

How many iterations are sufficient for semiparametric estimation?

Guang Cheng

Department of Statistics, Purdue University

Statistics Department Colloquium
Northwestern University
March 9th, 2011

Outline

Introduction

Semiparametric Models

General Iterative Estimation Procedure

Grid Search of the Initial Estimate

Semiparametric Maximum Likelihood Estimation

Example I: Cox Model with Current Status Data

Semiparametric Estimation under Regularization

Example II: Conditionally Normal Model

Example III: Sparse and Efficient Est. of Partial Spline Model

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.

- ▶ Example I: The Cox regression model with survival data

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.

- ▶ Example I: The Cox regression model with survival data
 - θ : regression covariate;
 - η : cumulative hazard function.

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.
- ▶ Example I: The Cox regression model with survival data
 - θ : regression covariate;
 - η : cumulative hazard function.
- ▶ Example II: The conditional normal model

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.
- ▶ Example I: The Cox regression model with survival data
 - θ : regression covariate;
 - η : cumulative hazard function.
- ▶ Example II: The conditional normal model
 - Conditional distribution $Y|(W = w, Z = z) \sim N(\theta'w, \eta(z))$

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.
- ▶ Example I: The Cox regression model with survival data
 - θ : regression covariate;
 - η : cumulative hazard function.
- ▶ Example II: The conditional normal model
 - Conditional distribution $Y|(W = w, Z = z) \sim N(\theta'w, \eta(z))$
- ▶ Example III: The partly linear model

Semiparametric Models

- ▶ Random Variable $X \sim \{P_{\theta, \eta} : \theta \in \Theta \subset \mathbb{R}^k, \eta \in \mathcal{H}\}$
 - θ : an Euclidean parameter of interest;
 - η : a *possibly* infinite dimensional nuisance parameter.
- ▶ Example I: The Cox regression model with survival data
 - θ : regression covariate;
 - η : cumulative hazard function.
- ▶ Example II: The conditional normal model
 - Conditional distribution $Y|(W = w, Z = z) \sim N(\theta'w, \eta(z))$
- ▶ Example III: The partly linear model
 - $Y = W'\theta + \eta(Z) + \epsilon.$

- ▶ Under regularity conditions, the semiparametric MLE $\hat{\theta}$ is shown to be efficient in the sense that it achieves the minimal asymptotic variance over all regular estimates.

- ▶ Under regularity conditions, the semiparametric MLE $\hat{\theta}$ is shown to be efficient in the sense that it achieves the minimal asymptotic variance over all regular estimates.
- ▶ A common practice to obtain the MLE $\hat{\theta}$ is through maximizing its log-profile likelihood

$$\log pl_n(\theta) = \sup_{\eta \in \mathcal{H}} \log lik_n(\theta, \eta)$$

via some optimization algorithm.

- ▶ Under regularity conditions, the semiparametric MLE $\hat{\theta}$ is shown to be efficient in the sense that it achieves the minimal asymptotic variance over all regular estimates.
- ▶ A common practice to obtain the MLE $\hat{\theta}$ is through maximizing its log-profile likelihood

$$\log pl_n(\theta) = \sup_{\eta \in \mathcal{H}} \log lik_n(\theta, \eta)$$

via some optimization algorithm.

- ▶ For example, the Newton-Raphson (NR) algorithm is applied to the partial likelihood of the Cox model to obtain $\hat{\theta}$ in \mathbf{R} .

A **general iterative** procedure is implemented as follows

A **general iterative** procedure is implemented as follows

(I) Identify an initial estimate $\hat{\theta}^{(0)}$;

A **general iterative** procedure is implemented as follows

- (I) Identify an initial estimate $\hat{\theta}^{(0)}$;
- (II) Construct the corresponding nuisance estimate $\hat{\eta}(\hat{\theta}^{(0)})$ either by pure nonparametric approach, e.g., isotonic estimation, or under some regularization, e.g., kernel or sieve estimation;

A **general iterative** procedure is implemented as follows

- (I) Identify an initial estimate $\hat{\theta}^{(0)}$;
- (II) Construct the corresponding nuisance estimate $\hat{\eta}(\hat{\theta}^{(0)})$ either by pure nonparametric approach, e.g., isotonic estimation, or under some regularization, e.g., kernel or sieve estimation;
- (III) Apply NR (or other optimization) algorithm to the generalized profile likelihood (Severini and Wong, 1992)

$$\hat{S}(\theta) \equiv \log \text{lik}_n(\theta, \hat{\eta}(\theta))$$

at $\theta = \hat{\theta}^{(0)}$ to obtain $\hat{\theta}^{(1)}$;

A **general iterative** procedure is implemented as follows

- (I) Identify an initial estimate $\hat{\theta}^{(0)}$;
- (II) Construct the corresponding nuisance estimate $\hat{\eta}(\hat{\theta}^{(0)})$ either by pure nonparametric approach, e.g., isotonic estimation, or under some regularization, e.g., kernel or sieve estimation;
- (III) Apply NR (or other optimization) algorithm to the generalized profile likelihood (Severini and Wong, 1992)

$$\hat{S}(\theta) \equiv \log \text{lik}_n(\theta, \hat{\eta}(\theta))$$

at $\theta = \hat{\theta}^{(0)}$ to obtain $\hat{\theta}^{(1)}$;

- (IV) Repeat k^* iterations until $|\hat{S}(\hat{\theta}^{(k^*)}) - \hat{S}(\hat{\theta}^{(k^*-1)})| \leq \epsilon$ for some pre-determined sufficiently small ϵ .

