

# Semiparametric Identification and Estimation of Multinomial Discrete Choice Models using Error Symmetry

Arthur Lewbel<sup>y</sup>      Jin Yan<sup>z</sup>      Yu Zhou<sup>x</sup>

Original February 2019, revised December 2021

## Abstract

We provide a new method to point identify and estimate cross-sectional multinomial choice models, using conditional error symmetry. Our model nests common random coefficient specifications (without having to specify which regressors have random coefficients), and more generally allows for arbitrary heteroskedasticity on most regressors, unknown error distribution, and does not require a “large support” (such as identification at infinity) assumption. We propose an estimator that minimizes the squared differences of the estimated error density at pairs of symmetric points about the origin. Our estimator is root  $N$  consistent and asymptotically normal, making statistical inference straightforward.

## 1 Introduction

Traditional multinomial choice models, such as multinomial logit (MNL) and multinomial probit (MNP), e.g., McFadden (1974), assume homoskedastic errors. However, in reality substantial

---

We thank Songnian Chen, Jeremy Fox, Bruce Hansen, Shakeeb Khan, Brendan Kline, Rosa Matzkin, Byoung Park, Jack Porter, Andres Santos, Xiaoxia Shi, Matthew Shum, Christopher Taber, Yuanyuan Wan and Hanghui

unobserved heterogeneity is common, e.g., Heckman (2001). We provide a new method to point identify preference parameters in cross-sectional multinomial choice models in the presence of general unobserved individual heterogeneity. Our identification is semiparametric, in that we do not specify the joint distribution of the latent errors, and we allow for arbitrary heteroskedasticity with respect to most regressors, including possible random coefficients. We propose a corresponding

and Tang (2016) in contrast find that symmetry, when combined with conditional independence

Using the analogy principle, we construct a corresponding estimator that minimizes the squared differences of the estimated error densities at each data point with its corresponding symmetry point. We show this minimum distance estimator is root  $N$  consistent and asymptotically normal. Computing the objective function of our estimator does not entail either numerical integration or deconvolution techniques, which are often required by random coefficients models. Moreover, our estimator does not require specifying which covariates, if any, have random coefficients, and is no more or less complicated regardless of how many covariates have random coefficients, or any other more complicated forms of heteroskedasticity.

Many methods have been developed for identifying and estimating utility function parameters with cross-sectional multinomial choice data. Many of those methods assume independence between the covariates and error terms, ruling out the possibility of individual heterogeneity such as random coefficients (Ruuluydence

the general multinomial choice case, we show root-N consistency and asymptotic normality of our estimator, and we provide proofs for all of our theorems.

## 2 The Model and Identification

### 2.1 The Random Utility Framework

To simplify notation and presentation of our results, for the main text of this paper we restrict attention to the case of three choices, with the relative utility of the outside option, denoted  $j = 0$ , normalized to equal zero. General results for an arbitrary number of multinomial choices, and

single outcome like  $y_0$ , because the choice of any one outcome depends on the utilities of all of the outcomes.

an absolutely continuous density function,  $f_{\mathbf{t}_1, \mathbf{t}_2}(\mathbf{X})$ , which is centrally symmetric about the origin, i.e.,

$$f_{\mathbf{t}_1, \mathbf{t}_2}(\mathbf{X}) = f_{\mathbf{t}_1, \mathbf{t}_2}(-\mathbf{X});$$

for any vector  $(\mathbf{t}_1; \mathbf{t}_2)$

If, this yields the equations

$$\frac{\partial E(y_0 | z = z; X = X)}{\partial z} = \frac{\partial \Pr(y_0 = 1 | z = z; X = X)}{\partial z} \quad (4)$$

$$= f_{12} \left( z_1, x_1^0; z_2, x_2^0 \mid X = X \right) \quad (1)^2;$$

and

$$\frac{\partial E(y_0 | z = z + 2X; X = X)}{\partial z} = \frac{\partial \Pr(y_0 = 1 | z = z + 2X; X = X)}{\partial z} \quad (5)$$

$$= f_{12} \left( z_1 + 2x_1^0, x_1^0; z_2 + 2x_2^0, x_2^0 \mid X = X \right) \quad (1)^2;$$

The left sides of equations (4) and (5) are both identified, and can be estimated as nonparametric regression derivatives, given a value of  $z$ . If  $z = z^0$ , then by the symmetry Assumption I2, the right sides of equations (4) and (5) are equal to each other. Define the function  $d_0(z; X)$  as



have positive probability measure for any  $\theta$  in the parameter space other than  $\theta_0$ . Assumptions I4 and I5 give one set of conditions that suffice. Assumption I4 provides a subset of the support of covariates with positive measure on which the function  $d_0(\theta; z; X)$  can be identified, while Assumption I5 ensures that symmetry points are unique.

Given these assumptions we obtain identification as follows. All proofs are in the Supplementary Appendix.

**Theorem 2.1** If Assumption I hold, then the parameter vector  $\theta_0$  is point identified by Definition 2.1.

### 2.3.1 Discussion

Theorem 2.1 used expectations of  $y_0$ . Additional identifying information (resulting in more efficient associated estimators) can similarly be obtained from  $y_1$  and  $y_2$ . Details are in the supplemental appendix.

The conditional independence between  $z$  and  $\epsilon$  in Assumption I1 is known as a distributional exclusion restriction (Powell, 1994, p. 2484). This allows for interpersonal heteroskedasticity on a subset of covariates: Higher moments of  $\epsilon$  can depend (in unknown ways) on  $X$ , but not  $z$ . Assumption I2 is our error symmetry restriction. Without loss of generality we assume that the point of symmetry is the origin, because any nonzero term could be absorbed into the intercept of the utility index as discussed in equation (1).

Assumption I3 assumes a compact parameter space, which is a standard assumption for many nonlinear models, including semiparametric multinomial discrete choice models. Assumption I4(a)

Assumption I5(a) ensures that the error density functions in (4) and (5) are evaluated at interior points of their support. Assumption I5(b) requires that the error density function has a unique (local) symmetry point over a subset of its support,  $\mathcal{E}_i(X)$ . This does not rule out densities having flat sections, but it does limit the range of any such flat sections.

## 2.4 An Alternative Identification Strategy

Existing binary choice estimators that make use of latent error symmetry (e.g. Chen (2000) and Chen, Khan and Tang (2016)) are based on the error distribution function rather than on the error density function as in Theorem 2.1. To illustrate, take a simple binary choice model where  $y = I(z + a + v > 0)$ . If  $v$  is a symmetric random variable around zero and  $v \perp z$ , then

$$E(y | z = c) = \Pr(v > c - a) = \Pr(v < c + a) = E(1 - y | z = c - 2a) \quad (8)$$

The constant  $a$  is identified by equating the above two expectations, which only requires estimation of the conditional mean of  $y$  and not its derivatives. This immediately extends to identification of covariate coefficients instead of just a constant.

We could have similarly based identification and estimation of our multinomial model on the distribution instead of the density of the errors, and thereby only required nonparametric regressions and not their derivatives for estimation. However, unlike the binary choice case, identification and estimation using the distribution instead of the density of the errors becomes complicated and clumsy in the multinomial setting. This is because, in the binary choice case, error symmetry just equates two conditional expectations, corresponding to two error intervals, while for multinomial choice, one must equate error rectangles.

To see the issue, begin again from equation (3). Let  $[a; b]$  be a rectangle in the support of  $\theta$ . Point  $a = (a_1; a_2)$  is the lower left vertex of this rectangle and  $b = (b_1; b_2)$  is the upper right vertex. By central symmetry, the probability of  $\theta$  being in the rectangle  $[a; b] = [a_1; b_1] \times [a_2; b_2]$  is the same as the probability of  $\theta$  being in the rectangle  $[b; a] = [b_1; a_1] \times [b_2; a_2]$ . This

then implies

$$\int_{[a,b]} f_{\mu_1, \mu_2}(t_1; t_2 | X) dt = \int_{[a,b]} f_{\mu_1, \mu_2}(t_1; t_2 | X) dt = \int_{[b, a]} f_{\mu_1, \mu_2}(t_1; t_2 | X) dt; \quad (9)$$

where the first equality in (9) holds by Assumption A2 and the second one holds by changing of variables.<sup>7</sup>

The integrals on both sides of equation (9) can be computed using the conditional distribution function of  $\mu$ , which in turn is obtained from the conditional expectation of  $y_0$ . For example, consider the left-hand side integral:

$$\begin{aligned} \int_{[a,b]} f_{\mu}(t | X) dt &= \Pr(a \leq \mu \leq b | X) \\ &= \Pr(a_1 \leq \mu_1 \leq b_1; a_2 \leq \mu_2 \leq b_2 | X) \\ &= \Pr(\mu_1 \leq b_1; \mu_2 \leq b_2 | X) - \Pr(\mu_1 < a_1; \mu_2 \leq b_2 | X) \\ &\quad - \Pr(\mu_1 \leq b_1; \mu_2 < a_2 | X) + \Pr(\mu_1 < a_1; \mu_2 < a_2 | X): \end{aligned} \quad (10)$$

We prefer to identify and estimate  $\theta$  by matching each point in the data using densities, rather than by matching rectangles using distributions, for many reasons. First, equating error distribution rectangles involves more tuning parameters, since rectangles need to be chosen. Second, matching densities only requires finding enough points  $(z = z^*; X = X^*)$  in the data that have matches  $(z = z^* - 2X^*; X = X^*)$  that lie in the support of the covariates. In contrast, each matching rectangle requires finding an entire range of covariates that lie in the support and has a range of matches that also lie entirely in the support. Third, to gain efficiency we will later create more moments by replacing  $y_0$  with different choices  $y_j$ . When matching density points, the same covariate values (points) that work for any one choice  $j$  will also work for any other choice. The same is not true for matching distribution rectangles, because for rectangles each match entails pairs of observations rather than individual observations. Finally, the computation cost of estimation is lower for equating error densities than for distribution rectangles. For a sample of size  $N$ , we compute error densities at  $2N$  points, while in contrast, using rectangles would entail computing the error distribution at  $N(N-1)2^J$  points.

