



**Weierstrass Institute for
Applied Analysis and Stochastics**



Non-asymptotic confidence bounds via multiplier bootstrap

Mayya Zhilova (joint with Vladimir Spokoiny)

CRC 649 "ECONOMIC RISK" CONFERENCE

Motzen, 2014

-
- 1 Problem formulation**
 - 2 Likelihood-based confidence set**
 - 3 Multiplier bootstrap**
 - 4 Uniform confidence corridors**

1 Problem formulation

2 Likelihood-based confidence set

3 Multiplier bootstrap

4 Uniform confidence corridors

Generalized regression

- A random sample of a **fixed size** n from unknown distribution:

$$\mathbf{Y} = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n,$$

$$Y_i \sim P_i(f(X_i)), \text{ independent.}$$

- Deterministic design $X_1, \dots, X_n \in \mathbb{R}^d$.
- Unknown function $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$.

- Mean regression

$$Y_i \sim f(X_i) + \varepsilon_i, \quad \varepsilon_i \text{ are i.i.d.}$$

- Logistic regression

$$Y_i \sim \text{Bernoulli}(f(X_i)).$$

- Quantile regression

$$P(Y_i - f(X_i) < 0) = \tau$$

for $\tau \in (0, 1)$.

- Let $Y \sim \mathbb{P}$ (unknown).
- Fix some known parametric family:

$$\{\mathbb{P}(\boldsymbol{\theta}) \ll \mu_0, \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p\}.$$

- The target parameter $\boldsymbol{\theta}^* \in \Theta$ is defined by the projection:

$$\begin{aligned}\boldsymbol{\theta}^* &\stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmin}} KL(\mathbb{P}, \mathbb{P}(\boldsymbol{\theta})) \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \int \log \left\{ \frac{d\mathbb{P}}{d\mathbb{P}(\boldsymbol{\theta})} \right\} d\mathbb{P}.\end{aligned}$$

To find a **sharp confidence set** \mathcal{E} s.t. for a given coverage probability $1 - \alpha \in (0, 1)$ and for a **fixed sample size** n :

$$\mathbb{P} \{ \boldsymbol{\theta}^* \in \mathcal{E} \} \geq 1 - \alpha.$$

1 Problem formulation

2 Likelihood-based confidence set

3 Multiplier bootstrap

4 Uniform confidence corridors

Maximum likelihood method:

- Log-likelihood function

$$\begin{aligned}L(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} \log \frac{d\mathbb{P}(\boldsymbol{\theta})}{d\mu_0}(\mathbf{Y}) \\ &= \sum_{i=1}^n \log \frac{dP_i(\boldsymbol{\theta})}{d\mu_0}(Y_i).\end{aligned}$$

- The target parameter:

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} KL(\mathbb{P}, \mathbb{P}(\boldsymbol{\theta})).$$

- (Quasi) Maximum likelihood estimate (qMLE):

$$\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}).$$

Wilks phenomenon for the likelihood ratio:

For $\mathbb{P} \in \{P(\boldsymbol{\theta})\}$

$$L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) \approx \|\boldsymbol{\xi}\|^2/2, \quad \boldsymbol{\xi} \xrightarrow{w} \mathcal{N}(0, \mathbf{I}_p)$$

with $n \rightarrow \infty$.

p is the dimension of parametric set Θ ,

$\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*)$ normalized score,

$D_0^2 \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E} L(\boldsymbol{\theta}^*)$ full Fisher Information matrix.

Likelihood-based confidence set:

$$\mathcal{E}(\mathfrak{z}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}) \leq \mathfrak{z}\},$$

the level $\mathfrak{z} > 0$ can be found by the Wilks theorem:

$$\mathbb{P}\{\boldsymbol{\theta}^* \in \mathcal{E}(\mathfrak{z})\} \rightarrow P\{\chi_p^2 \leq 2\mathfrak{z}\}.$$

- The likelihood ratio is asymptotically pivotal,
- $\mathcal{E}(\mathfrak{z})$ does not require to estimate the variance,
- pivotality **fails** under model misspecification $\mathbb{P} \notin \{\mathbb{P}(\boldsymbol{\theta})\}$,
- the convergence speed is **slow**.

Example: Mean Gaussian regression

$$Y_i = \theta^* + \varepsilon_i, \quad Y_i \sim \mathcal{N}(\theta^*, 1), \text{ i.i.d.}$$

- The log-likelihood function is quadratic:

$$L(\theta) = -\frac{1}{2} \sum_{i=1}^n (Y_i - \theta)^2 + C$$

- The MLE:

$$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n Y_i \sim \mathcal{N}(\theta^*, \frac{1}{\sqrt{n}})$$

- Wilks Theorem:

$$2 \left\{ L(\tilde{\theta}) - L(\theta^*) \right\} = n(\tilde{\theta} - \theta^*)^2 \sim \chi_1^2.$$

If $Y_i \approx \mathcal{N}(\theta^*, 1)$, then the distribution of $L(\tilde{\theta}) - L(\theta^*)$ differs from χ_1^2 .

- It holds with probability $\geq 1 - 2e^{-y}$ (conditions come later)

$$\left| 2 \left\{ L(\tilde{\theta}) - L(\theta^*) \right\} - \|\xi\|^2 \right| \leq \Delta_{\text{sq}}(y) \asymp \sqrt{(p+y)^3/n},$$

$$\left| \sqrt{2 \left\{ L(\tilde{\theta}) - L(\theta^*) \right\}} - \|\xi\| \right| \leq \Delta(y) \asymp (p+y)/\sqrt{n}.$$

- Deviation bound for the approximating term:

$$\mathbb{P}(\|\xi\|^2 \geq p + 6y) \leq 2e^{-y}.$$

- Allows the model to be misspecified,
- finite sample result,
- distribution of $\|\xi\|$ in general depends on \mathbb{P} and θ^* ,
- the deviation bound is not sharp enough for a confidence set.

