

## Lec 8: Introduction to semiparametric models

Yanjin Han

Nov 7, 2023



Previous lectures: parametric models  $y_1, \dots, y_n \sim P_\theta$ ,  
 $\theta \in \mathbb{R}^p$  is finite dimensional

This lecture: semiparametric models  $y_1, \dots, y_n \sim P_{\theta, \eta}$ :

$\theta$ : target parameter (typically finite-dimensional)  
 $\eta$ : nuisance parameter (could be infinite-dimensional)

Historic remark: symmetric location family  
(by C. Stein, "efficient nonparametric testing & estimation", 1956)

Model:  $y_1, \dots, y_n \sim f(\cdot - \theta)$ , where

- $\theta \in \mathbb{R}$ : target location parameter
- $f$ : unknown density symmetric around zero  
(nuisance) (i.e.  $f(x) = f(-x)$ )

Estimators for  $\theta$ :

1. sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

$$\begin{aligned}\mathbb{E}_\theta[y_1] &= \int y f(y - \theta) dy = \theta + \int (y - \theta) f(y - \theta) dy = \theta \\ \Rightarrow \mathbb{E}_\theta[\bar{y}] &= \theta \\ \text{Var}_\theta(\bar{y}) &= \frac{1}{n} \int y^2 f(y) dy.\end{aligned}$$

2. efficient estimator with known  $f$ : MLE

$$\hat{\theta}^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log f(y_i - \theta) \Rightarrow \frac{1}{n} \sum_{i=1}^n \frac{f'(y_i - \hat{\theta}^{\text{MLE}})}{f(y_i - \hat{\theta}^{\text{MLE}})} = 0.$$

Using Fisher info, one can show that

$$\mathbb{E}_\theta[(\hat{\theta}^{\text{MLE}} - \theta)^2] = \frac{1 + o_p(1)}{n} \left( \int \frac{f'(y)^2}{f(y)} dy \right)^{-1}$$

↑  
asymptotically optimal MSE

3. What about unknown  $f$ ?

Stein (1956) showed that if we use some nonparametric procedures to estimate  $f$  by  $\hat{f}$ , then estimate  $\theta$  by

$$\hat{\theta} : \quad \frac{1}{n} \sum_{i=1}^n \frac{\hat{f}'(y_i - \hat{\theta})}{\hat{f}(y_i - \hat{\theta})} = 0 \quad (\text{plug-in approach})$$

then

$$\mathbb{E}_\theta[(\hat{\theta} - \theta)^2] = \frac{1 + o_p(1)}{n} \left( \int \frac{f'(y)^2}{f(y)} dy \right)^{-1}$$

semiparametric efficient!

(the same asymptotic efficiency can be achieved without knowing the nuisance; NOT all semiparametric problems admit semiparametric efficient estimators)

Key ideas behind semiparametric models:

- do not want to propose a restrictive model for the nuisance;
- hope that even if the nuisance estimation error is large, the target estimation error is still small;
- orthogonality will play a central role.

### Examples . 1. Linear regression.

$$Y = X\theta_0 + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0$$

Target :  $\theta_0 \in \mathbb{R}^p$

Nuisance : distribution of  $\varepsilon$

(Remark: we do not assume that  $\varepsilon \sim N(0, \sigma^2)$ ,  
nor the independence of  $(X, \varepsilon)$  )

### 2. Partial linear regression:

$$\begin{cases} Y = D\theta_0 + g_0(X) + \varepsilon_1, & \mathbb{E}[\varepsilon_1|X, D] = 0 \\ D = m_0(X) + \varepsilon_2, & \mathbb{E}[\varepsilon_2|X] = 0 \end{cases}$$

Data :  $(X_i, D_i, Y_i)$

Target :  $\theta_0$

Nuisance :  $(g_0, m_0, \text{distributions of } (\varepsilon_1, \varepsilon_2))$

(closely related to the potential outcome model in causal inference next lecture)

### 3. Errors in variables :

$$\begin{cases} Y = \alpha + \beta Z + \varepsilon_1, & \varepsilon_1 \sim N(0, \sigma_1^2) \\ X = Z + \varepsilon_2, & \varepsilon_2 \sim N(0, \sigma_2^2) \end{cases}$$

Data :  $(X_i, Y_i)$

Target :  $(\alpha, \beta)$

Nuisance : distribution of  $Z$ .