### 3 A Minimum Distance Estimator and its Asymptotic Properties

#### 3.1 Population Objective Functions for Estimation

Given the identification strategy described in Section 2, we develop a minimum distance estimator (hereafter, MD estimator) for  $\theta$  using the identifying restriction  $d_0(\theta; z; X) = 0$ , where  $d_0$  is defined by equation (6). Note that the function  $d_0(\theta; z; X)$  is well defined if both points  $(z; X)$  and  $(z - 2X; X)$  are in the interior of the support of covariates,  $(z; X)$ . For this reason, we only wish to evaluate the function  $d_0(\theta; z; X)$

the range of these values. Define functions  $\phi_0(\cdot)$  and  $\phi_1(\cdot)$  by

$$\phi_0(z; X; \bar{c}_1) = \phi_0(z; X) \phi_0(z - 2X\bar{c}_1; X) \quad \phi_1(z; X) = \phi_1(z - 2X\bar{c}_2; X); \quad (11)$$

and

$$\phi_0(z; X) = \phi_0(z - c_1) \phi_1(X - C_2). \quad (12)$$

Here the absolute value of a vector or matrix,  $\|\cdot\|$ , is defined as the corresponding vector or matrix

### 3.2 An Estimator

We now provide an estimator for function  $d_0(\cdot; z_n; X_n)$  in (13) as

$$\hat{d}_{0; n}(\cdot; z_n; X_n) = \hat{\lambda}_{o; n}^{(2)}(z_n; X_n) - \hat{\lambda}_{cs; n}^{(2)}(z_n; X_n; \cdot); \quad (14)$$

where  $\hat{\lambda}_{o; n}^{(2)}(z_n; X_n)$  and  $\hat{\lambda}_{cs; n}^{(2)}(z_n; X_n; \cdot)$  are leave-one-out, Nadaraya-Watson nonparametric regression kernel estimators for the derivatives on the right hand side of equation (6) (see the supplemental appendix for details). By replacing the expectation in  $Q_0(\cdot)$  with its sample mean and replacing the function  $d_0(\cdot; z_n; X_n)$



$Q_{Nj}(\cdot)$  over each choice. We also extend all our results to multinomial choice with an arbitrary number of choices, instead of just three as above.

## 4 Monte Carlo Experiments

In this section, we use Monte Carlo experiments to study the finite-sample properties of the minimum distance (MD) estimator proposed above. We consider four data generating processes (DGPs). In each DGP, individual  $n$ 's utility from alternative  $j$ ,  $u_{nj}$ , is specified as

$$u_{nj} = \beta_j + \alpha_j x_{nj} + \epsilon_{nj} \text{ for } n = 1, 2, \dots, N \text{ and } j = 0, 1, 2 \quad (17)$$

$n$   
 $n$



Table 1: Monte Carlo Results of estimating  $\theta$  (True Parameter  $\theta = 0.2$ )

DGP	N	MNP		MD ( $y_0$ )		MD ( $y_0; y_1; y_2$ )	
		Bias	RMSE	Bias	RMSE	Bias	RMSE
1	1000	-0.0012	0.0435	0.0216	0.2368	-0.0017	0.1337
	2000	-0.0010	0.0307	0.0055	0.1355	-0.0078	0.0788
2	1000	0.5656	0.5833	0.1047	0.3521	-0.0392	0.3048
	2000	0.5627	0.5714	0.0543	0.2308	-0.0289	0.1747
3	1000	-0.0013	0.0454	0.0317	0.2220	0.0015	0.1417
	2000	-0.0017	0.0319	0.0158	0.1301	-0.0051	0.0812
4	1000	-0.7512	0.7718	-0.0054	0.3765	-0.0748	0.3550
	2000	-0.7481	0.7585	0.0180	0.2616	-0.0343	0.2149

covariates. Under all four DGP's our MD estimator remains consistent.

Table 1 reports the bias and root mean square error (RMSE) of each estimator in our simulations. The first set of columns reports the MNP estimator, the second reports our MD estimator using only  $y_0$ , while the third uses observations of all choices  $y_0$ ,  $y_1$ , and  $y_2$  (MNP also uses observations of all choices).

Under DGP 1, the MD estimators have small finite sample bias, and RMSEs two to four times larger than that of the correctly specified efficient MNP estimator. Under DGP 2, the bias of the misspecified MNP estimator is around three times the true parameter value, and this bias remains as the sample size is doubled. In contrast, the bias and RMSE of the MD estimators are much smaller than the MNP estimator, and they decrease sharply as the sample size increases. In DGP 3, the random coefficients MNP is correctly specified, and so performs better than the MD estimators in terms of bias and RMSE. However, in DGP 4 where the random component is heterogeneous, the bias of MNP is almost four times the true parameter value and does not vanish as sample size grows. In contrast, the bias of the MD estimators is still relatively small.<sup>9</sup> In all the DGPs, in terms of RMSE, the MD estimator using  $y_0$ ,  $y_1$ , and  $y_2$  performs better than

<sup>9</sup>We speculate that the bias in the MD estimators might be further reduced by a bandwidth search, and/or using local linear estimation for the first stage choice probabilities.



## References

- [1] Ahn, H., Ichimura, H., Powell, J.L., and Ruud, P. (2018): "Simple Estimators for Invertible Index Models," *Journal of Business and Economic Statistics*, 36, 1-10.
- [2] Berry, S. and Haile P.A. (2010): "Nonparametric Identification of Multinomial Choice Demand Models with Heterogeneous Consumers," Cowles Foundation Discussion Paper #1718.
- [3] Berry, S., Levinsohn, J., and Pakes, A. (1995): "Automobile Prices in Market," *Econometrica*, 63, 841-890.
- [4] Berry, S., Levinsohn, J., and Pakes, A. (2004): "Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Car Market," *Journal of Political Economy*, 112(1), 68-105.
- [5] Blundell, R. and Powel, J.L. (2004): "Endogeneity in Semiparametric Binary Response Models," *Review of Economic Studies*, 71, 655–679.
- [6] Chen, S. (2000): "Efficient Estimation of Binary Choice Models under Symmetry," *Journal of Econometrics*, 96, 183-199.
- [7] Chen, S. and Z8):S.8(Y80.004(an7)-334(ln)28.018en,)-(onomicor)-334282.9ametric Ztudmation o(79n)-35(45

- [10] Delaigle, A. and Hall, P. (2016): "Methodology for Non-parametric Deconvolution When the Error Distribution is Unknown," *Journal of the Royal Statistical Society, Series B*, 78, 231–252.
- [11] Dong, Y. and Lewbel, A. (2011): "Nonparametric Identification of a Binary Random Factor in Cross Section Data," *Journal of Econometrics*, 163, 163-171.
- [12] Fox, J.T. (2007): "Semiparametric Estimation of Multinomial Discrete-choice Models Using a Subset of Choices," *RAND Journal of Economics*, 38, 1002-1019.
- [13] Fox, J.T. and Gandhi, A. (2016): "Nonparametric Identification and Estimation of Random Coefficients in Multinomial Choice Models," *RAND Journal of Economics*, 47, 118-139.
- [14] Goolsbee, A. and Petrin, A. (2004): "The Consumer Gains from Direct Broadcast Satellites and the Competition with Cable TV," *Econometrica*, 72, 351-381.
- [15] Hausman, J.A. and Wise, D.A. (1978): "A Conditional Probit Model for Qualitative Choice: Discrete Decisions Recognizing Interdependence and Heterogeneous Preferences," *Econometrica*, 46, 403-426.
- [16] Heckman, J. (2001): "Micro Data, Heterogeneity, and the Evaluation of Public Policy: Nobel Lecture" *Journal of Political Economy*, 109, 673-748.
- [17] Honoré, B. E., Kyriazidou, E., and Udry, C. (1997): "Estimation of Type 3 Tobit Models using Symmetric Trimming and Pairwise Comparisons," *Journal of Econometrics*, 76, 107-128.
- [18] Khan, S., Ouyang, F., and Tamer, E. (2019): "Inference in Semiparametric Multinomial Response Models," Working Paper.
- [19] Kumbhakar, S.C. and Lovell, C.A.K. (2000): *Stochastic Frontier Analysis*. Cambridge University Press, Cambridge.

