

# Semiparametric estimation in very high-dimensional models

James Robins

Lingling Li

Eric Tchetgen

Aad van der Vaart

Harvard University

Universiteit Leiden

Four facets of statistics and a Surprise, Zürich, November 2015

# Semiparametrics

# Semiparametric Inference

A **semiparametric model** is a model of **infinite dimension**

Interest is in estimating a finite-dimensional parameter defined through the structure of the model, e.g.

- a relative risk
- a coefficient of a particular regression variable
- a mean response

**Classical semiparametrics** (1980/90s): the combination of parameter and model is such that the “bias” is small relative to “variance” and estimation is possible at the “parametric rate”  $\sqrt{n}$ .

**Modern semiparametrics**: what if the model is too “large” for a parametric rate?

## Classical semiparametrics

$X_1, \dots, X_n$  i.i.d. with density  $p \in \mathcal{P}$

We want to estimate  $\chi(p)$ , for  $\chi: \mathcal{P} \rightarrow \mathbb{R}$ .

### META THEOREM

If  $\mathcal{P}$  and  $\chi$  are nice, then there exist  $T_n = T_n(X_1, \dots, X_n)$ , with

$$\sqrt{n}(T_n - \chi(p)) \rightsquigarrow N(0, \sigma_p^2).$$

**Classical semiparametrics** (1980/90s) was concerned with finding  $T_n$  with **minimal**  $\sigma_p^2$

*General methods such as (semiparametric, penalized, sieved) **maximum likelihood** or **Bayes** work well for many semiparametric models.*

## Example: symmetric location (Stein (1956), Stone (1975), Bickel (1981), ...)

Error  $\varepsilon$  with symmetric density  $\eta$ , Fisher information  $I_\eta < \infty$

Observe  $X = \theta + \varepsilon$

### THEOREM

There exists  $T_n = T_n(X_1, \dots, X_n)$  with, for all  $(\theta, \eta)$ ,

$$\sqrt{n}(T_n - \theta) \rightsquigarrow N(0, I_\eta^{-1})$$

## Example: Cox model (Cox (1975), Tsiatis, Gill, Wellner,...)

Covariate  $Z, \quad \sim f$

Survival time  $T$  with conditional hazard  $\lambda(t)e^{\theta^T z}$

Observe  $(Z, T)$

### THEOREM

There exists  $T_n = T_n(X_1, \dots, X_n)$  with, for all  $(\theta, \lambda, f)$ ,

$$\sqrt{n}(T_n - \theta) \rightsquigarrow N(0, \sigma_{\theta, \lambda, f}^2)$$

*Maximum likelihood estimator and certain Bayes estimators attain minimal  $\sigma_{\theta, \lambda, f}^2$ .*

## Example: semiparametric regression

Covariates  $(W, Z)$ , density  $f$

Error  $\varepsilon$ , such that  $\varepsilon | (W, Z)$  has density  $g(\cdot | W, Z)$  with  $E(\varepsilon | W, Z) = 0$

Outcome  $Y = \theta W + \eta(Z) + \varepsilon$

Observe  $X = (W, Z, Y)$

### THEOREM

There exists  $T_n = T_n(X_1, \dots, X_n)$  with, for all  $(\theta, f, g, \eta)$  such that  $g$  and  $\eta$  are sufficiently smooth,

$$\sqrt{n}(T_n - \theta) \rightsquigarrow N(0, \sigma_{\theta, f, g, \eta}^2)$$

## Example: missing data (Robins and Rotnitzky, ...)

Covariate  $Z$ ,  $\nu \sim f$

Response  $Y$ , with  $Y|Z \sim \text{binomial}(1, b(Z))$

Missingness indicator  $A$ , with  $A|Z \sim \text{binomial}(1, 1/a(Z))$

Missing at random:  $Y \perp\!\!\!\perp A|Z$

Observe  $X = (YA, A, Z) \in \{0, 1\} \times \{0, 1\} \times [0, 1]^d$

We wish to estimate mean response  $\chi(a, b, f) = \int b f d\nu = \mathbf{E}Y$ .



## Example: missing data (Robins and Rotnitzky, ...)

Covariate  $Z$ ,  $\nu \sim f$

Response  $Y$ , with  $Y|Z \sim \text{binomial}(1, b(Z))$

Missingness indicator  $A$ , with  $A|Z \sim \text{binomial}(1, 1/a(Z))$

Missing at random:  $Y \perp\!\!\!\perp A|Z$

Observe  $X = (YA, A, Z) \in \{0, 1\} \times \{0, 1\} \times [0, 1]^d$

We wish to estimate mean response  $\chi(a, b, f) = \int b f d\nu = \mathbf{E}Y$ .

*$Y$  is observed only if  $A = 1$*

*$Z$  is included to make the assumption  $Y \perp\!\!\!\perp A|Z$  realistic*

## Example: missing data (Robins and Rotnitzky, ...)

Covariate  $Z$ ,  $\nu \sim f$

Response  $Y$ , with  $Y|Z \sim \text{binomial}(1, b(Z))$

Missingness indicator  $A$ , with  $A|Z \sim \text{binomial}(1, 1/a(Z))$

Missing at random:  $Y \perp\!\!\!\perp A|Z$

Observe  $X = (YA, A, Z) \in \{0, 1\} \times \{0, 1\} \times [0, 1]^d$

We wish to estimate mean response  $\chi(a, b, f) = \int b f d\nu = \mathbf{E}Y$ .

