# Identification and Estimation of Average Marginal Effects in Fixed Effects Logit Models[*]

Laurent Davezies[†]    Xavier D'Haultfœuille[‡]    Louise Laage[§]

## Abstract

This article considers average marginal effects (AME) in a panel data fixed effects logit model. Relating the identified set of the AME to an extremal moment problem, we first show how to obtain sharp bounds on the AME straightforwardly, without any optimization. Then, we consider two strategies to build confidence intervals on the AME. In the first, we estimate the sharp bounds with a semiparametric two-step estimator. The second, very simple strategy estimates instead a quantity known to be at a bounded distance from the AME. It does not require any nonparametric estimation but may result in larger confidence intervals. Monte Carlo simulations suggest that both approaches work well in practice, the second being often very competitive. Finally, we show that our results also apply to average treatment effects, the average structural functions and ordered, fixed effects logit models.

**Keywords:** Fixed effects logit models, panel data, partial identification.

**JEL Codes:** C14, C23, C25.

# 1 Introduction

In this paper, we consider the identification and estimation of average marginal effects (AME) and average treatment effects (ATE) in the fixed effects (FE) binary logit model, with short panels. Estimation of the common slope parameters dates back to Rasch (1961) and Andersen (1970) (see also Chamberlain, 1980) but up to now, there has been no study of the identification and estimation of the AME and ATE in this model. These parameters are yet of more direct interest than the slope parameter, which only provides information on relative marginal effects. For this reason, and following the influential work of Angrist (2001) (see also Angrist and Pischke, 2008), many applied economists have turned to using FE linear probability models. Besides their simplicity, they allow one to identify the best linear approximation of the true model, if it is nonlinear. Then, one could argue that their slope parameter, which corresponds to the AME if the model is linear, is close to the true AME in practice.

However, FE linear models can be problematic in panel data if "stayers", namely units with constant covariates across the period, differ from "movers" in their unobserved characteristics. The reason is that the AME in FE linear models is identified using "movers" only. But the true AME may depend on the unobserved heterogeneity and thus be very different for the population of "stayers". Then, approximating the AME of the whole population using the linear model may be far from the truth. This is especially the case when the proportion of stayers is large, something we illustrate numerically in Appendix A. Also, it seems unfortunate that FE linear models only rely on movers given that in fact, the AME is nonparametrically identified for the population of stayers only (Hoderlein and White, 2012).[1]

Unlike the FE linear model, nonlinear models such as the FE logit models do allow for heterogeneity of treatment effects, in particular between stayers and movers. Moreover, we demonstrate in this paper that estimation and inference on the AME and the ATE in this model can be performed almost as simply as in the FE linear model. To this end, we first study in Section 2 the identification of the AME (the

---

[1]Nonlinearities may also cause trouble when using the FE linear model. In Appendix A, we give a simple example, with a difference-in-differences flavor, where the FE linear model identifies a negative ATE, even though the true ATE is positive.

study of the ATE, which is very similar, is postponed to Section 5). If this parameter is generally not point identified, sharp bounds can be obtained by solving an extremal moment problem, that is, maximizing a moment over probability distributions given the knowledge of some other moments. Using existing results on such problems, the bounds can then be obtained very simply, without any optimization. Other results on moment problems also highlight that that the bounds are very informative in practice, even if the panel is very short.

Next, we consider in Section 3 an estimator based on the theoretical expressions of the sharp bounds. This involves in particular the nonparametric estimation of a vector-valued function which, after a suitable transformation, corresponds at each point to a vector of raw moments $(1, E(U), ..., E(U^T))$ for some random variable $U \in [0, 1]$. One difficulty is that standard (e.g., local polynomial) nonparametric estimators, once transformed, may not be vectors of raw moments themselves. We show how to modify any initial estimator so as to satisfy this constraint. We then establish root-$n$ consistency of the estimators of the bounds under regularity conditions. The estimators of the bounds are asymptotically normal except if the corresponding slope coefficient is zero. We build confidence intervals of the true AME that are valid whether this is the case or not.

The previous estimator has the drawback of relying on a nonparametric first-step estimator. We then suggest in Section 4 an even simpler approach that avoids this issue. The idea is to estimate a simple approximation of the true AME and then make bias-aware inference, following the ideas of Donoho (1994) and Armstrong and Kolesár (2018). We can do so because the structure of the model allows us to consistently estimate an upper bound on the distance between the simple approximation and the true AME. The corresponding confidence intervals are asymptotically valid under mild conditions and, if slightly enlarged, even control the asymptotic size uniformly over a large set of data generating processes.

Section 5 shows that the same identification and estimation analysis can be applied to other parameters and models. Specifically, we study the ATE, the average structural function and FE ordered logit models. We also show that our method also applies when the number of observations varies per individual. Thus, our method easily accommodates unbalanced panels and hierarchical data. The R

package `MarginalFElogit`, developed with Christophe Gaillac and available `here`, estimates the AME and ATE (depending on whether $X$ is continuous or binary) in this model, and accommodates most of the aforementioned extensions.

Next, we study in Section 6 the finite sample properties of our two estimation and inference methods. In line with the theory, they show that the estimated bounds are very informative in practice. Also, the two confidence intervals have coverage close to their nominal level already for moderate sample sizes. Interestingly, we also find that the second inference method leads to confidence intervals often of the same size as, and sometimes even shorter than those obtained with the first method. This may seem surprising because as the sample size grows, such confidence intervals tend to an interval that strictly includes the true identified set. But for typical sample sizes and number of periods, it turns out that the distance between the simple approximation and the true AME is very small, leading to a tiny bias correction.

Our work is related to the literature on the identification of average marginal and treatment effects in panel data. Similar semiparametric approach looking at similar parameters in dynamic discrete choice models can be found in Honoré and Tamer (2006) and Aguirregabiria and Carro (2020). The latter point identifies a class of average marginal effects in dynamic FE logit models with four or more periods, exploiting the dynamic structure of the model. Another approach consists in using correlated random effects, following Mundlak (1978) and Chamberlain (1982); see e.g., for recent contributions Wooldridge (2019) for nonlinear models and Liu et al. (2021) for semiparametric binary response models. Compared to this approach, we do not impose any restriction on individual effects, which implies that average effects are only partially identified, though the bounds appear to be very informative in practice. An important part of this literature has also studied nonparametric identification (see in particular Altonji and Matzkin, 2005; Hoderlein and White, 2012; Chernozhukov et al., 2013, 2015). Our aim is different, as we consider a more constrained model, with the aim of providing a simple characterization and estimation of the bounds in this set-up.

Finally, our work is also related to moment problems, which have been studied extensively since Chebyshev and Markov. We refer to Karlin and Shapley (1953) and Krein and Nudelman (1977) for mathematical expositions and to Dette and Studden

(1997) for applications to various statistical problems. D'Haultfœuille and Rathelot (2017) use similar results on moment problems as here to obtain bounds on segregation measures with small units. Finally, Dobronyi et al. (2021) use other results on moment problems to characterize the identified set of common parameters in dynamic logit models, generalizing the work of Honoré and Weidner (2020).

## 2 Identification

We have a panel with $T$ periods and observe binary outcomes $Y_1, ..., Y_T$ and for each period $t$, a vector of covariates $X_t := (X_{t1}, ..., X_{tp})'$. We let $Y := (Y_1, ..., Y_T)'$, $X := (X_1', ..., X_T')'$ and make the following assumption.

**Assumption 1** *We have $Y_t = \mathbb{1}\{X_t'\beta_0 + \alpha + \varepsilon_t \geq 0\}$, where the $(\varepsilon_t)_{t=1,...,T}$ are i.i.d., independent of $(\alpha, X)$ and follow a logistic distribution.*

Importantly, the individual effect $\alpha$ is allowed to be correlated in an unspecified way with $X$. In this model, $S := \sum_{t=1}^T Y_t$ is a sufficient statistic for $\alpha$. As a result, identification of $\beta_0$ can be achieved by maximizing the theoretical conditional log-likelihood. For any $y = (y_1, ...y_T) \in \{0, 1\}^T$, let us define

$$C_k(x, \beta) := \sum_{(d_1, ..., d_T) \in \{0,1\}^T : \sum_{t=1}^T d_t = k} \exp\left(\sum_{t=1}^T d_t x_t'\beta\right),$$

$$\ell_c(y|x; \beta) := \sum_{t=1}^T y_t x_t'\beta - \ln\left[C_{\sum_{t=1}^T y_t}(x, \beta)\right].$$

To ensure that $\beta_0$ is identified as the unique maximizer of the theoretical conditional log-likelihood, we impose the following condition.

**Assumption 2** $E[\sum_{t,t'}(X_t - X_{t'})(X_t - X_{t'})']$ *is nonsingular.*

Assumption 2 is equivalent to the non-singularity of $E\left(\sum_{t=1}^T (X_t - \overline{X})(X_t - \overline{X})'\right)$ (with $\overline{X} = \sum_{t=1}^T X_t/T$) or $E\left(\sum_{t=2}^T (X_t - X_{t-1})(X_t - X_{t-1})'\right)$, which are necessary and sufficient conditions for the identification of the slope parameter in fixed effects linear models. The following proposition shows that this is also the case in FE logit models. It must be well-known, but we have not been able to find it in the literature. Its proof, as the other proofs of identification results, is presented in Appendix B.

**Proposition 1** *Suppose that Assumption 1 holds and for all $t \neq t'$ and $k \in \{1, ..., p\}$, $E[(X_{tk} - X_{t'k})^2] < \infty$. Then $\beta_0$ is identified if and only if Assumption 2 holds. In this case, $\beta_0 = \arg\max_\beta E\left(\ell_c(Y|X, \beta)\right)$ and $\mathcal{I}_0 = -E\left(\partial^2 \ell_c / \partial \beta \partial \beta'(Y_1, ..., Y_T|X; \beta_0)\right)$ is nonsingular.*

The second part of Proposition 1 shows that $\beta_0$ can be identified as the unique maximizer of the average conditional log-likelihood. Under mild regularity conditions, $\mathcal{I}_0^{-1}$ is the asymptotic variance of the conditional maximum likelihood estimator but also the semiparametric efficiency bound for $\beta_0$ (Hahn, 1997).

We now turn to average marginal effects. Without loss of generality, we focus hereafter on the last period. We first consider the case where $X_{Tk}$ is continuous; the case of a binary variable is deferred to Section 5.1. The average marginal effects is defined by

$$\Delta := E\left[\frac{\partial P(Y_T = 1|X, \alpha)}{\partial X_{Tk}}\right].$$

By Assumption 1, we have $P(Y_T = 1|X, \alpha) = \Lambda(X_T'\beta_0 + \alpha)$ with $\Lambda(x) := 1/(1 + \exp(-x))$. Therefore,

$$\Delta = \beta_{0k} E[\Lambda'(X_T'\beta_0 + \alpha)].$$

The identification of $\Delta$ is rendered difficult by the fact that $\alpha$ is unobserved and Assumption 1 imposes no restriction on $F_{\alpha|X}$, the cumulative distribution function (cdf) of $\alpha$ given $X$. The first point to note is that we can focus on $\Delta(x)$, defined by

$$\begin{aligned}\Delta(x) :=&\beta_{0k} E[\Lambda'(x_T'\beta_0 + \alpha)|X = x]\\=&\beta_{0k} \int \Lambda'(x_T'\beta_0 + a)dF_{\alpha|X}(a|x).\end{aligned} \tag{1}$$

As we shall see, the integral is only partially identified in general. Let $\underline{\Delta}(x)$ and $\overline{\Delta}(x)$ denote the sharp lower and upper bounds on $\Delta(x)$. In the absence of restrictions between the $(F_{\alpha|X}(\cdot|x))_{x \in \text{Supp}(X)}$, $\underline{\Delta} = E[\underline{\Delta}(X)]$ and $\overline{\Delta} = E[\overline{\Delta}(X)]$ are also the sharp lower and upper bounds on $\Delta$. Now, to determine $\underline{\Delta}(x)$ and $\overline{\Delta}(x)$, we seek the distributions of $\alpha|X = x$ minimizing and maximizing the average marginal effect, under the constraint that they are compatible with the distribution of $Y|X$. We exhibit the solutions of this problem below and show that their computation is very simple, leading also to simple estimators.

First, let us remark that because $S$ is a sufficient statistic for $\alpha$, its (conditional) distribution exhausts all the information available on $\alpha$. We have, for $k \in \{0, ..., T\}$,

$$P(S = k|X = x) = C_k(x, \beta_0) \int \frac{\exp(ka)}{\prod_{t=1}^{T}[1 + \exp(x_t'\beta_0 + a)]} dF_{\alpha|X}(a|x) \qquad (2)$$

Thus, sharp bounds on $\Delta(x)$ are obtained by finding the minimum and maximum value of the integral in (1), under the $T + 1$ constraints given by (2). One constraint is actually redundant since $\sum_{k=0}^{T} P(S = k|X = x) = 1$. Lemma 1 below shows that we can rewrite the optimization problem in a more convenient way. Before presenting it, we introduce additional notation. Let us define

$$\sum_{t=0}^{T+1} \lambda_t(x, \beta_0)u^t := u(1 - u) \prod_{t=1}^{T-1} (u(\exp((x_t - x_T)'\beta_0) - 1) + 1), \qquad (3)$$

$$c_t(x) := E\left[ \frac{\mathbb{1}\{S \geq t\} \binom{T-t}{S-t} \exp(Sx_T'\beta_0)}{C_S(x, \beta_0)} \bigg| X = x \right], \ t \in \{0, ..., T\} \qquad (4)$$

$$m_t(x) := \frac{c_t(x)}{c_0(x)}, \ t \in \{0, ..., T\}.$$

We then let $m(x) := (m_0(x), ..., m_T(x))'$. For any $m \in [0, 1]^{T+1}$ we denote by $\mathcal{D}(m)$ the set of positive measures $\mu$ on $[0, 1]$ whose vector of first $T + 1$ raw moments (starting from $\int_0^1 u^0 d\mu(u) = \mu([0, 1])$) is equal to $m$. We finally define

$$\underline{q}_T(m) := \inf_{\mu \in \mathcal{D}(m)} \int_0^1 u^{T+1} d\mu(u), \ \overline{q}_T(m) := \sup_{\mu \in \mathcal{D}(m)} \int_0^1 u^{T+1} d\mu(u). \qquad (5)$$

**Lemma 1** *Suppose that Assumptions 1-2 hold. Then, there exists $\mu \in \mathcal{D}(m(x))$ such that*

$$\Delta(x) = \beta_{0k} c_0(x) \int_0^1 \sum_{t=0}^{T+1} \lambda_t(x, \beta_0) u^t d\mu(u).$$

*Thus, the sharp identified set of $\Delta(x)$ is $[\underline{\Delta}(x), \overline{\Delta}(x)]$, with[2]*

$$\begin{cases} \underline{\Delta}(x) &= \beta_{0k} \left[ \sum_{t=1}^{T} \lambda_t(x, \beta_0) c_t(x) + c_0(x)\lambda_{T+1}(x, \beta_0)\underline{q}_T(m(x)) \right] \\ \overline{\Delta}(x) &= \beta_{0k} \left[ \sum_{t=1}^{T} \lambda_t(x, \beta_0) c_t(x) + c_0(x)\lambda_{T+1}(x, \beta_0)\overline{q}_T(m(x)) \right] \end{cases} \qquad (6)$$

---

[2]Technically, we define here the sharp identified set as the closure of all parameters that can rationalized by the data and the model. In some cases, the bounds $\underline{\Delta}(x)$ and $\overline{\Delta}(x)$ do not correspond to a valid probability distribution on $\alpha|X = x$, as they amount to put mass at plus or minus infinity. Nevertheless, these bounds can be approached arbitrarily well by sequences of appropriate probability distributions.

if $\beta_{0k}\lambda_{T+1}(x, \beta_0) \geq 0$. If $\beta_{0k}\lambda_{T+1}(x, \beta_0) < 0$, the same holds with $\overline{q}_T$ and $\underline{q}_T$ switched in the two bounds.

This lemma follows essentially by a change of measures in (1) and (2), and taking an appropriate linear transform of the constraints. With this result, the issue of computing $\underline{\Delta}(x)$ and $\overline{\Delta}(x)$ reduces to the computation of $\underline{q}_T(m(x))$ and $\overline{q}_T(m(x))$ defined by (5). This may seem difficult since the programs are infinite-dimensional. It turns out, however, that they can be computed easily and without any numerical optimization, using results on moment problems (see Karlin and Shapley, 1953; Krein and Nudelman, 1977; Dette and Studden, 1997). Specifically, let $\mathcal{D}$ denote the set of probability measures on $[0, 1]$ and for any $t \geq 1$, let

$$\mathcal{M}_t := \left\{ \left( \int_0^1 u^0 d\mu(u), \int_0^1 u^1 d\mu(u), ..., \int_0^1 u^t d\mu(u) \right)' : \mu \in \mathcal{D} \right\}.$$

The set $\mathcal{M}_t$ is the set of all possible vectors of first $t + 1$ raw moments, starting from the moment of order 0, of probability measures on $[0, 1]$. Also, for any $m = (m_0, ..., m_t) \in \mathbb{R}^{t+1}$, we define $\underline{\mathbb{H}}_t(m)$ and $\overline{\mathbb{H}}_t(m)$ as

If $t$ is even, $\quad \underline{\mathbb{H}}_t(m) = (m_{i+j-2})_{1 \leq i,j \leq t/2+1}, \quad \overline{\mathbb{H}}_t(m) = (m_{i+j-1} - m_{i+j})_{1 \leq i,j \leq t/2}.$

If $t$ is odd, $\quad \underline{\mathbb{H}}_t(m) = (m_{i+j-1})_{1 \leq i,j \leq (t+1)/2}, \quad \overline{\mathbb{H}}_t(m) = (m_{i+j-2} - m_{i+j-1})_{1 \leq i,j \leq (t+1)/2}.$

Then, we define $\underline{H}_t(m) = \det\left(\underline{\mathbb{H}}_t(m)\right)$ and $\overline{H}_t(m) = \det(\overline{\mathbb{H}}_t(m))$.

**Proposition 2** *For any $m = (m_0, ..., m_T) \in \mathcal{M}_T$, $\underline{H}_T(m) \times \overline{H}_T(m) \geq 0$ and:*

1. *If $\underline{H}_T(m) \times \overline{H}_T(m) > 0$, then $\underline{q}_T(m) < \overline{q}_T(m)$. Moreover, $q \mapsto \underline{H}_{T+1}(m, q)$ is strictly increasing, linear and $\underline{H}_{T+1}(m, \underline{q}_T(m)) = 0$. Similarly, $q \mapsto \overline{H}_{T+1}(m, q)$ is strictly decreasing, linear and $\overline{H}_{T+1}(m, \overline{q}_T(m)) = 0$.*

2. *If $\underline{H}_T(m) \times \overline{H}_T(m) = 0$, then $\underline{q}_T(m) = \overline{q}_T(m)$. Moreover, letting $T' = \min\{t \leq T : \underline{H}_t(m) \times \overline{H}_t(m) = 0\}$, $\underline{q}_T(m) = \overline{q}_T(m)$ is the unique solution of*

$$\begin{aligned} \underline{H}_{T'}(m_{T-T'+1}, ..., m_T, \underline{q}_T(m)) = 0 \quad &\text{if } \underline{H}_{T'}(m) = 0, \\ \overline{H}_{T'}(m_{T-T'+1}, ..., m_T, \underline{q}_T(m)) = 0 \quad &\text{if } \overline{H}_{T'}(m) = 0. \end{aligned}$$

The first point follows by classical results in moment theory, see e.g. Theorems 1.2.7 and 1.4.3 in Dette and Studden (1997). The first part of the second point is also well-known. On the other hand, to the best of our knowledge, the second part is new.

The first case $(\underline{H}_T(m) \times \overline{H}_T(m) > 0)$ corresponds to the situation where $m \in \text{Int } \mathcal{M}_T$, the interior of $\mathcal{M}_T$. Then, there are infinitely many distributions in $\mathcal{D}$ that rationalize $m$, and $\underline{q}_T(m) < \overline{q}_T(m)$. Moreover, the two bounds can be simply obtained by solving the one-dimensional linear equation $\underline{H}_{T+1}(m_0, ..., m_T, \underline{q}_T(m)) = 0$ and $\overline{H}_{T+1}(m_0, ..., m_T, \overline{q}_T(m)) = 0$. Though we do not pursue this here, we can also characterize the (unique) probability distributions $\mu$ corresponding to $\underline{q}_T(m)$ and $\overline{q}_T(m)$. These distributions have few points of support (around $T/2$) and are called principal representations. We refer to, e.g., Chapter 3 in Krein and Nudelman (1977) for more details on these distributions.

In the second case $(\underline{H}_T(m) \times \overline{H}_T(m) = 0)$, $m \in \partial \mathcal{M}_T$, the boundary of $\mathcal{M}_T$. Then, there is a single distribution in $\mathcal{D}$ rationalizing $m$, and $\underline{q}_T(m) = \overline{q}_T(m)$. We obtain this value by solving a similar linear equation as in the first case.

Using Lemma 1 and results related to Proposition 2, we obtain the following further properties of the identified set:

**Proposition 3** *Suppose that Assumptions 1-2 hold. Then:*

1. *The length of the identified set satisfies:*

$$\overline{\Delta}(x) - \underline{\Delta}(x) \leq \frac{|\beta_{0k}|}{2^{T+1}} \prod_{t=1}^{T-1} \frac{|\exp(x_t'\beta_0) - \exp(x_T'\beta_0)|}{\exp(x_t'\beta_0) + \exp(x_T'\beta_0)}$$
$$\leq \frac{|\beta_{0k}|}{2^{T+1}}.$$

2. *$\Delta(x)$ is point identified if and only if (i) $\beta_{0k} = 0$ or (ii) there exists $t < T$ such that $x_t'\beta_0 = x_T'\beta_0$ or (iii) $|Supp(\alpha|X = x)| \leq T/2$.*

The first point exploits in particular a result of the theory of moments, namely that $\overline{q}_T(m) - \underline{q}_T(m) \leq 1/4^T$ for any $m \in \mathcal{M}_T$ (see Karlin and Shapley, 1953). With the sole knowledge of $\beta_0$ (and assuming for instance that $\beta_{0k} > 0$), it follows from (1) and $\Lambda'(u) \in (0, 1/4)$ for all $u$ that the identified set of the AME is $[0, \beta_{0k}/4]$. By exploiting

9

the model and the data, we can therefore shrink this set by at least a factor $2^{T-1}$. A similar result on the length of the identified set of the average structural function $x \mapsto E[\Lambda(x'\beta_0 + \alpha)]$ was obtained by Chernozhukov et al. (2013), see their Theorem 4. Actually, their result imposed substantially weaker conditions on the distribution of $Y_t|X_t, \alpha$. On the other hand, their result only holds for discrete $X$ and imposes additional restriction on the distribution of $X$.[3] Also, the exponent of the rate could be arbitrarily close to 1 in their set-up. Proposition 3 shows that in the fixed effects logit model, such an exponential bound holds irrespective of the distribution of $X$, and with an exponent of at least 2. For specific distributions of $X$, the exponent could actually be larger than 2, if, roughly speaking, $(X_t - X_T)'\beta_0$ is small with a large enough probability. This can be often expected, as many variables exhibit strong persistence over time.

The second result of Proposition 3 characterizes the point identification of $\Delta(x)$. The cases $\beta_{0k} = 0$ and $x_t'\beta_0 = x_T'\beta_0$ for some $t < T$ can be directly deduced from the first result. That $\Delta(x)$ is point identified if $(x_t - x_T)'\beta_0 = 0$ for some $t < T$ could be expected. Indeed, Hoderlein and White (2012) show that average marginal effects are nonparametrically identified on "stayers", namely individuals for whom $X_{it}$ remains constant between two periods. The third possibility for point identification $(|\text{Supp}(\alpha|X = x)| \leq T/2)$ corresponds to the second case described in Proposition 2. Intuitively, if $\alpha|X = x$ has few points of supports, its full distribution is characterized by its first moments. Then, given these moments, the higher moments are fully determined. As an illustration, suppose $T = 2$ and $\alpha|X = x$ is degenerate and equal to $\alpha_0$. Then, some algebra shows that $m(x) = (1, \Lambda(\alpha_0 + x_T'\beta_0), \Lambda(\alpha_0 + x_T'\beta_0)^2)'$. In such a case, the variance of any distribution in $\mathcal{D}(m(x))$ is zero. Thus, $\mathcal{D}(m(x))$ is reduced to the Dirac at $\Lambda(\alpha_0 + x_T'\beta_0)$, which implies $\underline{q}_2(m) = \overline{q}_2(m) = \Lambda(\alpha_0 + x_T'\beta_0)^3$. Hence, $\Delta(x)$ is point identified in this case.

---

[3]Their parameter is also different. However, our analysis also applies to the average structural function (see Section 5.2 below), so we can obtain the same decrease in the length of the identified set for this parameter.

# 3    A first estimation and inference method

In this section, we estimate the sharp bounds on $\Delta$ and develop inference on this parameter based on these bounds, using a sample $(Y_i, X_i)_{i=1,...,n}$.

## 3.1    Definition of the estimators

First, let us define

$$U(x, s, \beta) := \beta_k \sum_{t=0}^{s} \binom{T-t}{s-t} \frac{\lambda_t(x; \beta_0) \exp(sx'_T \beta)}{C_s(x, \beta)}.$$

By Equation (6) and the law of iterated expectations, we have

$$\overline{\Delta} = E\left[U(X, S, \beta_0)\right] + \beta_{0k} E\Big[c_0(X)\lambda_{T+1}(X; \beta_0)\big(\overline{q}_T(m(X))\mathbb{1}\left\{\beta_{0k}\lambda_{T+1}(X, \beta_0) \geq 0\right\}$$
$$+ \underline{q}_T(m(X))\mathbb{1}\left\{\beta_{0k}\lambda_{T+1}(X, \beta_0) < 0\right\}\big)\Big],$$
$$\overline{\Delta} = E\left[U(X, S, \beta_0)\right] + \beta_{0k} E\Big[c_0(X)\lambda_{T+1}(X; \beta_0)\big(\underline{q}_T(m(X))\mathbb{1}\left\{\beta_{0k}\lambda_{T+1}(X, \beta_0) \geq 0\right\}$$
$$+ \overline{q}_T(m(X))\mathbb{1}\left\{\beta_{0k}\lambda_{T+1}(X, \beta_0) < 0\right\}\big)\Big].$$

We then estimate these bounds in three steps:

1. Estimation of $\beta_0$ by the conditional maximum likelihood estimator $\widehat{\beta}$.

2. Estimation of the functions $c_0, ..., c_T$ and $m$:

   (a) Nonparametric estimation of $c_0, ..., c_T$;

   (b) Nonparametric estimation of $m$ by an estimator $\widehat{m}$ satisfying $\widehat{m}(X_i) \in \mathcal{M}_T$ for all $i$.