- [20] Lee, L. (1995): "Semiparametric Maximum Likelihood Estimation of Polychotomous and Sequential Choice Models," *Journal of Econometrics*, 65, 381-428.
- [21] Lewbel, A. (1997): "Constructing Instruments for Regressions With Measurement Error When no Additional Data are Available, with An Application to Patents and R&D," *Econometrica*, 65, 1201-1213.
- [22] Lewbel, A. (2000): "Semiparametric Qualitative Response Model Estimation with Unknown Heteroskedasticity or Instrumental Variables," *Journal of Econometrics*, 97, 145-177.
- [23] Manski, C.F. (1975): "Maximum Score Estimation of the Stochastic Utility Model of Choice," *Journal of Econometrics*, 3, 205-228.
- [24] Manski CF. (1985): "Semiparametric Analysis of Discrete Response: Asymptotic Properties of the Maximum Score Estimator," *Journal of Econometrics*, 27, 313-333.
- [25] Manski, C.F. (1988): "Identification of Binary Response Models," *Journal of the American Statistical Association*, 83, 729-738.
- [26] McFadden, D. (1974): "Conditional Logit Analysis of Qualitative Choice Behavior," in Fron-

- [31] Powell, J.L. and Ruud, P.A. (2008): "Simple Estimators for Semiparametric Multinomial Choice Models," Working Paper.
- [32] Ruud, P.A. (1986): "Consistent Estimation of Limited Dependent Variable Models Despite Misspecification of Distribution," *Journal of Econometrics*, 32, 157-187. Cambridge University Press.
- [33] Serfling, R. (2006): "Multivariate Symmetry and Asymmetry," *Encyclopedia of Statistical Sciences*. John Wiley & Sons, Inc.
- [34] Shi, X., Shum, M., and Song, W. (2018): "Estimating Semi-Parametric Panel Multinomial Choice Models Using Cyclic Monotonicity," *Econometrica*, 86, 737-761.
- [35] Train, K. (2009): *Discrete Choice Methods with Simulation*.
- [36] Yan, J. (2013): "A Smoothed Maximum Score Estimator for Multinomial Discrete Choice Models," Working Paper.
- [37] Yan, J. and Yoo, H. (2019): "Semiparametric Estimation of the Random Utility Model with Rank-ordered Choice Data," *Journal of Econometrics*, 211, 414-438.
- [38] Zhou, Y. (2021): "Identification and Estimation of of Entry Games under the Symmetry of Unobservables," Working Paper, NYU, Shanghai.

Supplementary Appendix: Semiparametric Identification and  
Estimation of Multinomial Discrete Choice Models using Error  
Symmetry

Arthur Lewbel<sup>y</sup>

Jin Yan<sup>z</sup>

Yu Zhou<sup>x</sup>

Original February 2019, revised December 2021





we only require that differences  $z$  and  $X$  be observed, and regularity conditions (e.g., continuity of  $z$ ) are only be imposed on  $z_j$  and  $x_j$ , not on  $z_j$  and  $x_j$ . In addition to these covariates, our identification only requires that  $y_0$  be observed, not the entire vector of outcomes  $y$ . This is possible because  $z$  provides information about the other outcomes. Nevertheless, the associated estimators will be more efficient by observing and making use of more the elements of  $y$  since each additional outcome  $y_j$  one observes provides additional overidentifying information.

Assumption I1 immediately implies

$$\begin{aligned} \Pr(y_0 = 1 \mid z; X) &= F_{z_1, z_2, \dots, z_J}^0(z_1, x_1^0; z_2, x_2^0; \dots; z_J, x_J^0 \mid z; X) \\ &= F_{z_1, z_2, \dots, z_J}^0(z_1, x_1^0; z) \end{aligned} \quad (\text{S.A.5})$$

Based on Assumptions I1 and I2, we have that if  $\theta = \theta_0$ , then  $d_0(\theta; z; X) = 0$ . Given some regularity conditions, setting the function  $d_0$  equal to zero at a collection of values of  $z$  and  $X$  provides enough equations to point identify  $\theta_0$ . The proof of Theorem S.A.1 is provided in Section S.C.

**Theorem S.A.1** *If Assumptions I hold, then the parameter vector  $\theta_0$  is point identified by Definition 1.*

### S.A.2 Identification Using Multiple Choices

In Section A.1, we identified the parameter vector  $\theta_0$  using only derivatives of the conditional mean of  $y_0$ . Here we illustrate that identification can be achieved using the conditional mean of  $y_j$  for any  $j \geq 1$ . Later we will increase efficiency of estimation by combining the identifying moments based on each of the observed choices.

We now introduce some additional notation. For each  $j \geq 1$ , define  $X^{(j)}$  as the matrix that consists of differenced covariate vectors  $x_k - x_j$  for all  $k \geq 1$  and  $k \neq j$ . For example, when  $1 < j < J$ ,  $X^{(j)} = (x_0 - x_j; \dots; x_{j-1} - x_j; x_{j+1} - x_j; \dots; x_J - x_j) \in \mathbb{R}^{J \times q}$ . By this notation, we have  $X^{(0)} = (x_1 - x_0; \dots; x_J - x_0) = X$ . In the same fashion, define  $\eta^{(j)} \in \mathbb{R}^J$  as the vector



Proposition S.A.2 *If Assumption 12 holds, then for every  $j \in J$  and almost every  $X^{(j)} \in S_{X^{(j)}}$ , the conditional distribution function  $F_{n^{(j)}}(t^{(j)} | X^{(j)})$  admits an absolutely continuous density function,  $f_{n^{(j)}}(t^{(j)} | X^{(j)})$ , which is centrally symmetric about the origin, i.e.,*

$$f_{n^{(j)}}(t^{(j)} | X^{(j)}) = f_{n^{(j)}}(-t^{(j)} | X^{(j)}); \quad (\text{S.A.13})$$

*for any vector  $t^{(j)} \in S_{n^{(j)}}(X^{(j)})$  where  $S_{n^{(j)}}(X^{(j)}) \subset \mathbb{R}^J$ .*

of the left-hand sides of (S.A.15) and (S.A.16), that is,

$$d_j(z^{(i)}; X^{(i)}) = E(y_j | z^{(i)} = z^{(i)}; X^{(i)} = X^{(i)}) - \alpha_j^{(i)} \quad (S.A.17)$$

$$E(y_j | z^{(i)} = z^{(i)}; X^{(i)} = X^{(i)}) - \alpha_j^{(i)}$$

which always equals zero when  $\alpha_j = 0$  and may be non-zero when  $\alpha_j \neq 0$ .

Then, analogous to Definition 1, define

$$D_j(z^{(i)}; X^{(i)}) = \int_{z^{(i)} \in S_{(z^{(i)}; X^{(i)})}} (y_j - \alpha_j^{(i)})^2 dz^{(i)} + \int_{z^{(i)} \notin S_{(z^{(i)}; X^{(i)})}} (y_j - \alpha_j^{(i)})^2 dz^{(i)} \quad (S.A.18)$$

Recall that there is a one-to-one correspondence, respectively, between  $X^{(i)}$  and  $X$ ,  $z^{(i)}$  and  $z$ , and  $\alpha_j^{(i)}$  and  $\alpha_j$ . For every  $(z^{(i)}; X^{(i)}) \in \text{int}(S_{(z^{(i)}; X^{(i)})})$  such that  $(z^{(i)}; X^{(i)}) \in \text{int}(S_{(z; X)})$ , we immediately have  $(z^{(i)}; X^{(i)}) \in \text{int}(S_{(z^{(i)}; X^{(i)})})$  and  $(z^{(i)}; X^{(i)}) \in \text{int}(S_{(z^{(i)}; X^{(i)})})$ , as well as  $d_j(z^{(i)}; X^{(i)}) = 0$  if and only if  $d_j(z; X) = 0$ . Therefore, we can also use the choice probability of any alternative in the choice set to achieve identification.

### S.A.3 Individual Heterogeneity and Random Coefficient

Our identifying assumptions do not refer specifically to random coefficients. Here we provide sufficient conditions for our key identification assumptions I1 and I2 to hold when unobserved

$\beta_n = \beta_n$  Our symmetry assumption implies that  $\beta^0$  will also be the mean coefficients, as long as these coefficients exist, but we don't impose this existence.

We can rewrite the utility function (S.A.19) as  $u_{nj} = z_{nj} + x_{nj}^0 \beta^0 + \epsilon_{nj}$  where  $\epsilon_{nj} = (\beta_{nj} + x_{nj}^0 \beta_n)$  for  $j = 1, \dots, J$ . The vector  $\epsilon_n = (\epsilon_{n1}, \dots, \epsilon_{nJ})$  is often called the composite error in the presence of random coefficients. Theorem 2.1, if this composite error vector  $\epsilon_n$  satisfies Assumptions I1 (Exclusion Restriction) and I2 (Central Symmetry), then  $\beta^0$  is point identified under the regularity conditions given by Assumptions I3-I5. We now give sufficient conditions for I1 and I2 to hold with random coefficients.