*$Y$  is observed only if  $A = 1$*

*$Z$  is included to make the assumption  $Y \perp\!\!\!\perp A|Z$  realistic*

### THEOREM

There exists  $T_n = T_n(X_1, \dots, X_n)$  with, for all  $(a, b, f)$  such that  $a$  and  $b$  are sufficiently smooth,

$$\sqrt{n}(T_n - \chi(a, b, f)) \rightsquigarrow N(0, \sigma_{a,b,f}^2)$$

## Classical semiparametrics — minimal variance

$X_1, \dots, X_n$  i.i.d. with density  $p \in \mathcal{P}$

We want to estimate  $\chi(p)$ , for  $\chi: \mathcal{P} \rightarrow \mathbb{R}$ .

### META THEOREM

If  $\mathcal{P}$  and  $\chi$  are nice, then there exist  $T_n = T_n(X_1, \dots, X_n)$ , with

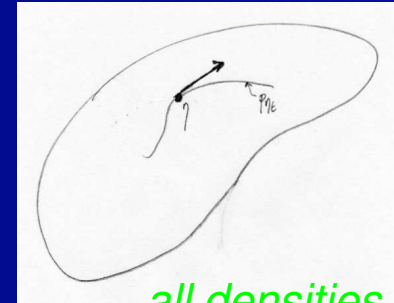
$$\sqrt{n}(T_n - \chi(p)) \rightsquigarrow N(0, \sigma_p^2).$$

What is the minimal variance  $\sigma_p^2$ ?

# First order tangent space and influence functions (Koshevnik and Levit (1976), Pfanzagl (1983),

vdV (1988).)

**Tangent set** (at  $p$ ): all score functions  $g = \left. \frac{d}{dt} \right|_{t=0} \log p_t$  of one-dimensional submodels  $t \mapsto p_t$  with  $p_0 = p$

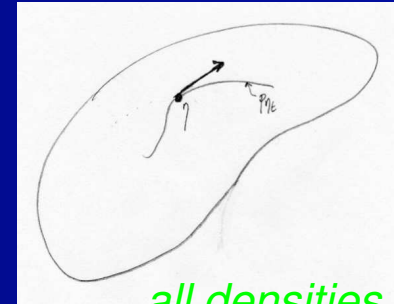


all densities  $p$

# First order tangent space and influence functions (Koshevnik and Levit (1976), Pfanzagl (1983),

vdV (1988).)

**Tangent set** (at  $p$ ): all score functions  $g = \frac{d}{dt}|_{t=0} \log p_t$  of one-dimensional submodels  $t \mapsto p_t$  with  $p_0 = p$



all densities  $p$

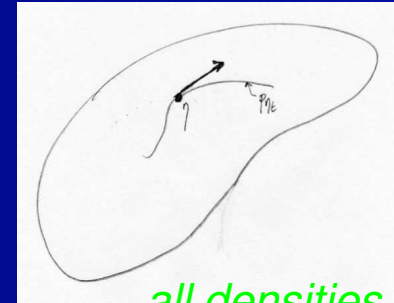
**Influence function** of  $p \mapsto \chi(p)$  is map  $x \mapsto \chi_p(x)$  with for all  $t \mapsto p_t$

$$\frac{d}{dt} \chi(p_t)|_{t=0} = P g \chi_p$$

# First order tangent space and influence functions (Koshevnik and Levit (1976), Pfanzagl (1983),

vdV (1988).)

**Tangent set** (at  $p$ ): all score functions  $g = \frac{d}{dt}|_{t=0} \log p_t$  of one-dimensional submodels  $t \mapsto p_t$  with  $p_0 = p$



**Influence function** of  $p \mapsto \chi(p)$  is map  $x \mapsto \chi_p(x)$  with for all  $t \mapsto p_t$

$$\frac{d}{dt} \chi(p_t)|_{t=0} = P g \chi_p$$

**THEOREM** If  $\sqrt{n}(T_n - \chi(p)) \rightsquigarrow L$ , locally uniformly in  $p$ , then for some  $M$

$$L = N(0, P(\Pi_p \chi_p)^2) * M,$$

for  $\Pi_p$  orthogonal projection onto closed linear span of tangent set.

## Example: missing data

Observe  $X = (Y A, A, Z)$

Parameter  $p \leftrightarrow (a, b, f)$

Likelihood  $f(Z)(1/a)(Z)^A(1 - 1/a(Z))^{1-A}b(Z)^{Y A}(1 - b(Z))^{(1-Y)A}$

$$\frac{Aa(Z) - 1}{a(Z)(a - 1)(Z)}\alpha(Z) \quad a\text{-score, } a_t = a + t\alpha$$

$$\frac{A(Y - b(Z))}{b(Z)(1 - b)(Z)}\beta(Z) \quad b\text{-score, } b_t = b + t\beta$$

$$\phi(Z) \quad f\text{-score, } f_t = f(1 + t\phi), \quad \int \phi f = 0$$

## Example: missing data

Observe  $X = (YA, A, Z)$

Parameter  $p \leftrightarrow (a, b, f)$

Likelihood  $f(Z)(1/a)(Z)^A(1 - 1/a(Z))^{1-A}b(Z)^{YA}(1 - b(Z))^{(1-Y)A}$

$$\frac{Aa(Z) - 1}{a(Z)(a - 1)(Z)}\alpha(Z) \quad a\text{-score, } a_t = a + t\alpha$$

