

## M- and Z- theorems

Wellner; 5/13/98, 1/26/07, 5/08/09

### Z-theorems: Notation and Context

Suppose that  $\Theta \subset R^k$ , and that

$$\begin{aligned}\Psi_n &: \Theta \rightarrow \mathbb{R}^k, \text{ random maps} \\ \Psi &: \Theta \rightarrow \mathbb{R}^k, \text{ deterministic maps.}\end{aligned}$$

Suppose that  $\hat{\theta}_n$  and  $\theta_0$  are the corresponding solutions (or approximate solutions) of

$$\begin{aligned}\Psi_n(\hat{\theta}_n) = 0 \quad \text{or} \quad \Psi_n(\hat{\theta}_n) = o_p(n^{-1/2}), \\ \Psi(\theta_0) = 0.\end{aligned}$$

In the simple case of i.i.d. data  $X_1, \dots, X_n$  i.i.d.  $P_0$  with empirical measure  $\mathbb{P}_n$ , and then, for the usual case of linear estimating equations, the functions  $\Psi_n$ ,  $\Psi$  are given by

$$\Psi_n(\theta) = \mathbb{P}_n \psi(\cdot, \theta), \quad \text{and} \quad \Psi(\theta) = P_0 \psi(\cdot, \theta)$$

for a vector of functions  $\psi : \mathcal{X} \times \Theta \rightarrow R^k$ ,  $\psi(x, \theta) = \underline{\psi}(x, \theta)$ ; often the functions  $\psi$  are score functions motivated by likelihood, pseudolikelihood, quasilikelihood, or some other “likelihood” for the data.

Here are the four basic conditions needed for Huber’s  $Z$ - theorem:

**A1**  $\Psi_n(\hat{\theta}_n) = o_p(n^{-1/2})$  and  $\Psi(\theta_0) = 0$ .

**A2**  $\sqrt{n}(\Psi_n - \Psi)(\theta_0) \rightarrow_d \mathbb{Z}_0$ .

**A3** For every sequence  $\delta_n \rightarrow 0$ ,

$$\sup_{|\theta - \theta_0| \leq \delta_n} \frac{|\sqrt{n}(\Psi_n - \Psi)(\theta) - \sqrt{n}(\Psi_n - \Psi)(\theta_0)|}{1 + \sqrt{n}|\theta - \theta_0|} = o_p(1).$$

**A4** The function  $\Psi$  is (Fréchet-)differentiable at  $\theta_0$  with nonsingular derivative  $\dot{\Psi}(\theta_0) \equiv \dot{\Psi}_0$ :

$$\Psi(\theta) - \Psi(\theta_0) - \dot{\Psi}_0(\theta - \theta_0) = o(|\theta - \theta_0|).$$

**Theorem.** (Huber (1967); Pollard (1985)). Suppose that A1 - A4 hold. Let  $\hat{\theta}_n$  be random maps into  $\Theta \subset R^k$  satisfying  $\hat{\theta}_n \rightarrow_p \theta_0$ . Then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d -\dot{\Psi}_0^{-1}(\mathbb{Z}_0);$$

if  $Z_0 \sim N_k(0, A)$ , then this yields, with  $\dot{\Psi}_0 \equiv B$ ,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d N_k(0, B^{-1}A(B^{-1})^T).$$

**Proof.** By definition of  $\hat{\theta}_n$  and  $\theta_0$ ,

$$\begin{aligned} \sqrt{n}(\Psi(\hat{\theta}_n) - \Psi(\theta_0)) &= \sqrt{n}(\Psi(\hat{\theta}_n) - \Psi_n(\hat{\theta}_n)) + o_p(1) \\ &= -\sqrt{n}(\Psi_n - \Psi)(\hat{\theta}_n) \\ &\quad - \left\{ \sqrt{n}(\Psi_n - \Psi)(\hat{\theta}_n) - \sqrt{n}(\Psi_n - \Psi)(\theta_0) \right\} + o_p(1) \\ &= -\sqrt{n}(\Psi_n - \Psi)(\theta_0) + o_p(1 + \sqrt{n}|\hat{\theta}_n - \theta_0|) + o_p(1); \end{aligned} \quad (1)$$

here the last equality holds by A3 and  $\hat{\theta}_n \rightarrow_p \theta_0$ . Since  $\dot{\Psi}_0$  is continuously invertible, there exists a constant  $c > 0$  such that

$$\|\dot{\Psi}_0(\theta - \theta_0)\| \geq c\|\theta - \theta_0\|$$

for every  $\theta$ ; this is just the basic property of a nonsingular matrix. By A4 (differentiability of  $\Psi$ ), this yields

$$|\Psi(\theta) - \Psi(\theta_0)| \geq c|\theta - \theta_0| + o(|\theta - \theta_0|).$$