Assumption RC

RC1: Condition on almost every  $X \in S_X$ , the covariate vector  $z$  is independent of  $(\epsilon_j)$ ;

and the conditional distribution function of



to know exactly which covariates have random coefficients and which do not. Last, our model does not require thin tails or unimodality, unlike, e.g., normal random coefficient MNP models<sup>4</sup>.

One restriction we do impose is that we require one covariate in each choice  $z_j$ , not have a random coefficient. Setting the coefficient of some covariate equal to one is often a natural, economically meaningful normalization. For example, utility of choices are typically modeled as benefits minus costs. Benefits may be subjective and so vary heterogeneously as in random coefficients, while costs are often objective and fixed. In these cases  $z_j$  would be a cost measure. Examples are willingness to pay studies where the benefits equal the willingness to pay, and consumer choice applications where  $z_j$  is the price of choice  $j$ . (See e.g., Bliemer and Rose 2013 for more discussion and examples<sup>5</sup>) Nevertheless, we could also assume that, before normalizing, the variable  $z$  has a random coefficient, provided that the random coefficient is the same for all choices and is positive (this latter restriction is a special case of the hemisphere condition required by semiparametric binary choice random coefficient estimators. See, e.g., Gautier and Kitamura 2013). This restriction is needed because we can't allow renormalizations that would change any individual's relative ranking of utilities. Note that in this case, we require our symmetry condition to hold after renormalization, not before.



for  $j = 0; \dots; J$ , where  $d_j$  is defined as the same as equation (S.A.17).

For each  $j$ , the function  $d_j(z^{(j)}; X^{(j)})$  is well defined if both points  $(z^{(j)}; X^{(j)})$  and  $(z^{(j)} - 2X^{(j)}; X^{(j)})$  are in the interior of the support of covariates,  $S_{(z^{(j)}; X^{(j)})}$ . For this reason, we only wish to evaluate the function  $d_j(z^{(j)}; X^{(j)})$  at such points. This can be achieved by multiplying each function  $d_j(z^{(j)}; X^{(j)})$  by a trimming function of the form

$$d_j(z^{(j)}; X^{(j)}; \bar{z}, \underline{z}) \otimes d_j(z^{(j)}; X^{(j)}) \otimes d_j(z^{(j)} - 2X^{(j)}; X^{(j)}) \otimes d_j(z^{(j)} - 2X^{(j)}; \bar{z}; X^{(j)});$$

where  $X^{(j)-}$  ( $X^{(j)-}$ ) gives the upper (lower) bound value that the index  $X^{(j)}$  can take. A simple choice for the function  $\otimes(\cdot)$  is  $\otimes(z^{(j)}; X^{(j)}) = 1 - |z^{(j)}|_j c_1^{(j)} - 1 - |X^{(j)}|_j C_2^{(j)}$ , where the absolute value of a vector or matrix,  $|z|_j$ , is defined as the corresponding vector or matrix of the absolute values of each element,  $c_1^{(j)} \in \mathbb{R}^J$  is a vector of trimming constants for the covariate vector  $z^{(j)}$ , and  $C_2^{(j)} \in \mathbb{R}^{J \times q}$  is a matrix of trimming constants for the covariate matrix  $X^{(j)}$  such that  $c_1^{(j)}; C_2^{(j)}$  is in the interior of the support of covariates  $S_{(z^{(j)}; X^{(j)})}$ . Denote  $S_{z^{(j)}; \bar{z}; \underline{z}}$  as the largest set of values  $z^{(j)}$  given  $\bar{z}, \underline{z}$ , and  $X^{(j)}$ , such that  $S_{z^{(j)}; \bar{z}; \underline{z}} \in \text{int } S_{z^{(j)}; X^{(j)}}$ .

We describe the regularity conditions on the trimming function in Assumption TR.

Assumption TR. The trimming function  $d_j(z^{(j)}; X^{(j)}; \bar{z}, \underline{z})$  is strictly positive and bounded on  $S_{z^{(j)}; \bar{z}; \underline{z}} \in \text{int } S_{X^{(j)}}$ , and is equal to zero on its complementary set for  $j = 0; \dots; J$ .

Theorem S.B.1 *If Assumptions I and TR hold, then (i)  $Q_j(\cdot) = 0$  for any  $\beta \in \mathbb{R}^J$  and (ii)  $Q_j(\cdot) = 0$  if and only if  $\beta = 0$ .*

Theorem S.B.1 shows identification based on the population objective function. Proofs is available at authors' webpage.

## S.B.2 MD Estimator and Regularity Conditions

Next, we derive the sample objective function based on population objective function and the asymptotic properties of the MD estimator. To ease notation, we denote the conditional means

$$E y_j | z^{(j)} = z_n^{(j)}; X^{(j)} = X_n^{(j)} \quad , \quad E z_n^{(j)} | X_n^{(j)} = z_{j;o}^{(j)}; X_n^{(j)} \quad ;$$

$$E y_j | z^{(j)} = z_n^{(j)} \quad 2X_n^{(j)}; X^{(j)} = X_n^{(j)} \quad , \quad E z_n^{(j)} \quad 2X_n^{(j)}; X_n^{(j)} \quad , \quad E z_n^{(j)} | X_n^{(j)}; \quad ;$$

and function

$$d_j ( ; z_n^{(j)}; X_n^{(j)} ) = \frac{E y_j | z^{(j)} = z_n^{(j)}; X^{(j)} = X_n^{(j)} \quad @ \quad E y_j | z^{(j)} = z_n^{(j)} \quad 2X_n^{(j)}; X^{(j)} = X_n^{(j)} \quad }{E z_n^{(j)} | X_n^{(j)} \quad E z_n^{(j)} | X_n^{(j)} \quad } \quad (S.B.2)$$

$$, \quad E z_n^{(j)} | X_n^{(j)} \quad , \quad E z_n^{(j)} | X_n^{(j)}; \quad ;$$

where  $E z_n^{(j)} | X_n^{(j)} \quad @ \quad E z_n^{(j)} | X_n^{(j)} = @ \quad @$  and  $E z_n^{(j)} | X_n^{(j)}; \quad$  is de..ned in

the similar way as  $E z_n^{(j)} | X_n^{(j)} \quad$ . Now, consider a leave-one-out (LOO) Nadaraya-Watson

(NW) estimator for  $E z_n^{(j)} | X_n^{(j)} \quad$  as  $\hat{\Lambda}_{j;o; n}^{(j)} ( ; ) = \frac{\frac{1}{N-1} \sum_{m=1; m \neq n}^N y_{mj} K_{h_z}^{(j)}(z_m^{(j)}) K_{h_x} X_m^{(j)}}{\frac{1}{N-1} \sum_{m=1; m \neq n}^N K_{h_z}^{(j)}(z_m^{(j)}) K_{h_x} X_m^{(j)}}$ , where

$$K_{h_z}^{(j)}(z_m^{(j)}) = \prod_{l=1}^J h_{z_l}^{-1} k(h_{z_l}^{-1} z_{ml}^{(j)}) \quad ; \quad \text{and} \quad K_{h_x} X_m^{(j)} = \prod_{l=1}^Q \prod_{r=1}^q h_{x_{lr}}^{-1} k(h_{x_{lr}}^{-1} x_{mlr}^{(j)}) \quad .$$

The properties of the kernel function  $k$  and those of the bandwidth  $h_z \quad (h_{z_1}; \quad ; h_{z_J})^0$  and

$h_x \quad ec-11s/T-2 /TT52z$

By replacing the expectation in  $Q_j(\cdot)$  with its sample mean and function  $d_j(\cdot; z_n^{(j)}; X_n^{(j)})$  with its LOO estimator  $\hat{d}_{j;n}(\cdot; z_n^{(j)}; X_n^{(j)})$ , we define the MD estimator

$$\hat{\theta}_j = \arg \min_{\theta} Q_{Nj}(\theta);$$

$$\text{where } Q_{Nj}(\theta) = \frac{1}{2N} \sum_{n=1}^N h_j(z_n^{(j)}; X_n^{(j)}) \hat{d}_{j;n}(\theta; z_n^{(j)}; X_n^{(j)})^2;$$

