

# Semiparametric Efficiency and Optimal Estimation, with Application to Auxiliary Outcome Problems

Jinbo Chen, Ph.D.

Biostatistics Branch

Division of Cancer Epidemiology and Genetics  
National Cancer Institute, Rockville, MD, USA

Norman Breslow, Ph.D.

Department of Biostatistics

University of Washington, Seattle, WA, USA

# Outline

- n Semiparametric efficiency and Godambe's optimality: Connection
- n Application: the auxiliary outcome problem with the conditional mean model
  - n MRC Cognitive Function and Aging Study
  - n The efficient estimation
- n Summary and Conclusions

# Semiparametric Model and Estimation

Model:  $P_{\theta,\eta}$        $\theta$ : finite dimensional  
                                  $\eta$ : infinite dimensional.

n The conditional mean model:

$$E(Y | X) = \mu(X; \theta)$$

n The Likelihood function:

$$f[\varepsilon | X]g(X), \quad \varepsilon = Y - \mu(X; \theta)$$

n The nuisance parameters  $\eta=(f, g)$ :

$$\{f(\varepsilon | X) : \int \varepsilon f(\varepsilon | x) d\varepsilon = 0, x \in X\}$$

# Estimation for the Semiparametric Model

Goal: estimation of  $\theta$

- n Semiparametric efficient estimation (SEE):  
finding a  $\sqrt{n}$  consistent estimator achieving the efficiency bound
- n Heuristic approach: solve an (optimal) estimating equation
  - n Optimal estimating function theory (Godambe, 1960)

# Motivation: Optimal Estimation and Semiparametric Efficiency

The conditional mean model:

n A Class of linear estimating functions

$$\{\sum_i h(X_i; \theta)[Y_i - \mu(X_i; \theta)], h \in H\}$$

n Optimal member / quasi-score (McCullagh and Nelder, 1993)

$$h^*(X)\varepsilon = \frac{\dot{\mu}(X; \theta)}{\text{Var}(Y | X)} \varepsilon$$

n  $h^*(X)\varepsilon$  is the efficient score (Chamberlain, 1987)

# Optimal Estimating Function Theory of Godambe and Heyde

- n A class of estimating functions indexed by  $H = \{h(Z; \theta)\}$ :

$$G = \{G(h; Z, \theta) : h \in H\}$$

- n Regular unbiased estimating function  $G(Z; \theta)$

$$E G(Z; \theta) = 0;$$

$$E [\partial G(Z; \theta) / \partial \theta] \text{ is nonsingular};$$

$$E [G(Z; \theta) G'(Z; \theta)] \text{ is nonsingular.}$$

# Criterion for Finding the Optimal Member $G^*$

n Optimality criterion:

Let  $\varepsilon(G) = (E \dot{G})^{-1} E G G' (E \dot{G}')^{-1}$ .  $G^*$  is optimal if  $\varepsilon(G^*) - \varepsilon(G)$  is nonpositive definite  $\forall G \in G$ ,  $\theta \in \Theta$  and  $\eta \in \Gamma$ .

n Solve  $G=0$  for  $\hat{\theta}$ :

$$\hat{\theta} - \theta \xrightarrow{D} [0, (E \dot{G})^{-1} E G G' (E \dot{G}')^{-1}]$$

where  $\dot{G} = \partial G / \partial \theta$

# Optimality Criterion: Continued

n Theorem (Corollary to Theorem 2.1 in Heyde, 1997)

$G^* \in G$  is a quasi-score estimating function  
if and only if

$$-E\dot{G} = EGG^{*'} \quad \forall G \in G.$$

n Geometric Interpretation:

$$-E\dot{G} = EGG^{*'} \Leftrightarrow E[G(\dot{l}_\theta - G^{*'})] = 0$$

# Semiparametric Efficient Estimation (BKRW, 1993)

n Notation:

$l_{\theta}^*$  : efficient score function

$[El_{\theta}^* l_{\theta}^*]^{-1}$  : semiparametric efficiency bound

$\tilde{l}_{\theta} = [El_{\theta}^* l_{\theta}^*]^{-1} l_{\theta}^*$  : efficient influence function

Data:  $\{Z_i, i = 1, \dots, n\}$  *i.i.d*

n Regular and asymptotically linear (RAL) estimator:

regular and  $\hat{\theta} = \theta + \frac{1}{n} \sum \varphi(Z_i; \theta, \eta) + o_p(n^{1/2})$

with  $E\varphi(Z_i; \theta, \eta) = 0$ ,  $\text{var}\varphi(Z_i; \theta, \eta) < \infty$ .

