor models where the unobservables are assumed to be given by the sum of a linear combination of mutually independent factors and an idiosyncratic shock.

In this paper we bring together the literature on factor models with the literature on the identification and estimation of triangular binary choice models (Chesher 2005, Vytlacil and Yildiz 2007, Shaikh and Vytlacil 2011, Han and Vytlacil 2017) by exploring the *informational content* of factor structures in this class of models. Focusing on this class can be well motivated from both an empirical and theoretical perspective. From the former, many treatment effect models fit into this framework as treatment is typically a binary and endogenous variable in the system, whose effect on outcomes is often a parameter the econometrician wishes to conduct inference on. From a theoretical perspective, inference on this type of system can be complicated, if not impossible without strong parametric assumptions, which may not be reflected in the observed data. A semiparametric approach to these models, while desirable from a theoretical point of view because of its generality, often fails to achieve identification of parameter, or at best only do so in sparse regions of the data, thus making inference impractical in practice. In this context, imposing a factor structure may be a useful "in-between" setting, which, at the very least, can be used to gauge the sensitivity of the parametric approach to their stringent assumptions.

We impose a particular factor structure to the two unobservables in this system and explore the informational content of this assumption. Specifically, we assume that the unobservables from the treatment equation (V) and the outcome equation (U) are related through the following factor model:

$$U = \gamma_0 V + \epsilon$$

where ϵ is an unobserved random variable assumed to be distributed independently of V .² Our main finding in this case is that there is indeed informational content of factor structures in the sense that, in contrast to prior literature - notably Vytlacil and Yildiz (2007) - one no longer requires an additional exclusion restriction nor the strong support conditions that are needed for

¹See Abbring and Heckman (2007) for an extensive discussion of factor structures and prior studies using these models in the context of treatment effect estimation.

²While this is our baseline specification, we also examine the informational content of more general factor structures involving nonparametric relationships between unobservables or multiple idiosyncratic errors.

identification in these models without the factor structure. Importantly, our identification results are constructive and translate directly into a rank based estimator of the coefficient associated with the binary endogenous variable.

The rest of the paper is organized as follows. In the next section we formally describe the triangular system with our factor structure, and discuss our main identification results for the parameters of interest in this model. Section 3 then proposes the estimation procedure and establishes its asymptotic properties. Section 4 explores identification in more complicated factor structure models which involve nonparametric relationships between unobservables or multiple idiosyncratic

restrictions on the unobserved variables in the model, $(U; V)$. Such parametric restrictions, such as the often assumed bivariate normality assumption, are not robust to misspecification in the sense that any estimator of β_0 based on these conditions will be inconsistent if $(U; V)$ have a different bivariate distribution.

The established difficulty of identifying β_0 in semi parametric, i.e., "distribution free" models, and the sensitivity of its identification to misspecification in parametric models is what motivates the factor structure we add to the above model in this paper. Specifically, to allow for endogeneity in the form of possible correlation between $U; V$, we augment the model and add the following equation:

$$U = \beta_0 V + \epsilon \tag{2.3}$$

where ϵ is an unobserved random variable, assumed to be distributed independently of $(V; Z_1; Z_2; Z_3)$, and β_0

nonparametrically identify the two marginals by assuming the existence of a full support regressor that is common to both equations. In contrast, our approach does not rely on the existence of a full support common regressor. Under the factor structure assumed in this paper, we bypass the nonparametric identification of the marginals as a whole and directly consider the identification of the structural parameters. It follows that our model cannot be nested by the one-parameter copula model considered by Han and Vytlacil (2017). On the other hand, there exists one-parameter copula models that cannot be decomposed into linear factor structures.⁴ This implies that our model does not nest Han and Vytlacil (2017) either.

To simplify the exposition of our strategy, in this and the following sections we will focus exclusively on the parameters $\beta_0; \beta_1$ and denote the linear indices by $X_1 = Z_1^0 \beta_0 + Z_3^0 \beta_1$ and $X_2 = Z_2^0 \beta_0$, where $Z = (Z_1; Z_2)$. In particular, we treat β_0 as known. In practice, β_0 can be identified and consistently estimated in a first step using a semi-parametric single index estimator such as the one proposed by Klein and Spady (1993). In addition, at the end of this section, we note that we can identify β_0 and β_1 simultaneously with β_0 . Then (2.1) and (2.2) are simplified to

$$Y_1 = 1f X_1 + \beta_1 Y_2 \quad U > 0g \tag{2.4}$$

and

$$Y_2 = 1f X_2 \quad V > 0g \tag{2.5}$$

Our proof will be based on the Assumptions A1 -A5 we state here:

- A1 The parameter $\beta_0 = (\beta_0; \beta_1)$ is an element of a compact subset of \mathbb{R}^2 .
- A2 The vector of unobserved variables, $(U; V)$ is continuously distributed with support on \mathbb{R}^2 and independently distributed of the vector $(Z_1; Z_2; Z_3)$. Furthermore, we assume the unobserved random variables $U; V$ are distributed independently of each other.
- A3 X is continuously distributed with absolute continuous density w.r.t. Lebesgue measure. The density is bounded and bounded away from zero on any compact subset of its support.
- A4 For any constant c , $P(X_2 = c | X_1 = \beta_0 X_1 + \beta_1 X_2) < 1$, where $(X; X_1)$ are an independent copy of $(X; X_1)$.
- A5 $\text{Supp}(X_1 + \beta_0 X_2) \setminus \text{Supp}(X_1 - \beta_0 X_2) \neq \emptyset$.

Before turning to our main identification result, a couple of remarks are in order.

⁴For instance, suppose that $(U; V)$ has a Gaussian copula with correlation ρ , and that the marginal distributions of U and V are uniform $[0; 1]$. It then follows that, denoting by $\Phi(\cdot)$ the standard normal cdf., $\Phi^{-1}(U); \Phi^{-1}(V)$ is bivariate normal with correlation ρ , which in turn yields the following non-linear relationship between U and V : $U = \Phi^{-1}(\rho V + W)$, where W is normally distributed and independent from V .

Remark 2.1.

Further assume we can identify β_0 , and thus, can treat X as an observable. Then, we can identify the choice probability

$$P^{ij}(z_1; z_3; x) = \text{Prob}(Y_1 = i; Y_2 = j | Z_1 = z_1; Z_3 = z_3; X = x)$$

and its derivative w.r.t. x , i.e., $\partial P^{ij}(z_1; z_3; x)$. Then, by the same argument, we can show that

$$\begin{aligned} \partial P^{11}(Z_1; Z_3; X) &= f_V(X) + \partial P^{10}(Z_1; Z_3; X) = f_V(X) = 0 \\ \partial (Z_1^0 \beta_0 + Z_3^0 \beta_0 + \beta_0 X + (Z_1^0 \beta_0 + Z_3^0 \beta_0 + \beta_0 X)) &= 0: \end{aligned}$$

Then, given sufficient variation in X , we can identify $(\beta_0; \beta_0)$ along with $(\beta_0; \beta_0)$ even when all elements of Z_1 and Z_3 are discrete.⁶

An important takeaway from this result is that imposing our factor structure yields point-identification under weaker support condition when compared to the existing literature, and does not require the second exclusion restriction either. In particular, our results yield point-identification of the parameters of interest even in situations where all of the regressors from the outcome equation are discrete. Interestingly, this indicates that, from the selection equation combined with the factor structure that we impose here, we can overturn the non-identification result of Bierens and Hartog (1988) which would apply to the outcome equation alone.