We denote the gradient of the objective function as  $q_{Nj}(\theta) = \nabla_{\theta} Q_{Nj}(\theta)$  and the Hessian matrix of the objective function as  $H_{Nj}(\theta) = \nabla_{\theta}^2 Q_{Nj}(\theta)$ : The smoothness of the objective function suggests the first-order condition (FOC):  $q_{Nj}(\hat{\theta}_j) = 0_q$ . Applying the standard first-order Taylor expansion to  $q_{Nj}(\hat{\theta}_j)$  around the true parameter vector  $\theta_0$  yields  $q_{Nj}(\hat{\theta}_j) = q_{Nj}(\theta_0) + H_{Nj}(\tilde{\theta}_j) (\hat{\theta}_j - \theta_0)$ , where  $\tilde{\theta}_j$  is a vector between the MD estimator  $\hat{\theta}_j$  and the true parameter vector  $\theta_0$ . Then the influence function will be given by

$$\hat{\theta}_j - \theta_0 = - [H_{Nj}(\tilde{\theta}_j)]^{-1} q_{Nj}(\theta_0); \quad (\text{S.B.4})$$

We will show that  $H_{Nj}(\tilde{\theta}_j) \xrightarrow{p} H_j(\theta_0)$ ; where

$$H_j(\theta_0) = E \left[ \sum_{j=1}^2 z_n^{(j)}; X_n^{(j)} \nabla_{\theta} d_j(\theta_0; z_n^{(j)}; X_n^{(j)}) \nabla_{\theta} d_j(\theta_0; z_n^{(j)}; X_n^{(j)})^i \right]; \quad (\text{S.B.5})$$

and  $\frac{1}{N} q_{Nj}(\theta_0) \xrightarrow{p} N(0_q; \Sigma_j)$ , where  $\Sigma_j$  is the probability limit of the variance-covariance matrix of  $q_{Nj}(\theta_0)$ . To obtain these properties, we assume the following regularity conditions.

Assumption E.

E1:  $f(y_n; z_n; X_n)$ , for  $n = 1, \dots, N$  is a random sample drawn from the infinite population distribution.

E2: The following smoothness conditions hold: (a) The density function  $f_j(z^{(j)}; X^{(j)})$  is continuous in the components of  $z^{(j)}$  for all  $z^{(j)} \in S_z^{\text{Tr}}(X^{(j)}; \underline{\cdot}); \underline{\cdot}$  and  $X^{(j)} \in \text{int}(S_X)$ .

In addition,  $f_j(z^{(j)}; X^{(j)})$  is bounded away from zero uniformly over its support. (b)

Functions  $f_j(z^{(j)}; X^{(j)})$ ,  $g_j(z^{(j)}; X^{(j)})$  and  $\rho_j(z^{(j)}; X^{(j)})$  are  $s$  ( $s = J + 1$ ) times continuously differentiable in the components of  $z^{(j)}$  for all  $z^{(j)} \in S_z^{\text{Tr}}(X^{(j)}; \underline{\cdot}); \underline{\cdot}$  and have bounded derivatives.

E3: The kernel function  $k$  is an  $l$ -th ( $l = 1$ ) order bias-reducing kernel that satisfies (a)

$k(u) = k(-u)$  for any  $u$  in the support of

Lipschitz conditions: for some  $L$  and  $K$ ,

$$\left| \frac{\partial f_j}{\partial z^{(i)}}(z^{(i)} + t; \cdot) - \frac{\partial f_j}{\partial z^{(i)}}(z^{(i)}; \cdot) \right| \leq L \|z^{(i)}\|; \quad \left| \frac{\partial f_j}{\partial t}(z^{(i)} + t; \cdot) - \frac{\partial f_j}{\partial t}(z^{(i)}; \cdot) \right| \leq K$$

(b) (Asymptotic Normality) *The MD estimator is asymptotically normal, i.e.,*

$$\sqrt{N}(\hat{\theta}_N - \theta_0) \xrightarrow{d} N(0, H_j^{-1} J_j H_j^{-1})$$

where matrix  $J_j = E[t_{nj} t_{nj}^0]$  and  $H_j$

would be equal, while under the alternative there must exist symmetric points where the densities are not equal. Also under the null, our estimator is consistent. So a test could be constructed based on the difference in error density estimates at many symmetry points (other than those used for estimation), using our estimated parameters to construct symmetry points. More general specification tests could also be constructed, using the fact that our parameters are over identified when more than one choice is observed.

### S.C Proof of Identification

Proof of Theorem S.A.1: First, we show that  $D_0(\theta)$  is a set of measure zero. If not, assume that there is a point  $(z; X)$  in set  $D_0(\theta)$ . By definition in equation (8), both points  $(z; X)$  and  $(z + 2X - \theta; X)$  are in set  $\text{int}(S_{(z; X)})$ . By Assumptions I1, I2, and equations (5)-(7), we have function

$$d_0(\theta; z; X) = (1)^{-J} [f_{\theta}(z - X - \theta; X = X) - f_{\theta}(z + X - \theta; X = X)] = 0;$$

which is a contradiction with definition in equation (8).

Next, we prove that  $\Pr[(z; X) \in D_0(\theta)] > 0$  for any  $\theta = \theta_0$ , where  $\theta_0$  and parameter space satisfies Assumption I3. Denote the set  $\mathcal{X}(\theta) = \{X \in S_X \mid X(\theta) = \theta\}$ , which is a collection of covariate values at which  $X = \theta$ . By Assumption I4(a) and the fact  $\theta_0 = \theta_0$ ,  $\mathcal{X}(\theta_0)$  is a subset in the support of  $S_X$  with positive measure, that is,

$$\Pr[X \in \mathcal{X}(\theta_0)] > 0; \tag{S.C.1}$$

Recall that we use  $X_c$  and  $X_d$ , respectively, to denote the continuous and discrete covariates in  $X$ . We define the interior of the support of  $X$  as  $\text{int}(S_X)$

$(X_c; X_d) \in S_{(X_c; X_d)} \mid X_c \in \text{int}(S_{X_c}(X_d)); X_d \in S_{X_d}$ . Define

$$\mathcal{S}_{(z; X)}(\theta) = \{(z; X) \in S_{(z; X)} \mid z \in \mathcal{S}_Z(X); X \in \mathcal{X}(\theta) \cap \text{int}(S_X)\}; \tag{S.C.2}$$

where  $\mathcal{S}_z(X)$  satisfies Assumption I4(c). By construction,  $\text{set}(\mathcal{S}_{(z;X)}(\cdot))$  is a Lebesgue measurable subset of  $\text{int}(\mathcal{S}_{(z;X)})$ . Next we construct a subset in the support of covariates  $(z; X)$  as follows:

$$\mathcal{D}_0(\cdot) = \bigcap_{i=1}^n (z; X)_i \in \mathcal{S}_{(z;X)}(\cdot)$$



## S.D Monte Carlo Details

As discussed in the paper, our Monte Carlo design includes 4 data generating processes (DGPs).

Details of the distribution of each DGP are provided in Table 1.

Table 1: Designs of the Data Generating Processes (DGPs)

DGP	Distribution of $\mathbf{z}_n$	Distribution of $\mathbf{u}_{nj}$
1	$\mathbf{z}_n = 0:2$	$\mathbf{u}_{nj} = \mathbf{u}_{nj}$
2	$\mathbf{z}_n = 0:2$	$\mathbf{u}_{nj} = \frac{1}{2} e^{2x_{nj}} \mathbf{u}_{nj}$ ,
3	$\mathbf{z}_n = 0:2 + \mathbf{z}_n$ where $\mathbf{z}_n = \frac{1}{2} \#_n$	$\mathbf{u}_{nj} = \frac{1}{2} \mathbf{u}_{nj}$
4	$\mathbf{z}_n = 0:2 + \mathbf{z}_n$ where $\mathbf{z}_n = (e^{x_{n1}} + e^{x_{n2}}) \#_n$	$\mathbf{u}_{nj} = \frac{1}{2} \mathbf{u}_{nj}$

*Note: both  $\#_n$  and  $\mathbf{u}_{nj}$  are standard normal random variables, and they are independent of each other and all the covariates, and i.i.d. across the subscripted dimension(s).*

For the MD estimator, we consider both the case where the researcher only observes whether the outside option (i.e., alternative 0) is chosen, and so just minimize  $Q_{N0}(\cdot)$ , and the case where the researcher also observes which alternative is chosen by each decision maker, and so minimizes the sum of  $Q_{Nj}(\cdot)$  for  $j = 0; 1; 2$ . In all DGPs, each covariate  $z_{nj}$  is a continuous uniform random variable over the interval  $[-9; 9]$  and  $x_{nj}$  is a binary variable that takes value of 0 or 1 with equal probability for  $j = 1; 2$ . The covariates of alternative 0 are  $z_{n0} = 0$  and  $x_{n0} = 0$ . All the observed covariates are independent of each other and are independent, identically distributed across the subscripted dimension(s).