$$\frac{A(Y - b(Z))}{b(Z)(1 - b)(Z)}\beta(Z) \quad b\text{-score, } b_t = b + t\beta$$

$$\phi(Z) \quad f\text{-score, } f_t = f(1 + t\phi)$$

Parameter of interest  $\chi(p) = \int bf = \mathbb{E}Y$

Influence function  $\chi_p(X) = Aa(Z)(Y - b(Z)) + b(Z) - \chi(p)$



## Example: missing data

Observe  $X = (YA, A, Z)$

Parameter  $p \leftrightarrow (a, b, f)$

Likelihood  $f(Z)(1/a)(Z)^A(1 - 1/a(Z))^{1-A}b(Z)^{YA}(1 - b(Z))^{(1-Y)A}$

$$\frac{Aa(Z) - 1}{a(Z)(a - 1)(Z)}\alpha(Z) \quad a\text{-score, } a_t = a + t\alpha$$

$$\frac{A(Y - b(Z))}{b(Z)(1 - b)(Z)}\beta(Z) \quad b\text{-score, } b_t = b + t\beta$$

$$\phi(Z) \quad f\text{-score, } f_t = f(1 + t\phi)$$

Parameter of interest  $\chi(p) = \int bf = \mathbf{E}Y$

Influence function  $\chi_p(X) = Aa(Z)(Y - b(Z)) + b(Z) - \chi(p)$

$$\mathbf{E}_p \chi_p(X)[a\text{-score}] = 0$$

$$\mathbf{E}_p \chi_p(X)[b\text{-score}] = \frac{\partial}{\partial t} \Big|_{t=0} \int b_t f \, d\nu = \int \beta f \, d\nu$$

$$\mathbf{E}_p \chi_p(X)[f\text{-score}] = \frac{\partial}{\partial t} \Big|_{t=0} \int b f_t \, d\nu = \int b \phi f \, d\nu$$

## Corrected plug-in estimators

## Heuristics — plug in and bias correction

Estimate  $\theta := \chi(p) \in \mathbb{R}$  from iid  $X_1, \dots, X_n \sim p$ .

Given  $\hat{p}$  and “influence function”  $(x_1, \dots, x_m) \mapsto \chi_p(x_1, \dots, x_m)$  use

$$\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi \hat{p},$$

for

$$\mathbb{U}_n f = \frac{(n-m)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_m \leq n} \dots \sum f(X_{i_1}, \dots, X_{i_m})$$

*(Classical semiparametrics:  $m = 1$  and  $\mathbb{U}_n = \mathbb{P}_n$ .)*

## Heuristics — plug in and bias correction

Estimate  $\theta := \chi(p) \in \mathbb{R}$  from iid  $X_1, \dots, X_n \sim p$ .

Given  $\hat{p}$  and “influence function”  $(x_1, \dots, x_m) \mapsto \chi_p(x_1, \dots, x_m)$  use

$$\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi \hat{p},$$

for

$$\mathbb{U}_n f = \frac{(n-m)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_m \leq n} \dots \sum f(X_{i_1}, \dots, X_{i_m})$$

*(Classical semiparametrics:  $m = 1$  and  $\mathbb{U}_n = \mathbb{P}_n$ .)*

What is a good influence function?

## Heuristics — plug in and bias correction

Estimate  $\theta := \chi(p) \in \mathbb{R}$  from iid  $X_1, \dots, X_n \sim p$ .

Given  $\hat{p}$  and “influence function”  $(x_1, \dots, x_m) \mapsto \chi_p(x_1, \dots, x_m)$  use

$$\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi \hat{p},$$

for

$$\mathbb{U}_n f = \frac{(n-m)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_m \leq n} \dots \sum f(X_{i_1}, \dots, X_{i_m})$$

*(Classical semiparametrics:  $m = 1$  and  $\mathbb{U}_n = \mathbb{P}_n$ .)*

What is a **good influence function** that works with a **general purpose  $p$** ?

## Heuristics — plug in and bias correction

Estimate  $\theta := \chi(p) \in \mathbb{R}$  from iid  $X_1, \dots, X_n \sim p$ .

Given  $\hat{p}$  and “influence function”  $(x_1, \dots, x_m) \mapsto \chi_p(x_1, \dots, x_m)$  use

$$\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi \hat{p},$$

for

$$\mathbb{U}_n f = \frac{(n-m)!}{n!} \sum_{1 \leq i_1 \neq \dots \neq i_m \leq n} \dots \sum f(X_{i_1}, \dots, X_{i_m})$$

*(Classical semiparametrics:  $m = 1$  and  $\mathbb{U}_n = \mathbb{P}_n$ .)*

What is a good influence function that works with a general purpose  $p$ ?