By (1) it follows that

$$\sqrt{n}|\hat{\theta}_n - \theta_0|(c + o_p(1)) \leq O_p(1) + o_p(1 + \sqrt{n}|\hat{\theta}_n - \theta_0|),$$

which implies

$$\sqrt{n}|\hat{\theta}_n - \theta_0| = O_p(1).$$

Hence from (1) again and A.4 it follows that

$$\dot{\Psi}_0(\sqrt{n}(\hat{\theta}_n - \theta_0)) + o_p(\sqrt{n}|\hat{\theta}_n - \theta_0|) = -\sqrt{n}(\Psi_n - \Psi)(\theta_0) + o_p(1)$$

and therefore

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d -\dot{\Psi}_0^{-1}(Z_0)$$

by A2 and A4. □

**Example.** Consider estimation of  $\theta$  in the simple triangular density family  $\mathcal{P} = \{P_\theta : \theta \in [0, 1]\}$  given by the densities

$$p(x; \theta) = 2 \left\{ \frac{x}{\theta} 1_{[0, \theta]}(x) + \frac{1-x}{1-\theta} 1_{(\theta, 1]}(x) \right\}, \quad \theta \in [0, 1].$$

Then the score function for  $\theta$  for one observation is

$$\dot{\mathbf{i}}_{\theta}(x; \theta) = -\frac{1}{\theta}1_{[0, \theta]}(x) + \frac{1}{1 - \theta}1_{(\theta, 1]}(x),$$

at least in the sense of a Hellinger derivative:

$$\int_0^1 \{\sqrt{p(x; \theta)} - \sqrt{p(x; \theta_0)} - 2^{-1}\dot{\mathbf{i}}_{\theta}(x; \theta_0)\sqrt{p(x; \theta_0)}\}^2 dx = o(|\theta - \theta_0|^2).$$

Our goal is to use Huber's theorem to study the behavior of the MLE of  $\theta$  in this model when (possibly)  $X, X_1, \dots, X_n$  are i.i.d.  $P$  on  $[0, 1]$  with  $P \notin \mathcal{P}$ . Let  $F(x) = P(X \leq x)$  be the distribution function corresponding to  $X$ . From the score calculation above, the score equation for estimation of  $\theta$  is equivalent to

$$\Psi_n(\theta) = \mathbb{F}_n(\theta) - \theta = 0.$$

The corresponding population version of  $\Psi_n$  is  $\Psi$  given by

$$\Psi(\theta) = F(\theta) - \theta, \quad \text{where } F(x) = \int_0^x p(y)dy.$$

It is reasonable to assume that  $\Psi(\theta_0) = 0$  has a unique solution  $\theta_0 = \theta_0(P)$  if  $P$  has a density  $p$  which is unimodal on  $[0, 1]$ . (Of course it is easy to construct examples in which this has only trivial solutions 0 or 1, or for which there are many solutions: for the former, consider the "anti-triangular density"  $p(x) = 2(1/2 - x)1_{[0, 1/2]}(x) + 2(x - 1/2)1_{(1/2, 1]}(x)$ ; for the latter consider  $p(x) = 1 + (1/2)\sin(6\pi x)$ .)

Let  $\theta_0 = \theta_0(P)$  satisfy  $\Psi(\theta_0) = 0 = F(\theta_0) - \theta_0$ . Then  $F(\theta_0) = \theta_0$ , and

$$\begin{aligned} \sqrt{n}\Psi_n(\theta_0) &= \sqrt{n}(\Psi_n(\theta_0) - \Psi(\theta_0)) \\ &= \sqrt{n}(\mathbb{F}_n(\theta_0) - F(\theta_0)) \\ &\rightarrow_d \mathbb{Z}_0 \sim N(0, F(\theta_0)(1 - F(\theta_0))) = N(0, \theta_0(1 - \theta_0)). \end{aligned}$$