3 Estimation and Asymptotic Properties

The previous section established a point identification result. The identification result is constructive in the sense that it motivates an estimator for the parameters of interest which we describe in detail here.

As we did in Section 2, to simplify exposition, in the following we focus exclusively on the parameters $\beta_0; \beta_0$. Recall the choice probabilities $P^{ij}(x_1; x) = \text{Prob}(Y_1 = i; Y_2 = j | X_1 = x_1; X = x)$ and its second derivative $\partial^2 P^{ij}(x_1; x)$, which can be estimated as we describe below. Another function needed for our identification result is the density function of the unobserved term V , denoted by $f_V(\cdot)$. This is also unknown, but from the structure of our model can be recovered from the derivative with respect to the instrument Z of $E[Y_2 | Z]$, and hence is estimable from the data. Note that the proof of Theorem 2.1 shows that the sign of the index evaluated at two different regressor values, which we denote here by X_1 and X' is determined by the choice probabilities via

$$\partial P^{11}(X_1; X) = f_V(X) + \partial P^{10}(X_1; X) = f_V(X) = 0 \quad \partial (X_1 + X + (X_1 - X)) = 0:$$

⁶An alternative approach to identifying this parameter can be found in Lewbel (2000). In his approach a second equation to model the endogenous variable is not needed, nor is the factor structure we impose. However, he imposes a strong support condition on a variable like Z_3 requiring that it exceeds the length of the unobservable U . As explained in Khan and Tamer (2010), such an approach precludes even bounding β_0 if the support condition on Z_3 is not satisfied.

This motivates us to use maximum rank correlation estimator proposed by Han (1987).

Implementation requires further details to pay attention to. The unknown choice probabilities, their derivatives, and the density of V will be estimated using nonparametric methods, and for this we adopt locally linear methods as they are particularly well suited for estimating derivatives of functions.

With functions and their derivatives estimated in the first stage of our procedure, the second stage plugs in these estimated values into an objective function to be optimized. Specifically, letting $\hat{\theta}$ denote $(\hat{\theta}; \hat{\theta})$, our estimator is of the form:

$$\hat{\theta} = \arg \max_{\theta} Q_{n,2}(\theta) = \sum_{i \in j} \hat{g}_{i,j}(\theta) \quad (3.1)$$

in which

$$\hat{g}_{i,j}(\theta) = [1f_{\theta_2} \hat{P}^{11}(X_{1,i}; X_i) = \hat{F}_V(X_i) + \theta_2 \hat{P}^{10}(X_{1,j}; X_j) = \hat{F}_V(X_j) \quad 0g1f_{\theta_1}(X_{1,i}; X_i; X_{1,j}; X_j) \quad 0g] + [1f_{\theta_2} \hat{P}^{11}(X_{1,i}; X_i) = \hat{F}_V(X_i) + \theta_2 \hat{P}^{10}(X_{1,j}; X_j) = \hat{F}_V(X_j) < 0g1f_{\theta_1}(X_{1,i}; X_i; X_{1,j}; X_j) < 0g];$$

with

$$(X_1; X; X_1; X) = X_1 + X (X_1 \quad X):$$

We note that this estimator falls into the class of those which optimize a nonsmooth U-process involving components estimated nonparametrically in a preliminary stage.⁷ Examples of other estimators in this class can be found in Khan (2001), Abrevaya, Hausman, and Khan (2010), Jochmans (2013), Chen, Khan, and Tang (2016), and our approach to deriving the limiting distribution theory of our estimator will follow along the steps used in those papers. Our limiting distribution theory for this estimator is based on the following regularity conditions:

RK1 θ_0 lies in the interior of Θ , a compact subset of R^2 .

RK2 The index X is continuously distributed with support on the real line, and has a density function which is twice continuously differentiable.

RK3 (Order of smoothness of probability functions and regressor density functions) The functions $P^{k;l;r}(\cdot)$ and $f_{X_1;X}(\cdot)$ (the density function of the random vector $(X_1; X)$) are continuously differentiable of order p_2 , where $p_2 > 5$.

⁷An alternative estimation procedure could be based on the exact relationship in (2.6). Note the equality on the left-hand side of (2.6) is a function of the data alone and not the unknown parameters. The right-hand side equality can then be regarded as a moment condition to estimate the unknown parameters. We describe this estimator and derive its asymptotic properties in the Online Supplement to the paper. While the two estimation approaches will have similar asymptotic properties (root- n consistent, asymptotically normal), we prefer the rank estimator in (3.1) which involves fewer tuning parameters. Furthermore rank type estimators in general are more robust to certain types of misspecification, as pointed out in Khan and Tamer (2018).

RK4 (First stage kernel function conditions) $K(\cdot)$, used to estimate the choice probabilities and their derivatives is an even function, integrating to 1 and is of order p_2 satisfying $p_2 > 5$.

RK5 (Rate condition on first stage bandwidth sequence) The first stage bandwidth sequence H_n used in the nonparametric estimator of the choice probability functions and their derivatives satisfies $\sqrt{n}H_n^{p_2-1} \rightarrow 0$ and $n^{-1}H_n^{-1} \rightarrow 0$.

Based on these conditions, we have the following theorem, whose proof is in Section B of the Supplementary Appendix which characterizes the rate of convergence and asymptotic distribution of the proposed estimator:

Theorem 3.1. *Under Assumptions RK1-RK5,*

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow N(0; V^{-1} \Sigma V^{-1}) \quad (3.2)$$

weighted semi linear regression model as in, e.g., Robinson (1988). Section C in the Supplementary Appendix provides details of how to construct such an estimator and outlines its large sample properties.

4.2 Model with Two Idiosyncratic Shocks and a Bounded Common Factor

We express this model as:

$$\begin{aligned} Y_1 &= \beta_1 X_1 + \beta_0 Y_2 + U \\ Y_2 &= \beta_2 X_2 + V \end{aligned} \tag{4.2}$$

where $U = \beta_0 W + \epsilon_1$, $V = \beta_2 W + \epsilon_2$, and $(W; \epsilon_1; \epsilon_2)$ are mutually independent. First we consider the case $\beta_0 = 1$ and X_1 is binary, because even in this context, for the baseline case with one idiosyncratic shock, we can identify β_0 . But identification of β_0 becomes more difficult in this model, as established in the following theorem

Theorem 4.1. *Suppose (4.2) holds, β_0 is known to be one, X_1 is binary, and W has a bounded support $[a; b]$ such that $0.5 > b - a$ and $1 - (b - a) > \beta_0 > b - a$, then β_0 is not point identified.*

This nonidentification result motivates imposing additional structure on W , and we consider the following model

B1 $U = \beta_0 W + \epsilon_1$ and $V = \beta_2 W + \epsilon_2$.

B2 W is standard normally distributed.

B3 W , ϵ_1 and ϵ_2 are mutually independent.

B4 X has full support.

B5 Denote the density of V as f_V , then f_V does not have a Gaussian component in the sense that

$$f_V \neq \int \phi(\cdot - g) g \, dG \text{ for some density } g \text{ implies that } f_V \neq \int \phi(\cdot - g) \, dG$$

where ϕ is the density for a normal distribution with zero mean and σ^2 variance.

