# Semiparametric Single Index Models

Li and Racine (2007, Chapter 8)

Yasuyuki Matsumura

December 8, 2024

Graduate School of Economics, Kyoto University

- A semiparametric single index model is given by

$$Y = g(X^T \beta_0) + u,$$

where

$Y \in \mathbb{R}$ : a dependent variable,

$X \in \mathbb{R}^q$ : a $q \times 1$ explanatory vector,

$\beta_0 \in \mathbb{R}^q$ : a $q \times 1$ vector of unknown parameters,

$u \in \mathbb{R}$ : an error term which satisfies $\mathbb{E}(u \mid X) = 0$,

$g(\cdot)$ : an unknown distribution function.

- Even though $x$ is a $q \times 1$ vector, $x^T \beta_0$ is a scalar of a single linear combination, which is called a single index.

- By the form of the single index model, we obtain

$$\mathbb{E}(Y \mid X) = g(X^T \beta_0),$$

  which means that the conditional expectation of $Y$ only depends on the vector $X$ through a single index $X^T \beta_0$.

- The model is semiparametric when $\beta \in \mathbb{R}^q$ is estimated with the parametric methods and $g(\cdot)$ with the nonparametric methods.

# Examples of Parametric Single Index Model

- If $g(\cdot)$ is the identity function, then the model turns out to be a linear regression model:

$$Y = g(X^T \beta_0) + u = X^T \beta_0 + u.$$

- If $g(\cdot)$ is the CDF of Normal$(0, 1)$, then the model turns out to be a probit model.
  - See the textbook for further discussions on a probit model.
- If $g(\cdot)$ is the CDF of logistic distribution, then the model turns out to be a logistic regression model.

## TOC

# Identification Conditions

## Identification Conditions

Proposition 8.1 (Identification of a Single Index Model)

For the semiparametric single index model $Y = g(x^T\beta_0) + u$, identification of $\beta_0$ and $g(\cdot)$ requires that

- (i) $x$ should not contain a constant/an intercept, and must contain at least one continuous variable. Moreover, $\|\beta_0\|=1$.

- (ii) $g(\cdot)$ is differentiable and is not a constant function on the support of $x^T\beta_0$.

- (iii) For the discrete components of $x$, varying the values of the discrete variables will not divide the support of $x^T\beta_0$ into disjoint subsets.

- Note that the location and the scale of $\beta_0$ are not identified.
- The vector $x$ cannot include an intercept because the function $g(\cdot)$ (which is to be estimated in nonparametric manners) includes any location and level shift.
  - That is, $\beta_0$ cannot contain a location parameter.

## Identification Condition (i)

- Some normalization criterion (scale restrictions) for $\beta_0$ are needed.
    - One approach is to set $\|\beta_0\| = 1$.
    - The second approach is to set one component of $\beta_0$ to equal one. This approach requires that the variable corresponding to the component set to equal one is continuously distributed and has a non-zero coefficient.
    - Then, $x$ must be dimension $2$ or larger. If $x$ is one-dimensional, then $\beta_0 \in \mathbb{R}^1$ is simply normalized to 1, and the model is the one-dimensional nonparametric regression $E(Y \mid x) = g(x)$ with no semiparametric component.

## Identification Conditions (ii) and (iii)

- The function $g(\cdot)$ cannot be a constant function and must be differentiable on the support of $x^T \beta_0$.
- $x$ must contain at least one continuously distributed variable and this continuous variable must have non-zero coefficient.
    - If not, $x^T \beta_0$ only takes a discrete set of values and it would be impossible to identify a continuous function $g(\cdot)$ on this discrete support.

# Ichimura's (1993) Method

- Textbook: Sections 8.2; 8.4.1; and 8.12.
- Suppose that the functional form of $g(\cdot)$ were known.
- Then we could estimate $\beta_0$ by minimizing the least-squares criterion:

$$\sum_{i=1} \left[ Y_i - g(X_i^T \beta) \right]^2$$

  with respect to $\beta$.
- We could think about replacing $g(\cdot)$ with a nonparametric estimator $\hat{g}(\cdot)$.
- However, since $g(z)$ is the conditional mean of $Y_i$ given $X_i^T \beta_0 = z$, $g(\cdot)$ depends on unknown $\beta_0$, so we cannot estimate $g(\cdot)$ here.