### 4. Cox model :

$$h(t|x) = e^{\beta^T x} h(t)$$

Target :  $\beta$

Nuisance : baseline hazard  $h$ .

Estimation. Joint/profile MLE: given  $y_1, \dots, y_n \sim P_{\theta, \eta}(y)$ , compute

$$(\hat{\theta}, \hat{\eta}) = \underset{(\theta, \eta)}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log P_{\theta, \eta}(y_i)$$

$$\text{or } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \left( \max_{\eta} \frac{1}{n} \sum_{i=1}^n \log P_{\theta, \eta}(y_i) \right).$$

Sometimes works (e.g. in Cox model), but in many cases computationally infeasible.

A simplified question:

Suppose we are given a (possibly coarse) estimator  $\hat{\eta}$  of  $\eta$ .  
How should we use  $\hat{\eta}$  to estimate  $\theta$ ?

Score function & estimating equation.

Score. For  $y \sim P_{\theta_0}$ , the score of  $y$  at  $\theta_0$  is

$$s_{\theta_0}(y) = \nabla_{\theta} \log P_{\theta}(y) \Big|_{\theta = \theta_0}.$$

Relationships between score and MLE.

For  $y_1, \dots, y_n \sim P_{\theta}$ , the MLE for  $\theta$  is

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y_i)$$

$$\xRightarrow{\text{F.O.C.}} 0 = \nabla_{\theta} \left[ \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y_i) \right] \Big|_{\theta = \hat{\theta}} = \frac{1}{n} \sum_{i=1}^n s_{\hat{\theta}}(y_i)$$

(estimating eqn. for MLE)

Another interpretation.

Lemma.  $\mathbb{E}_{\theta_0}[s_{\theta_0}(y)] = 0$  for all  $\theta_0$ .

Pf. 
$$\begin{aligned}\mathbb{E}_{\theta_0}[s_{\theta_0}(y)] &= \mathbb{E}_{\theta_0}[\nabla_{\theta} \log p_{\theta}(y) |_{\theta=\theta_0}] \\ &= \mathbb{E}_{\theta_0}\left[\frac{\nabla_{\theta} p_{\theta}(y) |_{\theta=\theta_0}}{p_{\theta_0}(y)}\right] \\ &= \int \cancel{p_{\theta_0}(y)} \frac{\nabla_{\theta} p_{\theta}(y) |_{\theta=\theta_0}}{\cancel{p_{\theta_0}(y)}} dy \\ &= \nabla_{\theta} \underbrace{\int p_{\theta}(y) dy}_{=1} |_{\theta=\theta_0} = 0 \quad \square\end{aligned}$$

View estimating equation in terms of score matching:

$$\underbrace{\frac{1}{n} \sum_{i=1}^n s_{\hat{\theta}}(y_i)}_{\text{empirical score at } \hat{\theta}} = 0 \quad \downarrow \\ = \mathbb{E}_{\theta_0}[s_{\theta_0}(y)] \text{ is true score at } \theta_0$$

(intuition: solve for  $\theta_0$  from

$$0 = \mathbb{E}_{\theta_0}[s_{\theta_0}(y)] \approx \frac{1}{n} \sum_{i=1}^n s_{\theta_0}(y_i). \quad )$$

General estimating equation.