$\chi_p = 0$  gives plug-in  $\chi(\hat{p})$ . *Not good!*

## Heuristics — plug in and bias correction

If  $\theta = \chi(p)$  is estimated by  $\hat{\theta}_n = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$ , then

$$\hat{\theta}_n - \chi(p) = [\chi(\hat{p}_n) - \chi(p) + P^m \chi_{\hat{p}_n}] + (\mathbb{U}_n - P^n) \chi_{\hat{p}_n}.$$

## Heuristics — plug in and bias correction

If  $\theta = \chi(p)$  is estimated by  $\hat{\theta}_n = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$ , then

$$\hat{\theta}_n - \chi(p) = [\chi(\hat{p}_n) - \chi(p) + P^m \chi_{\hat{p}_n}] + (\mathbb{U}_n - P^n) \chi_{\hat{p}_n}.$$

Construct  $\chi_p$  such that  $-P^m \chi_{\hat{p}_n}$  “represents” the first  $m$  terms of the Taylor expansion of  $\chi(\hat{p}_n) - \chi(p)$ :

$$P^m \chi_p = 0$$
$$\frac{d^j}{dt^j} \Big|_{t=0} \chi(p_t) = -\frac{d^j}{dt^j} \Big|_{t=0} P^m \chi_{p_t}, \quad j = 1, \dots, m$$

for “smooth” one-dimensional submodels  $t \mapsto p_t$  with  $p_0 = p$ .



## Heuristics — plug in and bias correction

If  $\theta = \chi(p)$  is estimated by  $\hat{\theta}_n = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$ , then

$$\hat{\theta}_n - \chi(p) = [\chi(\hat{p}_n) - \chi(p) + P^m \chi_{\hat{p}_n}] + (\mathbb{U}_n - P^n) \chi_{\hat{p}_n}.$$

Construct  $\chi_p$  such that  $-P^m \chi_{\hat{p}_n}$  “represents” the first  $m$  terms of the Taylor expansion of  $\chi(\hat{p}_n) - \chi(p)$ :

$$P^m \chi_p = 0$$
$$\frac{d^j}{dt^j} \Big|_{t=0} \chi(p_t) = - \frac{d^j}{dt^j} \Big|_{t=0} P^m \chi_{p_t}, \quad j = 1, \dots, m$$

for “smooth” one-dimensional submodels  $t \mapsto p_t$  with  $p_0 = p$ .

*This translates into inner products of influence function and scores.*

## First-order estimator

For  $m = 1$  we find the influence function  $\chi_p$  from classical semiparametrics.

### META THEOREM

First-order estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  satisfies

$$\begin{aligned}\hat{\theta} - \chi(p) &= (\mathbb{U}_n - P)\chi_{\hat{p}} + [\chi(\hat{p}) - \chi(p) - (\hat{P} - P)\chi_{\hat{p}}] \\ &= O_P\left(\frac{1}{\sqrt{n}}\right) + O_P(\|\hat{p} - p\|^2)\end{aligned}$$

(Worst case scenario for bias)

## Example: missing data

Observe  $X = (YA, A, Z)$

Parameter  $p \leftrightarrow (a, b, f)$

First-order estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  satisfies

$$\begin{aligned} \hat{\theta} - \chi(p) &= (\mathbb{U}_n - P)\chi_{\hat{p}} - \left[ \int (\hat{a} - a)(\hat{b} - b) \frac{f}{a} \right] \\ &= O_P\left(\frac{1}{\sqrt{n}}\right) + O_P(\|\hat{a} - a\| \|\hat{b} - b\|) \end{aligned}$$

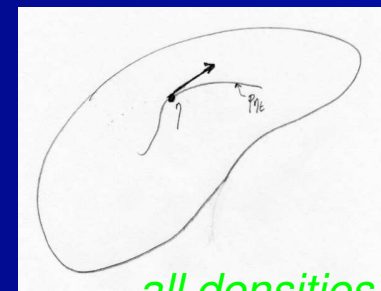
$a$	$b$	$f$	no bias if
$O_P(n^{-1/2})$	$O_P(1)$	—	$\dim(a) < \infty$
$O_P(1)$	$O_P(n^{-1/2})$	—	$\dim(b) < \infty$
$n^{-\alpha/(2\alpha+d)}$	$n^{-\alpha/(2\alpha+d)}$	—	$\alpha > d/2$
...	...	—	...

If  $Z$  of high dimension, then bias of linear estimator dominates variance.

## Higher-order tangent space and influence function

Tangent space of order  $m$  (at  $p$ ): all higher order score functions of one-dimensional submodels  $t \mapsto p_t$

$$g(x_1, \dots, x_m) = \frac{\frac{d^j}{dt^j} \Big|_{t=0} \prod_{i=1}^m p_t(x_i)}{\prod_{i=1}^m p(x_i)}, \quad j = 1, \dots, m$$



all densities  $p$

(These are  $U$ -statistics)

Influence function of order  $m$  of  $p \mapsto \chi(p)$  is map  $(x_1, \dots, x_m) \mapsto \chi_p(x_1, \dots, x_m)$  with for all submodels  $t \mapsto p_t$

$$\frac{d^j}{dt^j} \Big|_{t=0} \chi(p_t) = \frac{d^j}{dt^j} \Big|_{t=0} P^m \chi_p g, \quad j = 1, \dots, m$$

## Higher-order influence function — computation

Influence function  $\chi_p$  of parameter  $p \mapsto \chi(p)$  can be computed recursively from its **Hoeffding decomposition**

$$\mathbb{U}_n \chi_p = \mathbb{U}_n \chi_p^{(1)} + \frac{1}{2} \mathbb{U}_n \chi_p^{(2)} + \dots + \frac{1}{m!} \mathbb{U}_n \chi_p^{(m)}$$