Assumption **B5** effectively assumes that the distribution of ϵ_2 has tail properties different from those of a normal distribution. This type of assumption is made in the deconvolution literature as it is necessary for identification of the target density when the error distribution is not completely known- see, e.g., Butucea and Matias (2005).⁸ The importance of non-normality in factor models goes back to Geary (1942) and Reiersol (1950), who have shown that factor loadings are identified in a linear measurement error model if the factor is not Gaussian.

⁸In fact, based on the results in Butucea and Matias (2005), W can belong to a more general class of known distributions. Furthermore, we note that if β_0 is known, then Assumption **B5** is not necessary.

Theorem 4.2. *If Assumptions B1–B5 hold, then β_0 , α_0 and γ_0 are identified.*

Note that this identification result does not require any variation from X_1 , which is in spirit close to the one-factor model in our paper and is different from the identification result in Vytlacil and Yildiz (2007). We also note that this result does not contradict the counterexample in the paper. In the counterexample, we only assume that we know the support of W is bounded. Here we assume that the full density of W , and thus, the support of W is known.

5 Finite Sample Properties

In this section we explore the finite sample properties of the proposed estimation procedure via a simulation study. We will also see how sensitive the performance of the proposed estimator is to the factor structure assumption. As a base comparison, we also report results for the estimator proposed in Vytlacil and Yildiz (2007) to see how sensitive it is to their second instrument restriction.

Our data are simulated from base models of the form

$$Y_1 = \beta_0 X_1 + \alpha_0 Y_2 + U \quad (5.1)$$

$$Y_2 = \beta_1 X + V > 0 \quad (5.2)$$

where X_1 is binary with success probability 0.6, X has marginal distribution $N(0;1)$, X_1 and X are mutually independent, $(X_1; X) \perp (V; U)$, $U = \alpha_0 V + \epsilon$. $(V; \epsilon)$ are distributed independently of each other, where V is distributed following a standard normal distribution, and ϵ is distributed either standard normal, Laplace, or $T(3)$. The parameters $(\alpha_0; \beta_0) = (0.25; 1.2)$ or $(0.5; 1.2)$.

Since X_1 is discrete, Vytlacil and Yildiz (2007)'s identification condition does not hold. However, the identification condition in this paper becomes

$$j \neq j \quad \text{length of the support of } X;$$

which holds.

For each choice of sample size $n = 100; 200; 400; 800; 1;600$, we simulate 280 samples and report the bias, standard deviation (std), root mean squared error (RMSE), and median absolute deviation (MAD) for both Vytlacil and Yildiz (2007)'s estimator (VY051 fn0.9091) -the X

As results from the table indicate, the finite sample performance of our estimator generally agrees with the asymptotic theory. The RMSE for the estimator proposed here is decreasing as the sample size increases, as one could expect given the consistency property of our estimator. Besides, the decay rate of the RMSE and MAD is about $\frac{1}{\sqrt{2}}$ when $n = 400$ as sample sizes doubles, in line with the parametric rate of convergence of our estimator.

6 Conclusions

In this paper we explored the identifying power of factor structures in discrete simultaneous systems. We found that for a binary-binary system the factor structure we considered did indeed add informational content. Specifically, it enabled the relaxation of both the exclusion and support conditions typically employed in the identification of these models. As we then demonstrated factor structures then enabled the regular identification of parameters of interest, and we proposed a new rank based estimation procedure that converged at a parametric rate with a limiting normal distribution. Finite sample properties of the estimator were demonstrated through simulation studies.

Supplement to "Informational Content of Factor Structures in Simultaneous Binary Response Models"

Abstract

This paper gathers the supplementary material to the original paper. Section A proves

RK5 (Rate condition on first stage bandwidth sequence) The first stage bandwidth sequence H_n used in the nonparametric estimator of the choice probability functions and their derivatives satisfies $P \frac{1}{n} H_n^{p_2 - 1} \rightarrow 0$ and $n^{-1} H_n^{-1} \rightarrow 0$.

We first show consistency of the rank estimator. To do so we first define the objective function $Q_{n;2}^{if}(\cdot)$, defined as

$$Q_{n;2}^{if}(\cdot) = \sum_{i \in j} g_{i;j}(\cdot)$$

where

$$g_{i;j}(\cdot) = \begin{cases} 1 & \text{if } @P^{11}(X_{1;i}; X_i) = f_V(X_i) + @P^{10}(X_{1;j}; X_j) = f_V(X_j) \\ 0 & \text{if } (X_{1;i}; X_i; X_{1;j}; X_j) \\ 0 & \text{if } @P^{11}(X_{1;i}; X_i) = f_V(X_i) + @P^{10}(X_{1;j}; X_j) < 0 \\ 0 & \text{if } (X_{1;i}; X_i; X_{1;j}; X_j) < 0 \end{cases}$$

Since $g_{i;j}$ is bounded by $18i;j$, and our random sampling assumption, we have for each,

by Markov's inequality. But the expectation in the numerator on the right hand side is

$$P(\hat{m}(x_i) > 0; m(x_i) < 0) = P(\hat{m}(x_i) > 0; m(x_i) > \epsilon_n) + P(\hat{m}(x_i) > 0; m(x_i) \leq \epsilon_n)$$

where ϵ_n is a sequence of positive numbers converging to 0, at a slow rate, e.g. $(\log n)^{-1}$. The first term on the right hand side is bounded above by

$$P(\|\hat{m} - m\| > \epsilon_n) \leq P(\|\hat{m} - m\| > \epsilon_n)$$

where the notation $\|\hat{m} - m\|$ above denotes the sup norm over x_i . The right hand side probability above will be sufficiently small for n large enough by the rate of convergence of the nonparametric estimator. The second term, $P(\hat{m}(x_i) > 0; m(x_i) \leq \epsilon_n)$, is bounded above by $P(m(x_i) \leq \epsilon_n)$ which by the smoothness of $m(x_i)$ converges to 0, and hence can be made arbitrarily small.

To derive the rate of convergence and limiting distribution theory for the feasible estimator where we first estimate choice probability functions and their derivatives nonparametrically, we expand the nonparametric estimators around true functions that are inside the indicator function in Q_{n2} . Then we can follow the approach in Sherman (1994b). Having already established consistency of the estimator, we will first establish root- n consistency and then asymptotic normality. For root- n consistency we will apply Theorem 1 of Sherman (1994b) and so here we change notation to deliberately stay as close as possible to his. We will actually apply this theorem twice, first establishing a slower than root- n consistency result and then root- n consistency. Keeping our notation deliberately as close as possible to Sherman(1994b), here replacing our second stage rank objective function $\hat{Q}_{2;n}(\cdot)$ with $\hat{G}_n(\cdot)$, our infeasible objective function $Q_{n,2}^f(\cdot)$ with $G_n(\cdot)$, and denoting our limiting objective function, previously denoted by $Q_0(\cdot)$, by $G(\cdot)$. We have the following theorem:

Theorem B.1. (From Theorem 1 in Sherman (1994b)).

If ϵ_n and δ_n are sequences of positive numbers converging to 0, and

1. $\hat{Q}_0 = Q_0(\epsilon_n)$
2. There exists a neighborhood of Q_0 and a constant $\delta > 0$ such that $G(\cdot) \geq G(Q_0) - \delta k^2$ for all \cdot in this neighborhood.
3. Uniformly over $O_p(\epsilon_n)$ neighborhoods of Q_0

$$\hat{G}_n(\cdot) = G(\cdot) + O_p(k^{-1} \epsilon_n) + o_p(k^{-1} \epsilon_n^2) + O_p(\delta_n)$$

then $\hat{Q}_0 = O_p(\max(\epsilon_n^2; \delta_n^2))$.