# SEE: Continued

n Semiparametric efficient estimator:

$$\hat{\theta} = \theta + \frac{1}{n} \sum \tilde{l}_{\theta}(Z_i) + o_p(n^{1/2})$$

n  $\text{var}(\tilde{l}_{\theta}) \leq \text{var}(\varphi)$

$\Rightarrow$  The efficient estimator is the most precise RAL.

$\Rightarrow$  Sufficient to work with the class of RAL estimators to find the efficient estimator.

$\Rightarrow$  find  $\tilde{l}_{\theta} (l_{\theta}^*)$

$\Rightarrow$  construct efficient estimators

n Influence function  $\iff$  Regular estimating function

# SEE and Optimal Estimation: Connection (Chen 2002, Ph.D dissertation)

- n A corollary to Theorem 3.3.1, BKRW

Let  $G = \{G\}$  be the closed linear span of the influence functions for all RAL estimators. Then  $\dot{l}_\theta^*$  is uniquely identified by  $\dot{l}_\theta^* \in G$  and

$$E[G(\dot{l}_\theta - \dot{l}_\theta^*)] = 0 \quad \forall G \in G.$$

- n Recall Godambe's optimality criterion

$$-E\dot{G} = EGG^{*'} \Leftrightarrow E[G(\dot{l}_\theta - G^{*'})] = 0$$

# Connection: Continued

- The efficient score/influence function is the Godambe's optimal member in the closed linear span of influence functions for all RAL estimators!
- Useful to calculating the efficient score functions for the missing data problem (Robins, Rotnitzky and Zhao, 1994; RRZ)
- Maybe useful to obtain the efficient score function for non-iid data (Welfelmayer 1996)
  - Trial and verification approach

# An Example

- n The quasi-likelihood model:

$$E(Y | X) = \mu(X; \theta); \text{ var}(Y | X) = \phi v(\mu).$$

$\phi$  and  $v$  are known.

- n The Likelihood function:

$$f[\varepsilon_s | x]g(X) \text{ where } \varepsilon_s = (Y - \mu)/(\phi v)^{1/2}$$

- n The nuisance parameters  $\eta = (f, g)$ :

$$\{f(\varepsilon_s | X) : \varepsilon_s f(\varepsilon_s | x) d\varepsilon_s = 0, \\ \varepsilon_s^2 f(\varepsilon_s | x) d\varepsilon_s = 1; x \in X\}$$

# Example: Continued

n A class of linear estimating functions

$$\{G_1(h; \theta) : G_1(h; \theta) = \sum_i h(X_i; \theta) [Y_i - \mu(X_i; \theta)], h \in H\}$$

$$\text{Optimal member: } G_1^* = \sum_i \dot{\mu}_i v^{-1}(\mu_i) [Y_i - \mu(X_i; \theta)]$$

n A class of quadratic estimating functions

$$\{G_2(h_1, h_2; \theta) : G_2 = \sum_i h_1(X_i; \theta) \varepsilon_i + h_2(X_i; \theta) (\varepsilon_i^2 - \phi v_i)\}$$

Optimal member  $G_2^*$ : involves  $E(\varepsilon^3 | X)$  and  $E(\varepsilon^4 | X)$ .

## An Example: Continued

- n  $[\text{var}(G_1^*)]^{-1} \geq [\text{var}(G_2^*)]^{-1}$  .
- n  $\{G_2\}$  is the complete class of influence functions for all RAL estimators
- n  $G_2^*$  is the efficient score function
- n The same efficient score was obtained by Rotnitzky and Robins (1995) using a different approach

# Outline

- n Semiparametric efficiency and Godambe's optimality: Connection
- n Application: the auxiliary outcome problem with the conditional mean model (Chen and Breslow, Canadian J. Statistics, to appear)
  - n MRC Cognitive Function and Aging Study
  - n The efficient estimation
- n Summary and Conclusions

# MRC Cognitive Function and Aging Study (Clayton *et al.*, JRSS(B), 1999)

- n Goal: Estimate prevalence of dementia
  - n by sex and age ( $X$  = covariates)
- n Outcome: dementia assessed by
  - n Geriatric mental state exam ( $Y$ ) – gold standard
  - n Mini-mental state exam ( $S$ ) – auxiliary outcome
- n Design: two phase sampling
  - n 10,000 main sample known ( $X, S$ )
  - n 1,780 validation sample known ( $Y, X, S$ )
  - n Data simulated by Clayton *et al.*