- Estimation of the distribution of the likelihood ratio statistic with a **multiplier bootstrap procedure**.
- Theoretical justification of the multiplier bootstrap using:
 - Uniform **non-asymptotic local quadratic approximation** of the likelihood ratio statistic (Wilks-type bound) for the Y and bootstrap worlds.
 - **Non-asymptotic Gaussian approximation** to connect the Y and the bootstrap worlds.

1 Problem formulation

2 Likelihood-based confidence set

3 Multiplier bootstrap

4 Uniform confidence corridors

The idea:

To mimic the distribution of $L(\tilde{\theta}) - L(\theta^*)$ using a multiplier bootstrap:

- Take an i.i.d. sample u_1, \dots, u_n independent of the data \mathbf{Y} :
 $\mathbb{E}(u_i) = \text{Var}(u_i) = 1$, $\mathbb{E} \exp(u_i) < \infty$ (e.g. $\sim \exp(1)$, $\mathcal{N}(1, 1)$).
- Bootstrap the likelihood function:

$$L^\circ(\theta) \stackrel{\text{def}}{=} \sum_{i=1}^n \log \frac{dP_i^\circ(\theta)}{d\mu_0}(Y_i) u_i$$

- denotes the conditional probability with the fixed sample \mathbf{Y} .

“ Y world”	conditional “bootstrap world”
quasi MLE	
$\tilde{\theta} \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta} L(\theta)$	$\tilde{\theta}^{\circ} \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta} L^{\circ}(\theta)$
target	
$\theta^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta} \mathbb{E}L(\theta)$	$\tilde{\theta} \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta} \mathbb{E}^{\circ}L^{\circ}(\theta)$
likelihood ratio	
$L(\tilde{\theta}) - L(\theta^*)$	$L^{\circ}(\tilde{\theta}^{\circ}) - L^{\circ}(\tilde{\theta})$

- The true point in bootstrap world is exactly qMLE $\tilde{\theta}$.
- The bootstrap side is computable!
- The “bootstrap world” is inside of the parametric model, which, however, may be wrong.

- Wild bootstrap:
[Wu, 1986], [Beran, 1986], [Hardle and Mammen, 1993], [Mammen, 1993];
covers our approach for the case of a linear model and quadratic log-likelihood.
- Multiplier bootstrap for resampling of an objective function:
 - [Jin, Ying and Wei, 2001] use the multiplier bootstrap procedure for perturbing U-process of the k -th degree in order to estimate the covariance matrix of the U-process' minimiser.
 - [Lavergne and Patilea, 2013] use multiplier bootstrap for resampling of the smooth minimum distance statistic.
- Non-asymptotic results (very few):
 - [Arlot, Blanchard and Roquian, 2010] consider a sample of high dimensional i.i.d. Gaussian vectors, a non-asymptotical confidence set is constructed using exchangeable weighted bootstrap scheme.
 - [Chernozhukov, Chetverikov and Kato, 2013] approximate non-asymptotically the distribution of the maximum of a sum of independent high-dimensional vectors using Gaussian multiplier bootstrap.

We use the following approximating diagram:

$$\begin{array}{c} \sqrt{2\{L(\tilde{\theta}) - L(\theta^*)\}} \stackrel{(p+y)/\sqrt{n}}{\approx} \|\xi\| \\ \Downarrow \\ \sqrt{2\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}} \stackrel{(p+y)/\sqrt{n}}{\approx} \|\xi^\circ\| \end{array}$$

The approximation \approx holds between the distributions on two different probability spaces:

$$\mathcal{L}(\|\xi\|) \approx \mathcal{L}(\|\xi^\circ\| \mid \mathbf{Y}).$$

Validity of the bootstrap [Spokoiny and Zhilova, 2014a]

It holds for all $\alpha \in (0, 1)$ and the **computable** bootstrap $(1 - \alpha)$ -quantile

$$\mathfrak{z}_\alpha^\circ \stackrel{\text{def}}{=} \min \left\{ \mathfrak{z} \geq 0 : \mathbb{P}^\circ \left(\sqrt{2 \left\{ \sup_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}) - L^\circ(\tilde{\boldsymbol{\theta}}) \right\}} > \mathfrak{z} \right) \leq \alpha \right\}$$

$$\mathbb{P} \left(\sqrt{2 \{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z}_\alpha^\circ (1 + \delta) + \Delta_1 \right) - \alpha \leq \Delta_2,$$

$$\mathbb{P} \left(\sqrt{2 \{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z}_\alpha^\circ (1 - \delta) - \Delta_1 \right) - \alpha \geq -\Delta_2,$$

$\delta, \Delta_{1,2} \asymp \{\log(n)\}^{a_k} / n^{b_k}$, $a_k \geq 0, b_k > 0$

deterministic, explicit up to a generic constant.

- For a linear w.r.t. u_i function in the bootstrap world its deterministic part has an “unbiased” Y -random equivalent. For instance, for the likelihood process:

$$\mathbb{E}^\circ L^\circ(\boldsymbol{\theta}) = L(\boldsymbol{\theta}), \quad \mathbb{E}^\circ \{L^\circ(\boldsymbol{\theta}) - L^\circ(\tilde{\boldsymbol{\theta}})\} = L(\boldsymbol{\theta}) - L(\tilde{\boldsymbol{\theta}}),$$

and for the Hessian matrix:

$$\mathfrak{D}^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}^\circ L^\circ(\boldsymbol{\theta}) = -\sum_{i=1}^n \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}).$$

- The local geometry on Θ around target points is defined with matrices D_0^2 and $\mathfrak{D}_0^2 \stackrel{\text{def}}{=} \mathfrak{D}^2(\tilde{\boldsymbol{\theta}})$. The following bound holds with exponentially high probability:

$$\sup_{\boldsymbol{\theta} \in \text{Loc}(\boldsymbol{\theta}^*)} \|D_0^{-1} \mathfrak{D}^2(\boldsymbol{\theta}) D_0^{-1} - \mathbf{I}_p\|_\infty \leq \mathbf{C}(p + y) / \sqrt{n}.$$

It allows to use the same (deterministic) Fisher information matrix in the quadratic expansion of the likelihood processes.