- ▶ Extensively applied in the semiparametric literature including

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);
 - ▶ Semiparametric Mixture Models (Roeder et al 1996);

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);
 - ▶ Semiparametric Mixture Models (Roeder et al 1996);
 - ▶ Semiparametric Transformation Model (Linton et al 2008);

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);
 - ▶ Semiparametric Mixture Models (Roeder et al 1996);
 - ▶ Semiparametric Transformation Model (Linton et al 2008);
 - ▶ Generalized Partly Linear (Single Index) Model (Fan et al 1995);

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);
 - ▶ Semiparametric Mixture Models (Roeder et al 1996);
 - ▶ Semiparametric Transformation Model (Linton et al 2008);
 - ▶ Generalized Partly Linear (Single Index) Model (Fan et al 1995);
 - ▶ Survival Models (Huang 1996; Murphy et al., 1997);

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);
 - ▶ Semiparametric Mixture Models (Roeder et al 1996);
 - ▶ Semiparametric Transformation Model (Linton et al 2008);
 - ▶ Generalized Partly Linear (Single Index) Model (Fan et al 1995);
 - ▶ Survival Models (Huang 1996; Murphy et al., 1997);
- ▶ The above iterative procedure can also be adapted to the penalized estimation and selection of semiparametric models, e.g., Cheng and Zhang (2010).

- ▶ Extensively applied in the semiparametric literature including
 - ▶ Quasi Likelihood Estimation (Severini and Staniwalis 1994);
 - ▶ Semiparametric Mixture Models (Roeder et al 1996);
 - ▶ Semiparametric Transformation Model (Linton et al 2008);
 - ▶ Generalized Partly Linear (Single Index) Model (Fan et al 1995);
 - ▶ Survival Models (Huang 1996; Murphy et al., 1997);
- ▶ The above iterative procedure can also be adapted to the penalized estimation and selection of semiparametric models, e.g., Cheng and Zhang (2010).
- ▶ The choice of ϵ or k^* is quite **arbitrary** in the above papers.

In this talk, we will answer

“How many iterations do we really need?”

equivalently,

“ $k^* = ?$ for obtaining a semiparametric efficient $\hat{\theta}^{(k^*)}$?”

from a theoretical point of view.

In this talk, we will answer

“How many iterations do we really need?”

equivalently,

“ $k^* = ?$ for obtaining a semiparametric efficient $\hat{\theta}^{(k^*)}$?”

from a theoretical point of view.

- ▶ k^* depends on the convergence rates of $\hat{\theta}^{(0)}$ and $\hat{\eta}(\theta)$.

In this talk, we will answer

“How many iterations do we really need?”

equivalently,

“ $k^* = ?$ for obtaining a semiparametric efficient $\hat{\theta}^{(k^*)}$?”

from a theoretical point of view.

- ▶ k^* depends on the convergence rates of $\hat{\theta}^{(0)}$ and $\hat{\eta}(\theta)$.
- ▶ k^* depends on the bandwidth order if kernel approach is used;

In this talk, we will answer

“How many iterations do we really need?”

equivalently,

“ $k^* = ?$ for obtaining a semiparametric efficient $\hat{\theta}^{(k^*)}$?”

from a theoretical point of view.

- ▶ k^* depends on the convergence rates of $\hat{\theta}^{(0)}$ and $\hat{\eta}(\theta)$.
 - ▶ k^* depends on the bandwidth order if kernel approach is used;
 - ▶ k^* depends on the order of the smoothing parameter if the penalization approach is used;

Why are our results useful?

Why are our results useful?

- ▶ Knowing the minimal k^* for each bootstrap sample will significantly reduce the bootstrap computational cost for making semiparametric inferences.

Grid Search of the Initial Estimate

- ▶ It is critical to initiate the iterations in a suitable neighborhood of the true value θ_0 , i.e., $\hat{\theta}^{(0)}$ is assumed to be n^ψ -consistent.

Grid Search of the Initial Estimate

- ▶ It is critical to initiate the iterations in a suitable neighborhood of the true value θ_0 , i.e., $\hat{\theta}^{(0)}$ is assumed to be n^ψ -consistent.
- ▶ Without any prior model information, a natural way is to do grid search of $\hat{S}(\theta)$. Recall that $\hat{S}(\theta) = \log \text{lik}_n(\theta, \hat{\eta}(\theta))$.

