



**Weierstrass Institute for
Applied Analysis and Stochastics**



Bootstrap tuning in model choice problem

Vladimir Spokoiny,
(with Niklas Willrich)

WIAS, HU Berlin, MIPT, IITP Moscow

SFB 649 Motzen, 17.07.2015

1 Introduction

2 SmA procedure for known noise variance

3 Bootstrap tuning

4 Numerical results

Consider a linear model

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

for an unknown parameter vector $\boldsymbol{\theta}^* \in \mathbb{R}^p$ and a given $p \times n$ design matrix Ψ .

Suppose that a family of linear smoothers

$$\tilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y}$$

is given, where \mathcal{S}_m is for each $m \in \mathcal{M}$ a given $p \times n$ matrix.

We also assume that this family is ordered by the complexity of the method.

The task is to develop a data based model selector \hat{m} which performs nearly as good as the optimal choice which depends on the model and is not available.

We consider the following **linear Gaussian model**:

$$Y_i = \Psi_i^\top \boldsymbol{\theta}^* + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \text{ i.i.d.}, \quad i = 1, \dots, n. \quad (1)$$

We also write this equation in the vector form

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon} \in \mathbb{R}^n$$

where Ψ is $p \times n$ **deterministic** design matrix and $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$.

In what follows, we allow the model (1) to be completely **misspecified**:

- **True model:** Y_i are independent, the response $\mathbf{f}^* = \mathbb{E}\mathbf{Y} \in \mathbb{R}^n$ with entries f_i :

$$Y_i = f_i + \varepsilon_i. \quad (2)$$

- The linear parametric assumption $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$ can be violated;
- The noise $\boldsymbol{\varepsilon} = (\varepsilon_i)$ can be heterogeneous and non-Gaussian.

For the linear model (2), define $\boldsymbol{\theta}^* \in \mathbb{R}^p$ as the vector providing the best linear fit:

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathbb{E} \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 = (\Psi \Psi^\top)^{-1} \Psi \mathbf{f}^*.$$

Below we assume a family $\{\tilde{\theta}_m\}$ of linear estimators of θ^* to be given:

$$\tilde{\theta}_m = S_m \mathbf{Y}$$

Typical examples include

- **projection estimation** on a m -dimensional subspace;
- **regularized estimation** with a regularization parameter α_m ;
- **penalized estimators** with a quadratic penalty function;
- **kernel estimation** with a bandwidth h_m .

Introduce a weighting $q \times p$ -matrix W for some fixed $q \geq 1$ and define quadratic loss and risk with this weighting matrix W :

$$\varrho_m \stackrel{\text{def}}{=} \|W(\tilde{\theta}_m - \theta^*)\|^2, \quad \mathcal{R}_m \stackrel{\text{def}}{=} \mathbb{E}\|W(\tilde{\theta}_m - \theta^*)\|^2.$$

Typical examples of W are as follows:

■ **Estimation of the whole vector θ^* :**

Let W be the identity matrix $W = I_p$ with $q = p$. This means that the estimation loss is measured by the usual squared Euclidean norm $\|\tilde{\theta}_m - \theta^*\|^2$.

■ **Prediction:** Let W be the square root of the total Fisher information matrix, that is, $W^2 = \mathbb{F} = \sigma^{-2}\Psi\Psi^\top$. Usually referred to as *prediction loss*.

■ **Semiparametric estimation:** Let the target of estimation is not the whole vector θ^* but its subvector θ_0^* of dimension q . The matrix W can be defined as the projector Π_0 on the θ_0^* subspace.

■ **Linear functional estimation:** The choice of the weighting matrix W can be adjusted to the problem of estimating some functionals of the whole parameter θ^* .

In all cases, the most important feature of the estimators $\tilde{\theta}_m$ is *linearity*.

In all cases, the most important feature of the estimators $\tilde{\boldsymbol{\theta}}_m$ is *linearity*. It greatly simplifies the study of their properties including the prominent bias-variance decomposition of the risk of $\tilde{\boldsymbol{\theta}}_m$. Namely, for the model (2) with $\mathbb{E}\boldsymbol{\varepsilon} = 0$, it holds

$$\begin{aligned}\mathbb{E}\tilde{\boldsymbol{\theta}}_m &= \boldsymbol{\theta}_m^* = \mathcal{S}_m \mathbf{f}^*, \\ \mathcal{R}_m &= \|W(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\|^2 + \text{tr}\{W\mathcal{S}_m \text{Var}(\boldsymbol{\varepsilon}) \mathcal{S}_m^\top W^\top\} \\ &= \|W(\mathcal{S}_m - \mathcal{S})\mathbf{f}^*\|^2 + \text{tr}\{W\mathcal{S}_m \text{Var}(\boldsymbol{\varepsilon}) \mathcal{S}_m^\top W^\top\}.\end{aligned}\quad (3)$$

The optimal choice of the parameter m can be defined by risk minimization:

$$m^* \stackrel{\text{def}}{=} \underset{m \in \mathcal{M}}{\text{argmin}} \mathcal{R}_m.$$

The *model selection* problem can be described as the choice of m by data which *mimics the oracle*, that is, we aim at constructing a selector \hat{m} leading to the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$ with the properties similar to the oracle estimate $\tilde{\boldsymbol{\theta}}_{m^*}$.

