# Lec 8: Introduction to semiparametric models

Yanjun Han

Nov 5. 2024

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

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

$\theta$: 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: $y_1, \cdots, y_n \sim f(\cdot - \theta)$, 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 $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

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

2. efficient estimator with known $f$: $\boxed{\text{MLE}}$

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

Using Fisher info, one can show that

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

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

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

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.

<u>Examples</u>. 1. Linear regression.
$$Y = X\theta_0 + \varepsilon, \qquad \mathbb{E}[\varepsilon \mid x] = 0$$
Target : $\theta_0 \in \mathbb{R}^r$
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 :
$$\begin{cases} Y = D\theta_0 + g_0(x) + \varepsilon_1, & \mathbb{E}[\varepsilon_1 \mid x, D] = 0 \\ D = m_0(x) + \varepsilon_2, & \mathbb{E}[\varepsilon_2 \mid x] = 0 \end{cases}$$
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 :
$$h(t \mid x) = e^{\beta^T x} h(t)$$
Target : $\beta$
Nuisance : baseline hazard $h$.

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

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

or $\qquad \hat{\theta} = \underset{\theta}{\arg\max} \left( \underset{\eta}{\max} \ \frac{1}{n} \sum_{i=1}^{n} \log P_{\theta, \eta}(y_i) \right).$

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

<div style="border:1px solid red; padding:8px">

A simplified question:

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

</div>

<u>Score function & estimating equation</u>.

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

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

<u>Relationships between score and MLE</u>.

For $y_1, \cdots, y_n \sim P_\theta$, the MLE for $\theta$ is

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

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

<span style="color:red">(estimating eqn. for MLE)</span>

## Another interpretation.

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

**Pf.**
$$\mathbb{E}_{\theta_0}[s_{\theta_0}(y)] = \mathbb{E}_{\theta_0}\left[ \nabla_\theta \log p_\theta(y) \big|_{\theta=\theta_0} \right]$$
$$= \mathbb{E}_{\theta_0}\left[ \frac{\nabla_\theta p_\theta(y)|_{\theta=\theta_0}}{p_{\theta_0}(y)} \right]$$
$$= \int p_{\theta_0}(y) \frac{\nabla_\theta p_\theta(y)|_{\theta=\theta_0}}{p_{\theta_0}(y)} \, dy$$
$$= \nabla_\theta \underbrace{\int p_\theta(y) \, dy}_{=1} \bigg|_{\theta=\theta_0} = 0 \qquad \square$$

View estimating equation in terms of score matching:

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

( intuition: solve for $\theta_0$ from
$$0 = \mathbb{E}_{\theta_0}[s_{\theta_0}(y)] \approx \frac{1}{n} \sum_{i=1}^{n} s_{\theta_0}(y_i). \qquad )$$

## 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.** $(x_1, y_1), \cdots, (x_n, y_n) \sim N\left(\begin{bmatrix} \theta_0 \\ \eta_0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$

$\qquad\qquad$ unknown: $(\theta_0, \eta_0)$ $\qquad$ 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:

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

**Example 2** $\qquad\qquad y_i = \langle \theta_0, x_i \rangle + \varepsilon_i. \qquad \mathbb{E}[\varepsilon_i | x_i] = 0, \quad i = 1, \cdots, n.$

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

$$\mathbb{E}_{\theta_0}[f(\theta_0, (x,y))] = \mathbb{E}_{\theta_0}[(y - \langle \theta_0, x \rangle)x]$$

$$= \mathbb{E}_{\theta_0}[\varepsilon x] = \mathbb{E}_{\theta_0}[\mathbb{E}[\varepsilon | x]x] = 0$$

$\implies$ estimating eqn:

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

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

---

Question in semiparametric models for example 1:
$\qquad$ 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 $\theta_0, \eta_0 \in \mathbb{R}$ )

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

Efficient score function for $\theta_0$:

$$s^{eff}_{\theta_0, \eta_0}(y) = s^{\theta}_{\theta_0, \eta_0}(y) - \frac{\mathbb{E}_{\theta_0, \eta_0}\left[ s^{\theta}_{\theta_0, \eta_0}(y) \, s^{\eta}_{\theta_0, \eta_0}(y) \right]}{\mathbb{E}_{\theta_0, \eta_0}\left[ s^{\eta}_{\theta_0, \eta_0}(y)^2 \right]} \, s^{\eta}_{\theta_0, \eta_0}(y)$$
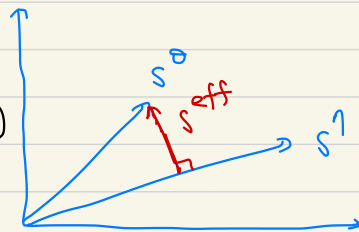
Estimating eqn. for $\theta_0$: given $\hat\eta$, solve

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

Geometric interpretation of $s^{eff}$:

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

( "orthogonalization" )