Thus A2 holds. Moreover if  $F$  is differentiable at  $\theta_0$  with derivative  $p(\theta_0)$ , then A4 holds with

$$\dot{\Psi}(\theta_0) = p(\theta_0) - 1.$$

Furthermore, the condition A3 holds since, for any  $\delta_n \rightarrow 0$ ,

$$\begin{aligned} &\sup_{\theta: |\theta - \theta_0| \leq \delta_n} |\sqrt{n}(\Psi_n - \Psi)(\theta) - \sqrt{n}(\Psi_n - \Psi)(\theta_0)| \\ &= \sup_{\theta: |\theta - \theta_0| \leq \delta_n} |\sqrt{n}(\mathbb{F}_n - F)(\theta) - \sqrt{n}(\mathbb{F}_n - F)(\theta_0)| \\ &\rightarrow_p 0 \end{aligned}$$

if  $F$  is continuous at  $\theta_0$  by the asymptotic equicontinuity of the empirical process. We conclude from Huber's theorem that any solution of  $\Psi_n(\hat{\theta}_n) = o_p(n^{-1/2})$  satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d -(p(\theta_0) - 1)^{-1}Z_0 \sim N\left(0, \frac{\theta_0(1 - \theta_0)}{(p(\theta_0) - 1)^2}\right).$$

When  $P \in \mathcal{P}$  holds,  $\theta_0(P_\theta) = \theta$  (so  $\theta_0(P_{\theta_0}) = \theta_0$ ), and  $p(\theta_0; \theta_0) = 2$ . Thus the asymptotic variance in the conclusion of Huber's theorem reduces to  $\theta_0(1 - \theta_0)$ , which agrees with the information bound calculation based on the score  $\dot{l}_\theta$ .

Now our goal is to extend this to an infinite-dimensional setting in which  $\Theta$  is a Banach space. A sufficiently general Banach space is the space

$$l^\infty(H) \equiv \{z : H \rightarrow R \mid \|z\| = \sup_{h \in H} |z(h)| < \infty\}$$

where  $H$  is a collection of functions. We suppose that

$$\Psi_n : \Theta \rightarrow L \equiv l^\infty(H'), \quad n = 1, 2, \dots$$

are random, and that

$$\Psi : \Theta \rightarrow L \equiv l^\infty(H'),$$

is deterministic. Suppose that either

$$\Psi_n(\hat{\theta}_n) = 0 \quad \text{in} \quad L;$$

(i.e.  $\Psi_n(\hat{\theta}_n)(h') = 0$  for all  $h' \in H'$ ), or

$$\Psi_n(\hat{\theta}_n) = o_p(n^{-1/2}) \quad \text{in} \quad L;$$

(i.e.  $\|\Psi_n(\hat{\theta}_n)\|_{H'} = o_p(n^{-1/2})$ ).

Here are the four basic conditions needed for the infinite-dimensional version of Huber's  $Z$ -theorem due to Van der Vaart (1995):

**B1**  $\Psi_n(\hat{\theta}_n) = o_p(n^{-1/2})$  in  $l^\infty(H')$  and  $\Psi(\theta_0) = 0$  in  $l^\infty(H')$ .

**B2**  $\sqrt{n}(\Psi_n - \Psi)(\theta_0) \Rightarrow Z_0$  in  $l^\infty(H')$ .

**B3** For every sequence  $\delta_n \rightarrow 0$ ,

$$\sup_{\|\theta - \theta_0\| \leq \delta_n} \frac{\|\sqrt{n}(\Psi_n - \Psi)(\theta) - \sqrt{n}(\Psi_n - \Psi)(\theta_0)\|}{1 + \sqrt{n}\|\theta - \theta_0\|} = o_p(1).$$

**B4** The function  $\Psi$  is (Fréchet-)differentiable at  $\theta_0$  with derivative  $\dot{\Psi}(\theta_0) \equiv \dot{\Psi}_0$  having a bounded (continuous) inverse:

$$\|\Psi(\theta) - \Psi(\theta_0) - \dot{\Psi}_0(\theta - \theta_0)\| = o(\|\theta - \theta_0\|).$$

**Theorem.** (van der Vaart, 1995). Suppose that B1 - B4 hold. Let  $\hat{\theta}_n$  be random maps into  $\Theta \subset l^\infty(H')$  satisfying  $\hat{\theta}_n \rightarrow_p \theta_0$ . Then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \Rightarrow -\dot{\Psi}_0^{-1}(Z_0) \quad \text{in} \quad l^\infty(H).$$