3. Estimation of the bounds by a plug-in estimator based on the formulas above.

Step 1 is straightforward. We now explain in details Steps 2a and 2b, before giving the formulas corresponding to Step 3.

### 3.1.1 Estimation of $(c_0, ..., c_T)$

Let $\gamma_{0j}(x) = P(S = j | X = x)$ for $j = 0, ..., T$. The functions $(c_t)_{t=0...T}$ and $(\gamma_{0j})_{j=0...T}$ are related through

$$(c_0(x), ..., c_T(x))' = \Gamma \left( \frac{\gamma_{00}(x) \exp(0 \times x_T' \beta_0)}{C_0(x, \beta_0)}, \ ... \ , \frac{\gamma_{0T}(x) \exp(T \times x_T' \beta_0)}{C_T(x, \beta_0)} \right)', \quad (7)$$

where $\Gamma$ is a square matrix of size $T + 1$ with coefficients $\Gamma_{ij} = \binom{T-i}{j-i} \mathbb{1}\{i \leq j\}$ for $i, j = 1, ..., T + 1$. We first estimate $\gamma_0 := (\gamma_{00}, ..., \gamma_{0T})$ nonparametrically. We use local polynomial estimators of order $\ell$ to avoid boundary effects. Let $K$ denote a kernel function and for a given $0 \leq j \leq T$, define

$$\widehat{a}^j(x) := \operatorname{argmin}_a \sum_{i=1}^n K \left( \frac{X_i - x}{h_n} \right) \left( \mathbb{1}\{S_i = j\} - \sum_{|b| \leq \ell} a_b (X_i - x)^b \right)^2, \quad (8)$$

where, in this definition, $b \in \mathbb{N}^{pT}$, $|b| = \sum_{j=1}^{pT} b_j$ and $x^b = x_1^{b_1}...x_{pT}^{b_{pT}}$. The estimator of $\gamma_{0j}(x)$ is then $\widehat{\gamma}_j(x) = \widehat{a}_0^j(x)$. Our estimator for $c_t(x)$, $\widehat{c}_t(x)$, uses (7), replacing $\gamma_0$ and $\beta_0$ with their estimators.

### 3.1.2 Estimation of $m$

Given its definition, a natural estimator of $m$ is

$$\widetilde{m}(x) = \left( 1, \ \frac{\widehat{c}_1(x)}{\widehat{c}_0(x)}, \ ..., \ \frac{\widehat{c}_T(x)}{\widehat{c}_0(x)} \right).$$

However, this estimator may not satisfy $\widetilde{m}(x) \in \mathcal{M}_T$. This is especially the case if $m(x)$ is at the boundary of $\mathcal{M}_T$, or for a "large" $T$, because the volume of $\mathcal{M}_T$ decreases very quickly with $T$ (Karlin and Shapley, 1953). In our simulations below, this already occurs with $T = 3$ and $n = 1,000$, even if $m(x)$ is in the interior of $\mathcal{M}_T$. That $\widetilde{m}(x) \notin \mathcal{M}_T$ is an issue because then, $\underline{q}_T(\widetilde{m}(x))$ and $\overline{q}_T(\widetilde{m}(x))$ are undefined. We thus consider another estimator $\widehat{m}$ such that $\widehat{m}(x) \in \mathcal{M}_T$.

To this end, we rely on Proposition 2. For any $(m_t)_{t \geq 0}$ and $t \in \{0, ..., T\}$, let $m_{\to t} = (m_0, ..., m_t)$. The idea of the estimator is to use the first elements of $\widetilde{m}(x)$, until $\widetilde{m}_t(x)$ falls too close to $\underline{q}_{t-1}(\widetilde{m}_{\to t-1}(x))$ or $\overline{q}_{t-1}(\widetilde{m}_{\to t-1}(x))$. In such a case, we simply replace $\widetilde{m}_t(x)$ by $\underline{q}_{t-1}(\widetilde{m}_{\to t-1}(x))$ or $\overline{q}_{t-1}(\widetilde{m}_{\to t-1}(x))$. We finally complete the vector using the second part of Proposition 2.

Specifically, let $c_n$ be a sequence tending to 0 at a rate specified later and define

$$\widehat{I}(x) := \max\left\{t \in \{1, ..., T\} : \underline{H}_t(\widetilde{m}_{\to t}(x)) \times \overline{H}_t(\widetilde{m}_{\to t}(x)) > c_n\right\}.$$

with the convention that $\max \emptyset = 0$. We then let

$$\widehat{m}_{\to \widehat{I}(x)}(x) := \widetilde{m}_{\to \widehat{I}(x)}(x).$$

If $\widehat{I}(x) = T$, $\widehat{m}(x)$ is fully defined. Otherwise, we complete $\widehat{m}(x)$ by first letting

$$\widehat{m}_{\widehat{I}(x)+1}(x) := \begin{vmatrix} \underline{q}_{\widehat{I}(x)}(\widetilde{m}_{\to \widehat{I}(x)}(x)) & \text{if } \underline{H}_{\widehat{I}(x)+1}(\widetilde{m}_{\to \widehat{I}(x)+1}(x)) < c_n^{1/2}, \\ \overline{q}_{\widehat{I}(x)}(\widetilde{m}_{\to \widehat{I}(x)}(x)) & \text{otherwise.} \end{vmatrix}$$

Next, if $\widehat{I}(x) + 1 < T$, by construction, we have

$$\underline{H}_{\widehat{I}(x)+1}(\widehat{m}_{\to \widehat{I}(x)+1}(x)) \times \overline{H}_{\widehat{I}(x)+1}(\widehat{m}_{\to \widehat{I}(x)+1}(x)) = 0.$$

Then, applying Part 2 of Proposition 2, we construct by induction the unique possible moments $\widehat{m}_{\widehat{I}(x)+2}, ..., \widehat{m}_T$ that are compatible with $\widehat{m}_{\to \widehat{I}(x)+1}(x)$. By construction, the corresponding vector $\widehat{m}(x)$ belongs to $\mathcal{M}_T$.

### 3.1.3  Estimation of the bounds

Given Steps 1 and 2 above, the estimators of the bounds of the AME are defined by

$$\widehat{\overline{\Delta}} = \frac{1}{n}\sum_{i=1}^{n} U(X_i, S_i, \widehat{\beta}) + \widehat{\beta}_k \widehat{c}_0(X_i)\lambda_{T+1}(X_i, \widehat{\beta})\left[\overline{q}_T(\widehat{m}(X_i))\mathbb{1}\left\{\widehat{\beta}_k\lambda_{T+1}(X_i, \widehat{\beta}) \geq 0\right\}\right.$$

$$\left. + \underline{q}_T(\widehat{m}(X_i))\mathbb{1}\left\{\widehat{\beta}_k\lambda_{T+1}(X_i, \widehat{\beta}) < 0\right\}\right],$$

$$\widehat{\underline{\Delta}} = \frac{1}{n}\sum_{i=1}^{n} U(X_i, S_i, \widehat{\beta}) + \widehat{\beta}_k \widehat{c}_0(X_i)\lambda_{T+1}(X_i, \widehat{\beta})\left[\underline{q}_T(\widehat{m}(X_i))\mathbb{1}\left\{\widehat{\beta}_k\lambda_{T+1}(X_i, \widehat{\beta}) \geq 0\right\}\right.$$

$$\left. + \overline{q}_T(\widehat{m}(X_i))\mathbb{1}\left\{\widehat{\beta}_k\lambda_{T+1}(X_i, \widehat{\beta}) < 0\right\}\right]. \tag{9}$$

## 3.2  Asymptotic properties

### 3.2.1  Consistency

We first establish consistency of the estimated bounds under the following conditions.

## Assumption 3

1. *The variables $(X_i, \alpha_i, \varepsilon_{i1}, ..., \varepsilon_{iT})$ are i.i.d across i.*

2. *$Supp(X)$ is a compact set and $\beta_0 \in \Theta$, where $\Theta$ is a compact set.*

## Assumption 4

1. *$X$ admits a density $f_X$ with respect to the Lebesgue measure on $\mathbb{R}^{pT}$. $f_X$ is $C^1$ and bounded away from $0$ on $Supp(X)$,*

2. *$\gamma_0$ is $C^{\ell+2}$ on $Supp(X)$,*

3. *$K$ is a Lipschitz density on $\mathbb{R}^{pT}$ with compact support including a neighborhood of 0,*

4. *$h_n \to 0$ and $nh_n^{pT}/\ln n \to \infty$ as $n \to \infty$.*

Assumption 3 is sufficient for $\widehat{\beta}$ to be consistent. Also, Assumptions 3 and 4 guarantee that $\widehat{\gamma}$ converges uniformly to $\gamma_0$ over the support of $X$ at a rate at least $\delta_n$, with $\delta_n := (\ln n/(nh_n^{pT}))^{1/2} + h_n^{\ell+1}$. Note that Assumption 4.2 is in fact a smoothness condition on the distribution of $\alpha$ given $X$. For instance, if this distribution is discrete and both support points and weighting probabilities are $C^{\ell+2}$ as functions of $X$ on $Supp(X)$, then Assumption 4.2 holds.

**Theorem 1** *Suppose that Assumptions 1-4 hold and $\delta_n/c_n \to 0$. Then*

$$(\widehat{\underline{\Delta}}, \widehat{\overline{\Delta}}) \xrightarrow{P} (\underline{\Delta}, \overline{\Delta}).$$

The proofs of this theorem and other asymptotic results are given in the Online Appendix. The key step therein is to show that $\widehat{m}$ is uniformly consistent, which is not straightforward because $\widehat{m}$ is a complicated function of $\widehat{\gamma}$.

### 3.2.2 Root-n consistency and inference

We now establish the asymptotic distribution of $(\widehat{\underline{\Delta}}, \widehat{\overline{\Delta}})$. To this end, we impose additional regularity conditions. For any $d \geq 1, u \in \mathbb{R}^d$ and $\epsilon > 0$, we let $\mathcal{B}(u, \epsilon)$ denote the closed ball centered at $u$ and with radius $\epsilon$.

**Assumption 5**

1. $\ell \geq pT/2$,

2. $nh_n^{2(\ell+1)} \to 0$ and $n[h_n^{pT}/\ln n]^3 \to \infty$ as $n \to \infty$,

**Assumption 6** *There exists $\epsilon > 0$ and $I \in \{1, ..., T\}$ such that for all $x \in Supp(X)$; (i) $\mathcal{B}(m_{\to I}(x), \epsilon) \subset Int \, \mathcal{M}_I$; (ii) if $I < T$, $m_{\to I+1}(x) \in \partial \mathcal{M}_{I+1}$.*

Assumption 5 is standard in semiparametric estimation and complements Assumption 4. Assumption 6, on the other hand, is specific to our context. The index $I$ therein relates to $|\text{Supp}(\alpha|X = x)|$. Specifically, one can prove that if $I < T$, then $I$ is odd and $|\text{Supp}(\alpha|X = x)| = (I+1)/2$ for all $x \in \text{Supp}(X)$. Also, if $I = T$, then $|\text{Supp}(\alpha|X = x)| > T/2$ for all $x \in \text{Supp}(X)$, in which case $|\text{Supp}(\alpha|X = x)|$ may vary with $x$. This means that Assumption 6 is violated when there exists $(x, x') \in \text{Supp}(X)^2$ such that

$$|\text{Supp}(\alpha|X = x)| \neq |\text{Supp}(\alpha|X = x')|,$$
$$\min\left(|\text{Supp}(\alpha|X = x)|, |\text{Supp}(\alpha|X = x')|\right) \leq T/2.$$

We impose this restriction because $\underline{q}_T$ and $\overline{q}_T$ are not regular everywhere for $T \geq 3$. Specifically, whereas these functions are continuous on $\mathcal{M}_T$ and infinitely differentiable on Int $\mathcal{M}_T$, they may not be even directionally differentiable at $m \in \partial \mathcal{M}_T$.[4]

Before presenting our asymptotic result, we introduce additional notation. First, for any vector of functions $\gamma = (\gamma_0, ..., \gamma_T)$, let

$$\left(c_0(\gamma, x, \beta), ..., c_T(\gamma, x, \beta)\right)' := \Gamma \left(\frac{\gamma_0(x)\exp(0 \times x'_T\beta)}{C_0(x, \beta)}, \, ... \, , \frac{\gamma_T(x)\exp(T \times x'_T\beta)}{C_T(x, \beta)}\right)'.$$

---

[4]See D'Haultfœuille and Rathelot (2017) for proofs of the first two statements. Regarding the third, one can show that, e.g., $m_1 \mapsto \underline{q}_3(m_0, m_1, m_2, m_3)$ is not differentiable at $m = (1, m_1, m_1^2, m_1^3)$.

Note that $\widehat{c}_t(x) = c_t(\widehat{\gamma}, x, \widehat{\beta})$. Then, with $I$ defined in Assumption 6, let

$$m(\gamma, x, \beta) := \left(1, \frac{c_1(\gamma, x, \beta)}{c_0(\gamma, x, \beta)}, ..., \frac{c_I(\gamma, x, \beta)}{c_0(\gamma, x, \beta)}\right),$$

so that $m(\gamma_0, x, \beta_0) = m_{\rightarrow I}(x)$ and $m(\widehat{\gamma}, x, \widehat{\beta}) = \widetilde{m}_{\rightarrow I}(x)$. Now, if $I = T$, we let, with a slight abuse of notation, $\underline{q}_T(\gamma, x, \beta) = \underline{q}_T(m(\gamma, x, \beta))$. If $I < T$, by Assumption 6 and Proposition 2, $m_{I+1}(x) = \underline{q}_I(m_{\rightarrow I}(x))$ or $m_{I+1}(x) = \overline{q}_I(m_{\rightarrow I}(x))$. Then, by Proposition 2 again and a straightforward induction, we can define $m_t(x)$ for $t \in \{I+1, ..., T\}$ as a function of $m_{\rightarrow I}(x)$. We let $E(.)$ denote the corresponding extension function. Then $m(x) = E(m_{\rightarrow I}(x))$. Finally, we let

$$\underline{q}_T(\gamma, x, \beta) := \underline{q}_T(E(m(\gamma, x, \beta))).$$

We define similarly $\overline{q}_T(\gamma, x, \beta)$. Note that $\underline{q}_T(\cdot, \cdot, \cdot)$ and $\overline{q}_T(\cdot, \cdot, \cdot)$ depend on the unknown $I$, and when $I < T$, on the true function $m$, since the definition of $E$ involves this true function. However, we show in the proof of Theorem 2 below that with probability approaching one, $\underline{q}_T(\widehat{m}(x)) = \underline{q}_T(\widehat{\gamma}, x, \widehat{\beta})$.

Then, we also define

$$\underline{h}(x, s, \gamma, \beta) = U(x, s, \beta) + \beta_k c_0(\gamma, x, \beta)\lambda_{T+1}(x, \beta)\Big[\underline{q}_T(\gamma, x, \beta)\mathbb{1}\{\lambda_{T+1}(x, \beta_0) > 0\}$$
$$+ \overline{q}_T(\gamma, x, \beta)\mathbb{1}\{\lambda_{T+1}(x, \beta_0) < 0\}\Big],$$

$$\overline{h}(x, s, \gamma, \beta) = U(x, s, \beta) + \beta_k c_0(\gamma, x, \beta)\lambda_{T+1}(x, \beta)\Big[\overline{q}_T(\gamma, x, \beta)\mathbb{1}\{\lambda_{T+1}(x, \beta_0) > 0\}$$
$$+ \underline{q}_T(\gamma, x, \beta)\mathbb{1}\{\lambda_{T+1}(x, \beta_0) < 0\}\Big].$$

Note that $\underline{h}(x, s, \gamma, \beta)$ (and similarly $\overline{h}(x, s, \gamma, \beta)$) depends on $\gamma$ only through $\gamma(x)$. Also, $\underline{h}$ is differentiable with respect to $\beta$ and the vector $\gamma(x)$. We denote its corresponding partial derivatives as $D_\beta\underline{h}(x, s, \gamma, \beta)$ and $D_\gamma\underline{h}(x, s, \gamma, \beta)$.

As Theorem 2 below shows, the influence functions of $\widehat{\underline{\Delta}}$ and $\widehat{\overline{\Delta}}$ are:

$$\underline{\psi}_i = \underline{h}(X_i, S_i, \gamma_0, \beta_0) - E[\underline{h}(X, S, \gamma_0, \beta_0)] + E\left[D_\beta\underline{h}(X, S, \gamma_0, \beta_0)\right]' \phi_i$$
$$+ D_\gamma\underline{h}(X_i, S_i, \gamma_0, \beta_0)'[Z_i - \gamma_0(X_i)],$$

$$\overline{\psi}_i = \overline{h}(X_i, S_i, \gamma_0, \beta_0) - E[\overline{h}(X, S, \gamma_0, \beta_0)] + E\left[D_\beta\overline{h}(X, S, \gamma_0, \beta_0)\right]' \phi_i$$
$$+ D_\gamma\overline{h}(X_i, S_i, \gamma_0, \beta_0)'[Z_i - \gamma_0(X_i)],$$

16

where $Z_i = (\mathbb{1}\{S_i = 0\}, ..., \mathbb{1}\{S_i = T\})'$ and $\phi_i = \mathcal{I}_0^{-1}\partial\ell_c/\partial\beta(Y_i|X_i;\beta_0)$ is the influence function of $\widehat{\beta}$. We let $\Sigma$ denote the variance-covariance matrix of $(\underline{\psi}, \overline{\psi})$. We introduce $\widehat{\phi}_i$ as the sample analog of $\phi_i$ and similarly, sample analogs of $\underline{\psi}_i$ and $\overline{\psi}_i$ are

$$\widehat{\underline{\psi}}_i = \underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) - \frac{1}{n}\sum_{j=1}^{n}\underline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta}) + \left(\frac{1}{n}\sum_{j=1}^{n}D_\beta\underline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta})\right)'\widehat{\phi}_i$$
$$+ D_\gamma\underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta})'[Z_i - \widehat{\gamma}(X_i)]$$

$$\widehat{\overline{\psi}}_i = \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) - \frac{1}{n}\sum_{j=1}^{n}\overline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta}) + \left(\frac{1}{n}\sum_{j=1}^{n}D_\beta\overline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta})\right)'\widehat{\phi}_i$$
$$+ D_\gamma\overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta})'[Z_i - \widehat{\gamma}(X_i)]$$

We finally estimate $\Sigma$ by $\widehat{\Sigma} = \frac{1}{n}\sum_{i=1}^{n}(\widehat{\underline{\psi}}_i, \widehat{\overline{\psi}}_i)'(\widehat{\underline{\psi}}_i, \widehat{\overline{\psi}}_i)$.

**Theorem 2** *Suppose that Assumptions 1-6 hold and $\delta_n/c_n \to 0$. Then:*

1. *If $\beta_{0k} > 0$,*

$$\sqrt{n}\left(\widehat{\underline{\Delta}} - \underline{\Delta}, \widehat{\overline{\Delta}} - \overline{\Delta}\right) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\left(\underline{\psi}_i, \overline{\psi}_i\right) + o_P(1) \xrightarrow{d} (\underline{Z}, \overline{Z}),$$

   *with $(\underline{Z}, \overline{Z}) \sim \mathcal{N}(0, \Sigma)$. If $\beta_{0k} < 0$, the same result holds by just exchanging the roles of $\underline{\psi}_i$ (resp. $\underline{Z}$) and $\overline{\psi}_i$ (resp. $\overline{Z}$).*

2. *If $\beta_{0k} = 0$,*

$$\sqrt{n}\left(\widehat{\underline{\Delta}} - \underline{\Delta}, \widehat{\overline{\Delta}} - \overline{\Delta}\right) \xrightarrow{d} \left(\min(\underline{Z}, \overline{Z}), \max(\underline{Z}, \overline{Z})\right).$$

3. *We have $\widehat{\Sigma} \xrightarrow{P} \Sigma$.*

The key step for proving this theorem is to show that, $\widehat{\underline{\Delta}}$ and $\widehat{\overline{\Delta}}$ are asymptotically equivalent to a standard two-step semiparametric estimator with a nonparametric first-step estimator (neglecting here their additional dependence on $\widehat{\beta}$). Assumption 6 is key for this purpose.

The estimated bounds are asymptotically normal if $\beta_{0k} \neq 0$, but not in general if $\beta_{0k} = 0$. An exception is when the whole vector $\beta_0$ is equal to 0. Then $\underline{\psi} = \overline{\psi}$, which implies that $\underline{Z} = \overline{Z}$. In this case $\sqrt{n}(\widehat{\overline{\Delta}} - \widehat{\underline{\Delta}}) = o_P(1)$, and both bounds are asymptotically normal.

With Theorem 2 at hand, we can construct confidence intervals on $\Delta$ that are asymptotically valid whether or not $\beta_{0k} = 0$, at least in a pointwise sense. To this end, let $\varphi_\alpha$ denote a consistent test with asymptotic level $\alpha$ of $\beta_{0k} = 0$, e.g. a $t$-test. Following Imbens and Manski (2004), let $c_\alpha$ denote the unique solution to

$$\Phi\left(c_\alpha + \frac{n^{1/2}\left(\widehat{\overline{\Delta}} - \widehat{\underline{\Delta}}\right)}{\max\left(\widehat{\Sigma}_{11}^{1/2}, \widehat{\Sigma}_{22}^{1/2}\right)}\right) - \Phi(-c_\alpha) = 1 - \alpha,$$

with $\Phi$ the cdf of a standard normal distribution and $\Sigma_{ij}$ is the $(i, j)$ term of $\Sigma$. Then, we define $\text{CI}_{1-\alpha}^1$ as

$$\text{CI}_{1-\alpha}^1 := \left|\begin{array}{ll} \left[\widehat{\underline{\Delta}} - c_\alpha(\widehat{\Sigma}_{11}/n)^{1/2},\ \widehat{\overline{\Delta}} + c_\alpha(\widehat{\Sigma}_{22}/n)^{1/2}\right] & \text{if } \varphi_\alpha = 1, \\ \left[\min\left(0, \widehat{\underline{\Delta}} - c_\alpha(\widehat{\Sigma}_{11}/n)^{1/2}\right),\ \max\left(0, \widehat{\overline{\Delta}} + c_\alpha(\widehat{\Sigma}_{22}/n)^{1/2}\right)\right] & \text{if } \varphi_\alpha = 0. \end{array}\right.$$

The following proposition shows that $\text{CI}_{1-\alpha}^1$ is pointwise valid as $n \to \infty$.

**Proposition 4** *Suppose that Assumptions 1-6 hold, $\delta_n/c_n \to 0$ and $\min(\Sigma_{11}, \Sigma_{22}) > 0$. Then $\liminf_n \inf_{\Delta \in [\underline{\Delta}, \overline{\Delta}]} P(\Delta \in CI_{1-\alpha}^1) \geq 1 - \alpha$, with equality when $\beta_{0k} \neq 0$.*

Intuitively, $\text{CI}_{1-\alpha}^1$ asymptotically reaches its nominal level when $\beta_{0k} \neq 0$ because it includes $[\widehat{\underline{\Delta}} - c_\alpha(\widehat{\Sigma}_{11}/n)^{1/2},\ \widehat{\overline{\Delta}} + c_\alpha(\widehat{\Sigma}_{22}/n)^{1/2}]$, and the latter interval has asymptotic coverage $1 - \alpha$, by Theorem 2. When $\beta_{0k} = 0$, the asymptotic coverage of $\text{CI}_{1-\alpha}^1$ is also at least $1 - \alpha$, because $\Delta = 0 \in \text{CI}_{1-\alpha}^1$ as soon as $\varphi_\alpha = 0$.

The interval $\text{CI}_{1-\alpha}^1$ may have a uniform coverage over an appropriate set of data generating processes (DGPs), even if $\beta_0$ varies over $\Theta$. Establishing this formally would however require to establish the uniform convergence in distribution of $(\widehat{\underline{\Delta}}, \widehat{\overline{\Delta}})$, a multistep estimator with a nonparametric first step. We are not aware of such uniformity results in the literature, and thus leave this issue for future research. Note, on the other hand, that we consider below other confidence intervals that are uniformly conservative.

# 4  An alternative, simple estimator and inference method

## 4.1  The estimator

The first result of Lemma 1 reveals the source of non-identification of $\Delta(x)$: the fact that $\Omega(u,x) := \sum_{t=0}^{T+1} \lambda_t(x,\beta_0)u^t$ is a polynomial of degree $T+1$ rather than $T$. Another idea, then, is to approximate $\Omega$ by a polynomial of degree $T$. Following this idea, we construct a very simple estimator that, in particular, does not require any first-step nonparametric estimator.

Specifically, note that among polynomials of degree $T+1$ with leading coefficient equal to 1, the (renormalized) Chebyshev polynomial $\mathbb{T}_{T+1}^c$ has the lowest supremum norm over $[-1,1]$. Thus, the same holds on $[0,1]$ for $\mathbb{T}_{T+1}(u) := 2^{-T-1}\mathbb{T}_{T+1}^c(2u-1)$. Then, the best approximation of $\Omega$ by a polynomial of degree $T$ for the supremum norm is

$$P_T^*(u,x) := \Omega(u,x) - \lambda_{T+1}(x,\beta_0)\mathbb{T}_{T+1}(u).$$

Figure 1 displays $P_T^*(.,x)$ and $\Omega(.,x)$ with $T \in \{2,3,4\}$ and $x_t'\beta_0 = -1/2 + (t-1)/(T-1)$. As we can see, the approximation is already good for $T=2$, and the two functions become indistinguishable for $T=4$.