- Local quadratic approximations hold locally around the target points:

$\forall \boldsymbol{\theta} : \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq r_0$ and some $0 < r_0 \leq C\sqrt{p+y}$:

$$L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*) \stackrel{\Delta_{\text{sq}}(y)}{\approx} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2$$

$$\stackrel{\boldsymbol{\theta} := \tilde{\boldsymbol{\theta}}}{\approx} \|\boldsymbol{\xi}\|^2/2,$$

$$L^\circ(\boldsymbol{\theta}) - L^\circ(\tilde{\boldsymbol{\theta}}) \stackrel{\Delta_{\text{sq}}(y)}{\approx} \nabla_{\boldsymbol{\theta}} L^\circ(\tilde{\boldsymbol{\theta}})^\top (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) - \|D_0(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})\|^2/2$$

$$\stackrel{\boldsymbol{\theta} := \tilde{\boldsymbol{\theta}}^\circ}{\approx} \|\boldsymbol{\xi}^\circ\|^2/2,$$

where $\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*)$, $\boldsymbol{\xi}^\circ \stackrel{\text{def}}{=} D_0^{-1} \nabla_{\boldsymbol{\theta}} L^\circ(\tilde{\boldsymbol{\theta}})$.

- The approximating term $\|\xi^\circ\|^2$ is of the 2-d order w.r.t. u_i :

$$\text{for } l_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \log \frac{dP_i(\boldsymbol{\theta})}{d\mu_0}(Y_i)$$

$$\|\xi^\circ\|^2 = \sum_{i,j=1}^n \nabla_{\boldsymbol{\theta}} l_i(\tilde{\boldsymbol{\theta}})^\top \nabla_{\boldsymbol{\theta}} l_j(\tilde{\boldsymbol{\theta}}) (u_i - 1)(u_j - 1),$$

$$\|\xi\|^2 = \sum_{i,j=1}^n \{\nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*) - \nabla_{\boldsymbol{\theta}} \mathbb{E} l_i(\boldsymbol{\theta}^*)\}^\top \{\nabla_{\boldsymbol{\theta}} l_j(\boldsymbol{\theta}^*) - \nabla_{\boldsymbol{\theta}} \mathbb{E} l_j(\boldsymbol{\theta}^*)\}.$$

Let $\check{\xi} \sim \mathcal{N}(0, \text{Var } \xi)$, $\check{\xi}^\circ \sim \mathcal{N}(0, \text{Var } \xi^\circ)$

Wilks
approximation

Gaussian
approximation

$$\sqrt{2\{L(\tilde{\theta}) - L(\theta^*)\}} \approx \|\xi\| \stackrel{\mathcal{L}}{\approx} \|\check{\xi}\|$$

$$\sqrt{2\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}} \approx \|\xi^\circ\| \stackrel{\mathcal{L}^\circ}{\approx} \|\check{\xi}^\circ\|$$

For $\check{\xi} \sim \mathcal{N}(0, \text{Var } \xi)$ $\mathcal{L}(\|\xi\|) \approx \mathcal{L}(\|\check{\xi}\|)$.

$$\blacksquare \|\xi\| = \sup_{\gamma \in \mathbb{R}^p, \|\gamma\|=1} \gamma^\top \xi \approx \sup_{\gamma \in G_\gamma} \gamma^\top \xi$$

for a finite grid G_γ on the sphere $\{\gamma \in \mathbb{R}^p : \|\gamma\| = 1\}$.

$$\blacksquare \mathbb{P}(\|\xi\| > t) = \mathbb{E} \mathbb{I}(\|\xi\| > t).$$

$$\blacksquare \mathbb{I}\left(\sup_{\gamma \in G_\gamma} \gamma^\top \xi > t\right) \approx f_\beta(\xi, t),$$

$$f_\beta(\xi, t) \stackrel{\text{def}}{=} \frac{\sum_{\gamma \in G_S(\varepsilon)} \exp\{\beta(\gamma^\top \xi - t)\}}{1 + \sum_{\gamma \in G_S(\varepsilon)} \exp\{\beta(\gamma^\top \xi - t)\}}, \beta > 0.$$

$f_\beta(\xi, t)$ is a composition of logistic function and smooth max function:

$$\mathbb{I}(\|\xi\| > t) \approx \frac{\exp\{\beta(\|\xi\| - t)\}}{1 + \exp\{\beta(\|\xi\| - t)\}} \quad \sup_{\gamma \in G_\gamma} \gamma^\top \xi \approx \frac{1}{\beta} \log \left\{ \sum_{\gamma \in G_\gamma} \exp(\beta \gamma^\top \xi) \right\}$$

For $\check{\xi} \sim \mathcal{N}(0, \text{Var } \xi)$ $\mathcal{L}(\|\xi\|) \approx \mathcal{L}(\|\check{\xi}\|)$.

- $\|\xi\| = \sup_{\gamma \in \mathbb{R}^p, \|\gamma\|=1} \gamma^\top \xi \approx \sup_{\gamma \in G_\gamma} \gamma^\top \xi$
for a finite grid G_γ on the sphere $\{\gamma \in \mathbb{R}^p : \|\gamma\| = 1\}$.

- $\mathbb{P}(\|\xi\| > t) = \mathbb{E} \mathbb{I}(\|\xi\| > t)$.