Once we use this theorem to establish the rate of convergence of our rank estimator, we can attain limiting distribution theory, which will follow from the following theorem:

Theorem B.2. (From Theorem 2 in Sherman (1994b)). Suppose $\hat{\theta}$ is $P_{\bar{n}}$ -consistent for θ_0 , an interior point of Θ : Suppose also that uniformly over $O_p(n^{-1/2})$ neighborhoods of θ_0 ,

$$\hat{G}_n(\theta) = \frac{1}{2}(\theta - \theta_0)'V(\theta - \theta_0) + \frac{1}{\sqrt{n}}(\theta - \theta_0)'W_n + o_p(1/n) \quad (B.1)$$

where V is a negative definite matrix, and W_n converges in distribution to a $N(0; \Sigma)$ random vector. Then

$$P_{\bar{n}}(\hat{\theta} - \theta_0) \rightarrow N(0; V^{-1} \Sigma V^{-1}) \quad (B.2)$$

We first turn attention to applying Theorem B.1 to derive the rate of convergence of our estimator. Having already established consistency of our rank estimator, we turn attention to the second condition in Theorem B.1. To show the second condition, we will first derive an expansion for $G(\theta)$ around $G(\theta_0)$. We denote that even though $G_n(\theta)$ is not differentiable in θ , $G(\theta)$ is sufficiently smooth for Taylor expansions to apply as the expectation operator is a smoothing operator and the smoothness conditions in Assumptions RK2, RK3. Taking a second order expansion of $G(\theta)$ around $G(\theta_0)$, we obtain

$$G(\theta) = G(\theta_0) + r_1(\theta - \theta_0) + \frac{1}{2}(\theta - \theta_0)'r_2(\theta - \theta_0) + o_p(\|\theta - \theta_0\|^2) \quad (B.3)$$

where r_1 and r_2 denote first and second derivative operators and θ denotes an intermediate value. We note that the first two terms of the right hand side of the above equation are 0, the first by how we defined the objective function, and the second by our identification result in Theorem 2.1. Define

$$V = r_2(\theta_0) \quad (B.4)$$

and V is positive definite by Assumption A3, so we have

$$(\theta - \theta_0)'r_2(\theta_0)(\theta - \theta_0) > 0 \quad (B.5)$$

$r_2(\theta)$ is also continuous at $\theta = \theta_0$ by Assumptions RK2 and RK3, so there exists a neighborhood of θ_0 such that for all θ in this neighborhood, we have

$$(\theta - \theta_0)'r_2(\theta)(\theta - \theta_0) > 0 \quad (B.6)$$

which suffices for the second condition to hold.

To show the third condition in Theorem B.1, we next establish the form of the remainder term when we replace nonparametric estimators with the true functions they are estimating. Specifically we wish to evaluate the difference between

$$\begin{aligned} & [E\{\hat{g}^{11}(X_{1j}; X_i) - \hat{g}^V(X_i) + E\{\hat{g}^{10}(X_{1j}; X_j) - \hat{g}^V(X_j) \mid X_i, X_j\} - 0] \\ & + [E\{f^{11}(X_{1j}; X_i) - f^V(X_i) + E\{f^{10}(X_{1j}; X_j) - f^V(X_j) \mid X_i, X_j\} - 0] \end{aligned} \quad (B.7)$$

and

$$\begin{aligned}
 & [f @ P^{11}(X_{1i}; X_i) = f_v(X_i) + @ P^{10}(X_{1j}; X_j) = f_v(X_j) < 0] g_1 f(X_{1i}; X_i; X_{1j}; X_j) < 0 & (B.9) \\
 + & [f @ P^{11}(X_{1i}; X_i) = f_v(X_i) + @ P^{10}(X_{1j}; X_j) = f_v(X_j) < 0] g_1 f(X_{1i}; X_i; X_{1j}; X_j) < 0 & (B.10)
 \end{aligned}$$

To establish a representation for this difference, we first simplify notation we write the expressions as:

$$\begin{aligned}
 & I[\hat{m}_1(x_i) + \hat{m}_2(x_j) < 0] / I[x_{ij}^0 < 0] & (B.11) \\
 + & I[\hat{m}_1(x_i) + \hat{m}_2(x_j) < 0] / I[x_{ij}^0 < 0] & (B.12)
 \end{aligned}$$

and

$$\begin{aligned}
 & I[m_1(x_i) + m_2(x_j) < 0] / I[x_{ij}^0 < 0] & (B.13) \\
 + & I[m_1(x_i) + m_2(x_j) < 0] / I[x_{ij}^0 < 0] & (B.14)
 \end{aligned}$$

respectively, where here x_i denotes the separate components of $x_{1i}; x_i$, and analogous for x_j . We first explore

$$(I[\hat{m}_1(x_i) + \hat{m}_2(x_j) < 0] - I[m_1(x_i) + m_2(x_j) < 0]) / I[x_{ij}^0 < 0]$$

for each $i; j$ inside the double summation:

$$\frac{1}{n(n-1)} \sum_{i \neq j} \dots$$

RK3, RK4, RK5, is $o_p(n^{-1/4})$. Thus by repeated application of Theorem B.1, we can conclude that the estimator is root- n consistent. To show that the estimator is also asymptotically normal, we will first derive a linear representation for the term:

$$\frac{1}{n(n-1)} \sum_{i \neq j}^X (0) \hat{f}_{m_{ij}}(0) (\hat{m}_1(x_i) - m_1(x_i)) / [x_{ij}^0 = 0] \quad (\text{B.17})$$

As this term is linear in the nonparametric estimator $\hat{m}_1(x_i)$, the desired linear representation follows from arguments used in Khan (2001). One slight difference here compared to Khan (2001) is that here our nonparametric estimators and estimands are each ratios of derivatives. Nonetheless, after linearizing these ratios as done in, e.g. Newey and McFadden (1994). Specifically, we have that B.17 can be expressed as:

$$\frac{1}{n(n-1)} \sum_{i \neq j}^X (0) \hat{f}_{m_{ij}}(0) \frac{1}{m_{1\text{den}}(x_i)} (\hat{m}_{1\text{num}}(x_i) - m_{1\text{num}}(x_i)) / [x_{ij}^0 = 0] \quad (\text{B.18})$$

$$\frac{1}{n(n-1)} \sum_{i \neq j}^X (0) \hat{f}_{m_{ij}}(0) \frac{m_{1\text{num}}(x_i)}{m_{1\text{den}}(x_i)^2} (\hat{m}_{1\text{den}}(x_i) - m_{1\text{den}}(x_i)) / [x_{ij}^0 = 0] \quad (\text{B.19})$$

where $\hat{m}_{1\text{num}}(x_i)$ denotes the numerator $f @_2 \hat{P}^{11}(X_{1i}; X_i)g$, the estimator of $m_{1\text{num}}(x_i)$ which denotes $f @_2 P^{11}(X_{1i}; X_i)g$, and $\hat{m}_{1\text{den}}(x_i)$ denotes the denominator $\hat{f}_V(X_i)$, the estimator of $m_{1\text{den}}(x_i)$ which denotes $f_V(X_i)$.