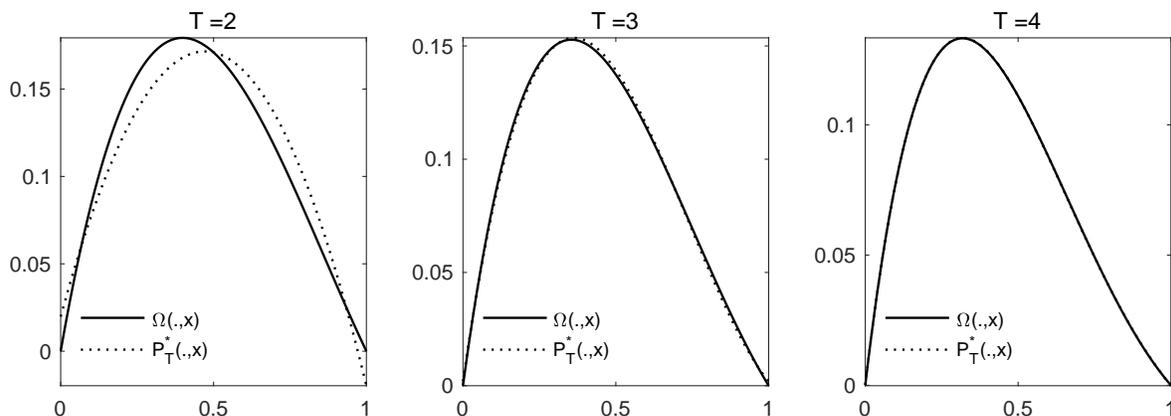Figure 1: Approximation of $\Omega(.,x)$ by $P_T^*(.,x)$ for different $T$.

Next, let $a_t(x,\beta_0)$ denote the coefficient of $u^t$ in $P_T^*(u,x)$. Then, we simply approxi-

mate $\Delta(x)$ by

$$\tilde{\Delta}(x) = \beta_{0k}c_0(x)\int_0^1 P_T^*(u,x)dG(u|x) \tag{10}$$

$$= \beta_{0k}c_0(x)\sum_{t=0}^{T} a_t(x,\beta_0)m_t(x)$$

$$= \beta_{0k}\sum_{t=0}^{T} a_t(x,\beta_0)c_t(x).$$

This leads to the following approximation of $\Delta$:

$$\tilde{\Delta} = \beta_{0k}E\left[\sum_{t=0}^{T} a_t(X,\beta_0)c_t(X)\right]$$

$$= \beta_{0k}E\left[\sum_{t=0}^{S}\frac{\exp(SX_T'\beta_0)a_t(X,\beta_0)\binom{T-t}{S-t}}{C_S(X,\beta_0)}\right].$$

We then estimate $\Delta$ by a plug-in estimator of $\tilde{\Delta}$:

$$\widehat{\Delta} = \frac{\widehat{\beta}_k}{n}\sum_{i=1}^{n}\sum_{t=0}^{S_i}\frac{\exp(S_iX_{iT}'\widehat{\beta})a_t(X_i,\widehat{\beta})\binom{T-t}{S_i-t}}{C_{S_i}(X_i,\widehat{\beta})}. \tag{11}$$

## 4.2 Inference on $\Delta$

Even if $\widehat{\Delta}$ is not a consistent estimator of $\Delta$ when $T$ is fixed, we now show that we can build asymptotically valid confidence intervals for $\Delta$ using $\widehat{\Delta}$. By a slight adjustment, we can even control their asymptotic size over a large class of DGPs. The confidence intervals shrink to $\{\Delta\}$ if $T \to \infty$ and remain of positive length otherwise. To construct these confidence intervals, we rely on two results. The first is the root-$n$ asymptotic normality of $\widehat{\Delta}$. Before displaying this result, we introduce some notation. Recall that $\phi_i = (\phi_{i1}, ..., \phi_{iK})'$ is the influence function of $\widehat{\beta}$ and $\widehat{\phi}_i$ its estimator. Then let

$$\psi_i = E\left[\sum_{t=0}^{S}\frac{a_t(X,\beta_0)\binom{T-t}{S-t}\exp(SX_T'\beta_0)}{C_S(X,\beta_0)}\right]\phi_{ik} + \beta_{0k}\left\{\sum_{t=0}^{S}\frac{a_t(X_i,\beta_0)\binom{T-t}{S_i-t}\exp(S_iX_{iT}'\beta_0)}{C_{S_i}(X_i,\beta_0)}\right.$$

$$\left. + E\left[\sum_{t=0}^{S}\binom{T-t}{S-t}\frac{\partial}{\partial\beta}\left(\frac{a_t(X,\beta_0)\exp(SX_T'\beta_0)}{C_S(X,\beta_0)}\right)\right]'\phi\right\},$$

$$\widehat{\psi}_i = \left[\frac{1}{n}\sum_{j=1}^{n}\sum_{t=0}^{S_i}\frac{a_t(X_i,\widehat{\beta})\binom{T-t}{S_i-t}\exp(S_iX_{iT}'\widehat{\beta})}{C_{S_i}(X_i,\widehat{\beta})}\right]\widehat{\phi}_{ik} + \widehat{\beta}_k\left\{\sum_{t=0}^{S}\frac{a_t(X_i,\widehat{\beta})\binom{T-t}{S_i-t}\exp(S_iX_{iT}'\widehat{\beta})}{C_{S_i}(X_i,\widehat{\beta})}\right.$$

20

$$+\left[\frac{1}{n}\sum_{j=1}^{n}\sum_{t=0}^{S_i}\binom{T-t}{S_i-t}\frac{\partial}{\partial\beta}\left(\frac{a_t(X_i,\widehat{\beta})\exp(S_iX'_{iT}\widehat{\beta})}{C_{S_i}(X_i,\widehat{\beta})}\right)\right]'\widehat{\phi}_i\Bigg\}.$$

Finally, we define $\sigma^2 = V(\psi)$ and $\widehat{\sigma}^2 = \sum_{i=1}^{n}\widehat{\psi}_i^2/n$.

**Lemma 2** *Suppose that Assumptions 1-3 hold. Then*

$$n^{1/2}\left(\widehat{\Delta}-\tilde{\Delta}\right) \xrightarrow{d} \mathcal{N}(0,\sigma^2). \tag{12}$$

*Moreover, $\widehat{\sigma} \xrightarrow{P} \sigma$.*

The second result is a bound on $\tilde{\Delta} - \Delta$, which essentially follows from the fact that the Chebyshev polynomial $\mathbb{T}_{T+1}$ satisfies $\sup_{u\in[0,1]}|\mathbb{T}_{T+1}(u)| \leq 1/[2\times 4^T]$. Below, we let $\mathcal{M}^+ = \arg\max_{u\in[0,1]}\mathbb{T}_{T+1}(u)$ and $\mathcal{M}^- = \arg\min_{u\in[0,1]}\mathbb{T}_{T+1}(u)$.

**Lemma 3** *Suppose that Assumption 1 holds. Then*

$$|\tilde{\Delta}-\Delta| \leq \bar{b} := |\beta_{0k}|E\left[\frac{|\lambda_{T+1}(X,\beta_0)|\binom{T}{S}\exp(SX'_T\beta_0)}{2\times 4^T\times C_S(X,\beta_0)}\right].$$

*Moreover, equality holds if and only if:*

1. *$\beta_{0k} = 0$;*

2. *Or, conditional on $X = x$, $\Lambda(x'_T\beta_0 + \alpha)$ is supported on $\mathcal{M}^+$ for almost all $x$ such that $\lambda_{T+1}(x,\beta_0) > 0$ and on $\mathcal{M}^-$ for almost all $x$ such that $\lambda_{T+1}(x,\beta_0) < 0$;*

3. *Or, conditional on $X = x$, $\Lambda(x'_T\beta_0 + \alpha)$ is supported on $\mathcal{M}^-$ for almost all $x$ such that $\lambda_{T+1}(x,\beta_0) > 0$ and on $\mathcal{M}^+$ for almost all $x$ such that $\lambda_{T+1}(x,\beta_0) < 0$.*

Interestingly, the maximal distance between $\tilde{\Delta}$ and $\Delta$ is equal to half the maximal length of the identified set (since $\bar{q}_T(m) - \underline{q}_T(m) \leq 1/4^T$). This implies that in the worst-case scenario where $\bar{q}_T(m) - \underline{q}_T(m) = 1/4^T$, $\tilde{\Delta}$ is at the middle of the identified set, which is optimal in the sense that $\tilde{\Delta} = \arg\min_d \max_{\Delta\in[\underline{\Delta},\overline{\Delta}]}|d - \Delta|$. Otherwise, on the other hand, $\tilde{\Delta}$ may not even belong to the identified set.

To build a confidence interval on $\Delta$, we first estimate $\bar{b}$ by

$$\widehat{\bar{b}} = \frac{|\widehat{\beta}_k|}{2\times 4^T}\frac{1}{n}\sum_{i=1}^{n}|\lambda_{T+1}(X_i,\widehat{\beta})|\frac{\binom{T}{S_i}\exp(S_iX'_{iT}\widehat{\beta})}{C_{S_i}(X_i,\widehat{\beta})}.$$

Let $Z_n := n^{1/2}\left(\widehat{\Delta} - \Delta\right)/\widehat{\sigma}$ and $b_n := n^{1/2}|b|/\sigma$. To motivate the construction of the confidence intervals, let us first assume that $\widehat{\sigma} = \sigma$, $\widehat{\overline{b}} = \overline{b}$ and the asymptotic approximation (12) is exact. Then

$$Z_n \sim \mathcal{N}\left(n^{1/2}\frac{\tilde{\Delta} - \Delta}{\widehat{\sigma}}, 1\right). \tag{13}$$

Let $q_\alpha(b)$ denote the quantile of order $1 - \alpha$ of a $|\mathcal{N}(b, 1)|$. It is not difficult to show that $b \mapsto q_\alpha(b)$ is symmetric and increasing on $[0, \infty)$. Then, by Lemma 3, if $\widehat{\overline{b}} = \overline{b}$ and (13) holds,

$$P\left(|Z_n| \leq q_\alpha\left(\frac{n^{1/2}\widehat{\overline{b}}}{\widehat{\sigma}}\right)\right) \geq 1 - \alpha. \tag{14}$$

We then define

$$\mathrm{CI}_{1-\alpha}^2 = \left[\widehat{\Delta} \pm q_\alpha\left(\frac{n^{1/2}\widehat{\overline{b}}}{\widehat{\sigma}}\right)\frac{\widehat{\sigma}}{n^{1/2}}\right].$$

The only difference between this confidence interval and a standard one is that because of the possible bias, we consider $q_\alpha\left(n^{1/2}\widehat{\overline{b}}/\widehat{\sigma}\right)$ instead of the usual normal quantile $q_\alpha(0)$. Inequality (14) implies that if (13) holds and $(\widehat{\sigma}, \widehat{\overline{b}}) = (\sigma, \overline{b})$, $\mathrm{CI}_{1-\alpha}$ has a level greater than $1 - \alpha$. Theorem 3 below shows that actually, the same property holds asymptotically without these conditions, as long as $|\tilde{\Delta} - \Delta| < \overline{b}$ or $\beta_{0k} = 0$.

We consider a second class of confidence intervals that handle the other equality cases $|\tilde{\Delta} - \Delta| = \overline{b}$ and is uniformly valid over a large set of DGPs. Specifically, fix $\Theta$ a compact subset of $\mathbb{R}^p$, $\overline{M}$, $\underline{\sigma} \geq 0$ and let $\underline{A}$ be a symmetric positive definite matrix. We consider the following subset of probability distributions:

$$\mathcal{P} = \left\{P : \text{Assumption 1 holds}, \beta_0 \in \Theta, \ P(\|X\| \leq \overline{M}) = 1, \mathcal{I}_{0P} >> \underline{A} \text{ and } \sigma_P^2 \geq \underline{\sigma}^2\right\},$$

where $B >> A$ means that $B - A$ is symmetric positive definite and we index $\mathcal{I}_0$ and $\sigma^2$ by $P$ to underline their dependence in $P$. We then consider the following simple modification of $\mathrm{CI}_{1-\alpha}^1$:

$$\mathrm{CI}_{1-\alpha}^3 = \left[\widehat{\Delta} \pm q_\alpha\left(\frac{n^{1/2}\widehat{\overline{b}} + \varepsilon_n}{\widehat{\sigma}}\right)\frac{\widehat{\sigma}}{n^{1/2}}\right],$$

where $\varepsilon_n \to \infty$. Note that $\varepsilon_n$ may tend to infinity very slowly, for instance we may have $\varepsilon_n = [2\ln\ln n]^{1/2}$.

22

**Theorem 3**

1.  *Suppose that Assumptions 1-3 hold, $\sigma^2 > 0$ and either $|\tilde{\Delta} - \Delta| < \bar{b}$ or $\beta_{0k} = 0$. Then:*

$$\liminf_{n \to \infty} P\left(\Delta \in CI_{1-\alpha}^2\right) \geq 1 - \alpha.$$

2.  *Suppose that Assumption 3.1 holds and $\underline{\sigma} > 0$ in the definition of $\mathcal{P}$. Then:*

$$\liminf_{n \to \infty} \inf_{P \in \mathcal{P}} P\left(\Delta \in CI_{1-\alpha}^3\right) \geq 1 - \alpha.$$

Given Lemma 3, the condition that either $|\tilde{\Delta} - \Delta| < \bar{b}$ or $\beta_{0k} = 0$ is very weak: it is violated only for very peculiar $F_{\alpha|X}$. The first point shows that under this condition, $CI_{1-\alpha}^2$ is pointwise conservative. One potential issue, however, is that $CI_{1-\alpha}^2$ may not be uniformly valid over DGPs if $|\tilde{\Delta} - \Delta|$ becomes very close to $\bar{b}$, or if $\beta_{0k}$ is close but not equal to 0. $CI_{1-\alpha}^3$, on the other hand, is uniformly valid on $\mathcal{P}$, at the price of only slightly enlarging the confidence interval.

## 5 Extensions

We mention in this section other parameters and models for which very similar identification and estimation strategies apply. Specifically, we study the ATE, the average structural function and FE ordered logit models. We also consider the case where $T$ varies per individual. We mostly focus on identification below, but also discuss in some cases how the inference methods above adapt to these set-ups.

### 5.1 Average treatment effects

When the regressor $X_k$ is a binary treatment, we usually consider other parameters than the average marginal effect. Let $X_T^0$ (resp. $X_T^1$) be as $X_T$ but with a 0 (resp. 1) in its $k$-th component. One usual parameter is the average treatment on the treated at period $T$:

$$\Delta^{ATT} = E\left[\Lambda\left(X_T^{1\prime}\beta_0 + \alpha\right) - \Lambda\left(X_T^{0\prime}\beta_0 + \alpha\right)|X_{kT} = 1\right].$$

This is the average effect of a ceteris paribus change of $X_{Tk}$ from 0 to 1, for all individuals satisfying $X_{Tk} = 1$. Because

$$E\left(\Lambda\left(X_T^{1\prime}\beta_0 + \alpha\right)|X_{kT} = 1\right) = E\left(\Lambda\left(X_T^{\prime}\beta_0 + \alpha\right)|X_{kT} = 1\right) = E(Y|X_{kT} = 1)$$

is identified, we just have to focus on the bounds of

$$\Delta^{(1)}(x) = E\left[\Lambda(X_T^{0\prime}\beta_0 + \alpha)|X = x, X_{kT} = 1\right].$$

The functions $c_0, ..., c_T, m_0, ..., m_T$ and $\lambda_0, ..., \lambda_{T+1}$ used for the average marginal effect have to be slightly adapted to $\Delta^{(1)}(x)$. Specifically, for all $t = 0, ..., T$, let

$$\sum_{t=1}^{T+1} \lambda_t^{(1)}(x, \beta_0)u^t := u\prod_{t=1}^{T}\left[1 + u(\exp((x_t - x_T^0)^{\prime}\beta_0) - 1)\right],$$

$$c_t^{(1)}(x) := E\left[\mathbb{1}\{S \geq t\}\binom{T-t}{S-t}\exp(Sx_T^{0\prime}\beta_0)/C_S(x, \beta_0)|X = x, X_{kT} = 1\right],$$

$$m_t^{(1)}(x) := c_t^{(1)}(x)/c_0^{(1)}(x).$$

The following result mimics Lemma 1 for $\Delta^{(1)}(x)$. Its proof is very similar and is thus omitted.

**Lemma 4** *Suppose that Assumptions 1-2 hold, $X_{kT} \in \{0; 1\}$ and $\lambda_{T+1}^{(1)}(x) \geq 0$. Then the sharp identified set of $\Delta^{(1)}(x)$ is $[\underline{\Delta}^{(1)}(x), \overline{\Delta}^{(1)}(x)]$, with*

$$\begin{cases} \underline{\Delta}^{(1)}(x) &= \sum_{t=1}^{T}\lambda_t^{(1)}(x, \beta_0)c_t^{(1)}(x) + c_0^{(1)}(x)\lambda_{T+1}^{(1)}(x, \beta_0)\underline{q}_T(m^{(1)}(x)) \\ \overline{\Delta}^{(1)}(x) &= \sum_{t=1}^{T}\lambda_t^{(1)}(x, \beta_0)c_t^{(1)}(x) + c_0^{(1)}(x)\lambda_{T+1}^{(1)}(x, \beta_0)\overline{q}_T(m^{(1)}(x)) \end{cases}. \quad (15)$$

*If $\lambda_{T+1}^{(1)}(x, \beta_0) < 0$, the same holds with $\overline{q}_T$ and $\underline{q}_T$ switched in the two bounds.*

A similar result holds for the average treatment on the untreated, defined as

$$\Delta^{ATU} = E\left[\Lambda\left(X_T^{1\prime}\beta_0 + \alpha\right) - \Lambda\left(X_T^{0\prime}\beta_0 + \alpha\right)|X_{kT} = 0\right].$$

Finally, consider the average treatment effect

$$\Delta^{ATE} = E\left[\Lambda\left(X_T^{1\prime}\beta_0 + \alpha\right) - \Lambda\left(X_T^{0\prime}\beta_0 + \alpha\right)\right].$$

To compute the bounds on $\Delta^{ATU}$, remark that

$$\Delta^{ATE} = P(X_{kT} = 1)\Delta^{ATT} + P(X_{kT} = 0)\Delta^{ATU}.$$

Moreover, the distribution of $\alpha|X, X_{kT} = 1$ does not restrict the distributions of $\alpha|X, X_{kT} = 0$. As a result, the sharp lower bound $\underline{\Delta}^{ATE}$ on $\Delta^{ATE}$ simply satisfies

$$\underline{\Delta}^{ATE} = P(X_{kT} = 1)\underline{\Delta}^{ATT} + P(X_{kT} = 0)\underline{\Delta}^{ATU}.$$

The same holds for the sharp upper bound.

We can also simply estimate $\Delta^{ATE}$ using the second method. Following the same logic as in Section 4, we obtain the following approximation for $\Delta^{ATE}$:

$$\tilde{\Delta}^{ATE} = E[Y_T(2X_{kT} - 1)] + E\left[\sum_{t=0}^{S} \frac{a_t^{ATE}(X, S, \beta_0)\binom{T-t}{S-t}\exp(SX_T'\beta_0)}{C_S(X, \beta_0)}\right], \qquad (16)$$

where, for $t = 0, ..., T$, we define

$$a_t^{ATE}(x, s, \beta_0) = d_t(x, s, \beta_0) + b_t^* d_{T+1}(x, s, \beta_0),$$

$$d_t(x, s, \beta_0) = -\lambda_t^{(1)}(x, \beta_0)\exp(-s\beta_{0k})x_{kT} + \lambda_t^{(0)}(x, \beta_0)\exp(s\beta_{0k})(1 - x_{kT}),$$

$$c_t^{(0)}(x) = E\left[\mathbb{1}\{S \geq t\}\binom{T-t}{S-t}\exp(Sx_T^{1'}\beta_0)/C_S(x, \beta_0)|X = x, X_{kT} = 0\right],$$

$$m^{(0)}(x) = \left(c_0^{(0)}(x)/c_0^{(0)}(x), ..., c_T^{(0)}(x)/c_0^{(0)}(x)\right),$$

$$\sum_{t=1}^{T+1} \lambda_t^{(0)}(x, \beta_0)u^t = u\prod_{t=1}^{T}\left[1 + u(\exp((x_t - x_T^1)'\beta_0) - 1)\right],$$

and $-(b_0^*, ..., b_T^*)$ are the first $T$ coefficients of $\mathbb{T}_{T+1}$. The estimator $\hat{\Delta}^{ATE}$ of $\Delta^{ATE}$ is then a plug-in estimator based on (16). Also, with the same reasoning as for deriving the upper bound on $\tilde{\Delta} - \Delta$, we obtain

$$|\Delta^{ATE} - \tilde{\Delta}^{ATE}| \leq \bar{b}^{ATE} := E\left[\frac{\binom{T}{S}\exp(SX_T'\beta_0)}{2 \times 4^T \times C_S(X, \beta_0)}\left(|\lambda_{T+1}^{(1)}(X, \beta_0)|\exp(-S\beta_{0k})X_T\right.\right.$$

$$\left.\left. + |\lambda_{T+1}^{(0)}(X, \beta_0)|\exp(S\beta_{0k})(1 - X_T)\right)\right].$$

Then, we can build confidence intervals on $\Delta^{ATE}$ using $\hat{\Delta}^{ATE}$, a plug-in estimator of $\bar{b}^{ATE}$ and an estimator of the asymptotic variance of $\hat{\Delta}^{ATE}$, which is similar to $\hat{\sigma}$.

## 5.2   Average structural function

We now turn to the average structural function defined, for any $x_0 \in \mathbb{R}^p$, by:

$$\Delta_{x_0} := E\left(\Lambda(x_0'\beta + \alpha)\right).$$

As above, we focus below on its conditional counterpart $\Delta_{x_0}^{(2)}(x) := E\left(\Lambda(x_0'\beta_0 + \alpha)|X = x\right)$. We define, for all $t = 0, ..., T$,

$$\sum_{t=1}^{T+1} \lambda_t^{(2)}(x, \beta_0)u^t := u\prod_{t=1}^{T}\left[1 + u(\exp((x_t - x_0)'\beta_0) - 1)\right],$$

$$c_t^{(2)}(x) := E\left[\mathbb{1}\{S \geq t\}\binom{T-t}{S-t}\exp(Sx_0'\beta_0)/C_S(x, \beta_0)|X = x\right],$$

$$m_t^{(2)} := c_t^{(2)}(x)/c_0^{(2)}(x).$$

Again, we obtain a similar result as Lemma 1 on the sharp bounds of $\Delta_{x_0}^{(2)}(x)$:

**Lemma 5** *Suppose that Assumptions 1-2 hold, and $\lambda_{T+1}^{(2)}(x) \geq 0$. Then the sharp identified set of $\Delta_{x_0}^{(2)}(x)$ is $[\underline{\Delta}_{x_0}^{(2)}(x), \overline{\Delta}_{x_0}^{(2)}(x)]$, with*

$$\begin{cases} \underline{\Delta}_{x_0}^{(2)}(x) &= \sum_{t=1}^{T}\lambda_t^{(2)}(x, \beta_0)c_t^{(2)}(x) + c_0^{(2)}(x)\lambda_{T+1}^{(2)}(x, \beta_0)\underline{q}_T(m^{(2)}(x)) \\ \overline{\Delta}_{x_0}^{(2)}(x) &= \sum_{t=1}^{T}\lambda_t^{(2)}(x, \beta_0)c_t^{(2)}(x) + c_0^{(2)}(x)\lambda_{T+1}^{(2)}(x, \beta_0)\overline{q}_T(m^{(2)}(x)) \end{cases}. \quad (17)$$

*If $\lambda_{T+1}^{(2)}(x, \beta_0) < 0$, the same holds with the role of $\overline{q}_T$ and $\underline{q}_T$ switched in the two bounds.*

### 5.3 Average marginal effect in ordered logit models

We now consider a model where the outcome is ordered and takes $J \geq 2$ values.

**Assumption 7** *We have $Y_t = \sum_{k=1}^{J-1}k\mathbb{1}\{\gamma_k \leq X_t'\beta_0 + \alpha + \varepsilon_t < \gamma_{k+1}\}$ with $\gamma_1 = 0 < ... < \gamma_J = +\infty$ and $(\varepsilon_t)_{t=1,...,T}$ are iid, independent of $(\alpha, X)$ and follow a logistic distribution.*

The condition $\gamma_1 = 0$ is a mere normalization: only the differences $\gamma_j - \gamma_{j'}$ are identified since the location of the distribution of $\alpha$ is left unrestricted. In this model, we consider the following AME, for any $j_0 \in \{1, ..., J-1\}$:

$$\Delta^{(3)} = E\left[\frac{\partial P\left(Y_T \geq j_0|X, \alpha\right)}{\partial X_{Tk}}\right].$$

To identify $(\beta_0, \gamma_2, ..., \gamma_{J-1})$, we follow Muris (2017). Let $\Pi$ be the set of functions from $\{1, ..., T\}$ into $\{1, ..., J-1\}$ and for $\pi \in \Pi$, let $Y_t^\pi = \mathbb{1}\{Y_t \geq \pi(t)\}$. By conditioning

on $S^\pi = \sum_t Y_t^\pi$, we get the conditional log-likelihood

$$\ell_c^\pi(y|x; \beta, \gamma_2, ..., \gamma_{J-1}) := \sum_{t=1}^{T} y_t(x_t'\beta - \gamma_{\pi(t)}) - \ln\left[C_{\sum_{t=1}^{T} y_t}^\pi(x, \beta, \gamma)\right],$$

$$\text{with } C_k^\pi(x, \beta, \gamma) := \sum_{(d_1,...,d_T)\in\{0,1\}^T:\sum_{t=1}^{T} d_t=k} \exp\left(\sum_{t=1}^{T} d_t(x_t'\beta - \gamma_{\pi(t)})\right).$$

The parameters $\theta_0 = (\beta_0, \gamma_2, ..., \gamma_{J-1})$ are then identified by stacking, over all $\pi \in \Pi$, the first-order conditions $E[\partial\ell_c^\pi/\partial\theta(Y|X;\theta_0)] = 0$ of the conditional log-likelihood maximization.

Turning to $\Delta^{(3)}$, we consider $\Delta^{(3)}(x) = E\left[\partial P\left(Y_T \geq j_0|X,\alpha\right)/\partial X_{Tk}|X = x\right]$. For any $(j, t) \in \{1, ..., J-1\} \times \{1, ..., T\}$, let $\rho(j, t, x) = \exp((x_t - x_T)'\beta_0 - \gamma_j + \gamma_{j_0}) - 1$ and

$$w(u) = \frac{1}{\prod_{\substack{1\leq j\leq J-1 \\ 1\leq t\leq T}}(1 + u\rho(j, t, x))}.$$

Note that $w(u)$ is well-defined and positive on $[0, 1]$ since $\rho(j, t, x) > -1$. Finally, let $U = \Lambda(x_T'\beta_0 - \gamma_{j_0} + \alpha)$. We show in the proof of Lemma 6 below that:

$$\text{span}\left\{u \mapsto P((Y_1, ..., Y_T) = y|X = x, U = u), \ y \in \{1, ..., J\}^T\right\}$$
$$= \text{span}\left\{u \mapsto u^t w(u), \ t \in \{0, ..., (J-1)T\}\right\}.$$

This means that there exist identified $(c_0^{(3)}(x), ..., c_{(J-1)T}^{(3)}(x))$, such that the $(J-1)T+1$ equations $\int u^t w(u)dF_{U|X=x}(u) = c_t^{(3)}(x)$ exhaust the information provided by the knowledge of $(P((Y_1, ..., Y_T) = y|X = x))_{y\in\{1,...,J\}^T}$. Next, let

$$m^{(3)}(x) = (c_0^{(3)}(x)/c_0^{(3)}(x), ..., c_{(J-1)T}^{(3)}(x)/c_0^{(3)}(x))$$

and as previously, define

$$\sum_{t=0}^{(J-1)T+1} \lambda_t^{(3)}(x, \beta_0)u^t := \frac{u(1-u)}{w(u)}.$$

Again, sharp bounds on $\Delta^{(3)}(x)$ can be obtained as in Lemma 1.