- $\chi_p^{(1)}$  is a first order influence function of  $p \mapsto \chi(p)$
- $x_j \mapsto \chi_p^{(j)}(x_1, \dots, x_j)$  is a first order influence function of  $p \mapsto \chi_p^{(j-1)}(x_1, \dots, x_{j-1})$  ( $j = 2, \dots, m$ )

*(Optimal version may need projection in tangent space)*

## Higher-order estimator

Estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  with  $\chi_p$  an  $m$ th order influence function

$$\begin{aligned}\hat{\theta} - \chi(p) &= (\mathbb{U}_n - P^n) \chi_{\hat{p}} + [\chi(\hat{p}) - \chi(p) - (\hat{P}^m - P^m) \chi_{\hat{p}}] \\ &= O_P\left(\frac{1}{\sqrt{n}}\right) + O_P(\|\hat{p} - p\|^{m+1})\end{aligned}$$

## Higher-order estimator

Estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  with  $\chi_p$  an  $m$ th order influence function

$$\begin{aligned}\hat{\theta} - \chi(p) &= (\mathbb{U}_n - P^n) \chi_{\hat{p}} + [\chi(\hat{p}) - \chi(p) - (\hat{P}^m - P^m) \chi_{\hat{p}}] \\ &= O_P\left(\frac{1}{\sqrt{n}}\right) + O_P(\|\hat{p} - p\|^{m+1})\end{aligned}$$

Free lunch??

## Higher-order estimator

Estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  with  $\chi_p$  an  $m$ th order influence function

$$\begin{aligned}\hat{\theta} - \chi(p) &= (\mathbb{U}_n - P^n) \chi_{\hat{p}} + [\chi(\hat{p}) - \chi(p) - (\hat{P}^m - P^m) \chi_{\hat{p}}] \\ &= O_P\left(\frac{1}{\sqrt{n}}\right) + O_P(\|\hat{p} - p\|^{m+1})\end{aligned}$$

*Free lunch??*

*No! Higher order influence functions may not exist.*



## Higher-order estimator

Estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  with  $\chi_p$  an  $m$ th order influence function

$$\begin{aligned}\hat{\theta} - \chi(p) &= (\mathbb{U}_n - P^n) \chi_{\hat{p}} + [\chi(\hat{p}) - \chi(p) - (\hat{P}^m - P^m) \chi_{\hat{p}}] \\ &= O_P\left(\frac{1}{\sqrt{n}}\right) + O_P(\|\hat{p} - p\|^{m+1})\end{aligned}$$

*Free lunch??*

*No! Higher order influence functions may not exist.*

Must use *approximations* (representing derivative in selected directions).

### META META THEOREM

*m*th-order estimator  $\hat{\theta} = \chi(\hat{p}) + \mathbb{U}_n \chi_{\hat{p}}$  satisfies

$$\hat{\theta} - \chi(p) = (\mathbb{U}_n - P^n) \chi_{n,\hat{p}} + O_P(\|\hat{p} - p\|^{m+1}) + \text{approximation bias}$$

*One variance term and two bias terms.*

## Approximate functional

Choose a map  $p \mapsto \tilde{p}$  of the model onto a “smaller” model and consider

$$\tilde{\chi}(p) = \chi(\tilde{p}) + P\chi_{\tilde{p}}^{(1)}$$

Definition of  $\chi_p^{(1)}$  suggests

$$\tilde{\chi}(p) - \chi(p) = O(\|\tilde{p} - p\|^2)$$

Choose  $p \mapsto \tilde{p}$  such that, for any path  $t \mapsto p_t$ ,

$$\frac{d}{dt}\Big|_{t=0} \left( \chi(\tilde{p}_t) + P_0\chi_{\tilde{p}_t}^{(1)} \right) = 0$$

Then  $\tilde{\chi}_p^{(1)} = \chi_{\tilde{p}}^{(1)}$  and  $\tilde{\chi}$  ought to have influence functions to any order.

## Example: missing data — approximate functional

Define  $\tilde{\chi}(p) := \chi(\tilde{p})$ , for  $p \mapsto \tilde{p}$  given by projection onto  $L \subset L_2(g)$ :

$$(a, b, g) \mapsto (\tilde{a}, \tilde{b}, g) \in L \times L \times \{g\}$$

### THEOREM

For  $\Pi_{i,j} = \Pi_p(Z_i, Z_j)$  projection kernel on  $L \subset L_2(g)$ ,  $\tilde{Y} = A(Y - \tilde{b}(Z))$ ,  
 $\tilde{A} = A\tilde{a}(Z) - 1$ ,

$$\tilde{\chi}_p^{(1)}(X) = A\tilde{a}(Z)(Y - \tilde{b}(Z)) + \tilde{b}(Z) - \chi(\tilde{p})$$

$$\tilde{\chi}_p^{(2)}(X_1, X_2) = -2[\tilde{Y}_1 \Pi_{1,2} \tilde{A}_2]$$

$$\tilde{\chi}_p^{(3)}(X_1, X_2, X_3) = 6 \left[ \tilde{Y}_1 \Pi_{1,2} A_2 \Pi_{2,3} \tilde{A}_3 - \tilde{Y}_1 \Pi_{1,3} \tilde{A}_3 \right]$$

$$\begin{aligned} \tilde{\chi}_p^{(4)}(X_1, X_2, X_3, X_4) = & -24 \left[ \tilde{Y}_1 \Pi_{1,2} A_2 \Pi_{2,3} A_3 \Pi_{3,4} \tilde{A}_4 \right. \\ & \left. - \tilde{Y}_1 \Pi_{1,3} A_3 \Pi_{3,4} \tilde{A}_4 - \tilde{Y}_1 \Pi_{1,2} A_2 \Pi_{2,4} \tilde{A}_4 + \tilde{Y}_1 \Pi_{1,4} \tilde{A}_4 \right] \end{aligned}$$

etc.