We use a grid search to compute our MD estimator over a parameter space  $\theta \in [0; 8; 0; 8]$  with the bin width of 0.05. In the estimation of choice probabilities we apply a truncated normal density for the kernel function  $k_h(\cdot)$  with bandwidth  $h_j = \text{sd}(z_{nj}) N^{(1-22)}$ , where  $j = 1; 2$ . Our

are  $O(h^s)$  with  $s = J + 1$  and  $O(Nh^{-2J})$

Online Supplemental Appendix to: Semiparametric  
Identification and Estimation of Multinomial Discrete Choice  
Models using Error Symmetry

Arthur Lewbel<sup>y</sup>      Jin Yan<sup>z</sup>      Yu Zhou<sup>x</sup>

Original February 2019, revised December 2021

## S.D Proofs Regarding Estimation

In this section, we provide the proofs of Theorems S.B.1-S.B.3 in Section S.B of the Supplementary Appendix and their related lemmas. Specifically, Section S.D.1 provides the proof of Theorem S.B.1 on the population sample objective function; Section S.D.2 collects preliminary lemmas needed for the asymptotic properties of the MD estimator defined in Section S.B.2; Section S.D.3 provides the proofs of Theorem S.B.2, the consistency of the MD estimator; and Section S.D.4 gives the proofs of Theorem S.B.3, the asymptotic linearity and normality of the estimator and related lemmas. Throughout this appendix, we use the same notations and acronyms defined in the main text.

---

## S.D.1 Proof of the Population Objective Function

Proof of Theorem S.B.1: Part (i) can be shown directly from the quadratic form of the population objective function. We will explicitly prove that Part (ii) holds. To show the existence of a minimizer, recall the population objective function

$$Q_j(\cdot)$$

mentary Appendix Section S.A.1; and the second term equals to zero since  $\phi^{(j)}(z_n^{(j)}; X_n^{(j)}) = 0$ .

Q.E.D.

Proof of Lemma S.D.2: The proofs for three terms are similar. We will focus on the proof for  $g_j(z_n^{(j)}; X_n^{(j)})$ . Other terms can be done in a similar fashion. First, by the fact that the outcome variables are binary and function  $f_j$  is bounded away from zero, applying the results of Lemma B.1 and Lemma B.2 in Newey (1994) gives the first equality in each equation. Second, the second equality follows from Assumption 10 using Lemma 8.10 in Newey and McFadden (1994) Q.E.D.

Define  $f_j^{(t)}(z_n^{(j)}; X_n^{(j)}) = \frac{\partial^t f_j}{\partial z_{r(t)}^{(j)}}$  be the derivative with respect to  $z_{r(t)}^{(j)}$ , where  $z_{r(t)}^{(j)} = z_{1:(t)}^{(j)}$ ;  $z_{r(t)}^{(j)}$  be any t-element of  $z^{(j)}$ . Similarly, we can define  $g_j^{(t)}(z_n^{(j)}; X_n^{(j)})$ ,  $h_j^{(t)}(z_n^{(j)}; X_n^{(j)})$  and  $v_j^{(t)}(z_n^{(j)}; X_n^{(j)})$ :

Lemma S.D.3 Under Assumptions E2-E5, for  $t = 1; \dots; J$ ,

$$\sup_{z_n^{(j)}; X_n^{(j)}} \frac{1}{2S_{(z^{(j)}; X^{(j)})}^{Tr}} \left( f_j^{(t)}(z_n^{(j)}; X_n^{(j)}) - f_j^{(t)}(z_n^{(j)}; X_n^{(j)}) \right) = O_p \left( \frac{\ln N}{Nh_N^{J+2t+Jq}} + h_N^s \right) = o_p(N^{-1/4})$$

$$\sup_{z_n^{(j)}; X_n^{(j)}} \frac{1}{2S_{(z^{(j)}; X^{(j)})}^{Tr}} \left( g_j^{(t)}(z_n^{(j)}; X_n^{(j)}) - g_j^{(t)}(z_n^{(j)}; X_n^{(j)}) \right) = O_p \left( \frac{\ln N}{Nh_N^{J+2t+Jq}} + h_N^s \right) = o_p(N^{-1/4})$$

$$\sup_{z_n^{(j)}; X_n^{(j)}} \frac{1}{2S_{(z^{(j)}; X^{(j)})}^{Tr}} \left( h_j^{(t)}(z_n^{(j)}; X_n^{(j)}) - h_j^{(t)}(z_n^{(j)}; X_n^{(j)}) \right) = O_p \left( \frac{\ln N}{Nh_N^{J+2t+Jq}} + h_N^s \right) = o_p(N^{-1/4})$$

Proof of Lemma S.D.3: The proof follows the same method used in Lemma S.D.2. Q.E.D.

Lemma S.D.4 Under Assumptions E2-E5, for  $t = 1; \dots; J$ ,

$$\sup_{z_n^{(j)}; X_n^{(j)}} r \hat{f}_j^{(t)}(z_n^{(j)}; X_n^{(j)}) - r f_j^{(t)}(z_n^{(j)}; X_n^{(j)})$$

$$= O_p \left( \frac{s}{\ln N} \right)$$

Condition (4) following Hong and Tamer (2003). We first introduce an infeasible sample objective function  $Q_{Nj}(\cdot)$ , defined as

$$Q_{Nj}(\cdot) = \frac{1}{2N} \sum_{n=1}^N \sum_j h_j(z_n^{(j)}; X_n^{(j)}) d_j(z_n^{(j)}; X_n^{(j)})^2$$

Following the triangle inequality, we have

$$|Q_{Nj}(\cdot) - Q_j(\cdot)| \leq Q_{Nj}(\cdot) + Q_j(\cdot) \tag{S.D.3}$$

Then, it is sufficient to show that the two terms on the right side of (S.D.3) go to zero uniformly, that is, (i)  $\sup_{\mathcal{Z}} Q_{Nj}(\cdot) - Q_{Nj}(\cdot) = o_p(1)$  and (ii)  $\sup_{\mathcal{Z}} Q_{Nj}(\cdot) - Q_j(\cdot) = o_p(1)$ .

For Part (i), we observe that

$$\sup_{\mathcal{Z}} Q_{Nj}(\cdot) - Q_{Nj}(\cdot) \tag{S.D.4}$$

$$= \sup_{\mathcal{Z}} \frac{1}{2N} \sum_{n=1}^N \sum_j h_j^2(z_n^{(j)}; X_n^{(j)}) \hat{d}_{j,n}^2(z_n^{(j)}; X_n^{(j)}) - d_j^2(z_n^{(j)}; X_n^{(j)})$$

$$= \sup_{\mathcal{Z}} \frac{1}{2N} \sum_{n=1}^N \sum_j h_j^2(z_n^{(j)}; X_n^{(j)}) \hat{d}_{j,n}(z_n^{(j)}; X_n^{(j)}) + d_j(z_n^{(j)}; X_n^{(j)})$$

$$\hat{d}_{j,n}(z_n^{(j)}; X_n^{(j)}) - d_j(z_n^{(j)}; X_n^{(j)})$$

$$\leq C \sup_{z_n^{(j)}; X_n^{(j)}} \sup_{\mathcal{S}_{Tr}(z^{(j)}, X^{(j)})} \hat{d}_{j,n}(z_n^{(j)}; X_n^{(j)}) - d_j(z_n^{(j)}; X_n^{(j)}) = o_p(1)$$

The first equality in (S.D.4) follows from definition and direct calculation. The second equality holds by factorization. The next inequality is satisfied by the fact that functions  $j$  and  $d_j$  are bounded QEXS. The last equality follows the fact that

$$\sup_{\mathcal{Z}} \sup_{z_n^{(j)}; X_n^{(j)}} \hat{d}_{j,n}(z_n^{(j)}; X_n^{(j)}) - d_j(z_n^{(j)}; X_n^{(j)})$$

is bounded by the product of a constant and the derivative functions shown by Lemma S.D.3.



Part (ii) holds by showing pointwise convergence and stochastic equicontinuity. By the Law of Large Numbers (LLN), we can directly obtain the pointwise convergence of  $Q_{Nj}^{(1)}$  to  $Q_j^{(1)}$ . Next we can conclude the uniformity by showing stochastic equicontinuity, that is,

$$\sup_{(1); (2) j; (1) (2) jj} Q_{Nj}^{(1)} - Q_{Nj}^{(2)} = o_p(1):$$

Following Andrews (1994), the stochastic equicontinuity can be shown by verifying that  $Q_{Nj}^{(1)}$  is the type II class of function, satisfying the Lipschitz condition  $Q_{Nj}^{(1)} - Q_{Nj}^{(2)} \leq C_{jj}^{(1)} - C_{jj}^{(2)}$ . We verify that this holds from the continuity of the quadratic form of the objective function and the continuity of the kernel derivative functions with bounded second derivatives. Q.E.D.