**Lemma 6** *Suppose that Assumptions 2 and 7 hold and* $\lambda_{(J-1)T+1}^{(3)}(x, \beta_0) \geq 0$. *Then the sharp identified set of* $\Delta^{(3)}(x)$ *is* $[\underline{\Delta}^{(3)}(x), \overline{\Delta}^{(3)}(x)]$, *with*

$$\begin{cases} \underline{\Delta}^{(3)}(x) = \beta_{0k}\left[\sum_{t=0}^{(J-1)T} \lambda_t^{(3)}(x, \beta_0)(x)c_t^{(3)}(x) + c_0^{(3)}(x)\lambda_{(J-1)T+1}^{(3)}(x, \beta_0)\underline{q}_{(J-1)T}(m^{(3)}(x))\right] \\ \overline{\Delta}^{(3)}(x) = \beta_{0k}\left[\sum_{t=1}^{(J-1)T} \lambda_t^{(3)}(x, \beta_0)(x)c_t^{(3)}(x) + c_0^{(3)}(x)\lambda_{(J-1)T+1}^{(3)}(x, \beta_0)\overline{q}_{(J-1)T}(m^{(3)}(x))\right] \end{cases}$$

27

If $\lambda^{(3)}_{(J-1)T+1}(x, \beta_0)) < 0$, *the same holds with* $\underline{q}_{(J-1)T}$ *and* $\overline{q}_{(J-1)T}$ *switched in the two bounds.*

The main difference between this result and Lemma 1 above is that the bounds are related to moments of order $(J-1)T+1$ of distributions for which the first $(J-1)T$ raw moments are known. Hence, not surprisingly, the bounds are tighter than in the binary case, and substantially more so given Proposition 3 above.

## 5.4   Varying number of periods

Missing data or attrition are common with panel data. A "panel" may also correspond to hierarchical data where $(i, t)$ corresponds to a unit $t$ belonging to a group $i$ (e.g. individuals within a household). In both cases, $T$ is a random variable varying from one individual (or group) to another. Our method still applies in this case, provided that $T$ is conditionally exogenous. Specifically, we assume that $(\varepsilon_1, ..., \varepsilon_T)$ is independent of $(T, X, \alpha)$. Note, on the other hand, that we remain agnostic on the dependence between $T$ and $(X, \alpha)$. Then, we consider the average marginal effects at period $\underline{T} = \min \text{Supp}(T)$. Other choices are of course possible but this parameter has the advantage of being easily interpretable in the panel case.[5]

Under the independence condition above, the identification and estimation of $\beta_0$ remains unchanged. Also, Lemma 1 holds for each subpopulation satisfying $T = \bar{t}$ and $(X_1, ..., X_{\bar{t}}) = (x_1, ..., x_{\bar{t}})$, for any $\bar{t} \in \text{Supp}(T)$. The only changes therein are that (i) the polynomial in (3) is now $u(1-u) \prod_{t \neq \underline{T}} (u(\exp((x_t - x_{\underline{T}})'\beta_0 - 1))$; (ii) one should replace $x_T$ by $x_{\underline{T}}$ in (4). Next, we obtain the sharp bounds on $\Delta$ by integrating over both $T$ and $(X_1, ..., X_T)$. Similarly, the first estimation method applies for each subpopulation satisfying $T = \bar{t}$, and then one can just sum over all $\bar{t} \in \text{Supp}(T)$.

The second method can also be easily adapted. An inspection of $\tilde{\Delta}$ reveals that the formula remains similar, with the following changes: (i) Equations (3) and (4) should be modified as above; (ii) the Chebyshev polynomials used for the approximation $P_T^*(u, x)$ now vary with $T$. The estimator $\widehat{\Delta}$ and the formulas of $\sigma^2$ and $\overline{b}$ should be adjusted in a similar way.

---

[5]With hierarchical data, the choice of the "period" does not matter anyway.

# 6   Monte Carlo simulations

We study in this section the finite sample performances of our two methods. We consider three DGPs that differ by their distributions of $\alpha|X$. In all cases, we assume that $(X_1, ..., X_T)$ are i.i.d., with $X_t \in \mathbb{R}$, uniformly distributed on $[-1/2, 1/2]$ and $\beta_0 = 1$. The three distributions of $\alpha|X$ are as follows:

1. $\alpha = 0$. Since $|\text{Supp}(\alpha|X)| = 1$ a.s., $\Delta_1 \simeq 0.2449$ is point identified for all $T \geq 2$.

2. Conditional on $X$, $\alpha$ is discrete and takes two values:

$$\alpha = X_T + \eta, \quad P(\eta = -1|X_1, ..., X_T) = P(\eta = 1|X_1, ..., X_T) = 1/2.$$

   Because $|\text{Supp}(\alpha|X)| = 2$, $\Delta_1 \simeq 0.1904$ is partially identified if $T < 4$ and point identified otherwise. The true bounds are $[0.1826, 0.1953]$ if $T = 2$ and $[0.1895, 0.1906]$ if $T = 3$.

3. Conditional on $X$, $\alpha$ is continuous:

$$\alpha = X_T + \eta, \quad \eta|X_1, ..., X_T \sim \mathcal{N}(0, 1).$$

   Because $|\text{Supp}(\alpha|X)| = \infty$, $\Delta_1 \simeq 0.1967$ is partially identified for all $T$. The true bounds are $[0.1905, 0.2015]$ for $T = 2$ and $[0.1961, 0.1970]$ for $T = 3$.

For each of the three DGPs above, we consider $T \in \{2, 3\}$ and $n \in \{250; 500; 1,000\}$. We then compute the estimators of the first and second methods, and $\text{CI}^1_{0.95}$ and $\text{CI}^2_{0.95}$. We use 500 simulations for the first method and 5,000 simulations for the second. To estimate nonparametrically $\gamma_0$, we use local linear estimators with a Gaussian product kernel. Further details on the computation of the estimated bounds are given in Section C of the Online Appendix.

Table 1 displays the properties of the estimators underlying the two methods. The estimators of the bounds appear to have a negligible bias in nearly all cases, even for $n = 250$. Strikingly, the bias of the estimator $\widehat{\Delta}_1$ is also very small compared to its standard deviation; even in the worst case scenarios with $T = 2$ and $n = 1,000$, it is still more than ten times smaller. The last column shows that as expected, in all but one case, $E(\widehat{\overline{b}})$ is larger than the bias of $\widehat{\Delta}_1$. But it is also much smaller than the

standard deviation of $\widehat{\Delta}_1$. We can thus expect the effect of the bias to be small on the length of $\text{CI}^2_{0.95}$.

| DGP | T | n | First method | | | | Second method | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | $\sigma(\widehat{\underline{\Delta}}_1)$ | $\text{Bias}(\widehat{\underline{\Delta}}_1)$ | $\sigma(\widehat{\widehat{\Delta}}_1)$ | $\text{Bias}(\widehat{\widehat{\Delta}}_1)$ | $\sigma(\widehat{\Delta}_1)$ | $\text{Bias}(\widehat{\Delta}_1)$ | $E(\widehat{\widehat{b}})$ |
| 1 | 2 | 250 | 0.133 | 0.005 | 0.138 | 0.011 | 0.115 | 0.0080 | 0.0136 |
| | | 500 | 0.087 | 0.000 | 0.090 | 0.006 | 0.080 | 0.0057 | 0.0119 |
| | | 1,000 | 0.066 | 0.007 | 0.068 | 0.013 | 0.056 | 0.0048 | 0.0112 |
| | 3 | 250 | 0.066 | -0.011 | 0.066 | -0.010 | 0.076 | 0.0027 | 0.0013 |
| | | 500 | 0.044 | -0.015 | 0.044 | -0.015 | 0.055 | 0.0003 | 0.0011 |
| | | 1,000 | 0.031 | -0.013 | 0.031 | -0.013 | 0.038 | 0.0007 | 0.0010 |
| 2 | 2 | 250 | 0.116 | -0.013 | 0.122 | -0.013 | 0.098 | 0.0058 | 0.0155 |
| | | 500 | 0.067 | -0.014 | 0.071 | -0.016 | 0.071 | 0.0035 | 0.0131 |
| | | 1,000 | 0.048 | -0.009 | 0.050 | -0.012 | 0.048 | 0.0034 | 0.0120 |
| | 3 | 250 | 0.065 | -0.003 | 0.066 | -0.004 | 0.066 | -0.0013 | 0.0016 |
| | | 500 | 0.047 | -0.002 | 0.047 | -0.003 | 0.047 | -0.0001 | 0.0013 |
| | | 1,000 | 0.032 | -0.006 | 0.032 | -0.007 | 0.033 | -0.0013 | 0.0012 |
| 3 | 2 | 250 | 0.097 | -0.018 | 0.102 | -0.018 | 0.103 | 0.0076 | 0.0155 |
| | | 500 | 0.061 | -0.020 | 0.065 | -0.022 | 0.072 | 0.0063 | 0.0132 |
| | | 1,000 | 0.048 | -0.009 | 0.051 | -0.011 | 0.051 | 0.0042 | 0.0120 |
| | 3 | 250 | 0.065 | -0.008 | 0.065 | -0.008 | 0.067 | -0.0002 | 0.0016 |
| | | 500 | 0.046 | -0.005 | 0.046 | -0.006 | 0.047 | 0.0001 | 0.0013 |
| | | 1,000 | 0.035 | -0.001 | 0.035 | -0.002 | 0.033 | 0.0008 | 0.0012 |

Notes: the three DGPs correspond respectively to $\alpha = 0$, $\alpha = X_T + \eta$ with $\eta|X \sim$ Rademacher and $\alpha = X_T + \eta$ with $\eta|X \sim \mathcal{N}(0,1)$. We ran 500 (resp. 5,000) simulations with the first (resp. second) method.

Table 1: Properties of the estimators

This is confirmed in Table 2, which presents the coverage rate and length of both confidence intervals. With $T = 3$, and also $T = 2$ in the first DGP, the second method is actually competitive, giving sometimes even slightly smaller confidence intervals, while maintaining a coverage very close to and, in line with the theory, almost always

larger than 95%. That the second method appears competitive with $T = 3$ could be expected: the maximal bias decreases very quickly with $T$, as emphasized in Lemma 3. With $T = 2$ and the second and third DGP, the first method displays better performance, but sometimes at the price of a slight undercoverage. Moreover, except with DGP 3 and $n = 500$, the second confidence interval is never more than 10% larger than the first.

| | | | $\text{CI}^1_{0.95}$ | | $\text{CI}^2_{0.95}$ | |
|---|---|---|---|---|---|---|
| DGP | T | n | Coverage | Avg. length | Coverage | Avg. length |
| 1 | 2 | 250 | 0.94 | 0.451 | 0.96 | 0.452 |
| | | 500 | 0.95 | 0.318 | 0.96 | 0.320 |
| | | 1,000 | 0.93 | 0.225 | 0.96 | 0.227 |
| | 3 | 250 | 0.98 | 0.296 | 0.95 | 0.297 |
| | | 500 | 0.98 | 0.208 | 0.94 | 0.210 |
| | | 1,000 | 0.97 | 0.146 | 0.95 | 0.149 |
| 2 | 2 | 250 | 0.93 | 0.365 | 0.97 | 0.395 |
| | | 500 | 0.94 | 0.255 | 0.96 | 0.280 |
| | | 1,000 | 0.94 | 0.182 | 0.97 | 0.199 |
| | 3 | 250 | 0.97 | 0.271 | 0.96 | 0.260 |
| | | 500 | 0.96 | 0.188 | 0.95 | 0.184 |
| | | 1,000 | 0.96 | 0.132 | 0.95 | 0.130 |
| 3 | 2 | 250 | 0.95 | 0.369 | 0.96 | 0.404 |
| | | 500 | 0.97 | 0.256 | 0.96 | 0.285 |
| | | 1,000 | 0.94 | 0.186 | 0.96 | 0.203 |
| | 3 | 250 | 0.96 | 0.271 | 0.95 | 0.261 |
| | | 500 | 0.97 | 0.191 | 0.95 | 0.185 |
| | | 1,000 | 0.95 | 0.134 | 0.95 | 0.130 |

Notes: the three DGPs correspond respectively to $\alpha = 0$, $\alpha = X_T + \eta$ with $\eta|X \sim$ Rademacher and $\alpha = X_T + \eta$ with $\eta|X \sim \mathcal{N}(0,1)$. We ran 500 (resp. 5,000) simulations with the first (resp. second) method.

Table 2: Coverage and average length of $\text{CI}^1_{0.95}$ and $\text{CI}^2_{0.95}$

# 7 Conclusion

In the FE logit model, the AME can be written as a function of the $(T + 1)$-th raw moment of an unknown distribution for which the first $T$ moments are known. By results in the theory of moments, this implies simple expressions for the sharp bounds of the AME. These bounds can be estimated consistently under weak conditions. Using instead the best uniform approximation of $u^{T+1}$ by a polynomial of degree $T$ yields an even simpler approach for inference on the AME. We expect both ideas to apply to other set-up involving latent variables, such as $\alpha$ in our context.[6]

The theory is simple here because only raw moments are involved; but similar results hold with other moments, provided that the corresponding functions form a so-called Chebyshev system (See, e.g., Krein and Nudelman, 1977, for a mathematical exposition). Results on these systems have already been applied to the optimal design of experiments (see Dette and Studden, 1997) and the measure of segregation with small units (D'Haultfœuille and Rathelot, 2017). By drawing attention on these tools, we hope that this paper will contribute to their use in econometrics.

---

[6]Noteworthy, Dobronyi et al. (2021) apply related results on moment problems to obtain a simple characterization of the identified set of slope parameters in a dynamic FE logit model.

# References

Aguirregabiria, V. and Carro, J. M. (2020), Identification of average marginal effects in fixed effects dynamic discrete choice models. Working paper.

Altonji, J. G. and Matzkin, R. L. (2005), 'Cross section and panel data estimators for nonseparable models with endogenous regressors', *Econometrica* **73**(4), 1053–1102.

Andersen, E. B. (1970), 'Asymptotic properties of conditional maximum-likelihood estimators', *Journal of the Royal Statistical Society. Series B (Methodological)* **32**(2), 283–301.

Angrist, J. D. (2001), 'Estimation of limited dependent variable models with dummy endogenous regressors: simple strategies for empirical practice', *Journal of business & economic statistics* **19**(1), 2–28.

Angrist, J. D. and Pischke, J.-S. (2008), *Mostly harmless econometrics: An empiricist's companion*, Princeton university press.

Armstrong, T. B. and Kolesár, M. (2018), 'Optimal inference in a class of regression models', *Econometrica* **86**(2), 655–683.

Chamberlain, G. (1980), 'Analysis of covariance with qualitative data', *The Review of Economic Studies* **47**(1), 225–238.

Chamberlain, G. (1982), 'Multivariate regression models for panel data', *Journal of econometrics* **18**(1), 5–46.

Chen, X., Linton, O. and Van Keilegom, I. (2003), 'Estimation of semiparametric models when the criterion function is not smooth', *Econometrica* **71**(5), 1591–1608.

Chernozhukov, V., Fernández-Val, I., Hahn, J. and Newey, W. (2013), 'Average and quantile effects in nonseparable panel models', *Econometrica* **81**(2), 535–580.

Chernozhukov, V., Fernandez-Val, I., Hoderlein, S., Holzmann, H. and Newey, W. (2015), 'Nonparametric identification in panels using quantiles', *Journal of Econometrics* **188**(2), 378–392.

Dette, H. and Studden, W. J. (1997), *The theory of canonical moments with applications in statistics, probability, and analysis*, Vol. 338, John Wiley & Sons.

D'Haultfœuille, X. and Rathelot, R. (2017), 'Measuring segregation on small units: A partial identification analysis', *Quantitative Economics* **8**(1), 39–73.

Dobronyi, C., Gu, J. and il Kim, K. (2021), Identification of dynamic panel logit models with fixed effects. arXiv preprint arXiv:2104.04590.

Donoho, D. L. (1994), 'Statistical estimation and optimal recovery', *The Annals of Statistics* pp. 238–270.

Gut, A. (1992), 'The weak law of large numbers for arrays', *Statistics & probability letters* **14**(1), 49–52.

Hahn, J. (1997), 'A note on the efficient semiparametric estimation of some exponential panel models', *Econometric Theory* **13**(4), 583–588.

Hoderlein, S. and White, H. (2012), 'Nonparametric identification in nonseparable panel data models with generalized fixed effects', *Journal of Econometrics* **168**(2), 300–314.

Honoré, B. E. and Tamer, E. (2006), 'Bounds on parameters in panel dynamic discrete choice models', *Econometrica* **74**(3), 611–629.

Honoré, B. E. and Weidner, M. (2020), Dynamic panel logit models with fixed effects. arXiv preprint arXiv:2005.05942.

Imbens, G. W. and Manski, C. F. (2004), 'Confidence intervals for partially identified parameters', *Econometrica* **72**(6), 1845–1857.

Karlin, S. and Shapley, L. S. (1953), *Geometry of moment spaces*, Vol. 12 of *Memoirs of the American Mathematical Society*, American Mathematical Society.

Kong, E., Linton, O. and Xia, Y. (2010), 'Uniform bahadur representation for local polynomial estimates of m-regression and its application to the additive model', *Econometric Theory* pp. 1529–1564.

Krein, M. and Nudelman, A. A. (1977), *The Markov Moment Problem and Extremal Problems*, American Mathematical Society.

Liu, L., Poirier, A. and Shiu, J.-L. (2021), Identification and estimation of average partial effects in semiparametric binary response panel models. arXiv preprint arXiv:2105.12891.

Mason, J. C. and Handscomb, D. C. (2002), *Chebyshev polynomials*, CRC press.

Masry, E. (1996), 'Multivariate local polynomial regression for time series: uniform strong consistency and rates', *Journal of Time Series Analysis* **17**(6), 571–599.

Mundlak, Y. (1978), 'On the pooling of time series and cross section data', *Econometrica: journal of the Econometric Society* pp. 69–85.

Muris, C. (2017), 'Estimation in the fixed-effects ordered logit model', *The Review of Economics and Statistics* **99**(3), 465–477.

Newey, W. K. and McFadden, D. (1994), 'Large sample estimation and hypothesis testing', *Handbook of econometrics* **4**, 2111–2245.

Rasch, G. (1961), On general laws and the meaning of measurement in psychology, *in* 'Proceedings of the fourth Berkeley symposium on mathematical statistics and probability', Vol. 4, Berkeley, Calif., pp. 321–333.

Ruppert, D. and Wand, M. P. (1994), 'Multivariate locally weighted least squares regression', *The annals of statistics* pp. 1346–1370.

van der Vaart, A. W. (2000), *Asymptotic Statistics*, Cambridge University Press.

van der Vaart, A. W. and Wellner, J. A. (1996), *Weak convergence and empirical processes*, Springer.

Wooldridge, J. M. (2019), 'Correlated random effects models with unbalanced panels', *Journal of Econometrics* **211**(1), 137–150.

# A  Potential pitfalls of using FE linear models

We illustrate here two points made in the introduction on the use of FE linear models for binary outcomes. First, the FE linear model will only approximate the effect on the "movers" (in terms of covariates), and this effect may be very different from the effect on the "stayers", and thus also different from the average effect on the whole population. Second, the linear approximation of a true, nonlinear model may be so poor that the approximation of a true treatment effect is of the wrong sign.

To illustrate the first point, suppose the true model is the FE logit model and assume that $T = 2$. Suppose $X_t \in \mathbb{R}$, $\beta_0 = 1$ and we have a dummy variable $M$, with $M = 1$ if the individual is a mover, $M = 0$ otherwise. In the first case, $X_1$ and $X_2$ are i.i.d. and continuous (so that $P(X_1 = X_2|M = 1) = 0$) whereas in the second case, $X_1 = X_2$ a.s. Assume that the individual effect is such that $|\alpha|$ is very large when $M = 0$, whereas $\alpha = 0$ if $M = 1$. Then, the true AME is

$$\Delta = P(M = 1)E[\Lambda'(X_2)|M = 1] + P(M = 0)E[\Lambda'(X_2 + \alpha)|M = 0]$$
$$\simeq P(M = 1)E[\Lambda'(X_2)|M = 1].$$

On the other hand, the linear approximation $\Delta_{\text{lin}}$ of $\Delta$, equal to the slope parameter of the FE linear model, satisfies

$$\begin{aligned}
\Delta_{\text{lin}} &= \frac{E[(Y_2 - Y_1)(X_2 - X_1)]}{E[(X_2 - X_1)^2]} \\
&= \frac{E[(Y_2 - Y_1)(X_2 - X_1)|M = 1]}{E[(X_2 - X_1)^2|M = 1]} \\
&= \frac{E[(\Lambda(X_2) - \Lambda(X_1))(X_2 - X_1)|M = 1]}{E[(X_2 - X_1)^2|M = 1]} \\
&\simeq E[\Lambda'(X_2)|M = 1],
\end{aligned}$$

where the last approximation follows by a Taylor expansion, if $X_2 - X_1$ is small. Thus, in this example, $\Delta_{\text{lin}}$ will overestimate $\Delta$ by the factor $1/P(M = 1)$, which can be arbitrarily large. Note that the reverse holds true if, instead, $|\alpha|$ is very large when $M = 1$ and $\alpha = 0$ when $M = 0$.

To illustrate the second point, suppose that potential outcomes $Y_t(d)$ satisfy

$$Y_t(d) = \mathbb{1}\{\alpha + \mathbb{1}\{t = 2\} + d + \varepsilon_t \geq 0\}, \quad t \in \{1, 2\},$$

where $\varepsilon_2, \varepsilon_2$ are i.i.d. and follow a logistic distribution. We observe $Y_t := Y_t(D_t)$, where the binary treatment satisfies $D_1 = 0$ a.s., whereas $P(D_2 = 1) = 0.5$. Assume further that $\alpha = -0.5 + 1.5D_2$. Then, the true ATE at the second period satisfies

$$
\begin{aligned}
\Delta^{ATE} &= E\left[Y_2(1) - Y_2(0)\right] \\
&= 0.5\left[E\Lambda(1 + 1 + 1) - \Lambda(1 + 1)\right] + 0.5\left[\Lambda(-0.5 + 1 + 1) - \Lambda(-0.5 + 1)\right] \\
&\simeq 0.13.
\end{aligned}
$$

On the other hand, the FE linear model yields the simple difference-in-difference:

$$
\begin{aligned}
\Delta_{\text{lin}}^{ATE} &= E[Y_2 - Y_1|D_2 = 1] - E[Y_2 - Y_1|D_2 = 0] \\
&= \Lambda(1 + 1 + 1) - \Lambda(1) - [\Lambda(-0.5 + 1) - \Lambda(-0.5)] \\
&\simeq -0.02.
\end{aligned}
$$

This problem arises because, basically, the common trends condition is violated in this nonlinear model. One could argue that the difference-in-difference estimand identifies the ATT, not the ATE, under common trends. But note that the ATT is equal to 0.07 and thus also of opposite sign to $\Delta_{\text{lin}}^{ATE}$.

## B   Proofs of the identification results

### B.1   Proposition 1

First, suppose that Assumption 2 does not hold. Then, there exists $\lambda \neq 0$ such that $X_1'\lambda = ... = X_T'\lambda$ almost surely (a.s.). For any $v \in \mathbb{R}$, let $\alpha' = \alpha - vX_t'\lambda$ and $\beta = \beta_0 + v\lambda$. Then

$$
Y_t = \mathbb{1}\left\{X_t'\beta + \alpha' + \varepsilon \geq 0\right\}.
$$

This model satisfies Assumption 1. Thus, $\beta_0$ is not identified.

Now, assume that Assumption 2 holds. By the concavity of the logarithm and Jensen's inequality,

$$
E\left(\ell_c(Y_1, ..., Y_T|X; \beta)\right) \leq E\left(\ell_c(Y_1, ..., Y_T|X; \beta_0)\right)
$$

with equality if and only if $\ell_c(Y_1, ..., Y_T|X; \beta) = \ell_c(Y_1, ..., Y_T|X; \beta_0)$ a.s. Assume that the latter holds. Then, a.s.,

$$\exp[\ell_c(Y_1, ..., Y_T|X; \beta)]\mathbb{1}\{S = 1\} = \exp[\ell_c(Y_1, ..., Y_T|X; \beta_0)]\mathbb{1}\{S = 1\}. \qquad (18)$$

Let us define

$$P_t(\beta) := \frac{\exp(X_t'\beta)}{\sum_{s=1}^T \exp(X_s'\beta)}.$$

Equality (18) is equivalent to

$$\sum_{t=1}^T (P_t(\beta) - P_t(\beta_0))Y_t \prod_{s \neq t}(1 - Y_s) = 0 \quad \text{a.s.}$$

Because the variables $(Y_t \prod_{s \neq t}(1 - Y_s))_t$ are mutually exclusive, we have, for all $t$,

$$(P_t(\beta) - P_t(\beta_0))Y_t \prod_{s \neq t}(1 - Y_s) = 0 \quad \text{a.s.}$$

By taking the expectation with respect to $X$ and noting that $P(Y_t \prod_{s \neq t}(1 - Y_s)|X) > 0$ a.s., we get, a.s., $P_t(\beta) = P_t(\beta_0)$. This in turn, implies that $X_t'(\beta - \beta_0)$ does not depend on $t$. Hence, a.s.,

$$\sum_{t,s} [(X_t - X_s)'(\beta - \beta_0)]^2 = 0.$$

Taking the expectation, this implies that

$$(\beta - \beta_0)'E\left[\sum_{t,s}(X_t - X_s)(X_t - X_s)'\right](\beta - \beta_0) = 0.$$

By Assumption 2, $\beta = \beta_0$. Hence, $\beta_0$ is identified and $\beta_0 = \arg\max_\beta E(\ell_c(Y|X, \beta))$.