- $\mathbb{I}\left(\sup_{\gamma \in G_\gamma} \gamma^\top \xi > t\right) \approx f_\beta(\xi, t),$

$$f_\beta(\xi, t) \stackrel{\text{def}}{=} \frac{\sum_{\gamma \in G_S(\varepsilon)} \exp\{\beta(\gamma^\top \xi - t)\}}{1 + \sum_{\gamma \in G_S(\varepsilon)} \exp\{\beta(\gamma^\top \xi - t)\}}, \beta > 0.$$

For $\gamma \in \mathbb{R}^p$

$$\left| f_\beta(\gamma + \xi, t) - f_\beta(\gamma, t) - \xi^\top \nabla_\gamma f_\beta(\gamma, t) - \frac{1}{2} \xi^\top \nabla_\gamma^2 f_\beta(\gamma, t) \xi \right| \leq 0.7 \beta^3 \|\xi\|^3$$

- By Lindeberg's telescopic sum device [Lindeberg, 1922]:

$$\left| \mathbb{E} f_{\beta}(\xi, z) - \mathbb{E} f_{\beta}(\check{\xi}, z) \right| \leq \frac{\beta^3}{8} \sum_{i=1}^n \mathbb{E} \left(\|\xi_i\|^3 + \|\check{\xi}_i\|^3 \right),$$

where

$$\xi_i \stackrel{\text{def}}{=} D_0^{-1} \{ \nabla_{\theta} \ell_i(\theta^*) - \nabla_{\theta} \mathbb{E} \ell_i(\theta^*) \},$$

$$\check{\xi}_i \sim \mathcal{N}(0, \text{Var } \xi_i).$$

$$\begin{aligned} \mathbb{P}(\|\xi\| \geq z + \sqrt{p}\Delta) &\leq \mathbb{P}\left(\sup_{\gamma \in G_S(\varepsilon)} \gamma^\top (\xi - \bar{\Delta}) \geq (1 - \varepsilon)z\right) \\ &\leq \mathbb{P}\left(\|\check{\xi}\| \geq z \left(1 - \frac{2\varepsilon}{1 + \varepsilon}\right) - \sqrt{p}\Delta\right) \\ &\quad + [1 + \text{card}\{G_S(\varepsilon)\}] e^{-\beta\Delta} + \mathbf{C}_3 \frac{(\beta \mathbf{C}_2 \|\Sigma^{1/2}\|)^3}{4\sqrt{n}}. \end{aligned}$$

Similarly for the lower bound:

$$\begin{aligned} \mathbb{P}(\|\xi\| \geq z - \sqrt{p}\Delta) &\geq \mathbb{P}\left(\|\check{\xi}\| \geq z \left(1 + \frac{2\varepsilon}{1 - \varepsilon}\right) + \sqrt{p}\Delta\right) \\ &\quad - [1 + \text{card}\{G_S(\varepsilon)\}] e^{-\beta\Delta} - \mathbf{C}_3 \frac{(\beta \mathbf{C}_2 \|\Sigma^{1/2}\|)^3}{4\sqrt{n}}. \end{aligned}$$

- Finally we compare the norms of two Gaussian vectors: $\|\check{\xi}\|$ and $\|\check{\xi}^\circ\|$:

Gaussian approximation

$$\begin{array}{ccc} \|\check{\xi}\| & \stackrel{\mathcal{L}}{\approx} & \|\check{\xi}\| \\ & & \mathcal{L} \\ \|\check{\xi}^\circ\| & \stackrel{\mathcal{L}^\circ}{\approx} & \|\check{\xi}^\circ\| \end{array}$$

i.e. their covariance matrices:

$$\begin{aligned} \text{Var}(\check{\xi}) &= D_0^{-1} \text{Var}(\nabla_{\theta} L(\theta^*)) D_0^{-1}, \\ \text{Var}^\circ(\check{\xi}^\circ) &= D_0^{-1} \text{Var}^\circ(\nabla_{\theta} L^\circ(\tilde{\theta})) D_0^{-1}. \end{aligned}$$

Let $\Sigma \stackrel{\text{def}}{=} \text{Var}(\boldsymbol{\xi})$, $\Sigma^\circ \stackrel{\text{def}}{=} \text{Var}^\circ(\boldsymbol{\xi}^\circ)$.

$$D_0 \Sigma D_0 = \sum_{i=1}^n \mathbb{E} \{ \nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*)^\top \} - \sum_{i=1}^n \mathbb{E} \{ \nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*) \} \mathbb{E} \{ \nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*) \}^\top,$$

$$D_0 \Sigma^\circ D_0 = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} l_i(\boldsymbol{\theta}^*)^\top.$$

It holds:

$$\| \Sigma^{-1/2} \Sigma^\circ \Sigma^{-1/2} - \mathbf{I}_p \|_\infty \leq C_{\mathcal{I}}^2 \delta_\xi^2 \asymp (p + y)/n^{1/2},$$

the terms $C_{\mathcal{I}}^2, \delta_\xi^2$ are responsible for the modelling bias, they are proportional to the error terms δ and Δ_1 in the statement on the bootstrap validity above.

Δ_1, Δ_2 are also determined by $C_{\mathcal{I}}^2$ and the errors induced by

$$\mathbb{1}(\sup_{\mathbf{G}_\gamma} \boldsymbol{\gamma}^\top \boldsymbol{\xi} > t) \approx f_\beta(\boldsymbol{\xi}, t).$$

Conditions include

- Independency of the observations Y_i .
- Bounded exponential moments of gradient of $L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta})$ up to 3-d order.
- Small modelling bias:

$$\sum_{i=1}^n \|D_0^{-1} \mathbb{E} \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)\|^2 \leq \delta_{\xi}^2 = \mathfrak{C}(p + \mathfrak{y})/n^{1/2}$$

for $\ell_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \log \frac{dP_i(\boldsymbol{\theta})}{d\mu_0}(Y_i)$, $D_0^2 \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}L(\boldsymbol{\theta}^*)$.

- $\|D_0 \text{Var}(\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*))^{-1} D_0\|_{\infty} \leq \mathfrak{C}_{\mathcal{I}}^2$.
($\delta_{\xi}^2 = 0$, $\mathfrak{C}_{\mathcal{I}}^2 = 1$ if the parametric model is true)

I.i.d. errors, misspecified distribution

$$Y_i = 2 + \varepsilon_i, \quad \varepsilon_i \sim \text{Lap}(0, 2^{-1/2}), \text{ i.i.d.}, \quad \text{Var}(\varepsilon_i) = 1,$$

$$L(\theta) = - \sum_{i=1}^n (Y_i - \theta)^2 / 2 \quad \text{corresponds to } \varepsilon_i \sim \mathcal{N}(0, 1).$$

Quantiles of $\sqrt{2}\{L(\tilde{\theta}) - L(\theta^*)\}^{1/2}$ are estimated from 5000 \mathbf{Y} samples of length 50.