#### S.D.4 Asymptotic Linearity and Normality of the MD Estimator

In this section, we first show the lemmas that contribute to the proof of Theorem S.B.3

Lemma S.D.5 Under Assumptions I,TR and E,  $H_{Nj} \sim \sqrt{N} H_j$ , where

$$H_j = E \int_j^2 z_n^{(j)}; X_n^{(j)} r_{dj} \circ; z_n^{(j)}; X_n^{(j)} r_{dj} \circ; z_n^{(j)}; X_n^{(j)} i^0$$

Proof of Lemma S.D.5: To show the desired result, we first show that the following results hold:

(i)  $H_{Nj} \sim H_{Nj;1} \sim + H_{Nj;2} \sim$ ; where

$$H_{Nj;1} \sim = \frac{1}{N} \sum_{n=1}^N \int_j^2 z_n^{(j)}; X_n^{(j)} \hat{d}_{j;n} \sim; z_n^{(j)}; X_n^{(j)} r_{dj} \circ \hat{d}_{j;n} \sim; z_n^{(j)}; X_n^{(j)} i^0;$$

$$H_{Nj;2} \sim = \frac{1}{N} \sum_{n=1}^N \int_j^2 z_n^{(j)}; X_n^{(j)} r_{dj} \hat{d}_{j;n} \sim; z_n^{(j)}; X_n^{(j)} r_{dj} \hat{d}_{j;n} \sim; z_n^{(j)}; X_n^{(j)} i^0;$$

(ii)  $H_{Nj;1} \sim = o_p(1)$ ; and (iii)  $H_{Nj;2} \sim \sqrt{N} H_j$ .

The decomposition in Part (i) follows from direct calculation. For Part (ii), observe that

$$H_{N_j;1} \tilde{\sim} = H_{N_j;1}^h \tilde{\sim} H_{N_j;1}^i(\theta) + H_{N_j;1}(\theta) = o_p(1)$$

Given that  $\tilde{\sim}$  lies between  $\theta$  and  $\tilde{\sim}$ , we get that  $\tilde{\sim}$  is uniformly consistent, and by applying the

Delta method for the continuity of the choice probability, we obtain that  $H_{N_j;1} \tilde{\sim} H_{N_j;1}(\theta) =$

$o_p(1)$ . Next,  $H_{N_j;1}(\theta)$



Next we calculate

$$\begin{aligned}
 & q_{Nj;2}^{(0)} \tag{S.D.7} \\
 &= \frac{1}{N} \sum_{n=1}^N \sum_j^J z_n^{(j)}; X_n^{(j)} \left[ \sum_{i=0}^h \lambda_{j;o;n}^{(j)} z_n^{(j)}; X_n^{(j)} + \sum_{i=0}^i \lambda_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)} \right] \\
 &\quad + \sum_{i=0}^h r_{d_j} z_n^{(j)}; X_n^{(j)} + \sum_{i=0}^i r_{d_j} z_n^{(j)}; X_n^{(j)} \\
 &= A_1 + A_2 + A_3 + A_4
 \end{aligned}$$

where

$$A_1 = \frac{1}{N} \sum_{n=1}^N \sum_j^J z_n^{(j)}; X_n^{(j)} \left[ \sum_{i=0}^h \lambda_{j;o;n}^{(j)} z_n^{(j)}; X_n^{(j)} + \sum_{i=0}^i \lambda_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)} \right]$$

$$A_2 = \sum_{j=1}^J \sum_n^N z_n^{(j)}; X_n^{(j)} \left[ \sum_{i=0}^h \lambda_{j;o;n}^{(j)} z_n^{(j)}; X_n^{(j)} + \sum_{i=0}^i \lambda_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)} \right]$$

x.75 Td (0)Tj /TT21 Td (z(j))Tj /TT3-249 T3()Tj /TT2 7.97 Tf (j)Tj /TT3 7.97 Tf 7.177^



where  $N(! m; ! n) = N_{;o}(! m; ! n) + N_{;cs}(! m; ! n)$  with

$$N_{;o}(! m; ! n) = \sum_j z_n^{(j)}; X_n^{(j)}; \circ y_{mj} \quad ; \quad j; o \quad z_n^{(j)}; X_n^{(j)}$$

$$K_{h_z}^{(J)} z_m^{(j)} z_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)} f_j^{-1} z_n^{(j)}; X_n^{(j)} ;$$

$$N_{;cs}(! m; ! n) = \sum_j z_n^{(j)}; X_n^{(j)}; \circ y_{mj} \quad ; \quad j; cs \quad z_n^{(j)}; X_n^{(j)}; \circ$$

$$K_{h_z}^{(J)} z_m^{(j)} z_n^{(j)} 2X_n^{(j)} \circ K_{h_x} X_m^{(j)} X_n^{(j)} f_j^{-1} z_n^{(j)} 2X_n^{(j)} \circ; X_n^{(j)} :$$

where  $K_{h_z}^{(J)} z_m^{(j)} = \sum_{l=1}^J h_N^{2J} k^{(l)} h_{z_l}^{-1} z_{ml}^{(j)}$  where  $k^{(1)}$  is the ...rst derivative of kernel function.

Proof of Lemma S.D.7: We ...rst observe that

$$q_{Nj;1}(\circ) = \frac{1}{N} \sum_{n=1}^N \sum_j z_n^{(j)}; X_n^{(j)}; \circ \Delta_{j; n} \circ; z_n^{(j)}; X_n^{(j)} \quad (S.D.9)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_j z_n^{(j)}; X_n^{(j)}; \circ h \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)} \wedge_{j; cs; n}^{(J)} z_n^{(j)}; X_n^{(j)}; \circ^i :$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_j z_n^{(j)}; X_n^{(j)}; \circ nh \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)} \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)}^i$$

$$h \wedge_{j; cs; n}^{(J)} z_n^{(j)}; X_n^{(j)}; \circ \wedge_{j; cs; n}^{(J)} z_n^{(j)}; X_n^{(j)}; \circ^{io}$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_j z_n^{(j)}; X_n^{(j)}; \circ nh \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)} E \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)} z_n^{(j)}; X_n^{(j)}^i$$

$$+ E h \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)} z_n^{(j)}; X_n^{(j)}^i \wedge_{j; o; n}^{(J)} z_n^{(j)}; X_n^{(j)}^i$$

$$h \wedge_{j; cs; n}^{(J)} z_n^{(j)}; X_n^{(j)}; \circ E h \wedge_{j; cs; n}^{(J)} z_n^{(j)}; X_n^{(j)}; \circ z_n^{(j)}; X_n^{(j)}^i$$

$$+ E h \wedge_{j; cs; n}^{(J)} z_n^{(j)}; X_n^{(j)}; \circ X$$

The second, third and fourth equalities follows from adding and subtracting terms. The last equality holds by the fact that

$$\sup_{z_n^{(j)}; X_n^{(j)}} E \left[ \Lambda_{j;0;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) - \Lambda_{j;0;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) \right] = O(h^s)$$

$$\sup_{z_n^{(j)}; X_n^{(j)}} E \left[ \Lambda_{j;cs;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) - \Lambda_{j;cs;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) \right] = O(h^s)$$

Next, to derive  $\Lambda_{j;0;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) - \Lambda_{j;0;n}^{(j)}(z_n^{(j)}; X_n^{(j)})$ , we observe that

$$\Lambda_{j;0;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) - \Lambda_{j;0;n}^{(j)}(z_n^{(j)}; X_n^{(j)}) = \int \dots$$

Since each term is of order  $O_p N^{-1/2}$ , thus  $R_{o;1}$  is of order  $O_p N^{-1/2}$ . Denoting

$$N_{;o;1}(!m;!n) = \sum_j z_n^{(j)}; X_n^{(j)}; \circ$$

$$f_j^{-1} y_{mj} \cdot \sum_{j;o}^{(j)} z_n^{(j)}; X_n^{(j)} K_{h_z}^{(j)} z_m^{(j)} z_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)}$$

will give the first term of  $N_{;1}(!m;!n)$ .

In addition, to derive  $\sum_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)}; \circ E^h \sum_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)}; \circ z_n^{(j)}; X_n^{(j)}; \circ^i$ , we observe that

serve that

$$\begin{aligned} & \sum_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)}; \circ E^h \sum_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)}; \circ z_n^{(j)}; X_n^{(j)}; \circ^i \\ &= f_j^{-1} y_{mj} K_{h_z}^{(j)} z_m^{(j)} z_n^{(j)} z_n^{(j)} z_n^{(j)} \circ X_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)} \\ & E^h \sum_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)}; \circ \sum_{j;cs;n}^{(j)} z_n^{(j)}; X_n^{(j)}; \circ^i \\ &= f_j^{-1} \frac{1}{N} \sum_{m=1; m \leq n}^X \sum_{j;cs}^{(j)} y_{mj} z_n^{(j)}; X_n^{(j)}; \circ \\ & K_{h_z}^{(j)} z_m^{(j)} z_n^{(j)} z_n^{(j)} z_n^{(j)} \circ X_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)} + R_{cs;1} \end{aligned}$$

where the second equality holds by the same argument for  $f_j^{-1}$ , and  $R_{cs;1}$  collects the higher order terms from the decomposition of  $f_j^{-1}$ , with the order of  $O_p N^{-1/2}$ , by the same argument as above. Denoting

$$N_{;cs;1}(!m;!n) = \sum_j z_n^{(j)}; X_n^{(j)}; \circ$$

$$f_j^{-1} y_{mj} \cdot \sum_{j;cs}^{(j)} K_{h_z}^{(j)} z_m^{(j)} z_n^{(j)} z_n^{(j)} z_n^{(j)} \circ X_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)}$$

will give the second term of  $N_{;1}(!m;!n)$ .