Grid Search of the Initial Estimate

- ▶ It is critical to initiate the iterations in a suitable neighborhood of the true value θ_0 , i.e., $\hat{\theta}^{(0)}$ is assumed to be n^ψ -consistent.
- ▶ Without any prior model information, a natural way is to do grid search of $\hat{S}(\theta)$. Recall that $\hat{S}(\theta) = \log \text{lik}_n(\theta, \hat{\eta}(\theta))$.
- ▶ We next present two types of grid search algorithm:

Grid Search of the Initial Estimate

- ▶ It is critical to initiate the iterations in a suitable neighborhood of the true value θ_0 , i.e., $\hat{\theta}^{(0)}$ is assumed to be n^ψ -consistent.
- ▶ Without any prior model information, a natural way is to do grid search of $\hat{S}(\theta)$. Recall that $\hat{S}(\theta) = \log \text{lik}_n(\theta, \hat{\eta}(\theta))$.
- ▶ We next present two types of grid search algorithm:
 - ▶ Deterministic grid search;

Grid Search of the Initial Estimate

- ▶ It is critical to initiate the iterations in a suitable neighborhood of the true value θ_0 , i.e., $\hat{\theta}^{(0)}$ is assumed to be n^ψ -consistent.
- ▶ Without any prior model information, a natural way is to do grid search of $\hat{S}(\theta)$. Recall that $\hat{S}(\theta) = \log \text{lik}_n(\theta, \hat{\eta}(\theta))$.
- ▶ We next present two types of grid search algorithm:
 - ▶ Deterministic grid search;
 - ▶ Stochastic grid search.

Grid Search of the Initial Estimate

- ▶ It is critical to initiate the iterations in a suitable neighborhood of the true value θ_0 , i.e., $\hat{\theta}^{(0)}$ is assumed to be n^ψ -consistent.
- ▶ Without any prior model information, a natural way is to do grid search of $\hat{S}(\theta)$. Recall that $\hat{S}(\theta) = \log \text{lik}_n(\theta, \hat{\eta}(\theta))$.
- ▶ We next present two types of grid search algorithm:
 - ▶ Deterministic grid search;
 - ▶ Stochastic grid search.
- ▶ We will calculate the convergence rate of the above numerical outcome. The technical challenge is that $\hat{S}(\theta)$ usually has no explicit form and may not be continuous/smooth.

Primary Assumptions

S1. Θ is compact;

Primary Assumptions

- S1. Θ is compact;
- S2. [Asymptotic Uniqueness] For any random sequence $\{\tilde{\theta}_n\} \in \Theta$,
 $[\hat{S}(\tilde{\theta}) - \hat{S}(\hat{\theta})]/n = o_P(1)$ implies that $\tilde{\theta} - \theta_0 = o_P(1)$.

Primary Assumptions

- S1. Θ is compact;
- S2. [Asymptotic Uniqueness] For any random sequence $\{\tilde{\theta}_n\} \in \Theta$,
 $[\hat{S}(\tilde{\theta}) - \hat{S}(\hat{\theta})]/n = o_P(1)$ implies that $\tilde{\theta} - \theta_0 = o_P(1)$.
- S3. [Asymptotic Concavity] For any consistent $\tilde{\theta}$, $\hat{S}(\cdot)$ satisfies

$$\begin{aligned} \hat{S}(\tilde{\theta}) &= \hat{S}(\theta_0) + n(\tilde{\theta} - \theta_0)' \mathbb{P}_n \tilde{\ell}_0 - \frac{n}{2} (\tilde{\theta} - \theta_0)' \tilde{I}_0 (\tilde{\theta} - \theta_0) \\ &\quad + \Delta_n(\tilde{\theta}), \end{aligned}$$

where $\Delta_n(\theta) = n \|\theta - \theta_0\|^3 \vee n^{1-2\gamma} \|\theta - \theta_0\|$ and γ represents the convergence rate of $\hat{\eta}(\theta)$ given later.

Remark

- ▶ Condition S2 is usually satisfied if the model is identifiable.

Remark

- ▶ Condition S2 is usually satisfied if the model is identifiable.
- ▶ Condition S3 is very weak since we only require that $\widehat{S}(\theta)$ has such an asymptotic expansion, but not that $\widehat{S}(\theta)$ is continuous.

Remark

- ▶ Condition S2 is usually satisfied if the model is identifiable.
- ▶ Condition S3 is very weak since we only require that $\widehat{S}(\theta)$ has such an asymptotic expansion, but not that $\widehat{S}(\theta)$ is continuous.
- ▶ Condition S3 can be implied by some smoothness and empirical processes conditions (concerning about the least favorable submodel).