We finally turn to the last result. If $S = 1$, we have $\partial \ell_c/\partial \beta(Y_1, ..., Y_T|X; \beta_0) = \sum_{t=1}^T X_t (Y_t - P_t(\beta_0))$. Then, conditional on $S = 1$,

$$\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_1, ..., Y_T|X; \beta_0) = -\sum_{t=1}^T X_t P_t(\beta_0) \sum_{s=1}^T (X_t - X_s)'P_s(\beta_0)$$

$$= -\frac{1}{2}\sum_{s,t} P_s(\beta_0)P_t(\beta_0)(X_t - X_s)(X_t - X_s)'.$$

Let $\lambda$ be such that $\lambda'\mathcal{I}_0\lambda = 0$. Because $-\partial^2 \ell_c/\partial \beta \partial \beta'$ is positive semidefinite, we have

$$\lambda'\mathcal{I}_0\lambda \geq \lambda'E\left[\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_1, ..., Y_T|X; \beta_0)\mathbb{1}\{S = 1\}\right]\lambda$$

$$= \frac{1}{2} \sum_{s,t} E\left[P_s(\beta_0)P_t(\beta_0)\mathbb{1}\left\{S=1\right\}\lambda'\left(X_t - X_s\right)\left(X_t - X_s\right)'\lambda\right]$$

$$= \frac{1}{2} \sum_{s,t} E\left[P_s(\beta_0)P_t(\beta_0)P(S=1|X)\left[(X_t - X_s)'\lambda\right]^2\right].$$

Hence, for all $(s,t)$, $P_s(\beta_0)P_t(\beta_0)P(S=1|X)\left[(X_t - X_s)'\lambda\right]^2 = 0$ almost surely. Since $P(S=1|X) > 0$, we have $(X_t - X_s)'\lambda = 0$ almost surely. In turn, this implies that

$$\lambda'E\left[\sum_{s,t}(X_t - X_s)(X_t - X_s)'\right]\lambda = 0.$$

Thus, by Assumption 2, $\lambda = 0$, proving that $\mathcal{I}_0$ is nonsingular.

## B.2  Lemma 1

First, let us define $U = \Lambda(\alpha + x_T'\beta_0)$. Then, a change of variable in (1)-(2) shows that

$$\Delta(x) = \beta_{0k} \int_0^1 u(1-u)dF_{U|X}(u|x), \tag{19}$$

$$P(S=k|X=x) = C_k(x,\beta_0)\exp(-kx_T'\beta_0)$$
$$\times \int_0^1 \frac{u^k(1-u)^{T-k}}{\prod_{t=1}^{T-1}[1+u(\exp((x_t - x_T)'\beta_0)-1)]}dF_{U|X}(u|x). \tag{20}$$

Remark that for all $t \in \{0,...,T\}$, $\sum_{k=t}^{T}\binom{T-t}{k-t}u^{k-t}(1-u)^{T-k} = 1$. Then, for such $t$,

$$c_t(x) = E\left[\mathbb{1}\left\{S \geq t\right\}\binom{T-t}{S-t}\exp(Sx_T'\beta_0)/C_S(x,\beta_0)|X=x\right]$$
$$= \int_0^1 \frac{\sum_{k=t}^{T}\binom{T-t}{k-t}u^k(1-u)^{T-k}}{\prod_{t=1}^{T-1}[1+u(\exp((x_t - x_T)'\beta_0)-1)]}dF_{U|X}(u|x)$$
$$= \int_0^1 \frac{u^t}{\prod_{t=1}^{T-1}[1+u(\exp((x_t - x_T)'\beta_0)-1)]}dF_{U|X}(u|x).$$

Let $\mu$ be the measure having a density with respect to $F_{U|X}(\cdot|x)$ equal to

$$\frac{\prod_{t=1}^{T-1} 1/[1+u(\exp((x_t - x_T)'\beta_0)-1)]}{\int_0^1 \prod_{t=1}^{T-1} 1/[1+v(\exp((x_t - x_T)'\beta_0)-1)]dF_{U|X}(v|x)}.$$

Then, by definition of $m_t(x)$, we obtain, for all $t \in \{0,...,T\}$

$$m_t(x) = \int_0^1 u^t d\mu(u) \tag{21}$$

39

By a change of measure in (19) and by definition of $(\lambda_t(x, \beta_0))_{t=0,...,T+1}$, we also get

$$\Delta(x) = \beta_{0k}c_0(x) \int_0^1 \sum_{t=0}^{T+1} \lambda_t(x, \beta_0)u^t d\mu(u),$$

where $\mu \in \mathcal{D}(m(x))$. The second result follows using (21), $c_t(x) = c_0(x)m_t(x)$ and the fact that $\mathcal{D}(m(x))$, and thus $\{\int u^{T+1}d\mu(u) : \mu \in \mathcal{D}(m(x))\}$, are convex.

## B.3 Proposition 2

Theorem 1.4.3 in Dette and Studden (1997) ensures that $\underline{H}_T(m)\overline{H}_T(m) \geq 0$ for any $m \in \mathcal{M}_T$.

1. If $\underline{H}_T(m)\overline{H}_T(m) > 0$, again by Theorem 1.4.3 in Dette and Studden (1997), $m \in \text{Int } \mathcal{M}_T$. Then, by Theorem 1.2.7 in Dette and Studden (1997), $\underline{q}_T(m) < \overline{q}_T(m)$. Moreover, $(m_0, ..., m_{T-1}) \in \text{Int } \mathcal{M}_{T-1}$. Thus,

$$\underline{H}_{T-1}(m_0, ..., m_{T-1}) > 0.$$

By expanding the determinant $\underline{H}_{T+1}(m, q)$ along its last column, we get that $q \mapsto \underline{H}_{T+1}(m, q)$ is linear and strictly increasing. The same reasoning applies to $\overline{q}_T(m)$.

2. If $\underline{H}_T(m)\overline{H}_T(m) = 0$, Theorem 1.4.3 in Dette and Studden (1997) implies that $m \in \partial \mathcal{M}_T$. Then, by Theorem 1.2.5 in Dette and Studden (1997), there is a unique distribution corresponding to $m$, implying that $\underline{q}_T(m) = \overline{q}_T(m)$. Let $U$ denote a random variable with $T$ first moments equal to $m$. Suppose first that $T'$ is even and $\underline{H}_{T'}(m) = 0$. Then, there exists a vector $\lambda = (\lambda_1, ..., \lambda_{T'/2+1})'$ such that $\mathbb{H}_{T'}(m)\lambda = 0$. Hence, for all $i \in \{1, ..., T'/2 + 1\}$, $\sum_{j=1}^{T'/2+1} \lambda_j m_{i+j-2} = 0$. Thus, for all $i \in \{0, ..., T'/2\}$,

$$E\left[U^i \sum_{j=0}^{T'/2} \lambda_{j+1}U^j\right] = 0.$$

Hence, $E\left[\left(\sum_{j=0}^{T'/2} \lambda_{j+1}U^j\right)^2\right] = 0$, which implies that almost surely, $\sum_{j=0}^{T'/2} \lambda_{j+1}U^j = 0$. In particular, for all $k \geq 1$ and letting $m_k := E(U^k)$ for $k > T$, we have

$$\sum_{j=1}^{T'/2+1} \lambda_j m_{j+k-2} = 0.$$

40

Since this holds for $k \in \{T+2-T', ..., T+2-T'/2\}$, we have $\underline{\mathbb{H}}_{T'}(m_{T+1-T'}, ..., m_{T+1})\lambda = 0$. Therefore, $\underline{H}_{T'}(m_{T+1-T'}, ..., m_{T+1}) = 0$, with $\underline{q}_T(m) = \bar{q}_T(m) = m_{T+1}$.

The reasoning is the same if $T'$ is odd and still $\underline{H}_{T'}(m) = 0$, with just one difference. Instead of having $E\left[\left(\sum_{j=0}^{T'/2-1} \lambda_{j+1}U^j\right)^2\right] = 0$, we have

$$E\left[U\left(\sum_{j=0}^{(T'-1)/2-1} \lambda_{j+1}U^j\right)^2\right] = 0.$$

But since $U \geq 0$, this still implies $U\left(\sum_{j=0}^{(T'-1)/2-1} \lambda_{j+1}U^j\right)^2 = 0$, and the rest of the proof is similar as above. When instead $\overline{H}_{T'}(m) = 0$ and $T'$ is even, we have instead $E\left[U(1-U)\left(\sum_{j=0}^{T'/2-1} \lambda_{j+1}U^j\right)^2\right] = 0$, implying again $U(1-U)\left(\sum_{j=0}^{T'/2-1} \lambda_{j+1}U^j\right)^2 = 0$. Finally, when $\overline{H}_{T'}(m) = 0$ and $T'$ is odd, we have instead

$$E\left[(1-U)\left(\sum_{j=0}^{(T'-1)/2-1} \lambda_{j+1}U^j\right)^2\right] = 0,$$

implying again $(1-U)\left(\sum_{j=0}^{(T'-1)/2-1} \lambda_{j+1}U^j\right)^2 = 0$.

### B.4 Proposition 3

1. From Lemma 1, we have

$$\overline{\Delta}(x) - \underline{\Delta}(x) = |\beta_{0k}| \times c_0(x) \times |\lambda_{T+1}(x, \beta_0)| \times [\bar{q}_T(m) - \underline{q}_T(m)],$$

By a result of Karlin and Shapley (1953),

$$\bar{q}_T(m) - \underline{q}_T(m) \leq \frac{1}{4^T}.$$

Next, $\Lambda'(u) \in (0, 1/4)$ for all $u$ and therefore $E\left(\Lambda'(\alpha + x_T'\beta_0)|X = x\right) \in (0, 1/4)$ for any conditional distribution of $\alpha|X$. Hence, letting $\Omega(u, x) := \sum_{t=1}^{T+1} \lambda_t(x, \beta_0)u^t$, we have, for any distribution $G$ on $[0; 1]$,

$$c_0(x) \int \Omega(u, x)dG(u|x) \leq \frac{1}{4}.$$

In particular $c_0(x)\Omega(1/2, x) \leq 1/4$. Now,

$$\Omega(1/2, x) = \frac{1}{4}\prod_{t=1}^{T-1}\left(1 + \frac{\exp((x_t - x_T)'\beta_0) - 1}{2}\right),$$

which ensures that $c_0(x) \leq 2^{T-1}/\prod_{t=1}^{T-1}(\exp((x_t - x_T)'\beta_0) + 1)$. This proves the first inequality. Because $\frac{|a-1|}{a+1} \leq 1$ for $a \geq 0$, we have $\overline{\Delta}(x) - \underline{\Delta}(x) \leq |\beta_{0k}|/2^{T+1}$.

2. By Lemma 1, $\Delta(x)$ is point identified if and only if $\beta_{0k} = 0$ or $\lambda_{T+1}(x, \beta_0) = 0$ or $\underline{q}_T(m) = \overline{q}_T(m)$. We have $\lambda_{T+1}(x, \beta_0) = 0$ if and only if $(x_t - x_T)'\beta_0 = 0$ for some $t < T$. By Proposition 2, $\underline{q}_T(m) = \overline{q}_T(m)$ is equivalent to $\underline{H}_T(m) \times \overline{H}_T(m) = 0$. By the proof of Lemma 1 and, e.g., Theorem 1.2.5 in Dette and Studden (1997), this holds if and only if the index of the conditional distribution of $U = \Lambda(\alpha + x_T'\beta_0)$ is smaller than or equal to $T/2$. The index denotes here the number of support points, except that 0 and 1 are counted only as one-half. Now, because $U = \Lambda(\alpha + x_T'\beta_0)$, $P(U = 0|X = x) = P(U = 1|X = x) = 0$. Hence, the number of support points of $U|X = x$, or equivalently $\alpha|X = x$, is smaller than or equal to $T/2$.

## B.5  Lemma 3

Using (10) and $\Delta(x) = \beta_{0k}c_0(x)\lambda_{T+1}(x, \beta_0)\int_0^1 \Omega(u, x)dG(u|x)$, we obtain

$$
\begin{aligned}
\left|\tilde{\Delta}(x) - \Delta(x)\right| &\leq |\beta_{0k}\lambda_{T+1}(x, \beta_0)|c_0(x) \sup_{u \in [0,1]} |\mathbb{T}_{T+1}(u)| \qquad (22) \\
&= \frac{|\beta_{0k}\lambda_{T+1}(x, \beta_0)|c_0(x)}{2^{T+1}} \sup_{u \in [-1,1]} |\mathbb{T}_{T+1}^c(u)| \\
&= \frac{|\beta_{0k}\lambda_{T+1}(x, \beta_0)|c_0(x)}{2 \times 4^T}.
\end{aligned}
$$

The last equality follows by standard properties of Chebyshev polynomials, see, e.g., Mason and Handscomb (2002). The first result follows by integration, using

$$
|\tilde{\Delta} - \Delta| \leq E\left[\left|\tilde{\Delta}(X) - \Delta(X)\right|\right]. \qquad (23)
$$

By what precedes, $|\tilde{\Delta} - \Delta| = \bar{b}$ if and only if we have an equality in (22) for almost all $x$, and an equality in (23). The latter holds if and only if $\beta_{0k} = 0$, or the sign of $\lambda_{T+1}(X)\int_0^1 \mathbb{T}_{T+1}(u)dG(u|X)$ is constant. The former holds if and only if $\beta_{0k} = 0$ or the support of $G(\cdot|x)$ is either $\mathcal{M}^+$ or $\mathcal{M}^-$. The characterization of the equality $|\tilde{\Delta} - \Delta| = \bar{b}$ follows.

## B.6 Lemma 6

Let $\Pi_0$ be the set of functions from $\{1, ..., T\}$ into $\{0, 1, ..., J - 1\}$. First, we prove that the set of conditional probabilities $(P(S^\pi = s | X, U))_{s=0,...,T, \pi \in \Pi_0}$ is in one-to-one linear mapping with $P(Y = y | X, U)_{y \in \{0,1,...,J-1\}^T}$. First,

$$P(Y = (J - 1, ..., J - 1) | X, U) = P(S^{\bar{\pi}} = T | X, U)$$

with $\bar{\pi}$ such that $\bar{\pi}(t) = J - 1$ for all $t$. Next, for any $y = (y_1, ..., y_T) \in \{0, 1, ..., J-1\}^T$, let $\pi \in \Pi_0$ be such that $\pi(t) = y_t$. Then:

$$P(Y = y | X, U) = P(S^\pi = T | X, U) - \sum_{\substack{y':y' \neq y \\ \forall t, y'_t \geq y_t}} P(Y = y' | X, U).$$

Hence, by a decreasing induction on $y$, using the lexicographic order, $P(Y = y | X, U)$ is a linear combination of the $(P(S^\pi = T | X, U))_{\pi \in \Pi_0}$. Conversely, $P(S^\pi = s | X, U) = \sum_{y \in \mathcal{Y}_s^\pi} P(Y = y | X, U)$ with $\mathcal{Y}_s^\pi = \left\{ y \in \{0, 1, ..., J - 1\}^T : \sum_t \mathbb{1}\{y_t \geq \pi(t)\} = s \right\}$. This ensures that $(S^\pi)_{\pi \in \Pi_0}$ is exhaustive for $U$ and

$$\mathrm{span}\left\{ u \mapsto P(Y = y | X, U = u), \, y \in \{0, 1, ..., J - 1\}^T \right\}$$
$$= \mathrm{span}\left\{ u \mapsto P(S^\pi = s | X, U = u), \, s = 0, ..., T, \pi \in \Pi_0 \right\}.$$

Then the sharp lower bound (say) $\underline{\Delta}^{(3)}(x)$ satisfies:

$$\underline{\Delta}^{(3)}(x) = \arg\min_{F_{U|X}(.|x)} \int \frac{\partial P(Y_T \geq j_0 | X = x, U = u)}{\partial X_{Tk}} dF_{U|X}(u|x)$$
$$\text{s.t.} \int P(S^\pi = s | X = x, U = u) \, dF_{U|X}(u|x) = P(S^\pi = s | X = x),$$
$$\pi \in \Pi_0, \, s \in \{0, 1, ..., T\}.$$

Thus, to conclude the proof, it suffices to show

$$\mathrm{span}\left\{ u \mapsto P(S^\pi = s | X = x, U = u), \, \pi \in \Pi_0, s = 0, ..., T \right\}$$
$$= \mathrm{span}\left\{ u \mapsto u^t w(u), \, t = 0, ..., (J - 1)T \right\}. \tag{24}$$

For $\pi \in \Pi_0$, let $\mathcal{T}_+^\pi = \{t : \pi(t) > 0\}$ and for $k \leq |\mathcal{T}_+^\pi|$, let $\mathcal{D}_k^\pi = \{d \in \{0,1\}^{\mathcal{T}_+^\pi} : \sum_{t \in \mathcal{T}_+^\pi} d_t = k\}$ and

$$C_k^\pi(x, \beta, \gamma) := \sum_{d \in \mathcal{D}_k^\pi} \exp\left( \sum_{t \in \mathcal{T}_+^\pi} d_t(x_t' \beta - \gamma_{\pi(t)}) \right).$$

43

For any $\pi \in \Pi_0$, let $s_0^\pi = T - |\mathcal{T}_+^\pi|$. We have

$$
\begin{aligned}
&P\left(S^\pi = s | X = x, U = u\right) \\
&= \frac{C_{s-s_0^\pi}^\pi(x, \beta_0, \gamma) \exp(-(s - s_0^\pi)(x_T'\beta_0 - \gamma_{j_0})) u^{s-s_0^\pi}(1-u)^{T-s}}{\prod_{t \in \mathcal{T}_+^\pi} [1 + u\rho(\pi(t), t, x)]} \mathbb{1}\left\{s_0^\pi \le s \le T\right\}.
\end{aligned}
$$

The Bernstein polynomials $\{u \mapsto u^{s-s_0^\pi}(1-u)^{T-s}, \ s = s_0^\pi, ..., T\}$ are a basis of polynomials of degree lower than $|\mathcal{T}_+^\pi|$. Thus,

$$
\begin{aligned}
&\operatorname{span}\left\{u \mapsto P(S^\pi = s | X = x, U = u), \ \pi \in \Pi_0, \ s = 0, ..., T\right\} \\
&= \operatorname{span}\left\{u \mapsto \frac{u^t}{\prod_{t \in \mathcal{T}_+^\pi} [1 + u\rho(\pi(t), t, x)]}, \ \pi \in \Pi_0, \ t = 0, ..., |\mathcal{T}_+^\pi|\right\} \quad (25) \\
&\subset \operatorname{span}\left\{u \mapsto u^t w(u), \ t = 0, ..., (J-1)T\right\}.
\end{aligned}
$$

Conversely, let $\sim$ be the equivalence relation on $\{1, ..., J-1\} \times \{1, ..., T\}$ defined by:

$$(j, t) \sim (j', t') \Leftrightarrow \rho(j, t, x) = \rho(j', t', x).$$

Then let $[(j, t)]$ denote the equivalence class of $(j, t)$ and let $\mathcal{E}$ be a set of representatives of all the equivalence classes, except $[(j_0, T)]$. Let also $n(j, t) = |[(j, t)]|$. Using $\rho(j_0, T, x) = 0$ and partial fraction decompositions, we obtain

$$
\begin{aligned}
&\operatorname{span}\left\{u \mapsto u^t w(u), \ t = 0, ..., (J-1)T\right\} \\
&\subset \operatorname{span}\left\{u \mapsto u^d, \ d = 0, ..., n(j_0, T), \ u \mapsto (1 + u\rho(j, t', x))^{-d}, \right. \\
&\qquad\qquad \left. (j, t') \in \mathcal{E}, \ d = 1, ..., n(j, t')\right\}. \quad (26)
\end{aligned}
$$

Fix $(j, t') \in \mathcal{E}$, $d \in \{1, ..., n(j, t')\}$ and let $(j_1, t_1), ..., (j_d, t_d)$ denote $d$ distinct elements of $[(j, t')]$. By definition of $\rho$, $t_1, ..., t_d$ are all distinct. Then, define $\pi \in \Pi_0$ as $\pi(t_i) = j_i$ for $i = 1, ..., d$ and $\pi(t) = 0$ for $t \notin \{t_1, ..., t_d\}$. Then:

$$\frac{1}{(1 + u\rho(j, t', x))^d} = \frac{1}{\prod_{t \in \mathcal{T}_+^\pi} [1 + u\rho(\pi(t), t, x)]}. \quad (27)$$

Next, fix $d \in \{0, ..., n(j_0, T)\}$, and let $(j_1, t_1), ..., (j_d, t_d)$ denote $d$ distinct elements of $[(j_0, T)]$. Define $\pi \in \Pi_0$ exactly as above if $d > 0$ and $\pi(t) = 0$ for all $t$ if $d = 0$. Using $\rho(j_i, t_i, x) = 0$ for $i = 1, ..., d$ and the definition of $\mathcal{T}_+^\pi$, we obtain $d = |\mathcal{T}_+^\pi|$ and

$$u^d = \frac{u^d}{\prod_{t \in \mathcal{T}_+^\pi} [1 + u\rho(\pi(t), t, x)]}. \quad (28)$$

Using (26) (27) and (28) and then (25), we finally obtain

$$\text{span}\left\{u \mapsto u^t w(u), \, t = 0, ..., (J-1)T\right\}$$

$$\subset \text{span}\left\{u \mapsto \frac{u^t}{\prod_{t \in \mathcal{T}_+^\pi} [1 + u\rho(\pi(t), t, x)]} \,\, \pi \in \Pi_0, \, t = 0, ..., |\mathcal{T}_+^\pi|\right\}$$

$$= \text{span}\left\{u \mapsto P(S^\pi = s | X = x, U = u), \, \pi \in \Pi_0, \, s = 0, ..., T\right\}.$$

Equation (24) follows, and this concludes the proof.

# Online Appendix

## A  Proofs of the asymptotic results

### A.1  Theorem 1

We focus on $\widehat{\underline{\Delta}}$ hereafter, as the proof for the upper bound is the same. The proof proceeds in three steps. First, we show the uniform consistency of $\widetilde{m}$ over $\mathrm{Supp}(X)$. Second, we prove that $\widehat{m}$ is also uniformly consistent. Finally, we show the consistency of $\widehat{\underline{\Delta}}$. For any function $f$ from $\mathcal{D}$ to $\mathbb{R}^q$, we let $\|f\|_\infty = \sup_{x \in \mathcal{D}} \|f(x)\|$, where $\|\cdot\|$ denotes the Euclidean norm on $\mathbb{R}^q$. We denote by $C$, $\underline{C}$ and $\overline{C}$ generic constants subject to change from one line to the next.

*Step 1: Uniform consistency of $\widetilde{m}$*

Remark that the Under Assumptions 1-3, $P \in \mathcal{P}$ as defined in Subsection 4.2, with $\underline{\sigma} = 0$ and some appropriate $\overline{M}$ and $A$. Then, by Lemma 8, $\widehat{\beta} \xrightarrow{P} \beta_0$. Moreover, $\mathrm{Supp}(X)$ is compact. Then, for all $(k, x, \beta) \in \{0, ..., T\} \times \mathrm{Supp}(X) \times \{\beta_0, \widehat{\beta}\}$, with probability approaching one (wpao),

$$\overline{C} > C_k(x, \beta) \geq \underline{C} > 0, \quad \overline{C} > \exp(k x_T' \beta) \geq \underline{C}. \tag{29}$$

Moreover, by definition of $c_k(\gamma, x, \beta)$,

$$\left\| c(\widehat{\gamma}, x, \widehat{\beta}) - c(\gamma_0, x, \beta_0) \right\|$$
$$\leq C \left\| \left( \frac{\widehat{\gamma}_0(x) e^{0 \times x_T' \widehat{\beta}}}{C_0(x, \widehat{\beta})}, \ ... \ , \frac{\widehat{\gamma}_T(x) e^{T \times x_T' \widehat{\beta}}}{C_T(x, \widehat{\beta})} \right)' - \left( \frac{\gamma_{00}(x) e^{0 \times x_T' \beta_0}}{C_0(x, \beta_0)}, \ ... \ , \frac{\gamma_{0T}(x) e^{T \times x_T' \widehat{\beta}}}{C_T(x, \beta_0)} \right)' \right\|.$$

Fix $0 \leq k \leq T$. Then wpao,

$$\left| \frac{\widehat{\gamma}_k(x) e^{k \times x_T' \widehat{\beta}}}{C_k(x, \widehat{\beta})} - \frac{\gamma_{0k}(x) e^{k \times x_T' \beta_0}}{C_k(x, \beta_0)} \right| \leq \frac{|\widehat{\gamma}_k(x) - \gamma_{0k}(x)| e^{k \times x_T' \widehat{\beta}}}{C_k(x, \widehat{\beta})} + \gamma_{0k}(x) \left| \frac{e^{k \times x_T' \beta_0}}{C_k(x, \beta_0)} - \frac{e^{k \times x_T' \widehat{\beta}}}{C_k(x, \widehat{\beta})} \right|.$$

The derivatives of $\beta \mapsto e^{k \times x_T' \beta} / C_k(x, \beta)$ are uniformly bounded over $\beta$ and $x \in \mathrm{Supp}(X)$ wpao. Combined with (29), this implies that wpao,

$$\left| \frac{\widehat{\gamma}_k(x) e^{k \times x_T' \widehat{\beta}}}{C_k(x, \widehat{\beta})} - \frac{\gamma_{0k}(x) e^{k \times x_T' \beta_0}}{C_k(x, \beta_0)} \right| \leq C \left( |\widehat{\gamma}_k(x) - \gamma_{0k}(x)| + \|\beta_0 - \widehat{\beta}\| \right).$$

Therefore, recalling that $\widehat{c} = c(\widehat{\gamma}, x, \widehat{\beta})$,

$$\|\widehat{c} - c\|_\infty \leq C \left( \|\widehat{\gamma} - \gamma_0\|_\infty + \left\|\beta_0 - \widehat{\beta}\right\| \right). \tag{30}$$

Next, by (7), (29) and $\sum_{j=0}^{T} \gamma_{0j}(x) = 1$, for all $(x, \beta) \in \mathrm{Supp}(X) \times \{\beta_0, \widehat{\beta}\}$, wpao,

$$c_0(\gamma_0, x, \beta) > \sum_{j=0}^{T} \gamma_{0j}(x) \underline{C}/\overline{C} = \underline{C}/\overline{C}. \tag{31}$$

The conditions in Theorem 6 of Masry (1996) hold under Assumptions 1-4. Thus, $\widehat{\gamma}$ is uniformly consistent. Given (30) and (31), we then have $c_0(\widehat{\gamma}, x, \widehat{\beta}) > C$ wpao.