Quantiles of $\sqrt{2}\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}$ are estimated from 5000 $\{u_i\} \sim \text{exp}(1)$ samples and 50 \mathbf{Y} samples of length 50.

Confidence level	0.99	0.95	0.90	0.85	0.80	0.75
$\sqrt{2}\{L(\tilde{\theta}) - L(\theta^*)\}^{1/2}$	2.55	1.99	1.64	1.44	1.29	1.16
$\sqrt{2}\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}, \text{Mean}$	2.63	1.92	1.58	1.38	1.22	1.09
$\sqrt{2}\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}, \text{St. dev.}$	0.39	0.27	0.22	0.19	0.17	0.15

Comparison of the empirical quantiles for the sample size $n = 50$

Heteroscedastic errors, misspecified distribution

$$Y_i = 2 + \sigma_i \varepsilon_i, \quad \varepsilon_i \sim \text{Lap}(0, 2^{-1/2}), \quad \sigma_i = 2 - (i \bmod 4)/2,$$

$$L(\theta) = - \sum_{i=1}^n (Y_i - \theta)^2 / 2 \quad \text{corresponds to } \sigma_i \varepsilon_i \sim \mathcal{N}(0, 1).$$

Quantiles of $\sqrt{2}\{L(\tilde{\theta}) - L(\theta^*)\}^{1/2}$ are estimated from 5000 \mathbf{Y} samples of length 50.

Quantiles of $\sqrt{2}\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}$ are estimated from 5000 $\{u_i\} \sim \text{exp}(1)$ samples and 50 \mathbf{Y} samples of length 50.

Confidence level	0.99	0.95	0.90	0.85	0.80	0.75
$\sqrt{2}\{L(\tilde{\theta}) - L(\theta^*)\}^{1/2}$	5.06	3.86	3.87	2.85	2.52	2.27
$\sqrt{2}\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}$, Mean	5.26	3.87	3.19	2.78	2.46	2.20
$\sqrt{2}\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}$, St. dev.	0.88	0.65	0.54	0.46	0.41	0.37

Comparison of the empirical quantiles for the sample size $n = 50$

Biased model

$$Y_i = \beta \sin(X_i) + \varepsilon_i, \quad \varepsilon_i \sim Lap(0, 2^{-1/2}), \text{ i.i.d.},$$

X_i are equidistant in $[0, 2\pi]$,

$$L(\theta) = - \sum_{i=1}^n (Y_i - \theta)^2 / 2 \Rightarrow \theta^* = 0,$$

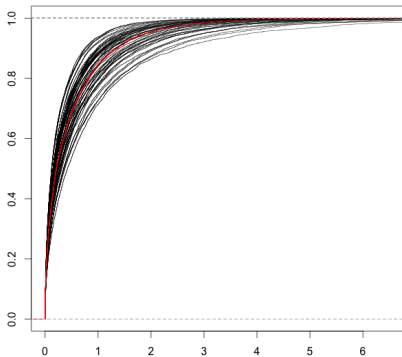
$\beta > 0$ is an amplitude.

Let us check how the confidence set for θ^* changes with β .

Numerical results (III). Biased model

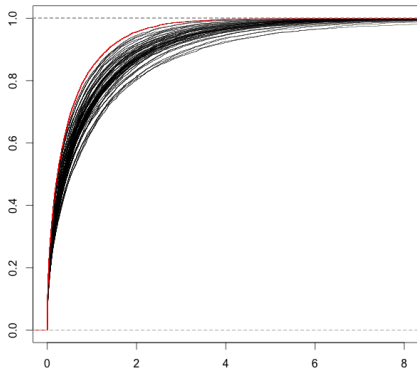
$n = 50$

- empirical distribution function of $L(\tilde{\theta}) - L(\theta^*)$ estimated with 5000 \mathbf{Y} samples
- 50 empirical distribution functions of $L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})$ estimated with 5000 $u_1, \dots, u_n \sim \exp(1)$ samples



$$Y_i = 0.25 \sin(X_i) + \text{Lap}(0, 2^{-1/2})$$

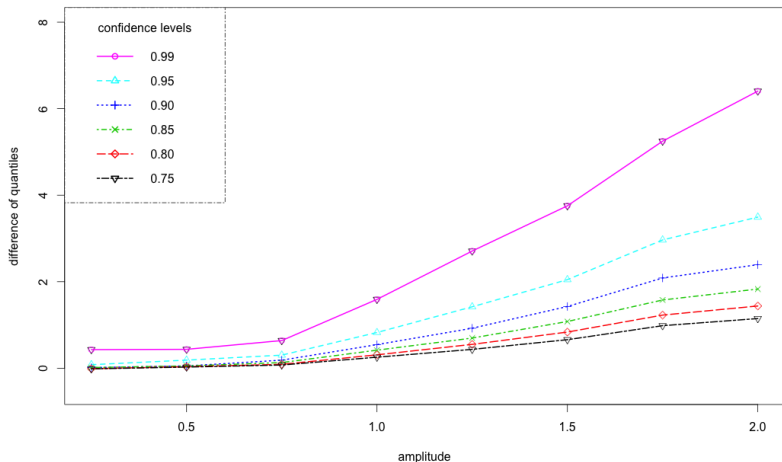
the “small modelling bias” case



$$Y_i = 1.25 \sin(X_i) + \text{Lap}(0, 2^{-1/2}),$$

“small modelling bias” condition is **not fulfilled**

The difference (“Bootstrap quantile” – “Y -quantile”) growing with modelling bias:



Logistic regression with bias

$$Y_i \sim \text{Bernoulli}(\beta X_i),$$

X_i are equidistant in $[0, 2]$,

$$\beta \in (0, 1/2],$$

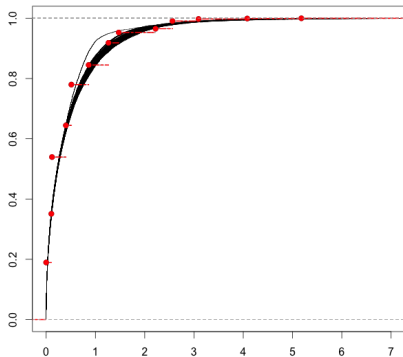
$$L(\theta) = \sum_{i=1}^n \{Y_i \theta - \log(1 + e^\theta)\},$$

$$\Rightarrow \theta^* = \log\{\beta/(1 - \beta)\}.$$

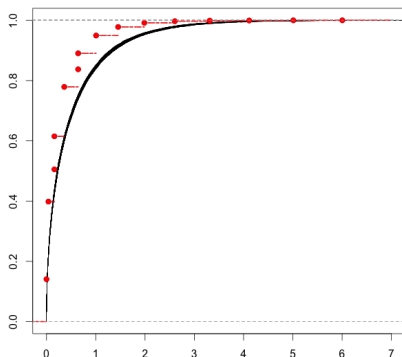
Numerical results (IV). Logistic regression

$n = 50$

- empirical distribution function of $L(\tilde{\theta}) - L(\theta^*)$ estimated with 10000 \mathbf{Y} samples
- 50 empirical distribution functions of $L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})$ estimated with 10000 $u_1, \dots, u_n \sim \exp(1)$ samples



$Y_i \sim \text{Bernoulli}(0.1X_i)$



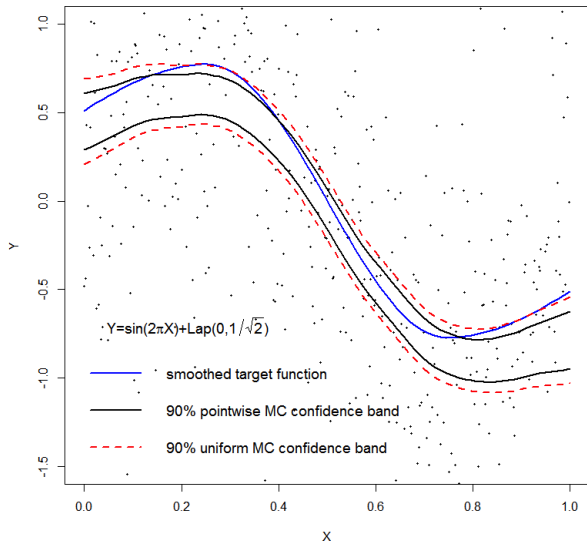
$Y_i \sim \text{Bernoulli}(0.5X_i)$

- 1 Problem formulation
- 2 Likelihood-based confidence set
- 3 Multiplier bootstrap
- 4 Uniform confidence corridors**

For a **fixed sample size** n construct a **simultaneous** confidence set for $f(x)$, when $x \in \mathcal{X} \subset \mathbb{R}^d$, with a given confidence level $1 - \alpha$ and a minimal possible diameter.

$$\mathbb{P} \left\{ \bigcap_{x \in \mathcal{X}} \{f(x) \in \mathcal{E}_x\} \right\} \geq 1 - \alpha.$$

Uniform vs. pointwise confidence corridors



Simulated data:

$$Y_i = \sin(2\pi X_i) + \varepsilon_i,$$

$$\varepsilon_i \sim \text{Lap}(0, 1/\sqrt{2}) \text{ i.i.d.},$$

X_i are equidistant on $[0, 1]$,

sample size = 500,

confidence sets are based on the Nadaraya-Watson estimator with parabolic weights and bandwidth = 0.25.

Local likelihood function for $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p$, $\boldsymbol{x} \in \mathcal{X} \subset \mathbb{R}^d$:

$$L(\boldsymbol{\theta}, \boldsymbol{x}) \stackrel{\text{def}}{=} \sum_{i=1}^n \ell(Y_i, f(X_i, \boldsymbol{\theta})) w_i(\boldsymbol{x}),$$

$$\ell(Y_i, f(X_i, \boldsymbol{\theta})) \stackrel{\text{def}}{=} \frac{dP_i(f(X_i, \boldsymbol{\theta}))}{d\mu_0}(Y_i),$$

$f(\boldsymbol{x}, \boldsymbol{\theta})$ is a known smooth function $\mathbb{R}^d \times \mathbb{R}^p \mapsto \mathbb{R}$,

$w_i(\boldsymbol{x}) \stackrel{\text{def}}{=} K\left(\frac{\boldsymbol{x} - X_i}{h}\right)$ are localising weights,

$K(\boldsymbol{x}) \geq 0$ is a symmetric kernel function,

$h > 0$ is a **fixed bandwidth** (however, the results here are easily extended to the uniform in bandwidth bounds, when h lies on a fine finite grid).

- Local maximum likelihood estimate:

$$\tilde{\boldsymbol{\theta}}(\mathbf{x}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}, \mathbf{x}),$$

$$\tilde{f}(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{x}, \tilde{\boldsymbol{\theta}}(\mathbf{x})).$$

- Target parameters (smoothed versions):

$$\boldsymbol{\theta}^*(\mathbf{x}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E}L(\boldsymbol{\theta}, \mathbf{x}),$$

$$f^*(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{x}, \boldsymbol{\theta}^*(\mathbf{x})).$$

- Pointwise confidence sets

$$\mathcal{E}(\mathfrak{z}, \mathbf{x}) \stackrel{\text{def}}{=} \left\{ \boldsymbol{\theta} : L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}) - L(\boldsymbol{\theta}, \mathbf{x}) \leq \mathfrak{z} \right\} \subset \mathbb{R}^p,$$

$$\mathcal{E}_f(\mathfrak{z}, \mathbf{x}) \stackrel{\text{def}}{=} \{f(\mathbf{x}, \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathcal{E}(\mathfrak{z}, \mathbf{x})\} \subset \mathbb{R}.$$

- Uniform (simultaneous) confidence corridor:

$$\mathcal{E}_f(\mathfrak{z}) \stackrel{\text{def}}{=} \bigcup_{\mathbf{x} \in \mathcal{X}} \{\mathbf{x} \times \mathcal{E}_f(\mathfrak{z}, \mathbf{x})\} \subset \mathbb{R}^d \times \mathbb{R}.$$