10

- Nevertheless, for a fixed value of $\beta$, we can estimate

$$G(X_i^T \beta) := \mathbb{E}(Y_i \mid X_i^T \beta) = \mathbb{E}(g(X_i^T \beta_0) \mid X_i^T \beta).$$

- In general $G(X_i^T \beta) \neq g(X_i^T \beta)$.
- When $\beta = \beta_0$, it holds that $G(X_i^T \beta_0) = g(X_i^T \beta_0)$. [1]

---

[1]一般の $X_i^T \beta$ を用いて条件付けると，$G$ と $g$ は通常は一致しないが，正しいインデックス $X_i^T \beta = X_i^T \beta_0$ のときだけ一致するということ．

- First, we estimate $G(X_i^T \beta)$ with the leave-one-out NW estimator:

$$
\hat{G}_{-i}(X_i^T \beta) := \hat{\mathbb{E}}_{-i}(Y_i \mid X_i^T \beta)
$$

$$
= \frac{\sum_{j \neq i} Y_j K\left(\frac{X_j^T \beta - X_i^T \beta}{h}\right)}{\sum_{j \neq i} K\left(\frac{X_j^T \beta - X_i^T \beta}{h}\right)}.
$$

- Second, using the leave-one-out NW estimator $\hat{G}_{-i}(X_i^T\beta)$, we estimate $\beta$ with

$$\hat{\beta} := \arg\min_{\beta} \sum_{i=1}^{n} \left[ Y_i - \hat{G}_{-i}(X_i^T\beta) \right]^2 w(X_i)\mathbf{1}(X_i \in A_n)$$

$$:= \arg\min_{\beta} S_n(\beta),$$

which is called Ichimura's estimator (the WSLS estimator).

- $w(X_i)$ is a nonnegative weight function.
- $\mathbf{1}(X_i \in A_n)$ is a trimming function to trim out small values of $\hat{p}(X_i^T\beta) = \frac{1}{nh} \sum_{j\neq i} K\left(\frac{X_j^T\beta - X_i^T\beta}{h}\right)$, so that we do not suffer the random denominator problem.
  - $A_\delta = \{x : p(x^T\beta) \geq \delta, \text{ for } {}^\forall\beta \in \mathcal{B}\}$.
  - $A_n = \{x : ||x - x^\star|| \leq 2h, \text{ for } {}^\exists x^\star \in A_\delta\}$, which shrinks to $A_\delta$ as $n \to \infty$ and $h \to 0$.

- Let $\hat{\beta}$ denote the semiparametric estimator of $\beta_0$ obtained from minimizing $S_n(\beta)$.

- To derive the asymptotic distribution of $\hat{\beta}$, the following conditions are neeeded:

## Asymptotic Distribution of Ichimura's Estimator

---

**Assumption 8.1**

The set $A_\delta$ is compact, and the weight function $w(\cdot)$ is bounded and posotive on $A_\delta$. Define the set

$$D_z = \{z : z = x^T\beta, \beta \in \mathcal{B}, x \in A_\delta\}.$$

Letting $p(\cdot)$ denote the PDF of $z \in D_z$, $p(\cdot)$ is bounded below by a positive constant for $^\forall z \in D_z$

---

**Assumption 8.2**

$g(\cdot)$ and $p(\cdot)$ are 3 times differentiable w.r.t. $z = x^\beta$. The third derivatives are Lipschitz continuous uniformly over $\mathcal{B}$ for $^\forall z \in D_z$.

# Asymptotic Distribution of Ichimura's Estimator

> **Assumption 8.3**
>
> The kernel function is a bounded second order kernel, which has bounded support; is twice differentiable; and its second derivative is Lipschitz continuous.

> **Assumption 8.4**
>
> $\mathbb{E}(|Y^m|) < \infty$ for $^{\exists}m \geq 3$. $\mathrm{var}(Y \mid x)$ is bounded and bounded away from zero for $^{\forall}x \in A_\delta$. $\frac{q\ln(h)}{nh^{3+\frac{3}{m-1}}} \to 0$ and $nh^8 \to 0$ as $n \to \infty$.