# Study Design

Age (yrs)	MMSE	No. of subjects		Sampling fraction
		Main	Validation	
65-74	0-21	291	291	1.000
	22-25	950	220	0.232
	26-30	3,759	386	0.103
75+	0-21	1,037	496	0.478
	22-25	1,486	208	0.140
	26-30	2,477	179	0.072
Total		10,000	1,780	0.178

# Notation

- $Y$  = true outcome
- $S$  = auxiliary (surrogate) outcome
- $X$  = covariates
- $R$  = sampling indicator ( $R=1 \rightarrow$  validation)
- Observed data

$$Z = (S, X, RY, R) \begin{cases} R=1: (S, X, Y) \\ R=0: (S, X) \end{cases}$$

- $\{Z_i, i = 1, \dots, n\}$  i.i.d.

# Conditional Mean Model

- n Conditional mean (parametric)

$$E(Y | X = x) = \mu(x; \theta), \quad \theta \in R^p$$

- n Joint distribution

$$p(y, s, x) = q(s | y, x) f(y - \mu(x; \theta)) g(x)$$

- u  $\eta = (q, f, g)$  (infinite dim.) "nuisance" parm  
finite variance,  $\int f(\varepsilon) d\varepsilon = 0$

# Validation Sampling: MAR

n Missing at random:

$$\Pr(R = 1 \mid y, s, x) = \Pr(R = 1 \mid s, x) = \pi(s, x)$$

n Positive validation probability:

$$\pi(s, x) \geq \sigma > 0$$

n Parametric missingness model:

$$\pi(s, x) = \pi(s, x; \alpha), \alpha \in R^q$$

# Previous Work

## n Semiparametric efficient estimators of $\theta$

- n Robins, Rotnitzky, Zhao (*JASA*, 1994)
- n Rotnitzky, Robins (*Scand J Statist*, 1995)
- n Holcroft, Rotnitzky, Robins (*J Stat Plan Inf*, 1997)
- n .....

## n Inefficient estimators of $\theta$

- n Pepe, Reilly, Fleming (*J Stat Plan Inf*, 1994)
  - n “Mean score” (Horwitz-Thompson) estimator
- n Y. Chen (*Biometrika*, 2000)
  - n “Robust imputation” estimator

# Rationale for More Work

- n “Simple” derivation of SEE via connection to Godambe optimality
- n Method may generalize to more complex problems
- n Explicit expression for efficient estimating equation
- n Useful pedagogical example

# Missing Data Problem (RRZ, 1994)

- n Space  $\Gamma$  of influence functions for RAL est

$$G(h; Z) = \frac{R}{\pi} h(X) \varepsilon - \frac{R - \pi}{\pi} h(X) E(\varepsilon | S, X)$$

- n Godambe optimal quasi score is  $G^* = h^*(X) \varepsilon^*$

$$\varepsilon^* = \frac{R}{\pi} Y - \frac{R - \pi}{\pi} E(Y | S, X) - \mu(X; \theta)$$

$$h^* = \dot{\mu}(X; \theta) \text{Var}^{-1}(\varepsilon^* | X)$$

- n Estimate  $E(Y | S, X)$ ,  $\text{Var}(\varepsilon^* | X)$  and  $\pi$  from data

# Semiparametric Efficiency Bound

n Efficiency bound is  $\text{Var}^{-1}(G^*)$  where

$$\text{Var}(G^*) = E \left\{ \dot{\mu}(X; \theta) [\text{Var}(\varepsilon | X) + A(X)]^{-1} \dot{\mu}^T(X; \theta) \right\}$$

$$A(X) = E \left( \frac{1 - \pi}{\pi} \text{Var}(\varepsilon | S, X) \middle| X \right)$$

n Perfect information:  $S = Y \Rightarrow A(X) = 0$

-Usual quasi likelihood bound

n No information:  $S \perp Y | X \Rightarrow$  Validation only

$$\text{Var } G^* = E \left[ \dot{\mu}(X; \theta) E^{-1} \left( \frac{1}{\pi} \middle| X \right) \text{Var}^{-1}(\varepsilon | X) \dot{\mu}^T(X; \theta) \right]$$

# Summary and Conclusions

- n SEE and optimal estimating equation theory
  - n Connection (Chen 2002, Dissertation)
- n The Auxiliary Outcome Problem (Chen and Breslow, Canadian J. Statistics):
  - n Efficient score has a simple closed form
- n Extensions of Chen's "robust imputation" method in progress for continuous  $S$  and  $X$