**Proof.** Exactly the same as in the finite-dimensional case: see van der Vaart (1995) or van der Vaart and Wellner (1996), pages 310-312.  $\square$

**M-theorems: Notation and context**

Suppose that  $\Theta \subset \mathbb{R}^k$  and that  $m : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$ . We often write  $m_\theta(x) = m(x, \theta)$  for  $x \in \mathcal{X}$ ,  $\theta \in \Theta$ . Suppose that  $\hat{\theta}_n$  and  $\theta_0$  are the corresponding maximizers (or approximate maximizers in the first case) of

$$\begin{aligned} \mathbb{M}_n(\theta) &\equiv \mathbb{P}_n m(X, \theta) = \mathbb{P}_n m_\theta(X), & \text{and} \\ M(\theta) &\equiv P_0 m(X, \theta) = P_0 m_\theta(X), \end{aligned}$$

respectively. A common choice for  $m(x, \theta)$  would be  $\log p(x; \theta) \equiv \log p_\theta(x)$  where  $p_\theta(\cdot)$  is the density of  $P_\theta$  with respect to some dominating measure  $\mu$  for a model  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ . Then  $\hat{\theta}_n$  is a Maximum Likelihood (or approximate maximum likelihood) estimator for the model  $\mathcal{P}$ .

**Theorem.** Suppose that for each  $\theta$  in an open subset of  $\Theta \subset \mathbb{R}^k$ ,  $x \mapsto m_\theta(x)$  is a measurable function such that  $\theta \mapsto m_\theta(x) = m(x, \theta)$  is differentiable at  $\theta_0$  for  $P_0$ -almost every  $x$  with derivative  $\dot{m}_{\theta_0}(x)$  and such that, for every  $\theta_1, \theta_2$  in a neighborhood of  $\theta_0$  and a measurable function  $\dot{m}$  with  $P_0 \dot{m}^2 < \infty$ ,

$$|m_{\theta_1}(x) - m_{\theta_2}(x)| \leq \dot{m} \|\theta_1 - \theta_2\|.$$

Furthermore, suppose that  $\theta \mapsto P_0 m_\theta$  has a second order Taylor expansion at a point of maximum  $\theta_0$  with nonsingular second derivative matrix  $V_{\theta_0}$ : i.e.

$$P_0 m_\theta = P_0 m_{\theta_0} + \frac{1}{2}(\theta - \theta_0)^T V_{\theta_0} (\theta - \theta_0) + o(\|\theta - \theta_0\|^2).$$

If  $\mathbb{P}_n m_{\hat{\theta}_n}(X) \geq \sup_\theta \mathbb{P}_n m_\theta(X) - o_p(n^{-1})$  and  $\hat{\theta}_n \rightarrow_p \theta_0$ , then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = -V_{\theta_0}^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{m}_{\theta_0}(X_i) + o_p(1).$$

In particular

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d N_k(0, V_{\theta_0}^{-1} P_0 (\dot{m}_{\theta_0} \dot{m}_{\theta_0}^T) V_{\theta_0}^{-1}).$$

**Proof.** See van der Vaart, *Asymptotic Statistics*, section 5.3, pages 51 - 60. □

## Applications and Extensions of Van der Vaart's Z-theorem:

- Gamma frailty model; Murphy (1995).
- Partially censored data; Van der Vaart (1995).
- Correlated gamma-frailty model; Parner (1998).
- Semiparametric biased sampling models; Gilbert (2000).
- Two-phase sampling models with data missing by design; Breslow, McNeney and Wellner (2003), Breslow and Wellner (2007), (2008).

However, in many statistical problems the parameter usually includes both a finite-dimensional parameter (e.g. regression parameters) and an infinite dimensional (nuisance) parameter. We now suppose that  $\theta = (\beta, \Lambda)$ , where  $\beta$  is a finite-dimensional parameter, say in  $\mathbb{R}^d$ , and  $\Lambda$  an infinite dimensional parameter (a function). The M-estimators of  $\beta$ ,  $\hat{\beta}_n$ , and of  $\Lambda$ ,  $\hat{\Lambda}_n$ , respectively, often have different convergence rates. The convergence rate for  $\hat{\Lambda}_n$  is often smaller than  $n^{1/2}$ , such as  $n^{1/3}$ , or  $n^{2/5}$  in some cases. Huang (1996) established a general theorem to show that under certain hypotheses, the maximum likelihood estimator of a finite dimensional parameter has  $n^{1/2}$  convergence rate and is asymptotically semiparametric efficient, even though the convergence rate for the maximum likelihood estimator of the infinite dimensional parameter is smaller than  $n^{1/2}$ . He also successfully applied his general theorem to the proportional hazards model with interval censored data.