## Example 1 (continued).

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

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

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

$$S^{eff}(x,y) = S^\theta(x,y) - \frac{-\frac{\rho}{1-\rho^2}}{1/(1-\rho^2)} S^\eta(x,y)$$

$$= S^\theta(x,y) + \rho S^\eta(x,y)$$

$$= \frac{1}{1-\rho^2}\left[(x-\theta_o) - \rho(y-\eta_o) + \rho(y-\eta_o) - \rho^2(x-\theta_o)\right]$$

$$= x - \theta_o$$

Estimating eqn. based on efficient score:

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

(independent of $\hat\eta$)

## Example 3 (Stein's symmetric location model)

$$y \sim f_{\eta_o}(\cdot - \theta_o) \quad (f \text{ symmetric around zero;}$$
assumed to be parametrized by $\eta_o$.)

$$S^\theta_{\theta_o,\eta_o}(y) = \frac{\partial}{\partial\theta}\log f_\eta(y-\theta)\Big|_{\substack{\theta=\theta_o\\\eta=\eta_o}} = -\frac{f'_{\eta_o}(y-\theta_o)}{f_{\eta_o}(y-\theta_o)}$$

$$S^\eta_{\theta_o,\eta_o}(y) = \frac{\partial}{\partial\eta}\log f_\eta(y-\theta)\Big|_{\substack{\theta=\theta_o\\\eta=\eta_o}} = \frac{1}{f_{\eta_o}(y-\theta_o)}\cdot\frac{\partial}{\partial\eta}f_\eta(y-\theta_o)\Big|_{\eta=\eta_o}$$

anti-symmetric around $\theta_o$

$$\mathbb{E}_{\theta_o,\eta_o}\left[S^\theta_{\theta_o,\eta_o}(y)S^\eta_{\theta_o,\eta_o}(y)\right] = \mathbb{E}_{\theta_o,\eta_o}\left[-\frac{f'_{\eta_o}(y-\theta_o)}{f_{\eta_o}(y-\theta_o)^2}\frac{\partial}{\partial\eta}f_\eta(y-\theta_o)\Big|_{\eta=\eta_o}\right]$$

symmetric around $\theta_o$

$$= 0$$

$$\implies S^{eff}(y) = S^\theta_{\theta_o,\eta_o}(y) = -\frac{f'_{\eta_o}(y-\theta_o)}{f_{\eta_o}(y-\theta_o)}$$

Estimating eqn: based on $\hat f = f_{\hat\eta}$, 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?

<div style="border:1px solid red; padding:6px;">

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

$$\mathbb{E}_{\theta_0, \eta_0}\left[\nabla_\eta f((\theta_0, \eta), y)\big|_{\eta = \eta_0}\right] = 0$$

</div>

**Insights**: Neyman orthogonal

Taylor expansion around $\hat{\eta} \approx \eta_0$

$$\Rightarrow \mathbb{E}_{\theta_0, \eta_0}\left[f((\theta_0, \hat{\eta}), y)\right] \approx \mathbb{E}_{\theta_0, \eta_0}\left[f((\theta_0, \eta_0), y)\right]$$
$$= 0$$

requirement of estimating eqn.

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

**Thm**: efficient scores are Neyman orthogonal.

**Pf (optional)**:
$$\nabla_\eta s_{\theta, \eta}^{eff}(y) = \nabla_\eta\left[s_{\theta, \eta}^0(y) - \alpha(\theta, \eta)\, s_{\theta, \eta}^1(y)\right]$$

$$\left(\alpha(\theta, \eta) = \frac{\mathbb{E}_{\theta, \eta}\left[s_{\theta, \eta}^0(y)\, s_{\theta, \eta}^1(y)\right]}{\mathbb{E}_{\theta, \eta}\left[s_{\theta, \eta}^1(y)^2\right]}\right)$$

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

can show:
$$\mathbb{E}_{\theta, \eta}\left[\nabla_\eta s_{\theta, \eta}^0(y)\right] = -\mathbb{E}_{\theta, \eta}\left[s_{\theta, \eta}^0(y)\, s_{\theta, \eta}^1(y)\right]$$
$$\mathbb{E}_{\theta, \eta}\left[\nabla_\eta s_{\theta, \eta}^1(y)\right] = -\mathbb{E}_{\theta, \eta}\left[s_{\theta, \eta}^1(y)^2\right]$$

$(*)$

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

does not depend on $y$ · has expectation zero

$$\mathbb{E}_{\theta, \eta}[\cdot] = 0$$

**Proof of $(*)$**:
$$0 = \nabla_\eta \mathbb{E}_{\theta, \eta}\left[s_{\theta, \eta}^0(y)\right]$$
$$= \nabla_\eta \int p_{\theta, \eta}(y)\, s_{\theta, \eta}^0(y)\, dy$$
$$= \int \left(\nabla_\eta p_{\theta, \eta}(y) \cdot s_{\theta, \eta}^0(y) + p_{\theta, \eta}(y) \cdot \nabla_\eta s_{\theta, \eta}^0(y)\right) dy$$
$$= \mathbb{E}_{\theta, \eta}\left[s_{\theta, \eta}^0(y)\, s_{\theta, \eta}^1(y) + \nabla_\eta s_{\theta, \eta}^0(y)\right]. \qquad \square$$