Combining all the terms gives the desired results. Q.E.D.



Lemma S.D.8 Under Assumptions E2-E5,

$$\frac{1}{N(N-1)} \sum_{m=1}^N \sum_{n=1; n \neq m}^N \mathbb{E}(\mathbb{1}_{m;n}) = \frac{1}{N} \sum_{m=1}^N t_{mj} + o_p(N^{-1/2})$$

and

$$N^{-1/2} \sum_{m=1}^N t_{mj} \xrightarrow{d} N(0, \sigma_j^2);$$

where  $t_{mj} = (z_m^{(j)}; X_m^{(j)})'$  and  $\sigma_j^2 = \text{Var}(z_m^{(j)}; X_m^{(j)})$ .

$$E \sum_{N;0}^h (\sum_{m;1}^i \sum_{n;1}^i) k^2 \tag{S.D.10}$$

$$= \sum_{h_z}^{(j)} K_{h_z}^{(j)} \sum_{z_m^{(j)}} \sum_{z_n^{(j)}} K_{h_x} \sum_{X_m^{(j)}} \sum_{X_n^{(j)}}^2$$

$$\sum_{j;1}^h \sum_{z_m^{(j)}; X_m^{(j)}} + \sum_{j;2}^2 \sum_{z_n^{(j)}; X_n^{(j)}} \sum_{j;1}^2 \sum_{z_m^{(j)}; X_m^{(j)}} \sum_{j;2}^2 \sum_{z_n^{(j)}; X_n^{(j)}}^i$$

$$f_{j;1} \sum_{z_m^{(j)}; X_m^{(j)}} f_{j;2} \sum_{z_n^{(j)}; X_n^{(j)}} f_{j;1} \sum_{z_m^{(j)}; X_m^{(j)}} \sum_{j;2}^2 \sum_{z_n^{(j)}; X_n^{(j)}}; \int dz_m^{(j)} dX_m^{(j)} dz_n^{(j)} dX_n^{(j)}$$

$$= \sum_{h_z}^{(j)} K_{h_z}^{(j)} \sum_{u_z^{(j)}} K_{h_x} \sum_{u_x^{(j)}}^2$$

$$\sum_{j;1}^h \sum_{z_m^{(j)}; X_m^{(j)}} + \sum_{j;2}^2 \sum_{z_m^{(j)}} u_z^{(j)} h_{z; X_m^{(j)}} \wedge u_x^{(j)} h_x$$

$$\sum_{j;1}^2 \sum_{z_m^{(j)}; X_m^{(j)}} \sum_{j;2}^2 \sum_{z_m^{(j)}} u_z^{(j)} h_{z; X_m^{(j)}} \wedge u_x^{(j)} h_x^i$$

$$f_{j;1} \sum_{z_m^{(j)}; X_m^{(j)}} \sum_{j;2}^2 \sum_{z_m^{(j)}} u_z^{(j)} h_{z; X_m^{(j)}} \wedge u_x^{(j)} h_x; \int dz_m^{(j)} dX_m^{(j)} du_z^{(j)} du_x^{(j)}$$

$$= O(h_N^{2J} N^{Jq}) = O(N^{-1} N h_N^{2J+J+rfJq}) = o(N);$$

where the first equality in (S.D.10) follows from definitions; the second equality holds using a change of variables; and the third equality is satisfied by Assumptions E3 and E4. The desired result then follows from Assumption E4. Similarly, we can show  $E \sum_{N;cs}^h (\sum_{m;1}^i \sum_{n;1}^i) k^2 = o(N)$ .

Next we show that the second term in  $\hat{U}_n$  contributes to the asymptotic linearity and normality, while the first and third terms are asymptotically negligible. In sum, we show that (i)

$$E[\sum_{N;1}^h (\sum_{m;1}^i \sum_{n;1}^i)] = E[r_{N1}(\sum_{m;1}^i)] = E[r_{N2}(\sum_{n;1}^i)] = o(N^{-1/2}), \text{ (ii) } \frac{1}{N} \sum_{n=1}^N (r_{N2}(\sum_{n;1}^i) - E[r_{N2}(\sum_{n;1}^i)])^2 = N^{-1} \sum_{n=1}^N t_{mj}^2; \text{ where } N^{-1/2} \sum_{n=1}^N t_{mj}^2 = 1.636 \text{ Td}(\cdot) - 2$$



Second, to show Part (ii) holds, by direct calculation, we have

$$\begin{aligned}
 r_{N2;0}(\mathbf{!}_n) &= E[ N_{;0}(\mathbf{!}_m; \mathbf{!}_n) \mathbf{!}_n ] && \text{(S.D.11)} \\
 &= E \int_j^h z_n^{(j)}; X_n^{(j)}; \circ y_{mj} \quad '_{j;0} z_n^{(j)}; X_n^{(j)} \\
 &K_{h_z}^{(j)} z_m^{(j)} z_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)} f_j^{-1} z_n^{(j)}; X_n^{(j)} \mathbf{!}_n^i \\
 &= E \int_j^h z_n^{(j)}; X_n^{(j)}; \circ \quad '_{j;0} z_m^{(j)}; X_m^{(j)} \quad '_{j;0} z_n^{(j)}; X_n^{(j)} \\
 &K_{h_z}^{(j)} z_m^{(j)} z_n^{(j)} K_{h_x} X_m^{(j)} X_n^{(j)} f_j^{-1} z_n^{(j)}; X_n^{(j)} \mathbf{!}_n^i \\
 &= \int_j^Z z_n^{(j)}; X_n^{(j)}; \circ \frac{ @ \quad '_{j;0} z_n^{(j)} + u_z^{(j)} h_z; X_n^{(j)} + \wedge u_x^{(j)} h_x \quad '_{j;0} z_n^{(j)}; X_n^{(j)} }{ @_z^{(j)} } \\
 &K_{h_z} u_z^{(j)} K_{h_x} u_x^{(j)} f_j^{-1} z_n^{(j)}; X_n^{(j)} \\
 &f_j z_n^{(j)} + u_z^{(j)} h_z; X_n^{(j)} + \wedge u_x^{(j)} h_x \quad du_z^{(j)} du_x^{(j)} \\
 &= O(h_N^s) = o N^{-1/2} :
 \end{aligned}$$

The last second equality follows from integration by parts and a Taylor expansion. We therefore

get that  $\frac{1}{N} \sum_{n=1}^N (r_{N2}(\mathbf{!}_n) - E[r_{N2}(\mathbf{!}_n)]) = o_p N^{-1/2}$ .



where

$$r_o(l_m) = \frac{Z @_{j, z_m^{(j)}; X_m^{(j)}} \circ '_{j, z_m^{(j)}; X_m^{(j)}}}{@_{f}^{(j)} @_{f}^{(j)}} K_{h_z} u$$

probability to zero. Note that

$$\begin{aligned}
 & \frac{1}{N} \sum_{m=1}^N (r_{o(j,m)} - E[r_{o(j,m)}]) \tag{S.D.15} \\
 &= \frac{1}{N} \sum_{m=1}^N \left( z_m^{(j)}; X_m^{(j)} - \frac{z_m^{(j)}; X_m^{(j)}}{z^{(j)}; X^{(j)}} \right) \\
 &= \frac{1}{N} \sum_{m=1}^N y_{mj} - \frac{z_m^{(j)}; X_m^{(j)}}{z^{(j)}; X^{(j)}} \\
 &+ E \left[ y_{mj} - \frac{z_m^{(j)}; X_m^{(j)}}{z^{(j)}; X^{(j)}} \right]
 \end{aligned}$$

and

$$\begin{aligned}
 & \frac{1}{N} \sum_{m=1}^N (r_{cs(j,m)} - E[r_{cs(j,m)}]) \tag{S.D.16} \\
 &= \frac{1}{N} \sum_{m=1}^N \left( z_m^{(j)}; X_m^{(j)} - \frac{z_m^{(j)}; X_m^{(j)}}{z^{(j)}; X^{(j)}} \right) \\
 &= \frac{1}{N} \sum_{m=1}^N \left( z_m^{(j)}; X_m^{(j)} - \frac{z_m^{(j)}; X_m^{(j)}}{z^{(j)}; X^{(j)}} \right)
 \end{aligned}$$

Then

$$\rho = \frac{1}{N} \sum_{m=1}^N (r_m)$$



Hong, H. and Tamer E. (2003): "Inference in Censored Models with Endogenous Regressors,"  
Econometrica, 71(3), 905-932.

Powell, J.L., Stock, J., and Stocker, T. (1989): "Semiparametric Estimation of Index Models,"  
Econometrica, 57, 1403-1430.

Newey, W.K. (1994): "Kernel Estimation of Partial Means and a General Variance Estimator,"  
Econometric Theory, 10(2), 233-253

Newey, W.K. and McFadden, D. (1994): "Large Sample Estimation and Hypothesis Testing,"  
Handbook of Econometrics Vol. 4, 2111-2245.