The following theorem due to Zhang (1998) generalizes the theorem of Huang (1996) to the case of inefficient M-estimators; it shows that under reasonable regularity hypotheses, the M-estimator of a finite-dimensional parameter  $\beta$  has  $n^{1/2}$  convergence rate, and that  $\hat{\beta}_n$  is asymptotically normal, even though the M-estimator of the corresponding infinite dimensional parameter  $\Lambda$  converges perhaps more slowly than  $n^{1/2}$ . The resulting asymptotic covariance matrix for the M-estimator of  $\beta$  has the well-known “sandwich” structure.

Here is the notation and conditions needed for the theorem. Let  $\theta = (\beta, \Lambda)$ , where  $\beta \in \mathbb{R}^d$ , and  $\Lambda$  is an infinite dimensional parameter in a class of functions  $\mathcal{F}$ .  $\Lambda_\eta$  is a parametric path in  $\mathcal{F}$  through  $\Lambda$ , i.e.  $\Lambda_\eta \in \mathcal{F}$ , and  $\Lambda_\eta|_{\eta=0} = \Lambda$ .

Let  $\mathbf{H} = \left\{ h : h = \frac{\partial \Lambda_\eta}{\partial \eta} \Big|_{\eta=0} \right\}$  and define

$$m_1(\beta, \Lambda; x) = \nabla_\beta m_{(\beta, \Lambda)}(x) \equiv \left( \frac{\partial}{\partial \beta_1} m_{(\beta, \Lambda)}(x), \dots, \frac{\partial}{\partial \beta_d} m_{(\beta, \Lambda)}(x) \right)'$$

$$m_2(\beta, \Lambda; x)[h] = \frac{\partial}{\partial \eta} m_{(\beta, \Lambda_\eta)}(x) \Big|_{\eta=0},$$

$$\begin{aligned}
m_{11}(\beta, \Lambda; x) &= \nabla_{\beta}^2 m_{(\beta, \Lambda)}(x), \\
m_{12}(\beta, \Lambda; x)[h] &= \left. \frac{\partial}{\partial \eta} m_1(\beta, \Lambda_{\eta}; x) \right|_{\eta=0}, \\
m_{21}(\beta, \Lambda; x)[h] &= \nabla_{\beta} m_2(\beta, \Lambda; x)[h],
\end{aligned}$$

and

$$m_{22}(\beta, \Lambda; x)[h, h] = \left. \frac{\partial^2}{\partial \eta^2} m(\beta, \Lambda_{\eta}; x) \right|_{\eta=0}.$$

We also define

$$\begin{aligned}
S_1(\beta, \Lambda) &= Pm_1(\beta, \Lambda; X), \\
S_2(\beta, \Lambda)[h] &= Pm_2(\beta, \Lambda; X)[h], \\
S_{1n}(\beta, \Lambda) &= \mathbb{P}_n m_1(\beta, \Lambda; X), \\
S_{2n}(\beta, \Lambda)[h] &= \mathbb{P}_n m_2(\beta, \Lambda; X)[h], \\
\dot{S}_{11}(\beta, \Lambda) &= Pm_{11}(\beta, \Lambda; X), \\
\dot{S}_{12}(\beta, \Lambda)[h] &= \dot{S}'_{21}(\beta, \Lambda)[h] = Pm_{12}(\beta, \Lambda; X)[h],
\end{aligned}$$

and

$$\dot{S}_{22}(\beta, \Lambda)[h, h] = Pm_{22}(\beta, \Lambda; X)[h, h].$$

Furthermore, for  $\mathbf{h} = (h_1, \dots, h_d)' \in \mathbf{H}^d$ , where  $h_j \in \mathbf{H}$  for  $j = 1, 2, \dots, d$ , and  $\mathbf{H}^d = \underbrace{\mathbf{H} \times \mathbf{H} \times \dots \times \mathbf{H}}_d$ , denote