By definition of $\widetilde{m}$, we have, for all $(k, x) \in \{0, ..., T\} \times \mathrm{Supp}(X)$, wpao,

$$
\begin{aligned}
|\widetilde{m}_k(x) - m_k(x)| \leq &\frac{1}{c_0(\gamma_0, x, \beta_0)} |c_k(\widehat{\gamma}, x, \widehat{\beta}) - c_k(\gamma, x, \beta_0)| \\
&+ \frac{1}{\widetilde{c}_0^2} |c_k(\widehat{\gamma}, x, \widehat{\beta})| \times |c_0(\widehat{\gamma}, x, \widehat{\beta}) - c_0(\gamma_0, x, \beta_0)|
\end{aligned} \tag{32}
$$

where $\widetilde{c}_0^2 \geq \min(c_0(\gamma_0, x, \beta)^2, c_0(\widehat{\gamma}, x, \widehat{\beta})^2) > C$ and $|c_k(\widehat{\gamma}, x, \widehat{\beta})|$ is bounded in probability in view of (30). Therefore, by (32) and, again, (30),

$$
\begin{aligned}
\|\widetilde{m} - m\|_\infty &\leq C \left( \|c - \widehat{c}\|_\infty + \|c_0 - \widehat{c}_0\|_\infty \right) \\
&\leq C \left( \|\widehat{\gamma} - \gamma_0\|_\infty + \left\|\beta_0 - \widehat{\beta}\right\| \right).
\end{aligned}
$$

The result follows by uniform consistency of $\widehat{\gamma}$ and consistency of $\widehat{\beta}$.

*Step 2: Uniform consistency of $\widehat{m}$*

We drop the dependence in $x$ and write $m, \widehat{m},...$ instead of $m(x), \widehat{m}(x),...$ to simplify notation as all the statements to follow hold uniformly over $x \in \mathrm{Supp}(X)$. We start by showing that for all $\epsilon > 0$ and for $n$ large enough, if $\widehat{I} = t$ then $|m_{t+1} - \widehat{m}_{t+1}| \leq 2\epsilon$. A first step is to notice that for all $\epsilon > 0$, there exists $N_0$ such that $n \geq N_0$, $m \in \mathcal{M}_T$ and $\underline{H}_{t+1}(m_1, ..., m_{t+1}) < 2c_n^{1/2}$ implies $|m_{t+1} - \underline{q}_t(m_{\to t})| = |m_{t+1} - \widehat{m}_{t+1}| \leq \epsilon$. To see this, suppose the contrary. Then there exists $\epsilon > 0$ and a subsequence $(m^{\phi(n)}) \in \mathcal{M}_T^{\mathbb{N}}$ such that for all $n \in \mathbb{N}$,

$$0 < \underline{H}_{t+1}(m_1^{\phi(n)}, ..., m_{t+1}^{\phi(n)}) < 2c_{\phi(n)}^{1/2} \quad \text{and} \quad |m_{t+1}^{\phi(n)} - \underline{q}_t(m_{\to t}^{\phi(n)})| > \epsilon.$$

The set $\mathcal{M}_T$ is compact, thus there exists a further subsequence $(m^{\phi'(n)})$ converging to some $m^0$. By continuity of the functions $\underline{q}_t$ and $\underline{H}_{t+1}$, we have $\underline{H}_{t+1}(m_1^0, ..., m_{t+1}^0) = 0$

and $|m_{t+1}^0 - \underline{q}_t(m_{\to t}^0)| \geq \epsilon > 0$. But this contradicts Proposition 2. The same result holds for $\overline{H}_{t+1}$.

Define $C'$ a Lipschitz constant valid for both $\overline{H}_t$ and $\underline{H}_t$ for all $t \leq T$. Take $\epsilon > 0$, $N_1$ larger than the corresponding $N_0$ and such that $n > N_1$ implies

$$\forall t \leq T, \ \|m_{\to t} - m'_{\to t}\| \leq \delta_n \Rightarrow \|q_t|(m_{\to t}) - q_t(m'_{\to t})| \leq \epsilon,$$

$$\delta_n \leq \epsilon \text{ and } \delta_n \leq c_n^{1/2}/C'.$$

Then for $n \geq N_1$, for all $t \leq T$, if $\widetilde{m}_{\to t} \in \mathcal{M}_t$ and $0 < \underline{H}_{t+1}(\widetilde{m}_1, ..., \widetilde{m}_t, \widetilde{m}_{t+1}) < c_n^{1/2}$ then wpao, $0 \leq \underline{H}_{t+1}(m_1, ..., m_{t+1}) \leq c_n^{1/2} + C' \times \delta_n \leq 2c_n^{1/2}$. Thus if $\widehat{I} = t$ and we are in the case $0 < \underline{H}_{t+1}(\widetilde{m}_1, ..., \widetilde{m}_t, \widetilde{m}_{t+1}) < c_n^{1/2}$ then wpao

$$|m_{t+1} - \widehat{m}_{t+1}| = |m_{t+1} - \underline{q}_t(\widetilde{m}_{\to t})| \leq |m_{t+1} - \underline{q}_t(m_{\to t})| + |\underline{q}_t(m_{\to t}) - \underline{q}_t(\widetilde{m}_{\to t})|$$

$$\leq 2\epsilon.$$

The same result holds for $\overline{H}_{t+1}$. We can then proceed by induction, as

$$|m_{t+2} - \widehat{m}_{t+2}| = |m_{t+2} - \underline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1})|$$

$$\leq |\overline{q}_{t+1}(m_{\to t+1}) - \underline{q}_{t+1}(m_{\to t+1})| + |\underline{q}_{t+1}(m_{\to t+1}) - \underline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1})|$$

$$\leq |\overline{q}_{t+1}(m_{\to t+1}) - \overline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1})| + |\underline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1}) - \underline{q}_{t+1}(m_{\to t+1})|$$

$$+ |\underline{q}_{t+1}(m_{\to t+1}) - \underline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1})|$$

where the last inequality follows from $\underline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1}) = \overline{q}_{t+1}(\widetilde{m}_{\to t}, \widehat{m}_{t+1})$. Using recursively the uniform continuity of $\overline{q}_{t'}$ and $\underline{q}_{t'}$ as functions of $m_{\to t}$ over $\mathcal{M}_t$ and properly adjusting recursive choices of the $\epsilon$'s, we then obtain the uniform convergence of $\widehat{m} - m$ to 0.

*Step 3: Consistency of the lower bound*

Let $\widehat{A}(x) := \widehat{\beta}_k \widehat{c}_0(x) \lambda_{T+1}(x, \widehat{\beta}) \widehat{\underline{q}}_T(\widehat{m}(x))$ and $\widehat{B}(x) := \widehat{\beta}_k \widehat{c}_0(x) \lambda_{T+1}(x, \widehat{\beta}) \widehat{\overline{q}}_T(\widehat{m}(x))$. By Equation (9), $\widehat{\underline{\Delta}}$ satisfies

$$\widehat{\underline{\Delta}} = \frac{1}{n} \sum_{i=1}^n U(X_i, S_i, \widehat{\beta}) + \frac{1}{n} \sum_{i=1}^n \min\left(\widehat{A}(X_i), \widehat{B}(X_i)\right). \tag{33}$$

Since $\lambda_t$ is infinitely differentiable for all $t \leq T$, $\text{Supp}(X)$ is compact and $\widehat{\beta}$ is consistent, $x \mapsto \lambda_t(x, \widehat{\beta})$ converges uniformly in probability to $x \mapsto \lambda_t(x, \beta_0)$. The same

holds for $(x, s) \mapsto C_s(x, \widehat{\beta})$ and $(x, s) \mapsto \exp(s X_T' \widehat{\beta})$. Because wpao $C_s(x, \widehat{\beta}) > \underline{C}$ for all $x \in \text{Supp}(X)$ and $s \leq T$, $(x, s) \mapsto U(x, s, \widehat{\beta})$ converges uniformly in probability to $(x, s) \mapsto U(x, s, \beta_0)$. Then, by the triangle inequality and the law of large numbers (LLN),

$$\frac{1}{n} \sum_{i=1}^{n} U(X_i, S_i, \widehat{\beta}) \xrightarrow{P} E\left(U(X, S, \beta_0)\right).$$

Next, let us show the convergence in probability of the second term in (33). The functions $\overline{q}_T$ and $\underline{q}_T$ are continuous and thus uniformly continuous over the compact set $\mathcal{M}_T$. Then, by Step 2 and since by construction $(m(x), \widehat{m}(x)) \in \mathcal{M}_T^2$, $x \mapsto \underline{q}_T(\widehat{m}(x))$ and $x \mapsto \overline{q}_T(\widehat{m}(x))$ converge uniformly in probability to $x \mapsto \underline{q}_T(m(x))$ and $x \mapsto \overline{q}_T(m(x))$ respectively. Thus, the functions $\widehat{A}$ and $\widehat{B}$ converge uniformly in probability to their corresponding limits, which we write $A$ and $B$. Since $\min(A, B) = (A + B - |A - B|)/2$, $x \mapsto \min(\widehat{A}(x), \widehat{B}(x))$ also converges uniformly in probability to $A$ and $B$. Then, by the triangle inequality and the LLN,

$$\frac{1}{n} \sum_{i=1}^{n} \min\left(\widehat{A}(X_i), \widehat{B}(X_i)\right) \xrightarrow{P} E\left(\min\left(A(X), B(X)\right)\right).$$

The result follows.


## A.2   Theorem 2

*Part 1: asymptotic approximation and normality when $\beta_{0k} \neq 0$.*

We show the linear approximation here. The convergence in distribution then follows directly from the central limit theorem (CLT). Also, we focus on $\widehat{\underline{\Delta}}$: the proofs for $\widehat{\overline{\Delta}}$ is similar. We prove the result in three steps. First, we show that wpao, $\widehat{I}(X_i) = I$ for all $i$. Second, we prove that

$$\widehat{\underline{\Delta}} = \frac{1}{n} \sum_{i=1}^{n} \underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) \mathbb{1}\left\{\widehat{\beta}_k \geq 0\right\} + \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) \mathbb{1}\left\{\widehat{\beta}_k < 0\right\} + o_P(n^{-1/2}). \quad (34)$$

These two steps are valid whatever the sign of $\beta_{0k}$. Finally, we show the result in the third step, assuming that $\beta_{0k} > 0$; the proof when $\beta_{0k} < 0$ follows similarly.

**Step 1: wpao, $\widehat{I}(X_i) = I$ for all $i$.**

First, let $T_t(m) := \underline{H}_t(m)\overline{H}_t(m)$. By definition of $\widehat{I}(x)$ and $I$,

$$\widehat{I}(x) > I \;\Rightarrow\; T_{I+1}(m(x)) = 0 \text{ and } T_{I+1}(\widetilde{m}(x)) > c_n.$$

Moreover,

$$T_{I+1}(\widetilde{m}(x)) > c_n \Rightarrow \underline{H}_{I+1}(\widetilde{m}(x))\overline{H}_{I+1}(\widetilde{m}(x)) - \underline{H}_{I+1}(m(x))\overline{H}_{I+1}(m(x)) > c_n,$$

The functions $\underline{H}_{I+1}$ and $\overline{H}_{I+1}$ are infinitely differentiable on the compact set $\mathcal{M}_{I+1}$. The product of these functions is thus Lipschitz on this set. By induction, $(\widetilde{m}_0(x), ..., \widetilde{m}_{I+1}(x))$ lies in $\mathcal{M}_{I+1}$. Indeed, otherwise we would not have $\widehat{I}(x) > I + 1$. This implies that for any given value $x \in \mathrm{Supp}(X)$, wpao

$$T_{I+1}(\widetilde{m}(x)) > c_n \Rightarrow \|\widetilde{m}(x) - m(x)\| > Cc_n.$$

Because $\delta_n/c_n \to 0$, this cannot occur for any $x \in \mathrm{Supp}(X)$, wpao. Hence, wpao, $\widehat{I}(X_i) \leq I$ for all $i$.

Now, assume that $\widehat{I}(x) < I$ for some $x \in \mathrm{Supp}(X)$. Then,

$$\exists k = \widehat{I}(x) < I, \; \underline{H}_k(\widetilde{m}(x))\overline{H}_k(\widetilde{m}(x)) \leq c_n \text{ and } \underline{H}_k(m(x))\overline{H}_k(m(x)) > 0.$$

We know $m_{\to k}(x) \in \mathcal{M}_k^\epsilon$ for any $k \leq t$. Thus if $\underline{H}_k(m(x))\overline{H}_k(m(x)) > 2c_n$ then

$$\underline{H}_k(\widetilde{m}(x))\overline{H}_k(\widetilde{m}(x)) > 2c_n - |\underline{H}_k(\widetilde{m}(x))\overline{H}_k(\widetilde{m}(x)) - \underline{H}_k(m(x))\overline{H}_k(m(x))|$$
$$> 2c_n - C\|m - \widetilde{m}\|_\infty,$$

using again the Lipschitz property of the product $\underline{H}_k\overline{H}_k$ on $\mathcal{M}_k$. By $\delta_n/c_n \to 0$ and by $\|\widetilde{m} - m\|_\infty = O_P(\delta_n)$, wpao and for $n$ large enough we have

$$\exists N_0, \; \forall x, \; n \geq N_0, \; \underline{H}_k(m(x))\overline{H}_k(m(x)) > 2c_n \;\Rightarrow\; \underline{H}_k(\widetilde{m}(x))\overline{H}_k(\widetilde{m}(x)) > c_n,$$

or alternatively

$$\underline{H}_k(\widetilde{m}(x))\overline{H}_k(\widetilde{m}(x)) \leq c_n \Rightarrow \underline{H}_k(m(x))\overline{H}_k(m(x)) \leq 2c_n$$

while $k < I$. Thus, wpao, for $n \geq N_0$

$$\widehat{I}(x) < I \;\Rightarrow\; \exists k = \widehat{I}(x) < I, \; \underline{H}_k(m(x))\overline{H}_k(m(x)) \leq 2c_n.$$

But since $m_{\to k}(x) \in \mathcal{M}_k^\epsilon$ for all $x \in \mathrm{Supp}(X)$, $\underline{H}_k\overline{H}_k$ is a continuous function and $\mathcal{M}_k$ is a compact set, we know that there exists $\epsilon'$ such that for all $x$, $\underline{H}_k(m(x))\overline{H}_k(m(x)) > \epsilon'$ is strictly positive. This makes it impossible to have for $n \geq N_0$, $\underline{H}_k(m(x))\overline{H}_k(m(x)) \leq 2c_n$ for any $x$. Thus we get $\{i \in \{1, ..., n\} : \widehat{I}(X_i) < I\} = \emptyset$ wpao.

In conclusion, wpao, we have $\widehat{I}(X_i) = I$ for all $i \in \{1, ..., n\}$.

**Step 2:** (34) **holds.**

$P\Big(\forall i \in \{1, ..., n\}, \ \underline{q}_T(\widehat{m}(X_i)) = \underline{q}_T(\widehat{\gamma}, X_i, \widehat{\beta})\Big) \to 1$ as $n \to \infty$. This in turn implies wpao

$$\widehat{\underline{\Delta}} = \frac{1}{n} \sum_{i=1}^{n} U(X_i, S_i, \widehat{\beta}) + \widehat{\beta}_k \widehat{c}_0(X_i) \lambda_{T+1}(X_i, \widehat{\beta}) \Big[ \underline{q}_T(\widehat{\gamma}, X_i, \widehat{\beta}) \mathbb{1} \left\{ \widehat{\beta}_k \lambda_{T+1}(X_i, \widehat{\beta}) \geq 0 \right\}$$

$$+ \overline{q}_T(\widehat{\gamma}, X_i, \widehat{\beta}) \mathbb{1} \left\{ \widehat{\beta}_k \lambda_{T+1}(X_i, \widehat{\beta}) < 0 \right\} \Big]. \tag{35}$$

To obtain (34), we define the set $\mathcal{V}_0 = \{ x \in \mathrm{Supp}(X) \,|\, \lambda_{T+1}(x, \beta_0) \geq 0 \}$ and $J_n :=$ $\widehat{\underline{\Delta}} - \frac{1}{n} \sum_{i=1}^{n} \underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) \mathbb{1} \left\{ \widehat{\beta}_k \geq 0 \right\} + \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) \mathbb{1} \left\{ \widehat{\beta}_k < 0 \right\}$. Note first that wpao

$$J_n = \frac{1}{n} \sum_{i=1}^{n} \widehat{\beta}_k \widehat{c}_0(X_i) \lambda_{T+1}(X_i, \widehat{\beta}) \Big[ \overline{q}_T(\widehat{m}(X_i)) \left( \mathbb{1} \left\{ \lambda_{T+1}(X_i, \widehat{\beta}) \geq 0 \right\} - \mathbb{1} \left\{ X_i \in \mathcal{V}_0 \right\} \right)$$

$$+ \underline{q}_T(\widehat{m}(X_i)) \left( \mathbb{1} \left\{ \lambda_{T+1}(X_i, \widehat{\beta}) < 0 \right\} - \mathbb{1} \left\{ X_i \in \mathcal{V}_0^c \right\} \right) \Big]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \widehat{\beta}_k \widehat{c}_0(X_i) \lambda_{T+1}(X_i, \widehat{\beta}) \Big[ \underline{q}_T(\widehat{\gamma}, X_i, \widehat{\beta}) - \overline{q}_T(\widehat{\gamma}, X_i, \widehat{\beta}) \Big]$$

$$\Big[ \mathbb{1} \left\{ \lambda_{T+1}(X_i, \widehat{\beta}) < 0 \right\} - \mathbb{1} \left\{ X_i \in \mathcal{V}_0^c \right\} \Big],$$

and we denote the right-hand side of the second equality as $I_n$. We prove now that $I_n = o_P(n^{-1/2})$ which will guarantee that (34) holds.

By definition, $\lambda_{T+1}(x, \beta) = -\Pi_{t=1}^{T-1} \left( e^{(x_t - x_T)'\beta} - 1 \right) = -e^{(T-1)x_T'\beta} \Pi_{t=1}^{T-1} \left( e^{x_t'\beta} - e^{x_T'\beta} \right)$. Define the random variables $V_{it} := e^{X_{it}'\beta_0}$ and $\widehat{V}_{it} = e^{X_{it}'\widehat{\beta}}$ for $i \leq n$ and $t \leq T$. Then $\lambda_{T+1}(X_i, \beta_0) = -V_{it}^{T-1} \Pi_{t=1}^{T-1} (V_{it} - V_{iT})$. The same equality holds replacing variables with their estimators. Define $L(X_i, \beta_0) = \Pi_{t=1}^{T-1} (V_{it} - V_{iT})$. Using previous results and Proposition 3, wpao,

$$I_n \leq C \frac{1}{n} \sum_{i=1}^{n} \left| L(X_i, \widehat{\beta}) \right| \left| \mathbb{1} \left\{ \lambda_{T+1}(X_i, \widehat{\beta}) < 0 \right\} - \mathbb{1} \left\{ X_i \in \mathcal{V}_0^c \right\} \right|.$$

Note that wpao,

$$\left| \mathbb{1} \left\{ \lambda_{T+1}(X_i, \widehat{\beta}) < 0 \right\} - \mathbb{1} \left\{ X_i \in \mathcal{V}_0^c \right\} \right|$$

$$\leq \mathbb{1} \left\{ \exists t < T : V_{it} - V_{iT} < 0 < \widehat{V}_{it} - \widehat{V}_{iT} \ \text{ou} \ V_{it} - V_{iT} > 0 > \widehat{V}_{it} - \widehat{V}_{iT} \right\}.$$

Moreover $\left| \widehat{V}_{it} - \widehat{V}_{iT} - (V_{it} - V_{iT}) \right| \leq |\widehat{V}_{it} - V_{it}| + |\widehat{V}_{iT} - V_{iT}|$. We use $|\exp(a) - \exp(b)| \leq \exp(b + |b - a|)|b - a|$, $\|X_{it}\| \leq C$ and the Cauchy-Schwarz inequality to obtain

$$|\widehat{V}_{it} - V_{it}| \leq C \|\widehat{\beta} - \beta_0\| \exp(C \|\widehat{\beta} - \beta_0\|). \tag{36}$$

51

Take $(r_n)_n$ a sequence such that $r_n \to \infty$ and $r_n = o(n^{1/4})$. The previous inequalities give

$$
\begin{aligned}
& |\, \mathbb{1}\left\{\lambda_{T+1}(X_i, \widehat{\beta}) < 0\right\} - \mathbb{1}\left\{x \in \mathcal{V}_{0k}^c\right\}\, | \\
& \quad \leq\ \mathbb{1}\left\{\sqrt{n}\|\widehat{\beta} - \beta_0\| \leq r_n\right\} \mathbb{1}\left\{\exists t : |V_{it} - V_{iT}| < 2C(r_n/\sqrt{n}) \exp(C(r_n/\sqrt{n}))\right\} \\
& \quad +\ \mathbb{1}\left\{\sqrt{n}\|\widehat{\beta} - \beta_0\| > r_n\right\}. \tag{37}
\end{aligned}
$$

Write $u_n = 2C(r_n/\sqrt{n}) \exp(Cr_n/\sqrt{n})$ where only here $C$ is fixed to be the constant in the previous inequality. Assume $\sqrt{n}\|\widehat{\beta} - \beta_0\| \leq r_n$ and $|V_{it^*} - V_{iT}| < u_n$ for some $t^*$. Then

$$
|\widehat{V}_{it^*} - \widehat{V}_{iT}| \leq 2u_n, \text{ and } |\widehat{V}_{it} - \widehat{V}_{iT}| \leq C + u_n.
$$

Thus we have

$$
\sqrt{n} I_n \leq \frac{C}{n^{1/2}} \sum_{i=1}^n I_{ni} + C\sqrt{n}\, \mathbb{1}\left\{\sqrt{n}\|\widehat{\beta} - \beta_0\| > r_n\right\}
$$

with $I_{ni} = \left|L(X_i, \widehat{\beta})\right| \mathbb{1}\{\exists t : |V_{it} - V_{iT}| \leq u_n\}$. We imposed $r_n \to \infty$, the second term is thus $o_P(1)$. Note that $\left|L(X_i, \widehat{\beta})\right| \leq C$. The second term will thus be an $o_P(1)$ as well if

$$
E\left[\frac{u_n}{n^{1/2}} \sum_{i=1}^n I_{ni}\right] \to 0. \tag{38}
$$

Note now that under Assumptions 3 and 4, if $\beta_0 \neq 0$,

$$
\exists u_0,\ u \leq u_0 \Rightarrow P(\exists t : |(X_{it} - X_{iT})'\beta_0| \leq u) \leq C'u.
$$

This is a consequence of $\text{Supp}(X)$ being compact and $f_X$ being absolutely continuous with respect to the Lebesgue measure on $\mathbb{R}^{pT}$. Since $u_n \to 0$, this implies that $E(I_{ni}) \leq C'u_n$ and

$$
E\left[\frac{u_n}{n^{1/2}} \sum_{i=1}^n I_{ni}\right] \lesssim \frac{r_n^2}{\sqrt{n}} = o(1),
$$

since $r_n = o(n^{1/4})$. This proves (38) and $I_n = o_P(n^{-1/2})$.

Finally, consider the case $\beta_0 = 0$. Using (36) and $\widehat{\beta} \in \Theta$ compact, we get

$$
I_{ni} \leq \left|L(X_i, \widehat{\beta})\right| \leq C\|\widehat{\beta}\|^{T-1}.
$$

By Lemma 8 (with $\mathcal{P}$ defined with $\underline{\sigma} = 0$ and some appropriate $\overline{M}$ and $A$) and the fact that $\widehat{\beta}$ is bounded imply that $n^{1/2}E\left[\|\widehat{\beta}\|^{T-1}\right]$ is bounded. Thus, $n^{1/2}E(I_{ni})$ is bounded, which implies that (38) holds. Thus, in all cases, (34) holds.

**Step 3: conclusion**

Define $H_n(\gamma, \beta) := \frac{1}{n}\sum_{i=1}^{n}\underline{h}(X_i, S_i, \gamma, \beta)$ and $H(\gamma, \beta) := E(\underline{h}(X, S, \gamma, \beta))$, so that $H(\gamma_0, \beta_0) = \underline{\Delta}$. Moreover, by Lemma 8, $\widehat{\beta}_k \xrightarrow{P} \beta_{0k} > 0$. Thus, $\widehat{\beta}_k \geq 0$ wpao. Then, by the previous step,

$$\widehat{\underline{\Delta}} = H_n(\widehat{\gamma}, \widehat{\beta}) + o_P(n^{-1/2}).$$

Now, $H_n(\widehat{\gamma}, \widehat{\beta})$ is a semiparametric estimator with a nonparametric first step. We then show the result by applying Chen et al. (2003). To this end, let $\alpha := \lceil pT/2 \rceil$ and following Chen et al. (2003), let us define, for any function $\gamma$ from $\mathrm{Supp}(X)$ to $\mathbb{R}^{T+1}$ admitting at least $\alpha$ derivatives,

$$\|\gamma\|_{\mathcal{G}} := \max_{|a| \leq \alpha} \|D^a \gamma\|_{\infty}.$$

For any $c > 0$, we let $\mathcal{C}_c^{\alpha}$ denote the set of functions $\gamma$ admitting at least $\alpha$ derivatives and such that $\|\gamma\|_{\mathcal{G}} \leq c$. By Assumption 5, there exists $C$ such that $\gamma_0 \in \mathcal{C}_C^{\alpha}$. Hereafter, we let $\mathcal{G} := \mathcal{C}_{C'}^{\alpha}$ for some $C' > C$. We prove in Lemma 7 below the five following conditions:

1. Condition 1: for all $(\epsilon_n)_{n \geq 1}$ such that $\epsilon_n \to 0$,

$$\sup_{\substack{\|\beta - \beta_0\| \leq \epsilon_n, \\ \|\gamma - \gamma_0\|_{\mathcal{G}} \leq \epsilon_n}} |\, [H_n(\gamma, \beta) - H(\gamma, \beta)] - [H_n(\gamma_0, \beta_0) - H(\gamma_0, \beta_0)]\,| = o_P(n^{-1/2}).$$

2. Condition 2: The functional pathwise and ordinary derivatives of $\underline{h}$ with respect to $\gamma$ and $\beta$ exist. Moreover, there exists $b(\cdot)$ such that $E(b(X_i)) < \infty$ and

$$\begin{aligned} &|\underline{h}(X_i, S_i, \gamma, \beta) - \underline{h}(X_i, S_i, \gamma_0, \beta_0) \\ &\quad - D_\gamma \underline{h}(X_i, S_i, \gamma_0, \beta_0)'[\gamma(X_i) - \gamma_0(X_i)] - D_\beta \underline{h}(X_i, S_i, \gamma_0, \beta_0)[\beta - \beta_0]| \\ &\leq b(X_i)\left(\|\gamma - \gamma_0\|_{\infty}^2 + |\beta - \beta_0|^2\right). \end{aligned}$$

3. Condition 3: We have $\sqrt{n}(\widehat{\beta} - \beta_0) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_i + o_P(1)$.