## Projection kernel — von Mises calculus

A projection  $\Pi_g: L_2(g) \rightarrow L$  onto a finite-dimensional space can be represented as

$$\Pi_g h(x) = \int h(y) \Pi_g(x, y) g(y) d\nu(y)$$

$y \mapsto \Pi_g(x, y)$  restricted to  $L$  works as Dirac kernel at  $x$ , because equivalent:

- $\Pi_g h = h$
- $h \in L$
- $h(x) = \int h(y) \Pi_g(x, y) d\nu(y)$  a.e.  $x$

Representation would be exact if  $L = L_2(g)$ , i.e.  $\Pi_g$  is the “Dirac kernel on the diagonal”, but this does not exist.

## Example: missing data – parametric rate – $m$ th-order

### THEOREM

For  $\sup_x \Pi_p(x, x) \lesssim k$ ,

$$\hat{\mathbb{E}}_p \hat{\theta}_n - \chi(p) = O\left(\|\hat{a} - a\|_r \|\hat{b} - b\|_r \|\hat{g} - g\|_{(m-1)r/(r-2)}^{m-1}\right) \\ + O\left(\left\| (I - \Pi_p)(\hat{a} - a) \right\|_2 \left\| (I - \Pi_p)(\hat{b} - b) \right\|_2\right),$$

$$\hat{\text{var}}_p \hat{\chi}_n \leq \sum_{j=1}^m \frac{1}{\binom{n}{j}} c^j k^{j-1}.$$

If  $(\alpha + \beta)/2 \geq d/4$  obtain  $\sqrt{n}$ -rate by choosing

- large enough order  $m$ .
- $L$  optimal for approximation in Hölder spaces, of dimension  $k = n/(\log n)^2$ .
- $\hat{a}, \hat{b}, \hat{g}$  that attain uniform minimax rates  $(\log n/n)^{-\delta/(2\delta+d)}$ .

If  $(\alpha + \beta)/2 > d/4$  obtain even efficiency  $\sqrt{n}(\hat{\chi}_n - \chi(p) - \mathbb{P}_n \chi_p^{(1)}) \xrightarrow{P} 0$ .

## Example: missing data – parametric rate – $m$ th-order

### THEOREM

For  $\sup_x \Pi_p(x, x) \lesssim k$ ,

$$\hat{\mathbb{E}}_p \hat{\theta}_n - \chi(p) = O\left(\|\hat{a} - a\|_r \|\hat{b} - b\|_r \|\hat{g} - g\|_{(m-1)r/(r-2)}^{m-1}\right) \\ + O\left(\left\| (I - \Pi_p)(\hat{a} - a) \right\|_2 \left\| (I - \Pi_p)(\hat{b} - b) \right\|_2\right),$$

$$\hat{\text{var}}_p \hat{\chi}_n \leq \sum_{j=1}^m \frac{1}{\binom{n}{j}} c^j k^{j-1}.$$

If  $(\alpha + \beta)/2 \geq d/4$  obtain  $\sqrt{n}$ -rate by choosing

- large enough order  $m$ .
- $L$  optimal for approximation in Hölder spaces, of dimension  $k = n/(\log n)^2$ .
- $\hat{a}, \hat{b}, \hat{g}$  that attain uniform minimax rates  $(\log n/n)^{-\delta/(2\delta+d)}$ .

If  $(\alpha + \beta)/2 > d/4$  obtain even efficiency  $\sqrt{n}(\hat{\chi}_n - \chi(p) - \mathbb{P}_n \chi_p^{(1)}) \xrightarrow{P} 0$ .

*Linear estimator ( $m = 1$ ) works only if  $(\alpha + \beta)/2 \geq d/2$ .*

## Example: missing data — lower smoothness — 3rd-order IF

Leading part of 3rd-order part of 3rd-order influence function of  $\tilde{\chi}$  is

$$6\tilde{A}_1\Pi_p(Z_1, Z_2)A_2\Pi_p(Z_2, Z_3)\tilde{Y}_3.$$

Decompose, for  $k_{-1} = l_{-1} = 1$  and  $k_R \sim l_S \sim k$ ,

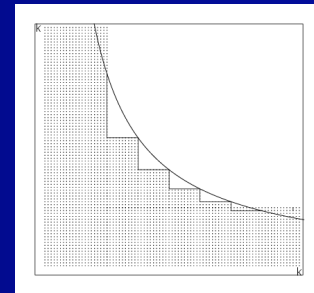
$$\Pi_p = \sum_{r=0}^R \Pi_p^{(k_{r-1}, k_r]}, \quad \Pi_p = \sum_{s=0}^S \Pi_p^{(l_{s-1}, l_s]}$$