## Asymptotic Distribution of Ichimura's Estimator

**Theorem 8.1.** Under assumptions 8.1 through 8.4,

$$\sqrt{n}(\hat{\beta} - \beta_0) \xrightarrow{d} \text{Normal}(0, \Omega_I),$$

with

$$
\begin{aligned}
\Omega_I &= V^{-1} \Sigma V^{-1}, \\
V &= \mathbb{E}\{w(X_i)(g_i^{(1)})^2 \\
&\quad \times (X_i - E_A(X_i \mid X_i^T \beta_0))(X_i - E_A(X_i \mid X_i^T \beta_0))^T\}, \\
\Sigma &= \mathbb{E}\{w(X_i)\sigma^2(X_i)(g_i^{(1)})^2 \\
&\quad \times (X_i - E_A(X_i \mid X_i^T \beta_0))(X_i - E_A(X_i \mid X_i^T \beta_0))^T\},
\end{aligned}
$$

where

- $(g_i^{(1)}) = \frac{\partial g(v)}{\partial v} \big|_{v = X_i^T \beta_0}$,
- $\mathbb{E}_A(X_i \mid v) = \mathbb{E}(X_i \mid x_A^T \beta_0 = v)$,
- $x_A$ has the distribution of $X_i$ conditional on $X_i \in A_\delta$. 17

## Asymptotic Distribution of Ichimura's Estimator

- See Ichimura (1993); and Hardle, Hall and Ichimura (1993) for the proof of **Theorem 8.1**.

- Horowitz (2009) provides an excellent heuristic outline for proving **Theorem 8.1**, using only the familiar Taylor series methods, the standard LLN, and the Lindeberg-Levy CLT.

# Optimal Weight under the Homoscedasticity Assumption

- We introduce the following homoscedasticity assumption:

$$\mathbb{E}(u_i^2 \mid X_i) = \sigma^2.$$

- Under this assumption, the optimal choice of $w(\cdot)$ is $w(X_i) = 1$.

- In this case,

$$\hat{\beta} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - \hat{G}_{-i}(X_i^T \beta)^2) \mathbf{1}(X_i \in A_n)$$

is semiparametrically efficient in the sense that $\Omega_I$ is the semiparametric variance lower bound (conditional on $X \in A_\delta$).

- In general, $\mathbb{E}(u_i^2 \mid X_i) = \sigma^2(X_i)$.

- An infeasible case: If one assues that $\mathbb{E}(u_i^2 \mid X_i) = \sigma^2(X_i^T \beta_0)$, that is, the conditional variance depends only on the single index $X_i^T \beta_0$, the choice of $w(X_i) = \frac{1}{\sigma^2(X_i^T \beta_0)}$ can lead to a semiparametrically efficient estimation.

- We could adopt a two-step procedure as follows.

# Two-Step Procedure to Choose Optimal Weight

- Suppose that the conditional variance is a function of $X_i^T \beta_0$ (Let $\sigma^2(X_i^T \beta_0)$ denote it).

- The first step: Use $w(X_i) = 1$ to obtain a $\sqrt{n}$-consistent estimator of $\beta_0$.

- Let $\tilde{\beta}_0$ denote the estimator of $\beta_0$, and $\tilde{u}_i = Y_i - \hat{g}(X_i^T \tilde{\beta}_0)$ denote the residual obtained from $\tilde{\beta}_0$.

- We can obtain a consistent nonparametric estimator of the conditional variance: $\hat{\sigma}^2(X_i^T \tilde{\beta}_0)$.

- The second step: Estimate $\beta_0$ again using $w(X_i) = \frac{1}{\hat{\sigma}^2(X_i^T \tilde{\beta}_0)}$:

$$\hat{\beta}_0 = \arg\min_{\beta} \sum_{i=1}^{n} \left[ Y_i - \hat{G}_{-i}(X_i^T \beta) \right]^2 \frac{1}{\hat{\sigma}^2(X_i^T \tilde{\beta}_0)} \mathbf{1}(X_i \in A_n).$$

- The estimator $\hat{\beta}_0$ is semiparametrically efficient because $\hat{\sigma}^2(v) - \sigma^2(v)$ converges to zero at a particular rate uniformly over $v \in D_v$ ($D_v$ is the support of $X_i^T \beta_0$). [2]

---

[2]$\hat{\sigma}^2(X_i^T \beta)$ を用いるケースもある.