Deterministic Search

- ▶ Form a grid of cubes with sides of length $sn^{-\psi}$ over \mathbb{R}^d for some $s > 0$ and $0 < \psi \leq 1/2$;

Deterministic Search

- ▶ Form a grid of cubes with sides of length $sn^{-\psi}$ over \mathbb{R}^d for some $s > 0$ and $0 < \psi \leq 1/2$;
- ▶ Obtain a set of points $\mathcal{D}_n = \{\theta_{iD}\}$ regularly spaced throughout Θ with cardinality $\text{card}(\mathcal{D}_n) \geq Cn^{d\psi}$ for some $C > 0$;

Deterministic Search

- ▶ Form a grid of cubes with sides of length $sn^{-\psi}$ over \mathbb{R}^d for some $s > 0$ and $0 < \psi \leq 1/2$;
- ▶ Obtain a set of points $\mathcal{D}_n = \{\theta_{iD}\}$ regularly spaced throughout Θ with cardinality $\text{card}(\mathcal{D}_n) \geq Cn^{d\psi}$ for some $C > 0$;
- ▶ Define $\widehat{\theta}_D^{(0)} = \arg \max_{\mathcal{D}_n} \widehat{S}(\theta)$.

Stochastic Search

However, the above deterministic search could be very slow if the dimension d of θ is high. This motivates the following computationally efficient stochastic search, i.e., of the order n^ψ .

Stochastic Search

However, the above deterministic search could be very slow if the dimension d of θ is high. This motivates the following computationally efficient stochastic search, i.e., of the order n^ψ .

- ▶ Assume $\bar{\theta}$ is independent of the data and admits a density having support Θ and bounded away from zero in some neighborhood of θ_0 ;

Stochastic Search

However, the above deterministic search could be very slow if the dimension d of θ is high. This motivates the following computationally efficient stochastic search, i.e., of the order n^ψ .

- ▶ Assume $\bar{\theta}$ is independent of the data and admits a density having support Θ and bounded away from zero in some neighborhood of θ_0 ;
- ▶ Let \mathcal{S}_n be a set of realizations of $\bar{\theta}$ with $\text{card}(\mathcal{S}_n) \geq \tilde{C}n^\psi$ for some $\tilde{C} > 0$;

Stochastic Search

However, the above deterministic search could be very slow if the dimension d of θ is high. This motivates the following computationally efficient stochastic search, i.e., of the order n^ψ .

- ▶ Assume $\bar{\theta}$ is independent of the data and admits a density having support Θ and bounded away from zero in some neighborhood of θ_0 ;
- ▶ Let \mathcal{S}_n be a set of realizations of $\bar{\theta}$ with $\text{card}(\mathcal{S}_n) \geq \tilde{C}n^\psi$ for some $\tilde{C} > 0$;
- ▶ Define $\hat{\theta}_S^{(0)} = \arg \max_{\mathcal{S}_n} \hat{S}(\theta)$.

Initial Estimate Theorem

Suppose Conditions S1-S3 hold. If $\hat{\theta}$ is consistent and the efficient information matrix \tilde{I}_0 is nonsingular, then we have

$$\theta_D^{(0)} - \theta_0 = O_P(n^{-\psi}), \quad (1)$$

$$\theta_S^{(0)} - \theta_0 = O_P(n^{-\psi}). \quad (2)$$

The above Theorem can be applied to a wide range of semiparametric models including the conditionally normal model, Cox model under survival data and semiparametric mixture model.

Semiparametric Maximum Likelihood Estimation

We first consider the maximum likelihood estimation of θ .

Semiparametric Maximum Likelihood Estimation

We first consider the maximum likelihood estimation of θ .

- ▶ For each fixed θ , η is estimated as a possibly nonsmooth NPMLE $\hat{\eta}(\theta)$ (usually under some shape constraints);

Semiparametric Maximum Likelihood Estimation

We first consider the maximum likelihood estimation of θ .

- ▶ For each fixed θ , η is estimated as a possibly nonsmooth NPMLE $\hat{\eta}(\theta)$ (usually under some shape constraints);
- ▶ $\hat{S}(\theta)$ becomes the log-profile likelihood $\log p'_n(\theta)$;

Semiparametric Maximum Likelihood Estimation

We first consider the maximum likelihood estimation of θ .

- ▶ For each fixed θ , η is estimated as a possibly nonsmooth NPMLE $\hat{\eta}(\theta)$ (usually under some shape constraints);
- ▶ $\hat{S}(\theta)$ becomes the log-profile likelihood $\log p'_n(\theta)$;
- ▶ **Challenge:** the profile likelihood is defined as a supremum over an infinite dimensional parameter space, and thus has no closed form and is possibly nonsmooth (although it can be computed numerically in practice).

The Construction of $\hat{\theta}^{(k)}$

We construct $\hat{\theta}^{(k)}$ as the following Newton-Raphson form:

$$\hat{\theta}^{(k)} = \hat{\theta}^{(k-1)} + \left[\hat{I}(\hat{\theta}^{(k-1)}, t_n^{(k-1)}) \right]^{-1} \hat{\ell}(\hat{\theta}^{(k-1)}, s_n^{(k-1)}), \quad (3)$$

where $\hat{\ell}(\theta, s_n)$ and $\hat{I}(\theta, t_n)$ are the first and second numerical derivatives of $\log pl_n(\theta)$ with step sizes $s_n, t_n \rightarrow 0$, respectively.