1. find  $f(\theta, y) \in \mathbb{R}^p$  s.t.  $\mathbb{E}_{\theta_0}[f(\theta_0, y)] = 0$
2. estimate  $\theta_0$  by  $\hat{\theta}$  from the estimating eqn.

$$\frac{1}{n} \sum_{i=1}^n f(\hat{\theta}, y_i) = 0.$$

Example 1.  $(x_1, y_1), \dots, (x_n, y_n) \sim N\left(\begin{bmatrix} \theta_0 \\ \eta_0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$   
 unknown:  $(\theta_0, \eta_0)$  known:  $\rho$ .

$$\log p_{\theta, \eta}(x, y) = \text{const} - \frac{(x - \theta)^2 + (y - \eta)^2 - 2\rho(x - \theta)(y - \eta)}{2(1 - \rho^2)}$$

$$S_{\theta_0, \eta_0}(x, y) = \begin{bmatrix} \nabla_{\theta} \log p_{\theta, \eta}(x, y) \Big|_{\substack{\theta = \theta_0 \\ \eta = \eta_0}} \\ \nabla_{\eta} \log p_{\theta, \eta}(x, y) \Big|_{\substack{\theta = \theta_0 \\ \eta = \eta_0}} \end{bmatrix} = \frac{1}{1 - \rho^2} \begin{bmatrix} x - \theta_0 - \rho(y - \eta_0) \\ y - \eta_0 - \rho(x - \theta_0) \end{bmatrix}$$

MLE estimating equation:

$$\begin{cases} \frac{1}{n} \sum_{i=1}^n [(x_i - \hat{\theta}) - \rho(y_i - \hat{\eta})] = 0 \\ \frac{1}{n} \sum_{i=1}^n [(y_i - \hat{\eta}) - \rho(x_i - \hat{\theta})] = 0 \end{cases}$$

Example 2

$$y_i = \langle \theta_0, x_i \rangle + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i | x_i] = 0, \quad i = 1, \dots, n.$$

Let  $f(\theta, (x, y)) = (y - \langle \theta, x \rangle)x \in \mathbb{R}^p$ , then

$$\begin{aligned} \mathbb{E}_{\theta_0}[f(\theta_0, (x, y))] &= \mathbb{E}_{\theta_0}[(y - \langle \theta_0, x \rangle)x] \\ &= \mathbb{E}_{\theta_0}[\varepsilon x] = \mathbb{E}_{\theta_0}[\mathbb{E}[\varepsilon | x]x] = 0 \end{aligned}$$

$\Rightarrow$  estimating eqn:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \langle \hat{\theta}, x_i \rangle) x_i = 0$$

$$\Rightarrow \hat{\theta} = (X^T X)^{-1} X^T Y \quad (\text{least squares})$$

Question in semiparametric models for example 1:

If  $\eta_0$  is a nuisance parameter and  $\hat{\eta}$  is given to us, which estimating eqn. should we use?

## Efficient score function

Let  $y \sim P_{\theta_0, \eta_0}$  in a semiparametric model with target  $\theta_0$  and nuisance  $\eta_0$  (for simplicity we assume  $\theta, \eta \in \mathbb{R}$ )

$$\text{Score function } s_{\theta, \eta}(y) = \begin{bmatrix} s_{\theta, \eta}^{\theta}(y) \\ s_{\theta, \eta}^{\eta}(y) \end{bmatrix} = \begin{bmatrix} \nabla_{\theta} \log p_{\theta, \eta}(y) \\ \nabla_{\eta} \log p_{\theta, \eta}(y) \end{bmatrix} \bigg|_{\substack{\theta = \theta_0 \\ \eta = \eta_0}}$$

Efficient score function for  $\theta_0$ :

$$s_{\theta, \eta}^{\text{eff}}(y) = s_{\theta, \eta}^{\theta}(y) - \frac{\mathbb{E}_{\theta_0, \eta_0}[s_{\theta, \eta}^{\theta}(y) s_{\theta, \eta}^{\eta}(y)]}{\mathbb{E}_{\theta_0, \eta_0}[s_{\theta, \eta}^{\eta}(y)^2]} s_{\theta, \eta}^{\eta}(y)$$