4. Condition 4: We have $\widehat{\gamma} \in \mathcal{G}$ wpao, $\|\widehat{\gamma} - \gamma_0\|_{\infty} = o_P(n^{-1/4})$ and $\|\widehat{\gamma} - \gamma_0\|_{\mathcal{G}} = O_P(\tilde{\epsilon}_n)$ for some $\tilde{\epsilon}_n \to 0$.

5. Condition 5: Holding fixed the nonparametric estimator $\widehat{\gamma}$ in the expectation,

$$\sqrt{n}E^*\Big(D_\gamma\underline{h}(X_i, S_i, \gamma_0, \beta_0)[\widehat{\gamma}(X_i) - \gamma_0(X_i)]\Big)$$
$$= \frac{1}{\sqrt{n}}\sum_{i=1}^n D_\gamma\underline{h}(X_i, S_i, \gamma_0, \beta_0)' \, [Z_i - E(Z_i|X_i)] + o_P(1).$$

Then, we have

$$\sqrt{n}[\widehat{\underline{\Delta}} - \underline{\Delta}] = \sqrt{n}[H_n(\widehat{\gamma}, \widehat{\beta}) - H(\gamma_0, \beta_0)]$$
$$= \sqrt{n}[H_n(\gamma_0, \beta_0) - H(\gamma_0, \beta_0)] + \sqrt{n}[H(\widehat{\gamma}, \widehat{\beta}) - H(\gamma_0, \beta_0)]$$
$$+ \sqrt{n}[H_n(\widehat{\gamma}, \widehat{\beta}) - H(\widehat{\gamma}, \widehat{\beta})] - \sqrt{n}[H_n(\gamma_0, \beta_0) - H(\gamma_0, \beta_0)]$$
$$= \frac{1}{\sqrt{n}}\sum_{i=1}^n \Bigg[\underline{h}(X_i, S_i, \gamma_0, \beta_0) - E(\underline{h}(X, S, \gamma_0, \beta_0)) + E\left[D_\beta\underline{h}(X, S, \gamma_0, \beta_0)\right]'$$
$$\times \phi_i + E\left[D_\gamma\underline{h}(X_i, S_i, \gamma_0, \beta_0)|X_i\right]' [Z_i - E(Z_i|X_i)]\Bigg] + o_P(1).$$

The result follows using $E\left[D_\gamma\underline{h}(X_i, S_i, \gamma_0, \beta_0)|X_i\right] = D_\gamma\underline{h}(X_i, S_i, \gamma_0, \beta_0)$ and the definition of $\underline{\psi}_k$.

*Part 2: case $\beta_{0k} = 0$.*

First, let us define

$$\underline{Z}_n = \frac{1}{n^{1/2}}\sum_{i=1}^n \underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}), \quad \overline{Z}_n = \frac{1}{n^{1/2}}\sum_{i=1}^n \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}).$$

Remark that when $\beta_{0k} = 0$, $\underline{\Delta} = \overline{\Delta} = 0$. Then, by Step 2 above (which holds regardless of the value of $\beta_{0k}$) and remarking that $\underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) \leq \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta})$ if and only if $\widehat{\beta}_k \geq 0$, we get

$$\sqrt{n}\left(\widehat{\underline{\Delta}} - \underline{\Delta}, \widehat{\overline{\Delta}} - \overline{\Delta}\right) = \left(\min\left(\underline{Z}_n, \overline{Z}_n\right), \max\left(\underline{Z}_n, \overline{Z}_n\right)\right) + o_P(1).$$

Now, the proof of asymptotic linearity of $\underline{Z}_n$ in Part 1 above also applies when $\beta_{0k} = 0$. Thus, $\left(\underline{Z}_n, \overline{Z}_n\right) \xrightarrow{d} \left(\underline{Z}, \overline{Z}\right)'$. The result follows by the continuous mapping theorem.

*Part 3: Consistency of $\widehat{\Sigma}$.*

Let us assume wlog that $\beta_{0k} > 0$. The estimator of the variance covariance matrix is $\widehat{\Sigma} = \frac{1}{n}\sum_{i=1}^n (\widehat{\underline{\psi}}_i, \widehat{\overline{\psi}}_i)(\widehat{\underline{\psi}}_i, \widehat{\overline{\psi}}_i)'$ where

$$\widehat{\underline{\psi}}_i = \underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) - \frac{1}{n}\sum_{j=1}^n \underline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta}) + \left(\frac{1}{n}\sum_{j=1}^n D_\beta\underline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta})\right)' \widehat{\phi}_i$$

54

$$+ D_\gamma \underline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta})'[Z_i - \widehat{\gamma}(X_i)]$$

$$\widehat{\overline{\psi}}_i = \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta}) - \frac{1}{n}\sum_{j=1}^n \overline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta}) + \left(\frac{1}{n}\sum_{j=1}^n D_\beta \overline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta})\right)' \widehat{\phi}_i$$

$$+ D_\gamma \overline{h}(X_i, S_i, \widehat{\gamma}, \widehat{\beta})'[Z_i - \widehat{\gamma}(X_i)]$$

We show that the functions of $(X_i, S_i)$ appearing in $\widehat{\underline{\psi}}_i$ and $\widehat{\overline{\psi}}_i$ converge uniformly to their pointwise limits. Similarly to what we argued in the proof of Theorem 1, $(x, s) \mapsto \underline{h}(x, s, \widehat{\gamma}, \widehat{\beta}) - \frac{1}{n}\sum_{j=1}^n \underline{h}(X_j, S_j, \widehat{\gamma}, \widehat{\beta})$ converges uniformly in probability to $(x, s) \mapsto \underline{h}(x, s, \gamma_0, \beta_0) - E(\underline{h}(X, S, \gamma_0, \beta_0))$.

The smoothness arguments given in Condition 1 (Part 2) implies in particular that the derivatives of $\underline{h}$ with respect to both the vector $\gamma(x)$ and $\beta$ are Lipschitz continuous on $\mathcal{C}_C^\alpha$ and $\Theta$, with Lipschitz constant uniform over $x \in \text{Supp}(X)$. This implies that $(x, s, z) \mapsto D_\gamma \underline{h}(x, s, \widehat{\gamma}, \widehat{\beta})'[z - \widehat{\gamma}(x)]$ converges uniformly in probability to $(x, s, z) \mapsto D_\gamma \underline{h}(x, s, \gamma_0, \beta_0)'[z - \gamma_0(x)]$. The same results follow for $\overline{h}$. By $\mathcal{I}_0$ nonsingular and $C_s(x, \beta)$ bounded away from 0 uniformly over $(s, x, \beta)$, the derivatives of $\beta \mapsto \left[\frac{1}{\sqrt{n}}\sum_{j=1}^n \frac{\partial^2 \ell_c}{\partial \beta^2}(Y_j|X_j; \beta)\right]^{-1} \frac{\partial \ell_c}{\partial \beta}(y|x; \beta)$ are uniformly bounded over $(y, x, \beta)$ wpao. Thus $\widehat{\phi}_i$ converges uniformly in probability to $\phi_i$.

In conclusion, the functions of $(X_i, S_i)$ appearing in $\widehat{\underline{\psi}}_i$ and $\widehat{\overline{\psi}}_i$ converge uniformly to their pointwise limits. This implies that $(\widehat{\underline{\psi}}_i, \widehat{\overline{\psi}}_i)(\widehat{\underline{\psi}}_i, \widehat{\overline{\psi}}_i)'$ converges uniformly to $(\underline{\psi}_i, \overline{\psi}_i)(\underline{\psi}_i, \overline{\psi}_i)'$. As in Theorem 1, we obtain using the LLN that $\widehat{\Sigma} \xrightarrow{P} \Sigma$.

### A.3  Proposition 4

First assume that $\beta_{0k} = 0$. Then $\Delta = 0$ and

$$P\left(\Delta \in \text{CI}_{1-\alpha}^1\right) \geq P(\varphi_\alpha = 0) \rightarrow 1 - \alpha,$$

where the latter follows since $\varphi_\alpha$ has asymptotic level $\alpha$. Now, assume $\beta_{0k} \neq 0$. Then $\varphi_\alpha \xrightarrow{P} 1$, so that $\text{CI}_{1-\alpha}^1$ takes the first form wpao. Suppose first that $\underline{\Delta} < \overline{\Delta}$. By consistency of the bounds, consistency of $\widehat{\Sigma}$ and $\min(\Sigma_{11}, \Sigma 12) > 0$, we have

$$\frac{n^{1/2}\left(\widehat{\overline{\Delta}} - \widehat{\underline{\Delta}}\right)}{\max\left(\widehat{\Sigma}_{11}^{1/2}, \widehat{\Sigma}_{12}^{1/2}\right)} \xrightarrow{P} \infty.$$

55

Then, by Lemma 5.10 of van der Vaart (2000), $c_\alpha \to \Phi^{-1}(1-\alpha)$. The result follows as in Lemma 2 of Imbens and Manski (2004). Next, assume $\underline{\Delta} = \overline{\Delta}$. Then, because $\widehat{\overline{\Delta}} \geq \widehat{\underline{\Delta}}$ a.s., $\overline{Z} - \underline{Z}$ must be degenerate, implying in turn $\overline{Z} = \underline{Z}$ a.s. Hence,

$$\frac{n^{1/2}\left(\widehat{\overline{\Delta}} - \widehat{\underline{\Delta}}\right)}{\max\left(\widehat{\Sigma}_{11}^{1/2}, \widehat{\Sigma}_{12}^{1/2}\right)} = o_P(1).$$

By, again, Lemma 5.10 of van der Vaart (2000), $c_\alpha \to \Phi^{-1}(1-\alpha/2)$. The result follows using standard arguments for this point identified case.


## A.4 Lemma 2

First, let $U_i = (X_i, S_i)$ and

$$g(U_i, \beta) = \beta_k \sum_{t=0}^{S_i} \frac{a_t(X_i, \beta)\binom{T-t}{S_i-t}\exp(S_i X'_{iT}\beta)}{C_{S_i}(X_i, \beta)},$$

so that $\tilde{\Delta} = E[g(U_1, \beta_0)]$ and $\widehat{\Delta} = \sum_{i=1}^{n} g(U_i, \widehat{\beta})/n$. By choosing appropriate $\overline{M}$, $\underline{\sigma}$ and $A$, $P \in \mathcal{P}$. Then, by Lemma 8, we have

$$\sqrt{n}\left(\widehat{\beta} - \beta_0\right) = \frac{1}{n^{1/2}}\sum_{i=1}^{n} \phi_i + o_P(1). \tag{39}$$

Since $\beta \mapsto g(u, \beta)$ is differentiable for all $u$, by the mean value theorem, there exists $\overline{\beta}_i = t_i\widehat{\beta} + (1-t_i)\beta_0$, with $t_i \in [0, 1]$, such that $g(U_i, \widehat{\beta}) - g(U_i, \beta_0) = \partial g/\partial\beta(U_i, \overline{\beta}_i)(\widehat{\beta} - \beta_0)$. Let

$$\widehat{G} = \frac{1}{n}\sum_{i=1}^{n} \frac{\partial g}{\partial\beta}(U_i, \overline{\beta}_i).$$

Then, by what precedes,

$$\sqrt{n}\left(\widehat{\Delta} - \tilde{\Delta}\right) = \frac{1}{n^{1/2}}\sum_{i=1}^{n}\left[\widehat{G}\phi_i + g(U_i, \beta_0) - \tilde{\Delta}\right] + o_P(1).$$

Because $\partial g/\partial\beta$ is continuous, $\mathrm{Supp}(U)$ is compact and $\widehat{\beta}$ is consistent, we have (see, e.g., Lemma 2.4 in Newey and McFadden, 1994)

$$\widehat{G} \xrightarrow{P} G := E\left[\frac{\partial g}{\partial\beta}(U_i, \beta_0)\right].$$

56

Thus,

$$\sqrt{n}\left(\widehat{\Delta} - \tilde{\Delta}\right) = \frac{1}{n^{1/2}}\sum_{i=1}^{n}\left[G\phi_i + g(U_i, \beta_0) - \tilde{\Delta}\right] + o_P(1).$$

The first result follows by the central limit theorem and a few algebra. The second result follows using the same reasoning as to prove $\widehat{G} \xrightarrow{P} G$.

## A.5 Theorem 3

1. Let $Z$ be a standard normal variable independent of the data. Let $\widehat{b}_n := n^{1/2}(\tilde{\Delta} - \Delta)/\widehat{\sigma}$, $\overline{\widehat{b}}_n := n^{1/2}\overline{b}/\widehat{\sigma}$ and $A_n = \{|\widehat{b}_n| \le \overline{\widehat{b}}_n\}$. By independence between $Z$ and $(\overline{\widehat{b}}_n, \widehat{b}_n)$ on the one hand, and $q_\alpha(\widehat{b}_n) \le q_\alpha(\overline{\widehat{b}}_n)$ under $A_n$ on the other hand, we have

$$P\left(|Z - \widehat{b}_n| \le q_\alpha(\overline{\widehat{b}}_n)|\overline{\widehat{b}}_n, \widehat{b}_n\right)\mathbb{1}_{A_n} \ge (1 - \alpha)\mathbb{1}_{A_n}.$$

As a result,

$$P\left(|Z - \widehat{b}_n| \le q_\alpha(\overline{\widehat{b}}_n)\right) \ge E\left[P\left(|Z - \widehat{b}_n| \le q_\alpha(\overline{\widehat{b}}_n)|\overline{\widehat{b}}_n, \widehat{b}_n\right)\mathbb{1}_{A_n}\right]$$

$$\ge (1 - \alpha)P(A_n). \tag{40}$$

Now, if $\beta_{0k} = 0$, then, by Lemma 3, $\widehat{b}_n = 0 \le \overline{\widehat{b}}_n$. Thus, $P(A_n) = 1$. If, instead, $|\tilde{\Delta} - \Delta| < \overline{b}$, then $A_n$ holds if and only if

$$n^{1/2}\left(\overline{\widehat{b}} - \overline{b}\right) \ge n^{1/2}\left(\left|\tilde{\Delta} - \Delta\right| - \overline{b}\right) \to -\infty. \tag{41}$$

We now prove that (41) holds wpao. To this end, define

$$f(U_i, \beta) = \frac{\beta_k \times \lambda_T(X_i, \beta)\binom{T}{S_i}\exp(S_i X'_{iT}\beta)}{2 \times 4^T \times C_{S_i}(X_i, \beta)}. \tag{42}$$

The function $f(u, .)$ is $C^1$ and $\partial f/\partial\beta$ is continuous over the compact $\mathrm{Supp}(U) \times \Theta$. Hence, there exists $M > 0$ such that $f(u, .)$ is Lipschitz with coefficient $M$ for all $u \in \mathrm{Supp}(U)$. The same property then holds for $|f(u, .)|$. As a result,

$$n^{1/2}|\overline{\widehat{b}} - \overline{b}| \le Mn^{1/2}|\widehat{\beta} - \beta_0| + \left|\frac{1}{n^{1/2}}\sum_{i=1}^{n}|f(U_i, \beta_0)| - E\left[|f(U_i, \beta_0)|\right]\right|. \tag{43}$$

By the proof of Theorem 2, $n^{1/2}(\widehat{\beta} - \beta_0) = O_P(1)$. The second term on the right-hand side of (43) is also an $O_P(1)$. As a result, $n^{1/2}\left(\overline{\widehat{b}} - \overline{b}\right) = O_P(1)$, and by (41), we have

$P(A_n) \to 1$. Hence, in both cases,

$$\liminf_{n\to\infty} P(|Z - \widehat{b}_n| \le q_\alpha(\widehat{\overline{b}}_n)) \ge 1 - \alpha.$$

Next, let $Z_n' := n^{1/2}(\widehat{\Delta} - \tilde{\Delta})/\widehat{\sigma}$, $\Phi$ be the standard normal cdf and $F_n$ of $Z_n'$. We have

$$\left| P\left(\Delta \in \mathrm{CI}_{1-\alpha}^2\right) - P(|Z - \widehat{b}_n| \le q_\alpha(\widehat{\overline{b}}_n)) \right|$$

$$= \left| P\left(|Z_n' - \widehat{b}_n| \le q_\alpha(\widehat{\overline{b}}_n)\right) - P\left(|Z - \widehat{b}_n| \le q_\alpha(\widehat{\overline{b}}_n)\right) \right|$$

$$\le \sup_{(x,y)\in\mathbb{R}\times\mathbb{R}^+} |P\left(|Z_n' - x| \le y\right) - P\left(|Z - x| \le y\right)|$$

$$= \sup_{(x,y)\in\mathbb{R}\times\mathbb{R}^+} |F_n(x+y) - F_n(x-y) - \Phi(x+y) + \Phi(x-y)|$$

$$\le 2\|F_n - \Phi\|_\infty.$$

Finally, Lemma 2 implies that for all $x$, $F_n(x) \to \Phi(x)$ with $\Phi$ continuous. By Lemma 2.11 in van der Vaart (2000), the convergence is uniform. The result follows.

2. The outline of the proof is the same as above, but (i) we have to make some statements uniform over $P \in \mathcal{P}$; (ii) we also have to account for the possibility that $\tilde{\Delta} - \Delta = \overline{b}$. Let $A_n' = \{|\widehat{b}_n| \le \widehat{\overline{b}}_n + \varepsilon_n/\widehat{\sigma}\}$. Then (40) still holds with $\widehat{\overline{b}}_n$ and $A_n$ respectively replaced by $\widehat{\overline{b}}_n + \varepsilon_n/\widehat{\sigma}$ and $A_n'$. Now, let us show that

$$\liminf_{n\to\infty} \inf_{P\in\mathcal{P}} P(A_n') = 1. \tag{44}$$

By definition of $A_n'$ and because $|\tilde{\Delta} - \Delta| \le \overline{b}$, we have

$$\left\{ n^{1/2}\left(\widehat{\overline{b}} - \overline{b}\right) \ge -\varepsilon_n \right\} \subset A_n'. \tag{45}$$

Next, using (43), (45), and an union bound, we get

$$P(A_n') \ge P\left(n^{1/2}\left(\widehat{\overline{b}} - \overline{b}\right) \ge -\varepsilon_n\right)$$

$$\ge P\left(n^{1/2}\left|\widehat{\overline{b}} - \overline{b}\right| \le \varepsilon_n\right)$$

$$\ge 1 - P\left(n^{1/2}\left|\widehat{\beta}_k - \beta_{0k}\right| > \frac{\varepsilon_n}{2M}\right)$$

$$- P\left(n^{-1/2}\left|\sum_{i=1}^n |f(U_i, \beta_0)| - E[|f(U_i, \beta_0)|]\right| > \frac{\varepsilon_n}{2}\right). \tag{46}$$

By Lemma 8, we have

$$\limsup_{n\to\infty} \sup_{P\in\mathcal{P}} P\left(n^{1/2}\left|\widehat{\beta}_k - \beta_{0k}\right| > \frac{\varepsilon_n}{2M}\right) = 0. \tag{47}$$

Turning to the second term in (46), we have, by Chebyshev's inequality

$$P\left(n^{-1/2}\left|\sum_{i=1}^{n}|f(U_i,\beta_0)| - E[|f(U_i,\beta_0)|]\right| > \frac{\varepsilon_n}{2}\right) \leq \frac{V(|f(U_i,\beta_0)|)}{(\varepsilon_n/2)^2}.$$

An inspection of (42) reveals that there exists $M' > 0$ such that for all $P \in \mathcal{P}$, $|f(U_i,\beta_0)| \leq M'$ a.s. Hence, $\sup_{P\in\mathcal{P}} V(f(U_i,\beta_0)) \leq M'^2$ and thus,

$$\limsup_{n\to\infty} \sup_{P\in\mathcal{P}} P\left(n^{-1/2}\left|\sum_{i=1}^{n}|f(U_i,\beta_0)| - E[|f(U_i,\beta_0)|]\right| \geq \frac{\varepsilon_n}{2}\right) = 0.$$

This, combined with (47), entails (44). As a result,

$$\liminf_{n\to\infty} \inf_{P\in\mathcal{P}} P\left(|Z - \widehat{b}_n| \leq q_\alpha\left(\widehat{\overline{b}}_n + \frac{\varepsilon_n}{\widehat{\sigma}}\right)\right) \geq 1 - \alpha.$$

Next, reasoning as in Point 1 above, we obtain

$$\left|P\left(\Delta \in \mathrm{CI}_{1-\alpha}^3\right) - P\left(|Z - \widehat{b}_n| \leq q_\alpha\left(\widehat{\overline{b}}_n + \frac{\varepsilon_n}{\widehat{\sigma}}\right)\right)\right| \leq \|F_{nP} - \Phi\|_\infty + \|G_{nP} - \Phi\|_\infty,$$

where we index $F_n$ and $G_n$ by $P$ to make the dependence explicit. By Lemma 8, we have, for all $x \in \mathbb{R}$,

$$\sup_{P\in\mathcal{P}} |F_{nP}(x) - \Phi(x)| \to 0, \quad . \tag{48}$$

Reasoning as in the proof of Lemma 2.11 in van der Vaart (2000), this entails $\sup_{P\in\mathcal{P}} \|F_{nP} - \Phi\|_\infty \to 0$ and $\sup_{P\in\mathcal{P}} \|G_{nP} - \Phi\|_\infty \to 0$. The result follows.

## B   Technical lemmas

**Lemma 7** *Suppose that Assumptions 1-5 and 6 hold and $\delta_n/c_n \to 0$. Then, the five conditions in Part 1, Step 3 of the proof of Theorem 2 hold.*

**Proof:** We use the same notation as that introduced in the proof of Theorem 2.

*Condition 1:* By Assumption 6 and uniform consistency of $\widetilde{m}$ (see Step 1 in the proof of Theorem 1), $\widetilde{m}(x) = m(\widehat{\gamma}, x, \widehat{\beta})$ lies in $\mathcal{M}_I^{\epsilon/2}$ for all $x \in \mathrm{Supp}(X)$ wpao, where, for any $\eta > 0$

$$\mathcal{M}_I^\eta := \mathrm{cl}\{m \in \mathcal{M}_I\,;\, \mathcal{B}(m,\eta) \subset \mathrm{Int}\,\mathcal{M}_I\}.$$

It is known that $\underline{q}_T$ and $\overline{q}_T$ are infinitely differentiable on $\mathcal{M}_I^{\epsilon/2}$. The function $(\gamma, x, \beta) \mapsto m(\gamma, x, \beta)$ depends on $\gamma$ only through its value when evaluated at $x$, $\gamma(x)$, and is infinitely differentiable with respect to the vector $\gamma(x)$ and with respect to $\beta$. Moreover, the set $\mathcal{V}_0$ is constructed using the known $\beta_0$ and thus does not depend on $\widehat{\beta}$. The function $\underline{h}$ is therefore infinitely differentiable in the vector $\gamma(x)$ and in $\beta$. It is in particular Fréchet differentiable in $\gamma$ and continuously differentiable in $\beta$ and we have

$$|\underline{h}(X_i, S_i, \gamma_1, \beta_1) - \underline{h}(X_i, S_i, \gamma_2, \beta_2)|$$
$$\leq \sup_{\beta \in \Theta} \|\partial_\beta \underline{h}(X_i, S_i, \gamma_2, \beta)\| \, \|\beta_1 - \beta_2\| + \sup_{\gamma(X_i), \, \gamma \in \mathcal{G}} \|\partial_\gamma \underline{h}(X_i, S_i, \gamma, \beta_2)\| \, \|\gamma_1(X_i) - \gamma_2(X_i)\|$$
$$\leq b(X_i, S_i) \left( \|\beta_1 - \beta_2\| + \|\gamma_1 - \gamma_2\|_{\mathcal{G}} \right),$$

where the suprema of the derivatives exist since $\Theta$ and $\{\gamma(X), \; \gamma \in \mathcal{G}\}$ are compact sets by Assumptions 3 and 5. Note additionally that by similar smoothness arguments and because indicator functions are bounded, $E(b(X_i, S_i)^r) < \infty$ for any $r \geq 2$ as $X$ and $S$ have bounded support. Moreover using the definitions of Chen et al. (2003) and the results they cite from van der Vaart and Wellner (1996), the covering number of $\mathcal{G}$ exists and is integrable if $\alpha > \dim(X)/2 = pT/2$. Thus by Theorem 3 of Chen et al. (2003), Condition 1 holds.

*Condition 2:* The difference

$$|\underline{h}(X_i, S_i, \gamma, \beta) - \underline{h}(X_i, S_i, \gamma_0, \beta_0)|$$
$$- D_\gamma \underline{h}(X_i, S_i, \gamma_0, \beta_0)'[\gamma(X_i) - \gamma_0(X_i)] - D_\beta \underline{h}(X_i, S_i, \gamma_0, \beta_0)'[\beta - \beta_0],$$

is equal to the second-order partial derivatives of $\underline{h}$ evaluated at some point $\tilde{\gamma}(X_i)$ and $\tilde{\beta}$, and applied to $\gamma(X_i) - \gamma_0(X_i)$ and $\beta - \beta_0$. By the same argument as in Condition 1, the second-order derivatives of $h$ can be bounded uniformly over $\beta$ and $\gamma(X)$ and these bounds have finite expectation over $(X_i, S_i)$. The residual can thus be bounded by a constant multiplied by $(\|\gamma - \gamma_0\|_\infty^2 + |\beta - \beta_0|^2)$.