and replace preceding display by

$$6 \sum_{\substack{(r,s): r+s \leq D \\ \vee r=0 \vee s=0}} \tilde{A}_1 \Pi_p^{(k_{r-1}, k_r]}(Z_1, Z_2) A_2 \Pi_p^{(l_{s-1}, l_s]}(Z_2, Z_3) \tilde{Y}_3.$$

$$k_r \sim n2^{r/\alpha}, \quad r = 0, \dots, R,$$

$$l_s \sim n2^{s/\beta}, \quad s = 0, \dots, S.$$



## Example: missing data — lower smoothness — higher order

Replace

$$\tilde{Y}_1 \Pi_{1,2} A_2 \Pi_{2,3} A_3 \times \cdots \times A_{j-1} \Pi_{j-1,j} \tilde{A}_j$$

by

$$\begin{aligned} & \sum_{i=1}^{j-2} j! (-1)^{j-1} \tilde{Y}_1 \Pi_{1,2}^{(0,n]} A_2 \times \cdots \times A_{i-1} \Pi_{i-1,i}^{(0,n]} A_i \times \\ & \quad \times \left[ \sum_{\substack{(r,s): r+s \leq D \\ \vee r=0 \vee s=0}} \Pi_{i,i+1}^{(k_{r-1}, k_r]} A_{i+1} \Pi_{i+1,i+2}^{(l_{s-1}, l_s]} \right] \times \\ & \quad \times A_{i+2} \Pi_{i+2,i+3}^{(0,n]} \times \cdots \times A_{j-1} \Pi_{j-1,j}^{(0,n]} \tilde{A}_j. \end{aligned}$$



## Example: missing data — lower smoothness — higher order

### THEOREM

For  $\sup_x \Pi_{\hat{p}}^{(0,l]}(x, x) \lesssim l$ , the pruned  $m$ th order estimator satisfies (for  $r \geq 2$ )

$$\begin{aligned} \hat{E}_p \hat{\chi}_n - \chi(p) &= O\left(\|\hat{a} - a\|_r \|\hat{b} - b\|_r \|\hat{g} - g\|_{\frac{mr}{r-2}}^{m-1}\right) \\ &\quad + O\left(\left\| (I - \Pi_p^{(0,k]}) (\hat{a} - a) \right\|_2 \left\| (I - \Pi_p^{(0,k]}) (\hat{b} - b) \right\|_2\right), \\ &\quad + O\left(\sum_{r=1}^R \left\| (I - \Pi_{\hat{p}}^{(0,k_{r-1}]}) (\hat{a} - a) \right\|_r \left\| (I - \Pi_{\hat{p}}^{(0,l_{D-r}]}) (\hat{b} - b) \right\|_r \|\hat{g} - g\|_{\frac{r}{r-2}}\right) \\ &\quad + O\left(R \left\| (I - \Pi_{\hat{p}}^{(0,n]}) (\hat{a} - a) \right\|_r \left\| (I - \Pi_{\hat{p}}^{(0,n]}) (\hat{b} - b) \right\|_r \|\hat{g} - g\|_{\frac{mr}{r-2}}^2\right), \\ \hat{\text{var}}_p \hat{\chi}_n &\lesssim \frac{1}{n} + \frac{k}{n^2} + \frac{D 2^{(\frac{1}{\alpha} \vee \frac{1}{\beta})D}}{n}. \end{aligned}$$

If  $\phi > \phi(\alpha, \beta)$  obtain rate  $n^{-(2\alpha+2\beta)/(2\alpha+2\beta+d)}$ ,

with sufficiently large  $m$ , suitable  $D$  and suitable initial estimators.

$$\phi(\alpha, \beta) = (\alpha/d \vee \beta/d)(d - 2\alpha - 2\beta)/(d + 2\alpha + 2\beta)$$

## Lower Bounds

## Classical semiparametrics

In classical semiparametrics the rate of estimation is  $\sqrt{n}$ .

The best limit distribution of  $\sqrt{n}(T_n - \chi(p))$  is normal, with variance equal to the **inverse efficient Fisher information**.

## Slow rates — testing argument (Le Cam)

$X_1, X_2, \dots, X_n$  i.i.d. sample from density  $p \in \mathcal{P}$ .

### THEOREM

If  $P_n$  and  $Q_n$  are in the convex hulls of the sets of measures  $\{P^n: p \in \mathcal{P}, \chi(p) \leq 0\}$  and  $\{P^n: p \in \mathcal{P}, \chi(p) \geq \varepsilon_n\}$ , and

$$\rho(P_n, Q_n) := \int \sqrt{dP_n} \sqrt{dQ_n} \gg 0$$

then the rate is not faster than  $\varepsilon_n$ .

Nontrivial details:

- find the **least favourable**  $P_n$  and  $Q_n$ .
- compute their **Hellinger affinity**  $\rho(P_n, Q_n)$ .

## Affinity bound (Birgé and Massart (1995), Robins and vdV (2008))

- Partition  $\mathcal{X} = \cup_{j=1}^k \mathcal{X}_j$ .
- Perturbation parameter  $\lambda = (\lambda_1, \dots, \lambda_k)$  with prior  $\pi = \pi_1 \otimes \dots \otimes \pi_k$ .
- $P_\lambda$  and  $Q_\lambda$  probability measures on  $\mathcal{X}$  such that restrictions  $P_{\lambda|\mathcal{X}_j}$  and  $Q_{\lambda|\mathcal{X}_j}$  depend on  $\lambda_j$  only and have equal mass  $p_j$ .