The Construction of $\hat{\theta}^{(k)}$

We construct $\hat{\theta}^{(k)}$ as the following Newton-Raphson form:

$$\hat{\theta}^{(k)} = \hat{\theta}^{(k-1)} + \left[\hat{I}(\hat{\theta}^{(k-1)}, t_n^{(k-1)}) \right]^{-1} \hat{\ell}(\hat{\theta}^{(k-1)}, s_n^{(k-1)}), \quad (3)$$

where $\hat{\ell}(\theta, s_n)$ and $\hat{I}(\theta, t_n)$ are the first and second numerical derivatives of $\log pl_n(\theta)$ with step sizes $s_n, t_n \rightarrow 0$, respectively.

- ▶ We use the **numerical derivatives** of $\log pl_n(\theta)$ since its differentiability is usually unknown;

The Construction of $\hat{\theta}^{(k)}$

We construct $\hat{\theta}^{(k)}$ as the following Newton-Raphson form:

$$\hat{\theta}^{(k)} = \hat{\theta}^{(k-1)} + \left[\hat{I}(\hat{\theta}^{(k-1)}, t_n^{(k-1)}) \right]^{-1} \hat{\ell}(\hat{\theta}^{(k-1)}, s_n^{(k-1)}), \quad (3)$$

where $\hat{\ell}(\theta, s_n)$ and $\hat{I}(\theta, t_n)$ are the first and second numerical derivatives of $\log pl_n(\theta)$ with step sizes $s_n, t_n \rightarrow 0$, respectively.

- ▶ We use the **numerical derivatives** of $\log pl_n(\theta)$ since its differentiability is usually unknown;
- ▶ A close inspection of (3) reveals that we have constructed $\hat{\theta}^{(k)}$ even **without** knowing the forms of efficient score function $\tilde{\ell}_0$ and efficient information matrix \tilde{I}_0 .

Intuition of calculating k^* :

Intuition of calculating k^* :

- ▶ We expect that $\hat{\theta}^{(k)}$ approaches to MLE $\hat{\theta}$, which is exactly the maximizer of $\log pI_n(\theta)$, asymptotically as $k \rightarrow \infty$ if $\hat{\ell}(\cdot)$ and $\hat{I}(\cdot)$ are consistent estimators of $\tilde{\ell}_0$ and \hat{I}_0 .

Intuition of calculating k^* :

- ▶ We expect that $\hat{\theta}^{(k)}$ approaches to MLE $\hat{\theta}$, which is exactly the maximizer of $\log p l_n(\theta)$, asymptotically as $k \rightarrow \infty$ if $\hat{\ell}(\cdot)$ and $\hat{I}(\cdot)$ are consistent estimators of $\tilde{\ell}_0$ and \tilde{I}_0 .
 - ▶ The above consistency can be guaranteed by the concave form of $\log p l_n(\theta)$ shown in Murphy and van der Vaart (2000).

Intuition of calculating k^* :

- ▶ We expect that $\hat{\theta}^{(k)}$ approaches to MLE $\hat{\theta}$, which is exactly the maximizer of $\log p l_n(\theta)$, asymptotically as $k \rightarrow \infty$ if $\hat{\ell}(\cdot)$ and $\hat{I}(\cdot)$ are consistent estimators of $\tilde{\ell}_0$ and \tilde{I}_0 .
 - ▶ The above consistency can be guaranteed by the concave form of $\log p l_n(\theta)$ shown in Murphy and van der Vaart (2000).
- ▶ We can further quantify how fast $\hat{\theta}^{(k)}$ converges to $\hat{\theta}$ if we know how fast $\hat{\ell}(\hat{I})$ converges to $\tilde{\ell}_0(\tilde{I}_0)$.

Intuition of calculating k^* :

- ▶ We expect that $\hat{\theta}^{(k)}$ approaches to MLE $\hat{\theta}$, which is exactly the maximizer of $\log p l_n(\theta)$, asymptotically as $k \rightarrow \infty$ if $\hat{\ell}(\cdot)$ and $\hat{I}(\cdot)$ are consistent estimators of $\tilde{\ell}_0$ and \tilde{I}_0 .
 - ▶ The above consistency can be guaranteed by the concave form of $\log p l_n(\theta)$ shown in Murphy and van der Vaart (2000).
- ▶ We can further quantify how fast $\hat{\theta}^{(k)}$ converges to $\hat{\theta}$ if we know how fast $\hat{\ell}(\hat{I})$ converges to $\tilde{\ell}_0(\tilde{I}_0)$.
 - ▶ The above convergence rates of $\hat{\ell}(\cdot)$ and $\hat{I}(\cdot)$ are derived based on a higher order quadratic expansion of $\log p l_n(\theta)$.

Primary Assumptions

1. Regularity conditions on the least favorable submodels (guarantee the valid **higher order** quadratic expansion of the log-profile likelihood);

Primary Assumptions

- I. Regularity conditions on the least favorable submodels (guarantee the valid **higher order** quadratic expansion of the log-profile likelihood);
- II. We assume that, for any random sequence $\tilde{\theta}_n \xrightarrow{P} \theta_0$,

$$\|\hat{\eta}(\tilde{\theta}_n) - \eta_0\| = O_P(\|\tilde{\theta}_n - \theta_0\| \vee n^{-\gamma}), \quad (4)$$

where $\|\cdot\|$ is some norm in \mathcal{H} , for some $1/4 < \gamma \leq 1/2$.