- **Aim:** to build for each $\mathbf{x} \in \mathcal{X}$ a **quantile function** $\mathfrak{z}(\alpha, \mathbf{x})$ s.t.

$$\mathbb{P} \left(\bigcap_{\mathbf{x} \in \mathcal{X}} \left\{ f^*(\mathbf{x}) \in \mathcal{E}_f(\mathfrak{z}(\alpha, \mathbf{x}), \mathbf{x}) \right\} \right) \geq 1 - \alpha.$$

For a smooth $f(\cdot, \boldsymbol{\theta})$

$$\begin{aligned} f^*(\mathbf{x}) \in \mathcal{E}_f(\mathfrak{z}, \mathbf{x}) &\Leftrightarrow \boldsymbol{\theta}^*(\mathbf{x}) \in \mathcal{E}(\mathfrak{z}, \mathbf{x}) \\ &\Leftrightarrow L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}) - L(\boldsymbol{\theta}^*(\mathbf{x}), \mathbf{x}) \leq \mathfrak{z}. \end{aligned}$$

Therefore, the initial problem translates into finding $\mathfrak{z}(\alpha, \mathbf{x})$ s.t.

$$\mathbb{P} \left(\sup_{\mathbf{x} \in \mathcal{X}} \left\{ L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}) - L(\boldsymbol{\theta}^*(\mathbf{x}), \mathbf{x}) - \mathfrak{z}(\alpha, \mathbf{x}) \right\} \leq 0 \right) \geq 1 - \alpha.$$

Similarly to the parametric approach define the bootstrap objects for the local likelihood conditioned on \mathbf{Y} :

$$L^\circ(\boldsymbol{\theta}, \mathbf{x}) \stackrel{\text{def}}{=} \sum_{i=1}^n \ell(Y_i, f(X_i, \boldsymbol{\theta})) w_i(\mathbf{x}) u_i.$$

The target parameter

$$\tilde{\boldsymbol{\theta}}(\mathbf{x}) = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E}^\circ L^\circ(\boldsymbol{\theta}, \mathbf{x}).$$

Maximum likelihood ratio for the bootstrap case:

$$L^\circ(\tilde{\boldsymbol{\theta}}^\circ, \mathbf{x}) - L^\circ(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}) \stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta} \in \Theta} L^\circ(\boldsymbol{\theta}, \mathbf{x}) - L^\circ(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}).$$

$$\sqrt{2\{L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}) - L(\boldsymbol{\theta}^*(\mathbf{x}), \mathbf{x})\}} \stackrel{\bigcap_{\mathbf{x} \in \mathcal{X}}}{\approx} \|\boldsymbol{\xi}(\mathbf{x})\|$$

$\gg w_X$

$$\sqrt{2\{L^\circ(\tilde{\boldsymbol{\theta}}^\circ(\mathbf{x}), \mathbf{x}) - L^\circ(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x})\}} \stackrel{\bigcap_{\mathbf{x} \in \mathcal{X}}}{\approx} \|\boldsymbol{\xi}^\circ(\mathbf{x})\|$$

- $\bigcap_{\mathbf{x} \in \mathcal{X}} \approx$ denotes the uniform Wilks-type approximation,
- \approx^{w_X} means approximation in distribution jointly in $\mathbf{x} \in \mathcal{X}$.
- The normalized scores and full Fisher Information matrices are:

$$\boldsymbol{\xi}(\mathbf{x}) \stackrel{\text{def}}{=} D_0^{-1}(\mathbf{x}_0) \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*(\mathbf{x}), \mathbf{x}) \qquad \boldsymbol{\xi}^\circ(\mathbf{x}) \stackrel{\text{def}}{=} \mathcal{D}_0^{-1}(\mathbf{x}) \nabla_{\boldsymbol{\theta}} L^\circ(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x})$$

$$D_0^2(\mathbf{x}) \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E} L(\boldsymbol{\theta}^*(\mathbf{x}), \mathbf{x}) \qquad \mathcal{D}_0^2(\mathbf{x}) \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}^\circ L^\circ(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x})$$

Let $G_{\mathcal{X}}(\varepsilon)$ denote a finite ε -net on \mathcal{X} : $\forall \mathbf{x} \in \mathcal{X} \exists \mathbf{x}_0 \in G_{\mathcal{X}}(\varepsilon) : \|\mathbf{x} - \mathbf{x}_0\| \leq \varepsilon$.

1. Pointwise step

For each $\mathbf{x} \in G_{\mathcal{X}}(\varepsilon)$ compute an empirical distribution of the ratio

$\sqrt{2\{L^\circ(\tilde{\theta}^\circ(\mathbf{x}), \mathbf{x}) - L^\circ(\tilde{\theta}(\mathbf{x}), \mathbf{x})\}}$, which yields the pointwise quantile function $\mathfrak{z}^\circ(\alpha, \mathbf{x})$ for each $\alpha \in (0, 1)$ and $\mathbf{x} \in G_{\mathcal{X}}(\varepsilon)$.

2. Correction for the uniformity

By an iterative procedure find a maximum value $\alpha(\varepsilon)$ s.t.

$$\mathbb{P} \left(\sup_{\mathbf{x} \in G_{\mathcal{X}}(\varepsilon)} \left[\sqrt{2\{L^\circ(\tilde{\theta}^\circ(\mathbf{x}), \mathbf{x}) - L^\circ(\tilde{\theta}(\mathbf{x}), \mathbf{x})\}} - \mathfrak{z}^\circ(\alpha(\varepsilon), \mathbf{x}) \right] \leq 0 \right) \geq 1 - \alpha.$$

The value $\alpha(\varepsilon)$ is a corrected confidence level for $\sup_{\mathbf{x} \in G_{\mathcal{X}}(\varepsilon)}$.

It holds $\alpha(\varepsilon)/\alpha \geq \mathfrak{C} |G_{\mathcal{X}}(\varepsilon)|^{-1}$.