$$\begin{aligned}
m_2(\beta, \Lambda; x)[\mathbf{h}] &= (m_2(\beta, \Lambda; x)[h_1], \dots, m_2(\beta, \Lambda; x)[h_d])', \\
m_{12}(\beta, \Lambda; x)[\mathbf{h}] &= (m_{12}(\beta, \Lambda; x)[h_1], \dots, m_{12}(\beta, \Lambda; x)[h_d]), \\
m_{21}(\beta, \Lambda; x)[\mathbf{h}] &= (m_{21}(\beta, \Lambda; x)[h_1], \dots, m_{21}(\beta, \Lambda; x)[h_d]), \\
m_{22}(\beta, \Lambda; x)[\mathbf{h}, h] &= (m_{22}(\beta, \Lambda; x)[h_1, h], \dots, m_{22}(\beta, \Lambda; x)[h_d, h])^T,
\end{aligned}$$

and define

$$\begin{aligned}
S_2(\beta, \Lambda)[\mathbf{h}] &= Pm_2(\beta, \Lambda; X)[\mathbf{h}], \\
S_{2n}(\beta, \Lambda)[\mathbf{h}] &= \mathbb{P}_n m_2(\beta, \Lambda; X)[\mathbf{h}], \\
\dot{S}_{12}(\beta, \Lambda)[\mathbf{h}] &= Pm_{12}(\beta, \Lambda; X)[\mathbf{h}], \\
\dot{S}_{21}(\beta, \Lambda)[\mathbf{h}] &= Pm_{21}(\beta, \Lambda; X)[\mathbf{h}],
\end{aligned}$$

and

$$\dot{S}_{22}(\beta, \Lambda)[\mathbf{h}, h] = Pm_{22}(\beta, \Lambda; X)[\mathbf{h}, h].$$

The following Assumptions will be used to formulate our general theorem:

A1. **(Consistency and rate of convergence):**

$$|\hat{\beta}_n - \beta_0| = o_p(1) \quad \text{and} \quad \|\hat{\Lambda}_n - \Lambda_0\| = O_p(n^{-\gamma})$$

for some  $\gamma > 0$ .

A2. **(Zero-mean structure):**

$$S_1(\beta_0, \Lambda_0) = 0, \quad \text{and} \quad S_2(\beta_0, \Lambda_0)[h] = 0, \quad \text{for all } h \in \mathbf{H}.$$

A3. **(Positive “pseudo-information”):** There exists an  $\mathbf{h}^* = (h_1^*, \dots, h_d^*)^T$ ,  $h_j^* \in \mathbf{H}$   $j = 1, \dots, d$ , such that

$$\dot{S}_{12}(\beta_0, \Lambda_0)[h] - \dot{S}_{22}(\beta_0, \Lambda_0)[\mathbf{h}^*, h] = 0, \quad (2)$$

for all  $h \in \mathbf{H}$ . Moreover, the matrix

$$A = -\dot{S}_{11}(\beta_0, \Lambda_0) + \dot{S}_{21}(\beta_0, \Lambda_0)[\mathbf{h}^*] = -P(m_{11}(\beta_0, \Lambda_0; X) - m_{21}(\beta_0, \Lambda_0; X)[\mathbf{h}^*])$$

is nonsingular.

A4. **(Approximate solution of pseudo-score equations):** The estimator  $(\hat{\beta}_n, \hat{\Lambda}_n)$  satisfies

$$S_{1n}(\hat{\beta}_n, \hat{\Lambda}_n) = o_{p^*}(n^{-1/2}),$$

and

$$S_{2n}(\hat{\beta}_n, \hat{\Lambda}_n)[\mathbf{h}^*] = o_{p^*}(n^{-1/2}).$$

A5. **(Stochastic equicontinuity):** For any  $\delta_n \downarrow 0$  and  $C > 0$ ,

$$\sup_{|\beta - \beta_0| \leq \delta_n, \|\Lambda - \Lambda_0\| \leq Cn^{-\gamma}} \left| \sqrt{n}(S_{1n} - S_1)(\beta, \Lambda) - \sqrt{n}(S_{1n} - S_1)(\beta_0, \Lambda_0) \right| = o_{p^*}(1),$$

and

$$\sup_{|\beta - \beta_0| \leq \delta_n, \|\Lambda - \Lambda_0\| \leq Cn^{-\gamma}} \left| \sqrt{n}(S_{2n} - S_2)(\beta, \Lambda)[\mathbf{h}^*] - \sqrt{n}(S_{2n} - S_2)(\beta_0, \Lambda_0)[\mathbf{h}^*] \right| = o_{p^*}(1).$$