Plugging in the definitions of the kernel estimators of $\hat{m}_{1\text{num}}(x_i)$, and $\hat{m}_{1\text{den}}(x_i)$, results in a third order process. Using arguments in Khan (2001) and Powell, Stock, and Stoker (1989) we can express the third order U process as a second order U process plus an asymptotically negligible remainder term. This is of the form:

$$\frac{1}{n} \sum_{i=1}^n (0) \frac{\hat{\psi}(x_i)}{m_{1\text{den}}(x_i)} (y_{1i} - m_{1\text{num}}(x_i)) E [f_{m_{ij}}(0) x_{ij}^0 = 0] j x_i \quad (\text{B.20})$$

where $\hat{\psi}(x_i) = \frac{f_X^0(x_i)}{f_X(x_i)}$. We note that the function $E [f_{m_{ij}}(0) / [x_{ij}^0 = 0] j x_i]$, which we denote here by $H(x_i; \cdot)$ is a smooth function in \cdot . We will use this feature to expand $H(x_i; \cdot)$ around $H(x_i; 0)$. Analogous arguments can be used to attain a linear representation of (B.19), which is of the form:

$$\frac{1}{n} \sum_{i=1}^n (0) \frac{\hat{\psi}_2(x_{1i}) m_{1\text{num}}(x_i)}{m_{1\text{den}}(x_i)^2} (y_{2i} - m_{1\text{den}}(x_i)) E [f_{m_{ij}}(0) x_{ij}^0 = 0] j x_i \quad (\text{B.21})$$

where $\hat{\psi}_2(x_{1i})$ that the f

(B.22)

Note that by Assumptions RK2 , RK3 , $H(x_i; \cdot)$ is smooth in \cdot implying the expansion

$$H(x_i; \cdot) = H(x_i; 0) + r \cdot H(x_i; 0)'(0)$$

Thus we can express (B.22) as the which we note is a mean 0 sum

$$\frac{1}{n} \sum_{i=1}^n 1_{rki} (0) \tag{B.23}$$

where

$$1_{rki} = (0) \frac{1}{m_{1den}(x_i)} \cdot (x_i)(y_{1i} - m_{1num}(x_i)) \frac{m_{1num}(x_i)}{m}$$

C Nonparametric Factor Structure

Here we describe an estimator for the case where we have a nonparametric factor structure. Recall for this model we had the following relationship between unobservable variables:

$$U = g_0(V) + \epsilon \quad (C.31)$$

where we assumed that $\epsilon \perp V$.

Our goal in this more general setup is to identify and estimate both α_0 and g_0 . Our identification is based on the condition that

$$x_1 + \alpha_0 + g_0(x) = x_1 + g_0(x):$$

if and only if

$$E_2 P^{11}(x_1^{(1)}; x^{(1)}) = f_V(x^{(1)}) + E_2 P^{10}(x_1^{(1)}; x^{(1)}) = f_V(x^{(1)}) = 0:$$

Using the same $i; j$ pair notation as before, this gives gives us, in the nonparametric case,

$$X_{1i} - X_{1j} = \alpha_0 + (g_0(X_i) - g_0(X_j)) \quad (C.32)$$

Note the above equation has a "semi parametric form", loosely related to the model considered in, for example, Robinson (1988). However, we point out crucial differences between what we have above and the standard semi linear model. Here we are trying to identify the intercept α_0 which is usually not identified in the semi linear model as it cannot be separately identified from the nonparametric function. However, note above on the right hand side, we do not just have a nonparametric function of $X_i; X_j$, but the difference of two *identical* and *additively separable* functions $g_0(\cdot)$. In fact it is this differencing of these functions which enables us to separately identify α_0 . Furthermore, as will now see when turning to our estimator of α_0 , the structure of the nonparametric component, specifically additive separability of two identical functions of $X_i; X_j$ respectively, can easily be incorporated into our approximation of each of them. From a theoretical perspective separable functions have the advantage of effectively being a one dimensional problem, as there are no interaction terms to have to deal with. It is well known that nonparametric

where y_i denotes the observed dependent variable, $x_i; z_i$ are observed regressors, $g(\cdot)$ is an unknown nuisance function, u_i is an unobserved disturbance term, and β_0 is the unknown regression coefficient vector which is the parameter of interest. There is a very extensive literature in both econometrics and statistics on estimation and inference methods for this model- see for example Powell (1994) for some references.

One popular way to estimate this model is to use an expansion of basis functions, for example polynomials or splines to approximate $g(\cdot)$, and from a random sample of n observations of $(y_i; x_i; z_i)$ regress y_i on $x_i; b(z_i)$ where $b(z_i)$ denotes the set of basis functions used to approximate $g(\cdot)$. As an illustrative example, assuming z_i were scalar, if one were to use polynomials as basis functions, one would estimate the approximate model,

$$y_i = \beta_0 + \beta_1 z_i + \beta_2 z_i^2 + \beta_3 z_i^3 + \dots + \beta_{k_n} z_i^{k_n} + u_{in}$$

where k_n is a positive integer smaller than the sample size n , and $\beta_1; \beta_2; \dots; \beta_{k_n}$ are additional unknown parameters. This has been done by regressing y_i on $x_i; z_i; z_i^2; \dots; z_i^{k_n}$, and our estimated coefficient of x_i would be the estimator of β_0 . The validity of this approach has been shown in, for example, Donald and Newey (1994). Now for our problem at hand, incorporating a nonparametric factor structure, we propose a kernel weighted least squares estimator. The weights are as they were before, assigning great weights to pairs of observations where the sum of derivatives of ratios of choice probabilities are closer to 0.

The dependent variable is identical to as before, the set of n choose 2 pairs $X_{1i} \quad X_{1j}$. The regressors now reflect the series approximation of $g_0(X_i) \quad g_0(X_j)$:

$$g_0(X_i) \quad g_0(X_j) \quad \beta_1(X_i \quad X_j) + \beta_2(X_i^2 \quad X_j^2) + \beta_3(X_i^3 \quad X_j^3) + \dots + \beta_{k_n}(X_i^{k_n} \quad X_j^{k_n})$$

is bounded away from 0 uniformly in k_n , where

$$P_{k_n} = (1; (X_i - X_j); (X_i - X_j)^2; \dots; (X_i - X_j)^{k_n})^0$$

Theorem C.1. Under Assumptions I, K, H, S, PS, FK, FH, BFC,

$$\hat{\Lambda}_{NP} \xrightarrow{P} \Lambda_0 \quad (C.33)$$

D Proof of Theorem 4.1

E Proof of Theorem 4.2

We first show that both f_0 and the density of Z are identified. Note X has full support. This implies the density of V denoted as $f_V(\cdot)$ is identified via

$$f_V(v) = \int E(Y_2 | X = v):$$

In addition, we have

$$f_V(\cdot) = f_W \circledast f_0(\cdot);$$

where \circledast denotes the convolution operator. Suppose $f_W(\cdot)$ and f_0 are not identified so that there exist $f_W^0(\cdot)$ and f_0^0 such that

$$f_V(\cdot) = f_W^0$$

and

$$F(\theta_x P^{10}(x_1;)) = F_0(F_1(x_1 - 0) f_W()) F($$

Theorem F.1. *Assumption 1 holds. (1) Then $j_{0j} > b - a$ is necessary and sufficient for θ_0 to be identified. (2) When $j_{0j} > b - a$, the sharp identified set for θ_0 is*

$$A = \{ \theta : j_{0j} > b - a \text{ if } \theta_0 > 0 \text{ and } j_{0j} < b - a \text{ if } \theta_0 < 0 \}.$$