*Condition 3:* This condition holds by Lemma 8 below, with $\mathcal{P} = \{P\}$ and

$$\phi(X_i, Y_i) = E \left[ \frac{\partial^2 \ell_c}{\partial \beta^2} (Y_i|X_i; \beta_0) \right]^{-1} \frac{\partial \ell_c}{\partial \beta} (Y_i|X_i; \beta_0).$$

*Condition 4:* We apply Theorem 6 of Masry (1996) on the convergence rate of local polynomial estimators. This theorem requires the conditional density $f_{X|Z}$ to exist and be bounded, which holds here as $Z = (\mathbb{1}\{S = 0\}, ..., \mathbb{1}\{S = T\})'$ and $X$ has compact support and bounded density. By Assumptions 3-5, the other conditions of the theorem hold. Thus, by Masry (1996)

$$\sup_{x \in \mathcal{D}} |\widehat{\gamma}_j(x) - \gamma_{0j}(x)| = O\left(\left(\frac{\ln n}{nh_n^{pT}}\right)^{1/2} + h_n^{\ell+1}\right) \text{ almost surely.} \qquad (49)$$

By Assumption 5, $\left(\ln n / \left(nh_n^{pT}\right)\right)^{1/2} + h_n^{\ell+1} = o\left(n^{-1/3}\right)$ thus $\|\widehat{\gamma} - \gamma_0\|_\infty = O_P(n^{-1/4})$. Theorem 6 of Masry (1996) also states that almost surely

$$\text{for } |a| \leq \ell, \ \sup_{x \in \mathcal{D}} \left|\frac{\partial^a \widehat{\gamma}_j(x)}{\partial x^a} - \frac{\partial^a \gamma_{0j}(x)}{\partial x^a}\right| = O\left(\left(\frac{\ln n}{nh_n^{pT+2|a|}}\right)^{1/2} + h_n^{\ell+1-|a|}\right). \qquad (50)$$

Define $\tilde{\epsilon}_n := \left[\ln n / (nh_n^{pT+2|\alpha|})\right]^{1/2} + h_n^{\ell+1-|\alpha|}$. Then by $\alpha \leq \ell$, $\|\widehat{\gamma} - \gamma_0\|_{\mathcal{G}} = O_P(\tilde{\epsilon}_n)$ and by $\alpha \leq \ell$, $\alpha \leq pT$ and Assumption 5.2, $\tilde{\epsilon}_n \to 0$. Moreover, note that $\widehat{\gamma}$ is continuous by construction and since $\gamma \in \mathcal{C}_{C'}^\alpha$, Equations (49) and (50) imply that $\widehat{\gamma} \in \mathcal{G}$ w.p.a.1.

*Condition 5:* First, we apply Corollary 1 of Kong et al. (2010). To this end, we check their Assumptions A1-A7. To avoid confusion with the notations of Kong et al. (2010), let $p' = pT$. For $u, e, \theta \in \mathbb{R}^3$, let $\rho(u, \theta) = \frac{1}{2}(u - \theta)^2$ and $\varphi(e) = -e$, we have $\rho(u, \theta) = \rho(u, 0) + \int_0^\theta \varphi(u - t)dt$ and $E(\varphi(\epsilon_j)|X) = 0$ for $\epsilon_j = \mathbb{1}\{S = j\} - \gamma_{0j}(X)$.

First, because $\epsilon_j$ has a bounded support and a bounded density, A1 and A2 in Kong et al. (2010) hold for any value of the parameter $\nu_1$ as defined in Kong et al. (2010).

Assumption 4.3 ensures that for any $\alpha = (\alpha_{j,t}) \in \mathbb{N}^{p'}$ such that $\sum_{j,t} \alpha_{j,t} \leq 2\ell + 1$, $u \mapsto u^\alpha K(u)$ is Lipschitz on any compact set (as a product of Lipschitz functions) and on $\mathbb{R}^{p'} \backslash \text{Supp}(K)$ (as the null function). If $u \in \text{Supp}(K)$ and $v \in \mathbb{R}^{p'} \backslash \text{Supp}(K)$ there exists $w \in \{\mu u + (1 - \mu)v : \mu \in [0; 1]\} \cap (2 \cdot \text{Supp}(K)) \cap \left(\mathbb{R}^{p'} \backslash \text{Supp}(K)\right)$. Because $2 \cdot \text{Supp}(K)$ is a compact containing $\text{Supp}(K)$, we have:

$$|u^\alpha K(u) - v^\alpha K(v)| \leq |u^\alpha K(u) - w^\alpha K(w)| + |w^\alpha K(w) - v^\alpha K(v)|$$
$$\leq C\left(|u - w| + |w - v|\right) = C|u - v|,$$

ensuring that $u \mapsto u^\alpha K(u)$ is Lipschitz on $\mathbb{R}^{p'}$. So, A3 holds.

Assumptions 3.1 and 4.1 (resp. 4.2) imply that A4 (resp. A5) holds.

To check A6, we fix the values $\lambda_1 = 3/4$, $\lambda_2 = 1/2$, $p = \ell$ and take any $\nu_2 > 12$, borrowing here the notation of Kong et al. (2010). Then, tedious algebra shows that under Assumption 5.2, the three conditions on the bandwidth $h_n$ in A6 hold.

Finally, the Bayes formula and Assumptions 1 and 4.1 ensure that $X|S$ admits a bounded density with respect to the Lebesgue measure. By independence across $i = 1, ..., n$, A7 holds.

Hence, by Corollary 1 in Kong et al. (2010), we have with probability 1 and uniformly in $x \in \mathcal{K}$, a compact subset of $\mathbb{R}^{pT}$,

$$\widehat{\gamma}_j(x) - \gamma_{0j}(x) = T_j(x) + O\left[\left(\frac{\log n}{n h_n^{pT}}\right)^{3/4}\right] + o\left(h_n^{\ell+1}\right), \tag{51}$$

where

$$T_j(x) := \alpha(x) h_n^{\ell+1} + \frac{1}{n} e' S_{h_n}(x)^{-1} \sum_{i=1}^{n} K\left(\frac{X_i - x}{h_n}\right) \epsilon_{ij} w\left(X_i - x\right).$$

In this expression, $e$ is the first element of the canonical basis of the corresponding vector space, the dimension of which depends on $\ell$ and $pT$ (it is the number of polynomials of $X$ of degree less than or equal to $\ell$). The function $\alpha(x)$ is a bounded function of $x$, $w(x)$ is the vector $(1, x, ..., x^k, ...)'$ for all $|k| \le \ell$ ordered by increasing degree and $S_{h_n}(x)$ is the matrix $E[w((X - x)/h_n)w((X - x)/h_n)'K((X - x)/h_n)]$.

Under Assumption 5.2, we have $\left(\ln n/(n h_n^{pT})\right)^{3/4} = o(n^{-1/2})$ and $h_n^{\ell+1} = o(n^{-1/2})$. Thus (51) implies

$$E^*\left(D_\gamma \underline{h}(X_i, S_i, \gamma_0, \beta_0)'[\widehat{\gamma}(X) - \gamma_0(X)]\right)$$
$$= \int_{x \in \mathbb{R}^{pT}} [D_\gamma \underline{h}(X_i, S_i, \gamma_0, \beta_0)]' T_j(x) f_X(x) \mathrm{d}x + o_P(n^{-1/2})$$
$$= \frac{1}{n} \sum_{i=1}^{n} \sum_{j \le T+1} \epsilon_{ij} \int_{x \in \mathbb{R}^{pT}} e' S_{h_n}(u)^{-1} w\left(X_i - x\right) \lambda_{h,j}(x) K\left(\frac{X_i - x}{h_n}\right) f_X(x) \mathrm{d}x$$
$$\quad + o_P(n^{-1/2}),$$

where $\lambda_{h,j}(x)$ is the derivative of $\underline{h}(x, s, \gamma_0, \beta_0)$ with respect to the $j^{\text{th}}$ component of $\gamma(x)$, written $D_{\gamma,j}\underline{h}(x, s, \gamma_0, \beta_0)$. Note that this derivative is not a function of $s$, as

$$\lambda_{h,j}(x) = \beta_k D_{\gamma,j} c_0(\gamma_0, x, \beta_0) \lambda_{T+1}(x, \beta_0) \left[\underline{q}_T(\gamma_0, x, \beta_0) \mathbb{1}\left\{\lambda_{T+1}(x, \beta_0) > 0\right\}\right.$$

$$+ \bar{q}_T(\gamma_0, x, \beta_0) \mathbb{1} \left\{ \lambda_{T+1}(x, \beta_0) < 0 \right\} \Big]$$

$$+ \beta_k c_0(\gamma_0, x, \beta_0) \lambda_{T+1}(x, \beta_0) \Big[ D_{\gamma,j} \underline{q}_T(\gamma_0, x, \beta_0) \mathbb{1} \left\{ \lambda_{T+1}(x, \beta_0) > 0 \right\}$$

$$+ D_{\gamma,j} \bar{q}_T(\gamma_0, x, \beta_0) \mathbb{1} \left\{ \lambda_{T+1}(x, \beta_0) < 0 \right\} \Big]. \quad (52)$$

Also, $\lambda_{h,j}(x)$ is a continuous function of $x$. Let $I_{i,j}$ denote the integral in the display above. After a change of variable, $I_{i,j}$ is equal to

$$I_{i,j} = h_n^{pT} \int_u e' S_{h_n}(X_i - h_n u)^{-1} w(h_n u) \, \lambda_{h,j}(X_i - h_n u) K(u) \, f_X(X_i - h_n u) \mathrm{d}u.$$

Assumptions A3-A6 of Kong et al. (2010) hold. Thus, by their Lemma 8,

$$\sup_{x \in \mathcal{D}} |S_{h_n}(x)/(h_n^{pT}) - f_X(x) S_\ell| = O(\nu_n),$$

with $\nu_n := h_n + \left[ n h_n^{pT} / \ln n \right]^{-1/2}$. Then, we have

$$I_{i,j} = \int_u e' S_\ell^{-1} w(h_n u) \, \lambda_{h,j}(X_i - h_n u) K(u) \, \mathrm{d}u + g_{j,n}(X_i),$$

where $g_{j,n}$ is a deterministic function and because $\mathcal{K}$ is compact, $\sup_{\mathcal{K}} |g_{j,n}| = O(\nu_n)$. While $\lambda_{h,j}$ is not differentiable, it is directionally differentiable and we can write $\lambda_{h,j}(X_i - h_n u) = \lambda_{h,j}(X_i) - h_n u' \nabla \lambda_{h,j}(\widetilde{X}, u)$, for some $\widetilde{X}$, where $\nabla \lambda_{h,j}(\widetilde{X}, u)$ is uniformly bounded over $\widetilde{X}$ and $u$. Including this new residual in the definition of $g_{j,n}$ and $\nu_n$ and noting that $w(h_n u)' e = 1$,

$$I_{i,j} = \int_u e' S_\ell^{-1} w(h_n u) \, w(h_n u)' e \lambda_{h,j}(X_i) K(u) \, \mathrm{d}u + g_{j,n}(X_i)$$

$$= e' S_\ell^{-1} \int_u H_n w(u) \, w(h_n u)' K(u) \, \mathrm{d}u \, e \lambda_{h,j}(X_i) + g_{j,n}(X_i)$$

$$= e' S_\ell^{-1} H_n S_\ell H_n e \lambda_{h,j}(X_i) + g_{j,n}(X_i),$$

where again $g_n$ is deterministic and such that $\sup_{\mathcal{K}} |g_{j,n}| = O(\nu_n)$, and $H_n$ is a diagonal matrix with diagonal entries $h_n^{|r|}$ for $|r| \leq \ell$. Their entries are ordered in the same order as for the polynomial terms in $w(x)$. One can show that

$$H_n S_\ell H_n = S_\ell + O(h_n)$$

where the $O(h_n)$ is an entry-wise bound. This gives $I_{i,j} = \lambda_{h,j}(X_i) + g_{j,n}(X_i)$, changing again the definition of $g_{j,n}$ to a function with the same properties. By Assumption

5.2, $\nu_n \to 0$. Then, by Chebyshev's inequality and since $E(\epsilon_{ij}|X_i) = 0$, we have $\frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_{ij} g_{j,n}(X_i) = o_P(1)$. Thus,

$$
\begin{aligned}
\sqrt{n} E^* \left( D_\gamma h(X_i, S_i, \gamma_0, \beta_0)'[\widehat{\gamma}(X_i) - \gamma_0(X_i)] \right) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \sum_{j \leq T+1} \epsilon_{ij} [\lambda_{h,j}(X_i) + g_{j,n}(X_i)] \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n \sum_{j \leq T+1} \epsilon_{ij} \lambda_{h,j}(X_i) + o_P(1).
\end{aligned}
$$

Hence, Condition 5 follows.

**Lemma 8** *Suppose that Assumption 3 holds. Then:*

$$
\limsup_{n \to \infty} \sup_{P \in \mathcal{P}} P \left( \left\| n^{1/2}(\widehat{\beta} - \beta_0) - \frac{1}{n^{1/2}} \sum_{i=1}^n \phi_i \right\| > \eta \right) = 0. \tag{53}
$$

*Moreover, if $\underline{\sigma} > 0$ in the definition of $\mathcal{P}$, Lemma 2 holds uniformly over $\mathcal{P}$.*

**Proof:** To show both results, it suffices to show that they hold along any sequence of probability distribution $(P_n)_{n \geq 1}$ in $\mathcal{P}$. We use the same notation as in the other proofs but index parameters, variables and the expectation operator by $n$ to underline their dependence on $P_n$ when deemed necessary. Relatedly, we use $o_{P_n}(1)$ as a shortcut for a sequence of random variable $\varepsilon_n$ satisfying $P_n(\|\varepsilon_n\| > \eta) \to 0$ for all $\eta > 0$.

To prove the first point, let us first prove that $\widehat{\beta} - \beta_{0n} = o_{P_n}(1)$. To that end, consider the class of functions $\mathcal{L} := \{(y, x) \mapsto \ell_c(y|x; \beta); \beta \in \Theta\}$. We apply a version of Glivenko-Cantelli theorem on $\mathcal{L}$ that is uniform over $P$. The functions $(y, x, \beta) \mapsto \ell_c(y|x; \beta)$ are $C^1$ on $\{0, 1\}^T \times \text{Supp}(X) \times \Theta$, which is a compact set. The class $\mathcal{L}$ thus satisfies the Lipschitz requirement of Theorem 2.7.11 of van der Vaart and Wellner (1996). Then, by that theorem and the fact that $\Theta$ is compact,

$$
N(\epsilon \|F\|_{Q,1}, \mathcal{L}, L_1(Q)) \leq N_{[\,]}(\epsilon \|F\|_{Q,1}, \mathcal{L}, L_1(Q)) \leq N(\epsilon/2, \Theta, \|.\|) < \infty,
$$

where $N_{[\,]}$ denotes bracketing numbers, $N$ denotes covering numbers and $F$ is the envelope function defined in the same theorem. Hence,

$$
\sup_Q \log N(\epsilon \|F\|_{Q,1}, \mathcal{L}, L_1(Q)) < \infty.
$$

In view of the comment after its proof, we can then apply Theorem 2.8.1 of van der Vaart and Wellner (1996). As a result,

$$\sup_{\beta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^{n} \ell_c(Y_i|X_i; \beta) - E_n[\ell_c(Y|X; \beta)] \right| = o_{P_n}(1). \tag{54}$$

We establish below a uniform version of the well-separation condition by proving that for all $\eta > 0$, there exists $\nu > 0$ such that for all $n \geq 1$,

$$\sup_{\beta: \|\beta - \beta_{0n}\| > \eta} M_n(\beta) < M_n(\beta_{0n}) - \nu, \tag{55}$$

where $M_n(\beta) = E_n[\ell_c(Y|X; \beta)]$. By suitably modifying the proof of Theorem 5.7 in van der Vaart (2000) to the sequence $(P_n)$, the result follows.

Now, we prove that for any $\eta > 0$, there exists $\nu > 0$ such that (55) holds. For any $\beta$ such that $\|\beta - \beta_{0n}\| > \eta$, let

$$\beta' = \frac{\eta}{\|\beta - \beta_{0n}\|} \beta + \left( 1 - \frac{\eta}{\|\beta - \beta_{0n}\|} \right) \beta_{0n}.$$

Then $\|\beta' - \beta_{0n}\| = \eta$. Moreover, by concavity of $M_n$,

$$M_n(\beta') \geq \frac{\eta}{\|\beta - \beta_{0n}\|} M_n(\beta) + \left( 1 - \frac{\eta}{\|\beta - \beta_{0n}\|} \right) M_n(\beta_{0n}) \geq M_n(\beta).$$

Thus,

$$\sup_{\beta: \|\beta - \beta_{0n}\| > \eta} M_n(\beta) \leq \sup_{\beta \in S_{n,\eta}} M_n(\beta),$$

where $S_{n,\eta} = \{\beta : \|\beta - \beta_{0n}\| = \eta\}$. Next, for any $\beta \in S_{n,\eta}$ by a Taylor expansion of $M_n$ at $\beta_{0n}$,

$$M_n(\beta) = M_n(\beta_{0n}) - \frac{1}{2}(\beta - \beta_{0n})' \mathcal{I}_{n,0} (\beta - \beta_{0n}) + \frac{\partial^3 M_n}{\partial \beta \partial \beta'}(\tilde{\beta})[\beta - \beta_{0n}],$$

where $\tilde{\beta} = t\beta + (1-t)\beta_{0n}$ for some $t \in (0,1)$ and $\frac{\partial^3 M_n}{\partial \beta \partial \beta'}(\tilde{\beta})[\beta - \beta_{0n}]$ is the third order differential of $M_n$ at $\tilde{\beta}$ evaluated at $\beta - \beta_{0n}$. We know that $\mathcal{I}_{n,0} >> A$, write $\rho$ the smallest eigenvalue of $A$. By Assumption 3, there exists $B > 0$ such that $\left| \frac{\partial^3 M_n}{\partial \beta \partial \beta'}(\tilde{\beta})[\beta - \beta_{0n}] \right| \leq B\eta^3$, which gives

$$M_n(\beta) \leq M_n(\beta_{0n}) + \eta^2 (B\eta - \frac{1}{2}\underline{\rho}) \leq M_n(\beta_{0n}) - \varepsilon\eta^2$$

if $\eta \leq (\frac{1}{2}\rho - \varepsilon)/B$ for some $\varepsilon > 0$. Taking $\eta$ small enough is without loss of generality, thus (55) follows.

Next, we prove (53). By a Taylor expansion, there exists $t_n \in (0, 1)$ such that

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \ell_c}{\partial \beta}(Y_i|X_i; \beta_{0n}) + \left[\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \tilde{\beta}_n)\right]\left(\widehat{\beta} - \beta_{0n}\right) = 0,$$

where $\tilde{\beta}_n = \widehat{\beta} + (1 - t_n)\beta_{0n}$. Thus, by definition of $\phi_{n,i}$,

$$\mathcal{I}_{n,0}^{-1}\left[\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \tilde{\beta}_n)\right]\sqrt{n}\left(\widehat{\beta} - \beta_{0n}\right) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_{n,i}. \tag{56}$$

Now, by the triangle inequality and the fact that the third derivatives of $\ell_c$ are uniformly bounded, there exists $C > 0$ such that

$$\left\|\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \tilde{\beta}_n) - \mathcal{I}_{n,0}\right\| \leq \left\|\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \tilde{\beta}_n) - \frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \beta_{0n})\right\|$$

$$+ \left\|\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \beta_{0n}) - \mathcal{I}_{n,0}\right\|$$

$$\leq C\left\|\widehat{\beta} - \beta_{0n}\right\| + \left\|\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \beta_{0n}) - \mathcal{I}_{n,0}\right\|.$$

By what precedes, the first term is an $o_{P_n}(1)$. Moreover, for all $i$ and $n$, each element of the matrix $\partial^2 \ell_c / \partial \beta \partial \beta'(Y_i|X_i; \beta_{0n})$ is bounded almost surely. Thus, the uniform integrability condition of Gut (1992) holds for this variable. Then, by his weak LLN, the second term of the right-hand side above is also an $o_{P_n}(1)$. Thus, because $\mathcal{I}_{n0}^{-1} << \underline{A}^{-1}$ (since $P_n \in \mathcal{P}$), we have

$$\mathcal{I}_{n,0}^{-1}\left[\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^2 \ell_c}{\partial \beta \partial \beta'}(Y_i|X_i; \tilde{\beta}_n)\right] = \mathrm{Id} + o_{P_n}(1).$$

Next, for all $i$ and $n$, we have

$$E_n[\phi_{n,i}] = 0, \quad V_n(\phi_{n,i}) = \mathcal{I}_{n0}^{-1} << \underline{A}^{-1}. \tag{57}$$

Hence, by Chebyshev's inequality, the right-hand side of (56) is bounded in probability uniformly over $n$. Thus, this is also the case of $\sqrt{n}\left(\widehat{\beta} - \beta_{0n}\right)$. Hence,

$$\sqrt{n}\left(\widehat{\beta} - \beta_0\right) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_{n,i} + o_{P_n}(1).$$

66

In other words, (53) holds.

We now show that Lemma 2 holds uniformly over $\mathcal{P}$. Reasoning as in the proof of Lemma 2 and using the first point above, we get

$$\sqrt{n}\left(\widehat{\Delta} - \widetilde{\Delta}\right) = \frac{1}{n^{1/2}} \sum_{i=1}^{n} \left[\widehat{G}\phi_{n,i} + g(U_i, \beta_0) - \widetilde{\Delta}\right] + o_{P_n}(1).$$

Note that $g$ is $C^2$ on the compact set $\text{Supp}(U) \times \Theta$. Moreover, $\overline{\beta}_i$ as defined in Lemma 2 satisfies $\left\|\overline{\beta}_i - \beta_0\right\| \leq \left\|\widehat{\beta} - \beta_0\right\|$. Hence, there exists $M > 0$ such that

$$\left\|\widehat{G} - G_n\right\| \leq M \left\|\widehat{\beta} - \beta_{0n}\right\| + \left\|\frac{1}{n}\sum_{i=1}^{n} \frac{\partial g}{\partial \beta}(U_i, \beta_{0n}) - G_n\right\|. \tag{58}$$

By the first part of the proof, $\widehat{\beta} - \beta_{0n} = o_{P_n}(1)$. Next, because $\partial g/\partial \beta(., \beta_{0n})$ is bounded on $\text{Supp}(U)$, the uniform integrability condition of Gut (1992) also holds for this variable. Then, by his weak LLN,, the second term of (58) is an $o_{P_n}(1)$. Thus, $\left\|\widehat{G} - G_n\right\| = o_{P_n}(1)$. As a result,

$$\sqrt{n}\frac{\widehat{\Delta} - \widetilde{\Delta}}{\sigma_n} = \frac{1}{n^{1/2}} \sum_{i=1}^{n} \frac{G_n\phi_{n,i} + g(U_i, \beta_{0n}) - E_n[g(U_i, \beta_{0n})]}{\sigma_n} + o_{P_n}(1).$$

Now, by the triangle and Cauchy-Schwarz inequalities, we have

$$\left|G_n\phi_{n,i} + g(U_i, \beta_{0n}) - E_n[g(U_i, \beta_{0n})]\right| \leq \|G_n\| \|\phi_{n,i}\| + |g(U_i, \beta_{0n}) - E_n[g(U_i, \beta_{0n})]| . \tag{59}$$

$\|G_n\|$ is bounded uniformly over $n$. The variables $|g(U_i, \beta_{0n}) - E_n[g(U_i, \beta_{0n})]|$ are also bounded. Next, $\phi_{n,i} = \mathcal{I}_{n0}^{-1} V_{n,i}$ where $V_{n,i}$ is a bounded vector (with $\|V_{n,i}\| \leq C$, say). Moreover, because $P_n \in \mathcal{P}$,

$$\left\|\mathcal{I}_{n0}^{-1} V_{n,i}\right\| \leq \left\|\underline{A}^{-1} V_{n,i}\right\| \leq \underline{\rho}^{-1} C,$$

where $\underline{\rho} > 0$ denotes the smallest eigenvalue of $\underline{A}$. Then, using (59) and $\sigma_n \geq \underline{\sigma}$, the variables $(G_n\phi_{n,i} + g(U_i, \beta_{0n})/\sigma_n$ are bounded by a constant independent of $n$. Thus, they satisfy the Lindeberg condition. Then, by the central limit theorem for triangular arrays,

$$\sqrt{n}\frac{\widehat{\Delta} - \widetilde{\Delta}}{\sigma_n} \xrightarrow{d} \mathcal{N}(0, 1).$$

# C  Further details on the simulations

First, to estimate $\gamma_{0t}(x)$, we use a local linear estimator with a common bandwidth $h_t$ for the $T$ components of $X$. To choose $h_t$, we aim at reaching a certain ratio between the (integrated) bias and standard deviation of the estimator. Specifically, let $B_t(x, h)$ and $\sigma_t^2(x, h)$ denote respectively the asymptotic bias and variance of $\widehat{\gamma}_t(x)$ with a bandwidth equal to $h$. Then (see, e.g. Ruppert and Wand, 1994),

$$B_t(x, h) = h^2 \left( \int u^2 K(u)du \right) \sum_{j=1}^{pT} \frac{\partial^2 \gamma_{0t}}{\partial x_j^2}(x),$$

$$\sigma_t^2(x, h) = \frac{1}{nh^T} \frac{\left( \int K(u)^2 du \right)^T \gamma_{0t}(x)(1 - \gamma_{0t}(x))}{f_X(x)}.$$

Then, define $B_t^2(h) := E[B_t^2(X, h)]$ and $\sigma_t^2(h) := E[\sigma_t^2(X, h)]$. Assuming first that $B_t^2(h)$ and $\sigma_t^2(h)$ are known, we would choose $h_t$ so that $\sigma_t^2(h_t) = R_n \times B_t^2(h_t)$, where $R_n > 0$ fixes to the degree of undersmoothing. For instance, $R_n = 1$ corresponds to the optimal bandwidth in terms of asymptotic mean integrated squared error. We use $R_n = 5(n/500)^2$ in our simulations. Now, $B_t^2(h)$ and $\sigma_t^2(h)$ are actually unknown. We estimate both assuming that $\alpha$ is constant. Then, we can estimate this constant by MLE (plugging the CMLE $\widehat{\beta}$ in the log-likelihood) and then estimate $\gamma_{0t}(x)$ by plug-in, using (2).

Finally, to obtain $\widehat{m}$, we must choose a threshold $c_n$. We actually slightly modify $\widehat{I}(x)$, by letting

$$\widehat{I}(x) := \max \left\{ t \in \{1, ..., T\} : \underline{H}_t(\widetilde{m}_{\to t}(x)) \geq \underline{c}_{nt}(x) \text{ and } \overline{H}_t(\widetilde{m}_{\to t}(x)) \geq \overline{c}_{nt}(x) \right\},$$

where $\underline{c}_{nt}(x) := \widehat{\underline{\sigma}}_t[2 \ln \ln(n)]^{1/2}$, $\overline{c}_{nt}(x) := \widehat{\overline{\sigma}}_t(x)[2 \ln \ln(n)]^{1/2}$ and $\widehat{\underline{\sigma}}_t^2(x)$ (resp. $\widehat{\overline{\sigma}}_t^2(x)$) is an estimator of the asymptotic variance of $\underline{H}_t(m_{\to t}(x))$ (resp. $\overline{H}_t(m_{\to t}(x))$).