Theorem 1

Suppose Conditions I&II hold and proper step sizes are chosen. Let $\|\widehat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\widehat{\theta}^{(k)} - \widehat{\theta}\| = O_P(n^{-r_k})$.

We show

Theorem 1

Suppose Conditions I&II hold and proper step sizes are chosen. Let $\|\widehat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\widehat{\theta}^{(k)} - \widehat{\theta}\| = O_P(n^{-r_k})$.

We show

- ▶ r_k continuously increases from ψ to $(\gamma + 1/4)$ as $k \rightarrow \infty$;

Theorem 1

Suppose Conditions I&II hold and proper step sizes are chosen. Let $\|\widehat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\widehat{\theta}^{(k)} - \widehat{\theta}\| = O_P(n^{-r_k})$.

We show

- ▶ r_k continuously increases from ψ to $(\gamma + 1/4)$ as $k \rightarrow \infty$;
- ▶ Specifically,

$$\|\widehat{\theta}^{(k)} - \widehat{\theta}\| = O_P(n^{-S(\psi, \gamma, k)}),$$

where $S(\psi, \gamma, k)$ has an easy-to-calculate explicit form.

Theorem 1

Suppose Conditions I&II hold and proper step sizes are chosen. Let $\|\widehat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\widehat{\theta}^{(k)} - \widehat{\theta}\| = O_P(n^{-r_k})$.

We show

- ▶ r_k continuously increases from ψ to $(\gamma + 1/4)$ as $k \rightarrow \infty$;
- ▶ Specifically,

$$\|\widehat{\theta}^{(k)} - \widehat{\theta}\| = O_P(n^{-S(\psi, \gamma, k)}),$$

where $S(\psi, \gamma, k)$ has an easy-to-calculate explicit form.

- ▶ $\|\widehat{\theta}^{(k^*)} - \widehat{\theta}\| = o_P(n^{-1/2})$ for $k^* \geq K(\psi, \gamma)$.

Remark 1

- ▶ It is well known that one-step estimate is efficient given that $\hat{\theta}^{(0)}$ is \sqrt{n} -consistent. However, this is not enough for the semiparametric estimation since $\hat{\theta}^{(0)}$ may have slower than \sqrt{n} rate as shown in the previous Theorem.

Remark 1

- ▶ It is well known that one-step estimate is efficient given that $\hat{\theta}^{(0)}$ is \sqrt{n} -consistent. However, this is not enough for the semiparametric estimation since $\hat{\theta}^{(0)}$ may have slower than \sqrt{n} rate as shown in the previous Theorem.
- ▶ More than k^* iterations, say k , will only improve the higher order asymptotic efficiency of $\hat{\theta}^{(k)}$.

Remark 1

- ▶ It is well known that one-step estimate is efficient given that $\hat{\theta}^{(0)}$ is \sqrt{n} -consistent. However, this is not enough for the semiparametric estimation since $\hat{\theta}^{(0)}$ may have slower than \sqrt{n} rate as shown in the previous Theorem.
- ▶ More than k^* iterations, say k , will only improve the higher order asymptotic efficiency of $\hat{\theta}^{(k)}$.
- ▶ Interestingly, the lower bound of $\|\hat{\theta}^{(k)} - \hat{\theta}\|$, i.e. $O_P(n^{-\gamma-1/4})$, is **intrinsically** decided, i.e., only dependent on the convergence rate of the nuisance parameter.

Example I: Cox Model with Current Status Data

The hazard function $\lambda(t|z)$ of the survival time T given the covariate Z is modeled as, with λ as the hazard function,

$$\lambda(t) \exp(\theta'z).$$

Current status data: observe $X = (Y, I\{T \leq Y\}, Z)$, where Y is the examination time.

We are interested in θ while treating the cumulative hazard function $\eta(y) = \int_0^y \lambda(t)dt$ as the nuisance parameter. The NPMLE $\hat{\eta}(\theta)$ and nonsmooth $\log pl_n(\theta)$ have no explicit forms, but can be calculated numerically via isotonic regression type algorithm.

The convergence rate of $\hat{\eta}(\theta)$ is shown to be $n^{-1/3}$. Thus, Theorem 1 implies the following table with $O_P(n^{-7/12})$ as the lower bound.

Table 1. *Cox Model under Current Status Data* ($\gamma = 1/3$)

$\psi = 1/2$	$\psi = 1/3$	$\psi = 1/4$
$r_1 = 7/12$	$r_1 = 1/2, r_2 = 7/12$	$r_1 = 3/8, r_2 = 25/48, r_3 = 7/12$
$k^* = 1$	$k^* = 2$	$k^* = 2$

Recall that $\|\hat{\theta}^{(k)} - \hat{\theta}\| = O_P(n^{-r_k})$.