Next, we assume, in addition to Assumption 1, the factor structure, i.e., (2.3) in Section 2. Recall in Section 3, under the factor structure, our rank estimator can be written as an M-estimator

$$\hat{\theta} = \arg \max_{\theta} Q_{n,2}(\theta) = \sum_{i \in j} \hat{g}_{i,j}(\theta)$$

in which

$$\hat{g}_{i,j}(\theta) = [1 \mathbf{f}_{\theta_2} \hat{P}^{11}(X_{1,i}; X_i) = \hat{F}_V(X_i) + \theta_2 \hat{P}^{10}(X_{1,j}; X_j) = \hat{F}_V(X_j) \geq 0 \mathbf{1} \mathbf{f}_{\theta_1} (X_{1,i}; X_i; X_{1,j}; X_j) \geq 0 \mathbf{g} + 1 \mathbf{f}_{\theta_2} \hat{P}^{11}(X_{1,i}; X_i) = \hat{F}_V(X_i) + \theta_2 \hat{P}^{10}(X_{1,j}; X_j) = \hat{F}_V(X_j) < 0 \mathbf{1} \mathbf{f}_{\theta_1} (X_{1,i}; X_i; X_{1,j}; X_j) < 0 \mathbf{g}];$$

with

$$(X_1; X; X_1; X) = x_1 + \theta_1 x \quad (x_1 \quad x):$$

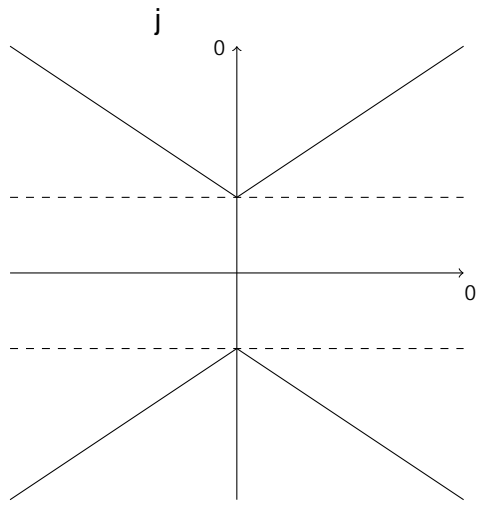
The information content explored by the M-estimator can be summarized as follows:

$$A_2(\theta) = \{ (X_1; X; X_1; X); (X_1; X; X_1; X; \theta) \geq 0 > (X_1; X; X_1; X; \theta) < 0 \text{ or } (X_1; X; X_1; X; \theta) < 0 < (X_1; X; X_1; X; \theta) \mathbf{g} \}.$$

Then we cannot distinguish, from the true parameter θ_0 , all impostors in

$$\bar{A}_2 = \{ \theta : P(A_2(\theta)) = 0 \}.$$

In a simple example, if $\text{Supp}(X_1; X) = [a; b] \times [c; d]$, then θ_0 is identified if $j_{0j} < b - a + j_{0j}(d - c)$. Recall Theorem F.1, without imposing factor structure, the necessary and sufficient condition for achieving identification is $j_{0j} > b - a$. Therefore, the blue area in the Figure below is the additional parts of parameter space that is identified with factor structure but not otherwise.



When $\epsilon_0 < a - b$, for any $\epsilon < \epsilon_0$, we can define

$$\begin{aligned} \mathcal{U} &= U + \epsilon_0 && \text{if } U \leq b + \epsilon_0 \\ \mathcal{U} &= U && \text{if } U > b + \epsilon_0 \end{aligned}$$

Then for any $x_1 \in [a; b]$,

$$\begin{aligned} P(\mathcal{U} \leq x_1 + \epsilon | V = v) &= P(\mathcal{U} \leq x_1 + \epsilon; U \leq b + \epsilon_0) + P(\mathcal{U} \leq x_1 + \epsilon; U > b + \epsilon_0 | V = v) \\ &= P(U \leq x_1 + \epsilon | V = v) \\ P(\mathcal{U} \leq x_1 | V = v) &= P(\mathcal{U} \leq x_1; U \leq b + \epsilon_0 | V = v) + P(\mathcal{U} \leq x_1; U > b + \epsilon_0 | V = v) \\ &= P(U \leq b + \epsilon_0; U \leq x_1 + \epsilon | V = v) + P(b + \epsilon_0 < U \leq x_1; | V = v) \\ &= P(U \leq b + \epsilon_0 | V = v) + P(b + \epsilon_0 < U \leq x_1; | V = v) \\ &= P(U \leq x_1 | V = v): \end{aligned}$$

Let $G_{U;V}$ and $G_{\mathcal{U};V}$ be the joint distribution of $(U; V)$ and $(\mathcal{U}; V)$ respectively. Then the above calculation with (G.38) imply that $(\epsilon_0; G_{U;V})$ and $(\epsilon; G_{\mathcal{U};V})$ are observationally equivalent.

When $\epsilon_0 > b - a$, for any $\epsilon > \epsilon_0$, we can define

$$\begin{aligned} \mathcal{U} &= U + \epsilon_0 && \text{if } U > a + \epsilon_0 \\ \mathcal{U} &= U && \text{if } U \leq a + \epsilon_0 \end{aligned}$$

Then for any $x_1 \in [a; b]$,

$$\begin{aligned} P(\mathcal{U} \leq x_1 + \epsilon | V = v) &= P(\mathcal{U} \leq x_1 + \epsilon; U \leq a + \epsilon_0) + P(\mathcal{U} \leq x_1 + \epsilon; U > a + \epsilon_0 | V = v) \\ &= P(U \leq a + \epsilon_0 | V = v) + P(a + \epsilon_0 < U \leq x_1 + \epsilon | V = v) \\ &= P(U \leq x_1 + \epsilon | V = v): \\ P(\mathcal{U} \leq x_1 | V = v) &= P(\mathcal{U} \leq x_1; U \leq a + \epsilon_0 | V = v) + P(\mathcal{U} \leq x_1; U > a + \epsilon_0 | V = v) \\ &= P(U \leq x_1 | V = v): \end{aligned}$$

So again, $(\epsilon_0; G_{U;V})$ and $(\epsilon; G_{\mathcal{U};V})$ are observationally equivalent.