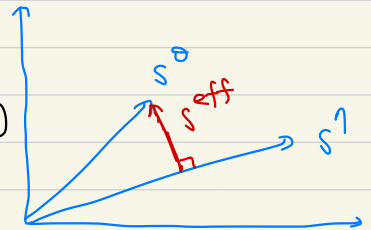
Estimating eqn. for  $\theta_0$ : given  $\hat{\eta}$ , solve

$$\frac{1}{n} \sum_{i=1}^n s_{\hat{\theta}, \hat{\eta}}^{\text{eff}}(y_i) = 0 \implies \hat{\theta}$$

Geometric interpretation of  $s^{\text{eff}}$ :

Gram-Schmidt orthogonalization of  $s^{\theta}$  with respect to  $s^{\eta}$  in  $L^2(P_{\theta_0, \eta_0})$

("orthogonalization")



Example 1 (continued)

$$s^{\theta}(x, y) = \frac{1}{1-\rho^2} [(x-\theta_0) - \rho(y-\eta_0)]$$

$$s^{\eta}(x, y) = \frac{1}{1-\rho^2} [(y-\eta_0) - \rho(x-\theta_0)]$$

$$\mathbb{E}_{\theta_0, \eta_0}[s^{\theta}(x, y) s^{\eta}(x, y)] = \frac{1}{(1-\rho^2)^2} [(1+\rho^2)\rho - 2\rho] = -\frac{\rho}{1-\rho^2}$$

$$\mathbb{E}_{\theta_0, \eta_0}[s^{\eta}(x, y)^2] = \frac{1}{(1-\rho^2)^2} [1+\rho^2-2\rho^2] = \frac{1}{1-\rho^2}$$



$$\begin{aligned}
 s^{\text{eff}}(x, y) &= s^{\theta}(x, y) - \frac{\rho}{1-\rho^2} s^{\eta}(x, y) \\
 &= s^{\theta}(x, y) + \rho s^{\eta}(x, y) \\
 &= \frac{1}{1-\rho^2} [(x-\theta_0) - \rho(y-\eta_0) + \rho(y-\eta_0) - \rho^2(x-\theta_0)] \\
 &= x - \theta_0.
 \end{aligned}$$

Estimating eqn. based on efficient score:

$$\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\theta}) = 0 \Rightarrow \hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i.$$

(independent of  $\eta$ )

Example 3 (Stein's symmetric location model)

$y \sim f_{\eta_0}(\cdot - \theta_0)$  ( $f$  symmetric around zero;  
assumed to be parametrized by  $\eta$ .)

$$s_{\theta_0, \eta_0}^{\theta}(y) = \frac{\partial}{\partial \theta} \log f_{\eta}(y - \theta) \Big|_{\theta = \theta_0, \eta = \eta_0} = - \frac{f'_{\eta}(y - \theta_0)}{f_{\eta}(y - \theta_0)}$$

$$s_{\theta_0, \eta_0}^{\eta}(y) = \frac{\partial}{\partial \eta} \log f_{\eta}(y - \theta) \Big|_{\theta = \theta_0, \eta = \eta_0} = \frac{1}{f_{\eta}(y - \theta_0)} \cdot \frac{\partial}{\partial \eta} f_{\eta}(y - \theta_0) \Big|_{\eta = \eta_0}$$

$$\mathbb{E}_{\theta_0, \eta_0} [s_{\theta_0, \eta_0}^{\theta}(y) s_{\theta_0, \eta_0}^{\eta}(y)] = \mathbb{E}_{\theta_0, \eta_0} \left[ - \frac{f'_{\eta}(y - \theta_0)}{f_{\eta}(y - \theta_0)^2} \frac{\partial}{\partial \eta} f_{\eta}(y - \theta_0) \Big|_{\eta = \eta_0} \right]$$