**THEOREM** If  $np_j(1 \vee a \vee b) \lesssim 1$  and  $0 \lesssim p_\lambda \lesssim 1$ , then

$$\rho\left(\int P_\lambda^n d\pi(\lambda), \int Q_\lambda^n d\pi(\lambda)\right) \geq 1 - Cn^2(\max_j p_j)(b^2 + ab) - Cnd.$$

$$p = \int p_\lambda d\pi(\lambda)$$

$$q = \int q_\lambda d\pi(\lambda)$$

$$a = \max_j \sup_\lambda \int_{\mathcal{X}_j} \frac{(p_\lambda - p)^2}{p_\lambda} \frac{d\nu}{p_j},$$

$$b = \max_j \sup_\lambda \int_{\mathcal{X}_j} \frac{(q_\lambda - p_\lambda)^2}{p_\lambda} \frac{d\nu}{p_j},$$

$$d = \max_j \sup_\lambda \int_{\mathcal{X}_j} \frac{(q - p)^2}{p_\lambda} \frac{d\nu}{p_j}.$$

## Example: missing data

Covariate  $Z$ ,  $\nu \sim f$

Response  $Y$ , with  $Y|Z \sim \text{binomial}(1, b(Z))$

Missingness indicator  $A$ , with  $A|Z \sim \text{binomial}(1, 1/a(Z))$

Missing at random:  $Y \perp\!\!\!\perp A|Z$

Observe  $X = (YA, A, Z) \in \{0, 1\} \times \{0, 1\} \times [0, 1]^d$

We wish to estimate mean response  $\chi(a, b, f) = \int b f d\nu = \mathbf{E}Y$ .

### THEOREM

If  $a$ ,  $b$ , and  $g$  belong to Hölder classes  $C^\alpha[0, 1]^d$ ,  $C^\beta[0, 1]^d$ ,  $C^\gamma[0, 1]^d$ , then the rate of estimation is not faster than  $n^{-(\alpha+\beta)/(2\alpha+2\beta+d)}$ .

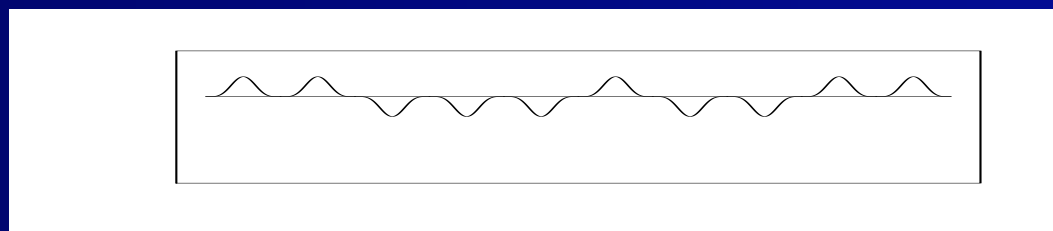
## Perturbations

- $H: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $C^\infty$ , support  $\subset [0, 1/2]^d$ ,  $\int H d\nu = 0$ .
- $k \sim \eta^{2d/(2\alpha+2\beta+d)}$ .
- $\mathcal{X}_j = \{0, 1\} \times \{0, 1\} \times \mathcal{Z}_j$ , for  $\mathcal{Z}_j$  disjoint translates of  $k^{-1/d}[0, 1/2]^d$ .
- $\pi$  uniform on  $\Lambda := \{0, 1\}^k$ .

For  $\lambda = (\lambda_1, \dots, \lambda_k) \in \Lambda$ :

$$a_\lambda(z) = 2 + \left(\frac{1}{k}\right)^{\alpha/d} \sum_{j=1}^k \lambda_j H((z - z_j)k^{1/d}),$$

$$b_\lambda(z) = \frac{1}{2} + \left(\frac{1}{k}\right)^{\beta/d} \sum_{j=1}^k \lambda_j H((z - z_j)k^{1/d}).$$



## Perturbations–2

- $\alpha \leq \beta$ :  $p_\lambda \leftrightarrow (a_\lambda, 1/2, 1/2)$  and  $q_\lambda \leftrightarrow (a_\lambda, b_\lambda, 1/2)$ .
- $\alpha \geq \beta$ :  $p_\lambda \leftrightarrow (2, b_\lambda, 1/2)$  and  $q_\lambda \leftrightarrow (a_\lambda, b_\lambda, 1/2)$ .

This leads to comparing the functional  $\chi(a, b, g)$  on two mixtures, where

- the first mixture  $\int P_\lambda^n d\pi(\lambda)$  perturbs only the **coarsest** of the two parameters  $a$  and  $b$ .
- the second mixture  $\int Q_\lambda^n d\pi(\lambda)$  perturbs **both** parameters.

*(The third parameter  $g$  is always taken 1/2. )*



## Concluding remarks

## Outlook

Adaptation to  $\alpha$  and  $\beta$ .

Implementation.

Prior parameter classes defined by sparsity.

Other models, e.g. semiparametric regression.