A6. **(Smoothness of the model):** For some  $\alpha > 1$  satisfying  $\alpha\gamma > 1/2$ , and for  $(\beta, \Lambda)$  in the neighborhood  $\{(\beta, \Lambda) : |\beta - \beta_0| \leq \delta_n, \|\Lambda - \Lambda_0\| \leq Cn^{-\gamma}\}$ ,

$$\begin{aligned} & \left| S_1(\beta, \Lambda) - S_1(\beta_0, \Lambda_0) - \dot{S}_{11}(\beta_0, \Lambda_0)(\beta - \beta_0) - \dot{S}_{12}(\beta_0, \Lambda_0)[\Lambda - \Lambda_0] \right| \\ & = o(|\beta - \beta_0|) + O(\|\Lambda - \Lambda_0\|^\alpha), \end{aligned}$$

$$\begin{aligned} & \left| S_2(\beta, \Lambda)[\mathbf{h}^*] - S_2(\beta_0, \Lambda_0)[\mathbf{h}^*] - \dot{S}_{21}(\beta_0, \Lambda_0)[\mathbf{h}^*](\beta - \beta_0) - (\dot{S}_{22}(\beta_0, \Lambda_0)[\mathbf{h}^*, \Lambda - \Lambda_0]) \right| \\ & = o(|\beta - \beta_0|) + O(\|\Lambda - \Lambda_0\|^\alpha). \end{aligned}$$

A7. (Asymptotic normality of projected pseudo-score): With

$$m^*(\beta_0, \Lambda_0; x) \equiv m_1(\beta_0, \Lambda_0; x) - m_2(\beta_0, \Lambda_0; x)[\mathbf{h}^*],$$

we have

$$\sqrt{n}\mathbb{P}_n m^*(\beta_0, \Lambda_0; X) \longrightarrow_d N(0, B),$$

where  $B = Em^*(\beta_0, \Lambda_0; X)^{\otimes 2} = Em^*(\beta_0, \Lambda_0; X)m^*(\beta_0, \Lambda_0; X)'$ .

**Theorem 2.3.5. (Asymptotic Normality)** Suppose that: assumptions A1-A7 hold. Then

$$\sqrt{n}(\hat{\beta}_n - \beta_0) = A^{-1}\sqrt{n}\mathbb{P}_n m^*(\beta_0, \Lambda_0; X) + o_{p^*}(1) \longrightarrow_d N\left(0, A^{-1}B(A^{-1})'\right).$$

*Proof* : A1 and A5 yield

$$\sqrt{n}(S_{1n} - S_1)(\hat{\beta}_n, \hat{\Lambda}_n) - \sqrt{n}(S_{1n} - S_1)(\beta_0, \Lambda_0) = o_{p^*}(1).$$

Since  $S_{1n}(\hat{\beta}_n, \hat{\Lambda}_n) = o_{p^*}(n^{-1/2})$  by A4 and  $S_1(\beta_0, \Lambda_0) = 0$  by A2, it follows that

$$\sqrt{n}S_1(\hat{\beta}_n, \hat{\Lambda}_n) + \sqrt{n}S_{1n}(\beta_0, \Lambda_0) = o_{p^*}(1).$$

Similarly, we have that

$$\sqrt{n}S_2(\hat{\beta}_n, \hat{\Lambda}_n)[\mathbf{h}^*] + \sqrt{n}S_{2n}(\beta_0, \Lambda_0)[\mathbf{h}^*] = o_{p^*}(1).$$

Combining these equalities and A6 yields

$$\begin{aligned} \dot{S}_{11}(\beta_0, \Lambda_0)[\hat{\beta}_n - \beta_0] + \dot{S}_{12}(\beta_0, \Lambda_0)[\hat{\Lambda}_n - \Lambda_0] + S_{1n}(\beta_0, \Lambda_0) \\ + o(|\hat{\beta}_n - \beta_0|) + O(\|\hat{\Lambda}_n - \Lambda_0\|^\alpha) = o_{p^*}(n^{-1/2}), \end{aligned} \quad (3)$$

$$\begin{aligned} \dot{S}_{21}(\beta_0, \Lambda_0)[\mathbf{h}^*][\hat{\beta}_n - \beta_0] + \dot{S}_{22}(\beta_0, \Lambda_0)[\mathbf{h}^*][\hat{\Lambda}_n - \Lambda_0] + S_{2n}(\beta_0, \Lambda_0)[\mathbf{h}^*] \\ + o(|\hat{\beta}_n - \beta_0|) + O(\|\hat{\Lambda}_n - \Lambda_0\|^\alpha) = o_{p^*}(n^{-1/2}). \end{aligned} \quad (4)$$