For the second result in the theorem, first note that, when $j - \epsilon_j > b - a$, the sign of ϵ_0 is identified by the data. We take $\epsilon_0 > b - a$ as an example. By the proof of Theorem F.1, we have already shown that all $\epsilon > \epsilon_0$ is in the identified set. Now we consider $\frac{b - a + \epsilon_0}{2} < \epsilon_0$.

$$\begin{aligned} \mathcal{U} &= U + \epsilon_0 && \text{if } U > a + \\ \mathcal{U} &= U && \text{if } U \leq a + \end{aligned}$$

Then for any $x_1 \in [a; b]$,

$$\begin{aligned}
 P(U \leq x_1 + \epsilon; jV = v) &= P(U \leq x_1 + \epsilon; U \leq a + \epsilon) + P(U \leq x_1 + \epsilon; U > a + \epsilon; jV = v) \\
 &= P(U \leq a + \epsilon; jV = v) + P(a + \epsilon < U \leq x_1 + \epsilon; jV = v) \\
 &= P(U \leq x_1 + \epsilon; jV = v): \\
 P(U \leq x_1; jV = v) &= P(U \leq x_1; U \leq a + \epsilon; jV = v) + P(U \leq x_1; U > a + \epsilon; jV = v) \\
 &= P(U \leq x_1; jV = v) + P(U \leq x_1 + \epsilon; U > a + \epsilon; jV = v): \\
 &= P(U \leq x_1; jV = v):
 \end{aligned}$$

Here note that the last equality is because $x_1 + \epsilon \leq b + \epsilon \leq a + \epsilon$ if $\epsilon = \frac{b - a + \epsilon}{2}$. Denote $\epsilon^{(1)} = \frac{b - a + \epsilon}{2}$. Then we have shown that there exists $U^{(1)}(\epsilon)$ which only depends on ϵ such that for any $x_1 \in [a; b]$, any v and any $\epsilon > 0$ $^{(1)}$

$$\begin{aligned}
 P(U^{(1)}(\epsilon) \leq x_1 + \epsilon; jV = v) &= P(U \leq x_1 + \epsilon; jV = v) \\
 P(U^{(1)}(\epsilon) \leq x_1; jV = v) &= P(U \leq x_1; jV = v):
 \end{aligned}$$

In particular, there exists $U^{(1)}(\epsilon^{(1)})$ such that

$$\begin{aligned}
 P(U^{(1)}(\epsilon^{(1)}) \leq x_1 + \epsilon^{(1)}; jV = v) &= P(U \leq x_1 + \epsilon^{(1)}; jV = v) \\
 P(U^{(1)}(\epsilon^{(1)}) \leq x_1; jV = v) &= P(U \leq x_1; jV = v):
 \end{aligned}$$

Now repeating the above construction but replacing U with $U^{(1)}$ and ϵ with $\epsilon^{(1)}$, we have for any $\epsilon^{(1)} > 0$ $^{(2)}$ $\frac{b - a + \epsilon^{(1)}}{2}$, there exists $U^{(2)}(\epsilon^{(1)})$

Online Supplement to "Informational Content of Factor Structures in Simultaneous Binary Response Models": Distribution Theory for Closed Form Estimator

H Distribution Theory for Closed Form Estimator

Many of the basic arguments follow those used in Chen and Khan (2008) and Chen, Khan, and Tang (2016). Recall what the key identification condition that motivated the weighted least squares estimator: For pairs of observations $(x_1; x)$ and $(x_1; x)$ in $\text{Supp}(X_1; X)$,

$$x_1 + \theta_0 x = x_1 + \theta_0 x:$$

if and only if

$$\partial_2 P^{11}(x_1; x) = f_V(x) + \partial_2 P^{10}(x_1; x) = f_V(x) = 0:$$

where recall ∂_2 denotes the partial derivative with respect to the second argument. Note that even though the random variable V is unobserved, the density function $f_V(\cdot)$ above can be recovered from the data from the partial derivative of the choice probability in the treatment equation with respect to the regressor in the treatment equation. Thus the above equation involves the sum of two ratios of derivatives of choice probabilities.

Recall $\theta_0 = (\theta_0; \theta_0)$. Our estimator of θ_0

$K(\cdot)$ and bandwidth H_n , whose properties are discussed below. The second problem can be dealt with through the use of "kernel weights" as has been frequently employed in the semiparametric literature.

$$4. \quad \hat{\beta}_{xxi} = E \left[\frac{1}{n} \sum_{i=1}^n \frac{X_i X_i^0 P_{0i}^{k;1;r}}{P_{0i}^{k;1;r}} \right]$$

$f_1(\rho_i^{1r}; \rho_j^{0r}) = E[x_{ij} | x_{ij}^0, \rho_i^{1r}; \rho_j^{0r}]$ where x_i denotes the 2×1 vector $(1; x_i)$, $f_0(\rho_j^{0r}) = E[x_j | \rho_j^{0r}]$, where x_j denotes the 2×1 vector $(1; x_j)$, $f_1(\cdot)$ denotes the density function of the random variable $P^{1;1;r}$, $f_0(\cdot)$ denotes the density function of the random variable $P^{1;0;r}$.

Our derivation of the asymptotic properties of this estimator are based on the following assumptions¹:

Assumption I (Identification) The 2×2 matrix:

$$M_1 = E \left[\frac{1}{n} \sum_{i=1}^n (\rho_i^{1r}; \rho_i^{1r})^0 f_0(\rho_i^{1r}) \right]$$

has full rank.

Assumption K (Second stage kernel function) The kernel function $k(\cdot)$ used in the second stage (to match the sum of ratios of derivatives to 0) is assumed to have the following properties:

K.1 $k(\cdot)$ is twice continuously differentiable, has compact support and integrates to 1.

K.2 $k(\cdot)$ is symmetric about 0.

K.3 $k(\cdot)$ is an eighth order kernel:

\int

$$\int_{-\infty}^{\infty} k(u) k(u) du = 0$$

The final set of assumptions involve restrictions for the first stage kernel estimator of the ratio of derivatives. This involves smoothness conditions on the choice probabilities $P_{0i}^{k;l;r}$, smoothness and moment conditions on the kernel function, and rate conditions on the first stage bandwidth sequence.

Assumption PS (Order of smoothness of probability functions and regressor density functions)

The functions $P^{k;l;r}(\cdot)$ and $f_{X_1;X}(\cdot)$

We will first derive a formula for the denominator term and then a linear representation for the numerator. For the denominator term here we aim to establish that the double sum $\sum_{i \neq j} \frac{1}{n(n-1)}$

Denoting a kernel estimator of the probability function of the outcome variable as a function of $x = (x_1; x)$, by $\hat{p}(x) = \frac{\sum_j y_{1j} K_H(x_j - x)}{\sum_j K_H(x_j - x)}$ where $K(\cdot)$ is our kernel function, H our bandwidth, and $K_H(\cdot) = \frac{1}{H} K(\frac{\cdot}{H})$, our estimator of the derivative of the probability function is

$$\hat{p}'(x) = \frac{\sum_k y_{1k} K_H^0(x_k - x) \frac{1}{H} - \sum_k K_H(x_k - x) \sum_k K_H^0(x_k - x) \frac{1}{H} \sum_k y_{1k} K_H(x_k - x)}{(\sum_k K_H(x_k - x))^2}$$

We plug in the first of the two terms in the above numerator into H.9 yielding

$$\frac{\frac{1}{n(n-1)(n-2)} \sum_{i \in j \in k} w_{ij}^0 f_{vi}^{-1} (y_{1k} K_H^0(x_k - x_i) \frac{1}{H} - \rho_i^1) x_{ij} (x_{1ij} - x_{ij}^0)}{\frac{1}{n} \sum_k K_H(x_k - x_i)}$$

In the above expression, we replace the denominator term with its plim², which is $f_x(x_i)$, which gives the expression:

$$\frac{1}{n(n-1)(n-2)} \sum_{i \in j \in k} \frac{y_{1k} K_H^0(x_k - x_i) \frac{1}{H} - \rho_i^1}{f_x(x_i)} x_{ij} (x_{1ij} - x_{ij}^0)$$