In contrast with the nonsmooth $\log p_l n(\theta)$, the regularized $\widehat{S}(\theta)$ is usually differentiable although its form may vary under different regularizations.

In contrast with the nonsmooth $\log p_l n(\theta)$, the regularized $\widehat{S}(\theta)$ is usually differentiable although its form may vary under different regularizations.

- A. We first present a unified framework for studying $\widehat{\theta}^{(k)}$ when $\widehat{S}(\theta)$ is third order differentiable.

In contrast with the nonsmooth $\log p_l n(\theta)$, the regularized $\widehat{S}(\theta)$ is usually differentiable although its form may vary under different regularizations.

- A. We first present a unified framework for studying $\widehat{\theta}^{(k)}$ when $\widehat{S}(\theta)$ is third order differentiable.
- B. We then discuss two special cases: kernel estimation of θ and penalized estimation of θ , in two examples.

In contrast with the nonsmooth $\log p_l(\theta)$, the regularized $\widehat{S}(\theta)$ is usually differentiable although its form may vary under different regularizations.

- A. We first present a unified framework for studying $\widehat{\theta}^{(k)}$ when $\widehat{S}(\theta)$ is third order differentiable.
- B. We then discuss two special cases: kernel estimation of θ and penalized estimation of θ , in two examples.
- C. In the end, we consider the semiparametric model selection as an extension of the penalized estimation.

Construction of $\widehat{\theta}^{(k)}$

Let $\widehat{\ell}(\cdot) = \widehat{S}^{(1)}(\cdot)/n$ and $\widehat{I}(\cdot) = -\widehat{S}^{(2)}(\cdot)/n$. We construct $\widehat{\theta}^{(k)}$ as

$$\widehat{\theta}^{(k)} = \widehat{\theta}^{(k-1)} + \left[\widehat{I}(\widehat{\theta}^{(k-1)}) \right]^{-1} \widehat{\ell}(\widehat{\theta}^{(k-1)}). \quad (5)$$

However, $\widehat{I}(\cdot)$ can also be constructed as the negative numerical derivative of $\widehat{\ell}(\cdot)$ when $\widehat{S}^{(2)}(\cdot)$ has no explicit form or is hard to compute.

Primary Conditions

The higher order quadratic expansion of $\widehat{S}(\theta)$ is valid under the following Condition G:

$$\frac{1}{n}\widehat{S}^{(1)}(\theta_0) - \frac{1}{n}S^{(1)}(\theta_0) = O_P(n^{-2g}), \quad (6)$$

$$\sup_{\theta \in \mathcal{N}(\theta_0)} \left| \frac{1}{n}\widehat{S}^{(2)}(\theta) - \frac{1}{n}S^{(2)}(\theta) \right| = O_P(n^{-g}), \quad (7)$$

$$\sup_{\theta \in \mathcal{N}(\theta_0)} \left| \frac{1}{n}\widehat{S}^{(3)}(\theta) \right| = O_P(1), \quad (8)$$

where $S(\theta) = \sup_{\eta \in \mathcal{H}} E \log \text{lik}(\theta, \eta)$ and $1/4 < g \leq 1/2$.

Theorem 2

Suppose Condition G holds and define $\hat{\theta}$ as the maximizer of $\hat{S}(\theta)$. Recall that $\|\hat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\hat{\theta}^{(k)} - \hat{\theta}\| = O_P(n^{-r_k})$.

Theorem 2

Suppose Condition G holds and define $\hat{\theta}$ as the maximizer of $\hat{S}(\theta)$. Recall that $\|\hat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\hat{\theta}^{(k)} - \hat{\theta}\| = O_P(n^{-r_k})$.

- ▶ $\hat{\theta}$ is semiparametric efficient;

Theorem 2

Suppose Condition G holds and define $\hat{\theta}$ as the maximizer of $\hat{S}(\theta)$. Recall that $\|\hat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\hat{\theta}^{(k)} - \hat{\theta}\| = O_P(n^{-r_k})$.

- ▶ $\hat{\theta}$ is semiparametric efficient;
- ▶ r_k continuously increases from ψ to $+\infty$ as $k \rightarrow \infty$;

Theorem 2

Suppose Condition G holds and define $\hat{\theta}$ as the maximizer of $\hat{S}(\theta)$. Recall that $\|\hat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\hat{\theta}^{(k)} - \hat{\theta}\| = O_P(n^{-r_k})$.

- ▶ $\hat{\theta}$ is semiparametric efficient;
- ▶ r_k continuously increases from ψ to $+\infty$ as $k \rightarrow \infty$;
- ▶ If $\hat{\theta}$ is constructed analytically, then

$$\begin{aligned} \|\hat{\theta}^{(k)} - \hat{\theta}\| &= O_P(n^{-2^k \psi}), \\ k^* &= L_1(\psi). \end{aligned}$$

Theorem 2

Suppose Condition G holds and define $\hat{\theta}$ as the maximizer of $\hat{S}(\theta)$. Recall that $\|\hat{\theta}^{(0)} - \theta_0\| = O_P(n^{-\psi})$ and $\|\hat{\theta}^{(k)} - \hat{\theta}\| = O_P(n^{-r_k})$.