- Recall that we assume in Assumption 8.4 that $\frac{q\ln(h)}{nh^{3+\frac{3}{m-1}}} \to 0$ and $nh^8 \to 0$ as $n \to \infty$, where $m \geq 3$ is a positive integer whose specific value depends on the existence of the number of finite moments of $Y$ along with the smoothness of the unknown function $g(\cdot)$. [3]

- The range of permissive smoothing parameters allows for optimal smoothing, i.e., $h = O(n^{-\frac{1}{5}})$. [4]

---

[3] Assumption 8.4 は，$g$ をノンパラメトリックに推定することがパラメトリックパートの収束レートに影響を与えないための十分条件になっている.
[4] このオーダーで選んだ $h$ は，Assumption 8.4 を満たしている.

## Bandwidth Selection for Ichimura's Estimator

- Our aim is to choose $\hat{\beta}$ close to $\beta_0$, and $h$ close to the value $h_0$, which minimize the average of

$$\mathbb{E}\{\hat{g}(X_i^T\beta_0 \mid X_i^T\beta_0) - g(X_i^T\beta_0)\}^2.$$

- Hardle, Hall and Ichimura (1993) suggest picking $\beta$ and the bandwidth $h$ jointly by minimization of $S_n(\beta)$.

- Further discussions follow in Section 8.4.

**Direct Semiparametric Estimators for $\beta$**

## Direct Semiparametric Estimators for $\beta$

- Textbook: Sections 8.3; and 8.4.2.
- Here we review:
    - Hardle and Stoker's (1989) Average Derivative Estimator,
    - Powell, Stock and Stoker's (1989) Density-Weighted Average Derivative Estimator, and
    - Li, Lu and Ullah's (2003) Estimator.
- The advantage of the direct estimation method is that we can estimate $\beta_0$ and $g(x^T \beta_0)$ directly without running the nonlinear least squares, which leads to the computational simplicity.
- We still suffer from a finite-sample problem.

- Suppose that $x$ is a $q \times 1$ vector of continuous variables.

- Then we obtain the average derivative of $\mathbb{E}(Y \mid X = x)$:

$$\mathbb{E}\left[\frac{\partial \mathbb{E}(Y \mid X = x)}{\partial x}\right] = \mathbb{E}\left[g^{(1)}(x^T \beta_0)\right]\beta_0$$

- Recall that the scale of $\beta_0$ is not identified, which means that the constant $\mathbb{E}\left[g^{(1)}(x^T \beta_0)\right]$ does not matter. That is, a normalized estimation of $\mathbb{E}\left[\frac{\partial \mathbb{E}(Y \mid X = x)}{\partial x}\right]$ is an estimation of normalized $\beta_0$.

## Hardle and Stoker's (1989) Average Derivative Estimator

- Let $\hat{\mathbb{E}}(Y_i \mid X_i)$ denote the NW estimator of $\mathbb{E}(Y_i \mid X_i)$:

$$\hat{\mathbb{E}}(Y_i \mid X_i) = \frac{\sum_{j=1}^{n} Y_j K\left(\frac{X_i - X_j}{a}\right)}{\sum_{j=1}^{n} K\left(\frac{X_i - X_j}{a}\right)}.$$

- Assuming that the kernel function is differentiable, we can estimate $\beta_0$, estimating $\mathbb{E}\left[\frac{\partial \mathbb{E}(Y \mid X = x)}{\partial x}\right]$ with its sample analogue:

$$\tilde{\beta}_{ave} = \frac{1}{n} \sum_{i=1}^{n} \frac{\partial \hat{\mathbb{E}}(Y_i \mid X_i)}{\partial X_i}.$$

- The scale normalization can also be implemented by $\frac{\tilde{\beta}_{ave}}{|\tilde{\beta}_{ave}|}$.

- An issue raised with this estimator is the random denominator problem, which leads to a difficulty in the derivation of the asymptotic properties.
- Rilstone (1991) establishes the $\sqrt{n}$-normality using a trimming function.

- As we obtain the average derivative above, we also obtain the weighted average derivative of $\mathbb{E}(Y \mid X = x)$:

$$\mathbb{E}\left[w(x)\frac{\partial \mathbb{E}(Y \mid X = x)}{\partial x}\right] = \mathbb{E}\left[w(x)g^{(1)}(x^T\beta_0)\right]\beta_0.$$

- Let $w(x)$ be the density function $f(x)$, and $\delta$ denote the density-weighted average derivative of $\mathbb{E}(Y \mid X = x)$.