Because  $\alpha\gamma > 1/2$ , then the rate of convergence assumption 1 implies

$$\sqrt{n}O(\|\hat{\Lambda}_n - \Lambda_0\|^\alpha) = o_{p^*}(1).$$

Thus by A4 and (2.3.4) minus (2.3.5), it follows that

$$\begin{aligned} (\dot{S}_{11}(\beta_0, \Lambda_0) - \dot{S}_{21}(\beta_0, \Lambda_0)[\mathbf{h}^*])(\hat{\beta}_n - \beta_0) + o(|\hat{\beta}_n - \beta_0|) \\ = -(S_{1n}(\beta_0, \Lambda_0) - S_{2n}(\beta_0, \Lambda_0)[\mathbf{h}^*]) + o_{p^*}(n^{-1/2}), \end{aligned}$$

i.e.

$$-(A + o(1))(\hat{\beta}_n - \beta_0) = -\mathbb{P}_n m^*(\beta_0, \Lambda_0; X) + o_p^*(n^{-1/2}).$$

Hence

$$\begin{aligned} \sqrt{n}(\hat{\beta}_n - \beta_0) &= (A + o(1))^{-1} \sqrt{n} \mathbb{P}_n m^*(\beta_0, \Lambda_0; X) + o_p^*(1) \\ &\rightarrow_d N\left(0, A^{-1} B (A^{-1})'\right). \end{aligned}$$

□

## References:

- Bickel, P. J., Klaassen, C. A. J., Ritov, Y., and Wellner, J. A. (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins University Press, Baltimore. Softcover edition published by Springer-Verlag, (1998).
- Breslow, N., McNeney, B., and Wellner, J. A. (2003). Large sample theory for semiparametric regression models with two-phase outcome dependent sampling. *Ann. Statist.* **31**, 1110-1139.
- Breslow, N., and Wellner, J. A. (2007). Weighted likelihood for semiparametric models and two-phase stratified samples, with application to Cox regression. *Scand. J. Statist.* **34**, 86 - 102.
- Breslow, N., and Wellner, J. A. (2008). A Z-theorem with estimated nuisance parameters and correction note for ‘Weighted likelihood for semiparametric models and two-phase stratified samples, with application to Cox regression’. *Scand. J. Statist.* **35**, 186-192.
- Gilbert, Peter B. (2000). Large sample theory of maximum likelihood estimates in semiparametric biased sampling models. *Ann. Statist.* **28**, 151-194.
- Hu, Hui-lin. (1998). Large sample theory for Pseudo-Maximum Likelihood Estimates in Semiparametric Models. *Unpublished Ph.D. dissertation*, University of Washington, Seattle.
- Huang, J. (1996). Efficient estimation for the Cox model with interval censoring. *Annals of Statistics* **24**, 540-568.
- Huber, P. J. (1967). The behavior of maximum likelihood estimates under non-standard conditions. *Proc. Fifth Berkeley Symp. Math. Statist. Prob.* **1**, 221 - 233. Univ. California Press.

- Murphy, S. A. (1995). Asymptotic theory for the frailty model. *Ann. Statist.* **23**, 182-198.
- Murphy, S. A. ; van der Vaart, A. W. (1997). Semiparametric likelihood ratio inference. *Ann. Statist.* **25**, 1471-1509.
- Parner, Erik (1998). Asymptotic theory for the correlated gamma-frailty model. *Ann. Statist.* **26**, 183-214.
- Pollard, D. (1985). New ways to prove central limit theorems. *Econometric Theory* **1**, 295 - 314.
- Van der Vaart, A. W. (1995a). Efficiency of infinite dimensional M-estimators. *Statistica Neerl.* **49**, 9 - 30.
- Van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer-Verlag, New York.
- Wellner, J. A. and Zhang, Y. (2007). Two likelihood-based semiparametric estimation methods for panel count data with covariates. *Ann. Statist.* **35**, 2106-2142.
- Zhang, Y. (1998). Estimation for Counting processes with Incomplete Data. Unpublished Ph.D. dissertation, University of Washington.