As an additional step we want a representation for w_{ij}^0 . By its definition,

$$\frac{1}{n(n-1)} \sum_{i \in j} X_{ij} = \frac{1}{n(n-1)} \sum_{i \in j} w_{ij}^0 x_{ij} (x_{1ij} \quad x_{ij}^0 \quad 0) = \frac{1}{n(n-1)} \sum_{i \in j} \frac{1}{h^2} k^0 \frac{\rho_i^{1r} + \rho_j^{0r}}{h} (x_i; x_j) \quad (H.12)$$

where $(x_i; x_j) = x_{ij} (x_{1ij} \quad x_{ij}^0 \quad 0)$. To attain this representation, we evaluate the expectation of the term inside the double summation. We express this as

$$\frac{1}{h^2} \int k^0 \frac{\rho_i^{1r} + \rho_j^{0r}}{h} (\rho_i^{1r}; \rho_j^{0r}) f_1(\rho_i^{1r}) f_0(\rho_j^{0r}) d\rho_i^{1r} d\rho_j^{0r}$$

where recall $f_1(\cdot)$ denotes the density function of the random variable $P^{1;1;r}$, $f_0(\cdot)$ denotes the density function of the random variable $P^{1;0;r}$, and here, $(\rho_i^{1r}; \rho_j^{0r}) = E[(x_i; x_j) | \rho_i^{1r}; \rho_j^{0r}]$. To evaluate the above integral we construct the change of variables $u = \frac{\rho_i^{1r} + \rho_j^{0r}}{h}$ and expand inside the integral. Before expanding the integral is of the form

$$\frac{1}{h} \int k^0(u) (\rho_i^{1r}; uh - \rho_i^{1r}) f_1(\rho_i^{1r}) f_0(uh - \rho_i^{1r}) du d\rho_i^{1r}$$

After expanding, the lead term is 0 because the function $k(\cdot)$ vanishes on the boundary of its support. The next term is of the form:

$$\int (\rho_i^{1r}; \rho_i^{1r}) f_1(\rho_i^{1r}) f_0(\rho_i^{1r}) + (\rho_i^{1r}; \rho_i^{1r}) f_1(\rho_i^{1r}) f_0'(\rho_i^{1r}) k^0(u) u du d\rho_i^{1r}$$

From our identification result the above integral simplifies to $E[(\rho_i^{1r}; \rho_i^{1r}) f_0'(\rho_i^{1r})]$ which we will denote by γ_1 . So plugging this result into H.8 we have the following result:

$$\frac{1}{n(n-1)} \sum_{i \in j} f_{vi}^{-1}(\rho_i^1 \quad \rho_i^1)_{ij} = \frac{1}{n} \sum_{i=1}^n f_{vi}^{-1} \left(y_{1i} \frac{f_x^0(x_i)}{f_x(x_i)} \right) p^1(x) f$$

The term \hat{f}_{vi} is our kernel estimator of the derivative of the probability function in the treatment equation: $\hat{f}_{vi} = \frac{1}{n} \sum_{j \in \mathcal{J}} E[Y_{2ij} | X_i]$. So we can use analogous arguments to attain a linear representation for this U -statistic in H.16 to conclude

$$\frac{1}{n(n-1)} \sum_{i \in \mathcal{J}} \sum_{j \in \mathcal{J}} \rho_{vi}^1 (\hat{f}_{vi} - f_{vi}) = \frac{1}{n} \sum_{i=1}^n \rho_{vi}^1 \rho_{vi}^1 y_{2i} \frac{f_X^0(x_i)}{f_X(x_i)} f_V(x_i) + o_p(n^{-1/2}) \quad (\text{H.17})$$

$$\frac{1}{n} \sum_{i=1}^n z_i + o_p(n^{-1/2}) \quad (\text{H.18})$$

where

$$z_i = \rho_{vi}^1 \rho_{vi}^1 y_{2i} \frac{f_X^0(x_i)}{f_X(x_i)} f_V(x_i) \quad (\text{H.19})$$

Next we can turn attention to the the second term in H.7,

$$\frac{1}{n(n-1)} \sum_{i \in \mathcal{J}} w_{ij}^0 (\rho_j^{0r} - \rho_j^{or}) x_{ij} (x_{1ij} - x_{ij}^0) \quad (\text{H.20})$$

The term $\rho_j^{or} - \rho_j^{0r}$ involves the ratio of two derivatives. So we can proceed as before by linearizing this ratio. This will yield the two expressions:

$$\frac{1}{n} \sum_{i=1}^n \rho_{vi}^1 \left(y_{1i} \frac{f_X^0(x_i)}{f_X(x_i)} - \rho^0(x_i) \right) + o_p(n^{-1/2}) = \frac{1}{n} \sum_{i=1}^n z_i + o_p(n^{-1/2}) \quad (\text{H.21})$$

where

$$z_i = \rho_{vi}^1 \left(y_{1i} \frac{f_X^0(x_i)}{f_X(x_i)} - \rho^0(x_i) \right) \quad (\text{H.22})$$

and

$$\frac{1}{n} \sum_{i=1}^n \rho_{vi}^1 \rho_{vi}^0 y_{2i} \frac{f_X^0(x_i)}{f_X(x_i)} f_V(x_i) + o_p(n^{-1/2}) = \frac{1}{n} \sum_{i=1}^n a_i + o_p(n^{-1/2}) \quad (\text{H.23})$$

where

$$a_i = \rho_{vi}^1 \rho_{vi}^0 y_{2i} \frac{f_X^0(x_i)}{f_X(x_i)} f_V(x_i) \quad (\text{H.24})$$

So collecting all results we can conclude that the estimator has the linear representation:

$$\hat{\theta}_0 = M_1 \frac{1}{n} \sum_{i=1}^n i + o_p(n^{-1/2}) \quad (\text{H.25})$$

where $i = 1i + 2i + 3i + 4i$.

Cunha, F., J. J. Heckman, and S. M. Schennach (2010): "Estimating the Technology of Cognitive and Noncognitive Skill Formation," *Econometrica*, 78, 883{931.

Donald, S., and W. Newey

Klein, R., C. Shan, and F. Vella (2015): "Semi-Parametric Estimation of Sample Selection Models with Binary Selection Rules and Binary Outcomes," *Journal of Econometrics*, 185, 82-94.

Klein, R., and R. Spady (1993): "An Efficient Semiparametric Estimator for Binary Response Model," *Econometrica*, pp. 387-421.

Lewbel, A. (2000): "Semiparametric Qualitative Response Model Estimation with Unknown Heteroscedasticity or Instrumental Variables," *Journal of Econometrics*, 97, 145-177.

Moneta, A., D. E. P. O. Hoyer, and A. Coad (2013): "Causal inference by independent component analysis: Theory and applications," *Oxford Bulletin of Economics and Statistics*, 75,

- Tamer, E. (2003): "Incomplete Bivariate Discrete Response Model with Multiple Equilibria," *Review of Economic Studies*, 70, 147{167.
- Vuong, Q., and H. Xu (2017): "Counterfactual mapping and individual treatment effects in nonseparable models with binary endogeneity," *Quantitative Economics*, pp. 589{610.
- Vytlacil, E. J., and N. Yildiz (2007): "Dummy Endogenous Variables in Weakly Separable Models," *Econometrica*, 75, 757{779.