- Then we obtain

$$\delta = \mathbb{E}\left[f(X)\frac{\partial \mathbb{E}(Y \mid X = x)}{\partial x}\right]$$

$$= \mathbb{E}\left[f(X)g^{(1)}(X^T\beta_0)\right]$$

$$= \int g^{(1)}(x^T\beta_0)f^2(x)dx$$

$$= g(x^T\beta_0)f^2(x) - 2\int g(x^T\beta_0)f^{(1)}(x)f(x)dx.$$

## Powell, Stock and Stoker (1989) Density-Weighted Average Derivative Estimator

- Assume that $f(x) = 0$ at the boundary of the support of $X$. Then we observe that $g(x^T \beta_0) f^2(x) = 0$, that is,

$$\delta = -2 \int g(x^T \beta_0) f^{(1)}(x) f(x) dx$$
$$= -2 \mathbb{E}[g(X^T \beta_0) f^{(1)}(X)]$$
$$= -2 \mathbb{E}[Y f^{(1)}(X)].$$

- We can estimate $\delta$ by its sample analogue:

$$\hat{\delta} = -\frac{2}{n} \sum_{i=1}^{n} Y_i \hat{f}_{-i}^{(1)}(X_i),$$

where $\hat{f}_{-i}(X_i)$ is the leave-one-out NW estimator of $f(X)$:

$$\hat{f}_{-i}(X_i) = \frac{1}{n-1} \sum_{j \neq i} \left(\frac{1}{h}\right)^q K\left(\frac{X_i - X_j}{h}\right).$$

**Powell, Stock and Stoker's (1989) Density-Weighted Average Derivative Estimator**

- The textbook uses the NW estimator $\hat{f}^{(1)}(X_i)$ in (8.17).
- However, Powell, Stock and Stoker (1989) define their estimator using the leave-one-out NW estimator $\hat{f}^{(1)}_{-i}(X_i)$.
- Here we proceed with Powell, Stock and Stoker (1989).

## Powell, Stock and Stoker's (1989) Density-Weighted Average Derivative Estimator

- A useful representation of $\hat{\delta}$ is given by

$$\hat{\delta} = \frac{-2}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \left(\frac{1}{h}\right)^{q+1} Y_i K^{(1)} \left(\frac{X_i - X_j}{h}\right).$$

- Under some assumptions, if $h \to 0$ and $nh^{q+2} \to \infty$ hold, then the density-weighted average derivative estimator $\hat{\delta}$ satisfies that

$$\sqrt{n}(\hat{\delta} - \mathbb{E}[\hat{\delta}]) \xrightarrow{d} \text{Normal}(0, \Sigma_\delta),$$

where
$\Sigma_\delta = 4\mathbb{E}[\sigma^2(X)f^{(1)}(X)f^{(1)}(X)^T] + 4\text{Var}(f(X)g^{(1)}(X)).$

- Recall that

# Bandwidth Selection

# Klein and Spady's (1993) Estimator

# Lewbel's (2000) Estimator

# Manski's (1975) Maximum Score Estimator

# Horowitz's (1992) Smoothed Maximum Score Estimator

# Han's (1987) Maximum Rank Estimator

# Multinomial Discrete Choice Models

# Ai's (1997) Semiparametric Maximum Likelihood Approach

# References

## References (1)

- Hardle, W, P. Hall and H. Ichimura (1993) "Optimal Smoothing in Single-Index Models," *Annals of Statistics,* 21, 157-178.
- Horowitz, J. L. (2009) *Semiparametric and Nonparametric Methods in Econometrics,* Springer.
- Ichimura, H. (1993) "Semiparametric Least Squares (SLS) and Weighted SLS Estimation of Single-Index Models," *Journal of Econometrics,* 3, 205-228.
- Li, Q. and J. S. Racine, (2007). *Nonparametric Econometrics: Theory and Practice,* Princeton University Press.
- 末石直也 (2024) 『データ駆動型回帰分析：計量経済学と機械学習の融合』日本評論社.
- 西山慶彦，人見光太郎 (2023) 『ノン・セミパラメトリック統計解析（理論統計学教程：数理統計の枠組み)』共立出版.

35

## References (2)

Useful references also include some lecture notes of the following topic courses:

- ECON 718 NonParametric Econometrics (Bruce Hansen, Spring 2009, University of Wisconsin-Madison),
- セミノンパラメトリック計量分析（末石直也，2014 年度後期，京都大学）.