- ▶ $\hat{\theta}$ is semiparametric efficient;
- ▶ r_k continuously increases from ψ to $+\infty$ as $k \rightarrow \infty$;
- ▶ If $\hat{\theta}$ is constructed analytically, then

$$\begin{aligned}\|\hat{\theta}^{(k)} - \hat{\theta}\| &= O_P(n^{-2^k \psi}), \\ k^* &= L_1(\psi).\end{aligned}$$

- ▶ If $\hat{\theta}$ is constructed numerically, then

$$\begin{aligned}\|\hat{\theta}^{(k)} - \hat{\theta}\| &= O_P(n^{-R(\psi, g, k)}), \\ k^* &= L_2(\psi, g).\end{aligned}$$

Four interesting applications of Theorem 2:

Four interesting applications of Theorem 2:

1. Parametric models where $\widehat{S}(\theta) = \log \text{lik}_n(\theta)$;

Four interesting applications of Theorem 2:

1. Parametric models where $\widehat{S}(\theta) = \log \text{lik}_n(\theta)$;
2. Kernel estimation of conditionally parametric models, i.e.,

$$\widehat{\eta}(\theta)(z) = \arg \sup_{\eta \in C^2(\mathcal{Z})} \sum_{i=1}^n \log \text{lik}(X_i; \theta, \eta(Z_i)) K\left(\frac{z - Z_i}{b_n}\right).$$

Condition G can be translated to the kernel conditions on $K(\cdot)$ and b_n . **Therefore, k^* is related to the order of b_n in this case.**

3. Penalized estimation of semiparametric models, i.e.,

$$\hat{\eta}_{\lambda_n}(\theta) = \arg \sup_{\eta \in \mathcal{H}_k} \left\{ \frac{1}{n} \sum_{i=1}^n \log \text{lik}(X_i; \theta, \eta) - \lambda_n^2 J^2(\eta) \right\}.$$

In this case, Condition G needs to be modified to take into account of λ_n . Therefore, k^* is related to the order of the smoothing parameter λ_n .

3. Penalized estimation of semiparametric models, i.e.,

$$\hat{\eta}_{\lambda_n}(\theta) = \arg \sup_{\eta \in \mathcal{H}_k} \left\{ \frac{1}{n} \sum_{i=1}^n \log \text{lik}(X_i; \theta, \eta) - \lambda_n^2 J^2(\eta) \right\}.$$

In this case, Condition G needs to be modified to take into account of λ_n . Therefore, k^* is related to the order of the smoothing parameter λ_n .

4. Semiparametric model selection in high dimensional data. We will show that k^* iterations are also sufficient to recover the estimation sparsity.

Example II: Conditionally Normal Model

We assume that $Y|(W = w, Z = z) \sim N(\theta'w, \eta(z))$ and thus

$$\hat{\eta}_\theta(z) = \frac{\sum_{i=1}^n (Y - \theta'W)^2 K((z - Z_i)/b_n)}{\sum_{i=1}^n K((z - Z_i)/b_n)}.$$

Given that $b_n \asymp n^{-1/5}$, Theorem 2 gives the following table.

Table 2. *Conditional Normal Model*

$$\psi = 1/4$$

$$r_1 = 151/600, r_2 = 153/600, r_3 = 157/600, r_4 = 165/600$$

$$r_5 = 181/600, r_6 = 213/600, r_7 = 277/600, r_8 = 405/600$$

$$k^* = 8$$

Example III: Sparse and Efficient Est. of Partial Spline Model

We consider the partial smoothing spline model:

$$Y = W'\theta + \eta(Z) + \epsilon, \quad (9)$$

where η belongs to the k -th order Sobolev space.

Under high dimensional data, it is common to assume that some components of θ_0 are exactly zero.

To achieve the estimation efficiency and sparsity of θ , Cheng and Zhang (2010) proposed the below regularization method

$$\min \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - W_i' \theta - \eta(Z_i))^2 + \lambda_n^2 J^2(\eta) + \tau_n^2 \sum_{j=1}^d \frac{|\theta_j|}{|\tilde{\theta}_j|} \right\},$$

where $\tilde{\theta} = (\tilde{\theta}_1, \dots, \tilde{\theta}_d)'$ is the consistent initial estimate. By incorporating LARS algorithm, the general iterative estimation procedure can also be adapted to this scenario.

Given that $\tilde{\theta}$ is \sqrt{n} -consistent, e.g., partial smoothing spline estimate, $\lambda_n \asymp n^{-k/(2k+1)}$ and $\tau_n \asymp n^{-k/(2k+1)}$, Theorem 2 (after adaptations) shows that $k^* = 1$ is sufficient to produce the efficient and sparse estimate of θ .

Thank you for your attention....

Guang Cheng
Department of Statistics, Purdue University
chengg@purdue.edu