- **unbiased risk estimation** [Kneip, 1994];
- **penalized model selection** [Barron et al., 1999], [Massart, 2007]);
- **Lepski's method** [Lepski, 1990], [Lepski, 1991], [Lepski, 1992], [Lepski and Spokoiny, 1997], [Lepski et al., 1997], [Birgé, 2001]
- **risk hull minimization** [Cavalier and Golubev, 2006].
- **Resampling for noise estimation:** For the penalized model selection, [Arlot, 2009] suggested the use of resampling methods for the choice of an optimal penalization, [Arlot and Bach, 2009] used the concept of minimal penalties from [Birgé and Massart, 2007].
- **Aggregation:**
An alternative approach to adaptive estimation is based on aggregation of different estimates; see [Goldenshluger, 2009] and [Dalalyan and Salmon, 2012];
- **Linear inverse problems:** [Tsybakov, 2000], [Cavalier et al., 2002].

Validity of a bootstrapping procedure for Lepski's method has been studied in [Chernozhukov et al., 2014] with applications to honest adaptive confidence bands.

[Spokoiny and Vial, 2009] offered a propagation approach to calibration of the Lepski's method in the case of the estimation of a one-dimensional quantity of interest.

A similar approach has been applied to local constant density estimation with sup-norm risk in [Gach et al., 2013] and to local quantile estimation in [Spokoiny et al., 2013].

Below we discuss the *ordered case*. The parameter $m \in \mathcal{M}$ is treated as complexity of the method $\tilde{\theta}_m = \mathcal{S}_m Y$. In some cases the set \mathcal{M} of possible m choices can be countable and/or continuous and even unbounded. For simplicity of presentation, we assume that \mathcal{M} is a finite set of positive numbers, $|\mathcal{M}|$ stands for its cardinality.

Typical examples:

- **the number of terms** in the Fourier expansion;
- **bandwidth** in the kernel smoothing.
- **regularization parameter**;
- **penalty coefficient**.

In general, complexity can be naturally expressed via the variance of the stochastic term of the estimator $\tilde{\theta}_m$: the larger m , the larger is the variance $\text{Var}(W\tilde{\theta}_m)$.

In the case of projection estimation with m -dimensional projectors, this variance is linear in m , $\text{Var}(\tilde{\theta}_m) = \sigma^2 m$.

In general, dependence of the variance term on m may be more complicated but the monotonicity of $\text{Var}(W\tilde{\theta}_m)$ has to be preserved.

1 Introduction

2 SmA procedure for known noise variance

3 Bootstrap tuning

4 Numerical results

Consider a family $\tilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y}$ and $\tilde{\boldsymbol{\phi}}_m = \mathcal{K}_m \mathbf{Y}$ with $\mathcal{K}_m = W \mathcal{S}_m : \mathbb{R}^n \rightarrow \mathbb{R}^q$, $m \in \mathcal{M}$, of linear estimators of the q -dimensional target parameter $\boldsymbol{\phi}^* = W \boldsymbol{\theta}^* = W \mathcal{S} \mathbf{f}^* = \mathcal{K} \mathbf{f}^*$ for $\mathcal{K} = W \mathcal{S}$.

Suppose that $\{\tilde{\boldsymbol{\phi}}_m, m \in \mathcal{M}\}$ is ordered by their complexity (variance):

$$\mathcal{K}_m \text{Var}(\boldsymbol{\varepsilon}) \mathcal{K}_m^\top \leq \mathcal{K}_{m'} \text{Var}(\boldsymbol{\varepsilon}) \mathcal{K}_{m'}^\top, \quad m' > m.$$

One would like to pick up a smallest possible index $m \in \mathcal{M}$ which still provides a reasonable fit. The latter means that the bias component

$$\|b_m\|^2 = \|\boldsymbol{\phi}_m^* - \boldsymbol{\phi}^*\|^2 = \|(\mathcal{K}_m - \mathcal{K}) \mathbf{f}^*\|^2$$

in the risk decomposition (3) is not significantly larger than the variance

$$\text{tr}\{\text{Var}(\tilde{\boldsymbol{\phi}}_m)\} = \text{tr}\{\mathcal{K}_m \text{Var}(\boldsymbol{\varepsilon}) \mathcal{K}_m^\top\}.$$

If $m^\circ \in \mathcal{M}$ is such a “good” choice, then our ordering assumption yields that a further increase of the index m over m° only increases the complexity (variance) of the method without real gain in the quality of approximation.

This latter fact can be interpreted in term of pairwise comparison: whatever $m \in \mathcal{M}$ with $m > m^\circ$ we take, there is no significant bias reduction in using a larger model m instead of m° .