Local constant regression:

$$Y_i = \sin(2\pi X_i) + \varepsilon_i, \quad \varepsilon_i \sim \text{Lap}(0, 2^{-1/2}), \text{ i.i.d.},$$

X_i are equidistant in $[0, 1]$,

$$L(\theta) = -\frac{1}{2} \sum_{i=1}^n (Y_i - \theta)^2 w_i(x),$$

$$w_i(x) = \begin{cases} 1 - \{(x - X_i)/h\}^2, & \text{if } |(x - X_i)/h| \leq 1, \\ 0, & \text{otherwise.} \end{cases}$$

Nadaraya-Watson estimate:

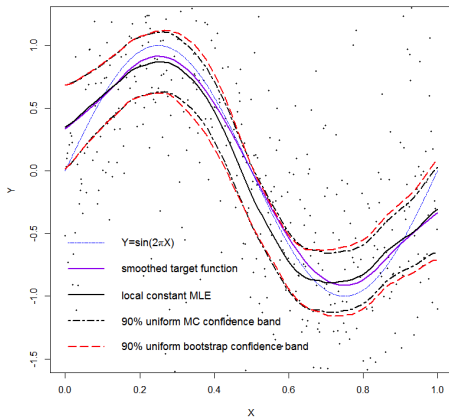
$$\tilde{\theta}(x) = \sum_{i=1}^n Y_i \frac{w_i(x)}{\sum_{i=1}^n w_i(x)},$$

Smoothed target function:

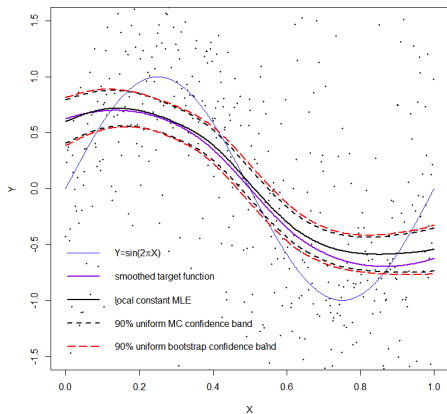
$$\theta^*(x) = \sum_{i=1}^n \mathbb{E}Y_i \frac{w_i(x)}{\sum_{i=1}^n w_i(x)}.$$

Numerical results. Uniform confidence corridors

$$Y_i = \sin(2\pi X_i) + \varepsilon_i, \quad \varepsilon_i \sim \text{Lap}(0, 2^{-1/2}), \text{ i.i.d.}, \quad n = 500.$$



$h = 0.15$



$h = 0.35$

Let \mathbf{x}_G denote the projection of a point $\mathbf{x} \in \mathcal{X}$ on the grid $G_{\mathcal{X}}(\varepsilon)$:





$$\mathbf{x}_G \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x}_1 \in G_{\mathcal{X}}(\varepsilon)} \|\mathbf{x} - \mathbf{x}_1\|.$$





It holds with probability $\geq 1 - \mathbf{C} |G_{\mathcal{X}}(\varepsilon)| e^{-y} = 1 - \mathbf{C} \exp\{-y + \log |G_{\mathcal{X}}(\varepsilon)|\}$

$$\sup_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}_G\| \leq \varepsilon} \left| \sqrt{2\{L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{x}) - L(\boldsymbol{\theta}^*(\mathbf{x}), \mathbf{x})\}} - \|\boldsymbol{\xi}(\mathbf{x}_G)\| \right| \leq \Delta(\varepsilon, y),$$

where $\Delta(\varepsilon, y) = \mathbf{C}(p + y)/\sqrt{(nh^d)}$.

- Adaptive choice of optimal bandwidth using a multiscale correction.
- Non-asymptotic bounds for Monte Carlo simulations (joint with Christian Bayer).
- Estimation of dynamic quantile/expectile curves (joint with project B1).
- Application to the multiple hypothesis testing: choice of a correction for multiplicity.

-  Arlot, S., Blanchard, G., and Roquain, E. (2010).
Some nonasymptotic results on resampling in high dimension. I. Confidence regions.
Ann. Statist., 38(1):51–82.
-  Beran, R. (1986).
Discussion: Jackknife, bootstrap and other resampling methods in regression analysis.
The Annals of Statistics, pages 1295–1298.
-  Chernozhukov, V., Chetverikov, D., and Kato, K. (2013).
Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors.
The Annals of Statistics, 41(6):2786–2819.
-  Hardle, W. and Mammen, E. (1993).
Comparing nonparametric versus parametric regression fits.
The Annals of Statistics, pages 1926–1947.

-  Jin, Z., Ying, Z., and Wei, L. J. (2001).
A simple resampling method by perturbing the minimand.
Biometrika, 88(2):381–390.
-  Lavergne, P. and Patilea, V. (2013).
Smooth minimum distance estimation and testing with conditional estimating equations: uniform in bandwidth theory.
Journal of Econometrics, 177(1):47–59.
-  Lindeberg, J. W. (1922).
Eine neue herleitung des exponentialgesetzes in der
wahrscheinlichkeitsrechnung.
Mathematische Zeitschrift, 15(1):211–225.
-  Mammen, E. (1993).
Bootstrap and wild bootstrap for high dimensional linear models.
The Annals of Statistics, pages 255–285.



Spokoiny, V. (2012).

Parametric estimation. finite sample theory.

The Annals of Statistics, 40(6):2877–2909.



Spokoiny, V. (2014).

Wilks theorem for penalized maximum likelihood estimators.

Manuscript. arXiv:1205.0498.



Spokoiny, V. and Zhilova, M. (2014a).

Parametric confidence sets via multiplier bootstrap.


to appear.



Spokoiny, V. and Zhilova, M. (2014b).

Uniform confidence sets via multiplier bootstrap.

to appear.

-  Wu, C. F. J. (1986).
Jackknife, Bootstrap and Other Resampling Methods in Regression Analysis.
The Annals of Statistics, 14(4):1261–1295+.

Thank you for your attention!