↖ anti-symmetric around  $\theta_0$   
↑ symmetric around  $\theta_0$

$$\Rightarrow s^{\text{eff}}(y) = s_{\theta_0, \eta_0}^{\theta}(y) \stackrel{=0}{=} - \frac{f'_{\eta}(y - \theta_0)}{f_{\eta_0}(y - \theta_0)}$$

Estimating eqn: based on  $\hat{f} = f_{\eta}$ , solve  $\hat{\theta}$  from

$$\frac{1}{n} \sum_{i=1}^n \frac{\hat{f}'(y_i - \hat{\theta})}{\hat{f}(y_i - \hat{\theta})} = 0 \quad (\text{Stein's estimator})$$

## Why efficient score?

Neyman orthogonality: an estimating eqn.  $f(\theta, \eta, y)$  is Neyman orthogonal iff

$$\mathbb{E}_{\theta, \eta} [\nabla_{\eta} f(\theta, \eta, y) |_{\eta=\eta_0}] = 0$$

Insights: Neyman orthogonal

$$\Rightarrow \mathbb{E}_{\theta, \eta} [f(\theta, \hat{\eta}, y)] \approx \mathbb{E}_{\theta, \eta} [f(\theta, \eta_0, y)]$$

↑ Taylor expansion around  $\hat{\eta} \approx \eta_0$   
= 0  
↳ requirement of estimating eqn.

(i.e. nuisance estimation error has second-order effects).

Thm: efficient scores are Neyman orthogonal.

Pf (optional):  $\nabla_{\eta} s_{\theta, \eta}^{\text{eff}}(y) = \nabla_{\eta} [s_{\theta, \eta}^0(y) - \alpha(\theta, \eta) s_{\theta, \eta}^1(y)]$   
 $(\alpha(\theta, \eta) = \frac{\mathbb{E}_{\theta, \eta} [s_{\theta, \eta}^0(y) s_{\theta, \eta}^1(y)]}{\mathbb{E}_{\theta, \eta} [s_{\theta, \eta}^1(y)^2]})$

$$= \nabla_{\eta} s_{\theta, \eta}^0(y) - \alpha(\theta, \eta) \nabla_{\eta} s_{\theta, \eta}^1(y)$$

can show:

$$\mathbb{E}_{\theta, \eta} [\nabla_{\eta} s_{\theta, \eta}^0(y)] = - \mathbb{E}_{\theta, \eta} [s_{\theta, \eta}^0(y) s_{\theta, \eta}^1(y)]$$

$$\mathbb{E}_{\theta, \eta} [\nabla_{\eta} s_{\theta, \eta}^1(y)] = - \mathbb{E}_{\theta, \eta} [s_{\theta, \eta}^1(y)^2] \quad (*)$$

$$\Rightarrow \mathbb{E}_{\theta, \eta} [\cdot] = 0 \text{ by defn. of } \alpha(\theta, \eta)$$

$$- \underbrace{\nabla_{\eta} \alpha(\theta, \eta)}_{\substack{\text{does not} \\ \text{depend on } \eta}} \cdot \underbrace{s_{\theta, \eta}^1(y)}_{\substack{\text{has expectation zero}}} \\ \mathbb{E}_{\theta, \eta} [\cdot] = 0$$

Proof of (\*):  $0 = \nabla_{\eta} \mathbb{E}_{\theta, \eta} [s_{\theta, \eta}^0(y)]$

$$= \nabla_{\eta} \int p_{\theta, \eta}(y) s_{\theta, \eta}^0(y) dy$$

$$= \int (\nabla_{\eta} p_{\theta, \eta}(y) \cdot s_{\theta, \eta}^0(y) + p_{\theta, \eta}(y) \cdot \nabla_{\eta} s_{\theta, \eta}^0(y)) dy$$

$$= \mathbb{E}_{\theta, \eta} [s_{\theta, \eta}^0(y) s_{\theta, \eta}^1(y) + \nabla_{\eta} s_{\theta, \eta}^0(y)].$$

□