Leads to a multiple test procedure: for each pair $m > m^\circ$ from \mathcal{M} , we consider a hypothesis of no significant bias between the models m° and m , and let τ_{m,m° be the corresponding test.

The model m° is accepted if $\tau_{m,m^\circ} = 0$ for all $m > m^\circ$. Finally, the selected model is the “smallest accepted”:

$$\hat{m} \stackrel{\text{def}}{=} \operatorname{argmin} \{ m^\circ \in \mathcal{M} : \tau_{m,m^\circ} = 0, \forall m > m^\circ \}.$$

Usually the test τ_{m,m° can be written in the form

$$\tau_{m,m^\circ} = \mathbb{I} \{ \mathbb{T}_{m,m^\circ} > \mathfrak{z}_{m,m^\circ} \}$$

for some *test statistics* \mathbb{T}_{m,m° and for *critical values* \mathfrak{z}_{m,m° .

The information-based criteria like AIC or BIC use the likelihood ratio test statistics

$\mathbb{T}_{m,m^\circ} = \sigma^{-2} \|\Psi^\top (\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})\|^2$. A great advantage of such tests is that the test statistic \mathbb{T}_{m,m° is pivotal (χ^2 with $m - m^\circ$ degrees of freedom) under the correct null hypothesis, this makes simple to compute the corresponding critical values.

Below we apply another choice based on the norm of differences $\tilde{\boldsymbol{\phi}}_m - \tilde{\boldsymbol{\phi}}_{m^\circ}$:

$$\mathbb{T}_{m,m^\circ} = \|\tilde{\boldsymbol{\phi}}_m - \tilde{\boldsymbol{\phi}}_{m^\circ}\| = \|\mathcal{K}_{m,m^\circ} \mathbf{Y}\|, \quad \mathcal{K}_{m,m^\circ} \stackrel{\text{def}}{=} \mathcal{K}_m - \mathcal{K}_{m^\circ}.$$

The main issue for such a method is a proper choice of the critical values \mathfrak{J}_{m,m° . One can say that the procedure is specified by a way of selecting these critical values.

Below we offer a novel way of doing this choice in a general situation by using a so called *propagation condition*: if a model m° is “good” it has to be accepted with a high probability.

This rule can be seen as analog of the family-wise level condition in a multiple test problem. Rejecting a “good” model is the family-wise error of first kind, and this error has to be controlled.

To specify precisely the meaning of a good model, consider for $m > m^\circ$ the decomposition

$$\mathbb{T}_{m,m^\circ} = \|\tilde{\phi}_m - \tilde{\phi}_{m^\circ}\| = \|\mathcal{K}_{m,m^\circ} \mathbf{Y}\| = \|\mathcal{K}_{m,m^\circ} (\mathbf{f}^* + \boldsymbol{\varepsilon})\| = \|b_{m,m^\circ} + \boldsymbol{\xi}_{m,m^\circ}\|,$$

where with $\mathcal{K}_{m,m^\circ} = \mathcal{K}_m - \mathcal{K}_{m^\circ}$

$$b_{m,m^\circ} \stackrel{\text{def}}{=} \mathcal{K}_{m,m^\circ} \mathbf{f}^*, \quad \boldsymbol{\xi}_{m,m^\circ} \stackrel{\text{def}}{=} \mathcal{K}_{m,m^\circ} \boldsymbol{\varepsilon}.$$

It obviously holds $\mathbb{E} \boldsymbol{\xi}_{m,m^\circ} = 0$. Introduce $q \times q$ -matrix \mathbb{V}_{m,m° as the variance of $\tilde{\phi}_m - \tilde{\phi}_{m^\circ}$:

$$\mathbb{V}_{m,m^\circ} \stackrel{\text{def}}{=} \text{Var}(\tilde{\phi}_m - \tilde{\phi}_{m^\circ}) = \text{Var}(\mathcal{K}_{m,m^\circ} \mathbf{Y}) = \mathcal{K}_{m,m^\circ} \text{Var}(\boldsymbol{\varepsilon}) \mathcal{K}_{m,m^\circ}^\top.$$

Further,

$$\begin{aligned} \mathbb{E} \mathbb{T}_{m,m^\circ}^2 &= \|b_{m,m^\circ}\|^2 + \mathbb{E} \|\boldsymbol{\xi}_{m,m^\circ}\|^2 = \|b_{m,m^\circ}\|^2 + \mathbf{p}_{m,m^\circ}, \\ \mathbf{p}_{m,m^\circ} &= \text{tr}(\mathbb{V}_{m,m^\circ}) = \mathbb{E} \|\boldsymbol{\xi}_{m,m^\circ}\|^2. \end{aligned}$$

The bias term $b_{m,m^\circ} \stackrel{\text{def}}{=} \mathcal{K}_{m,m^\circ} \mathbf{f}^*$ is significant if its squared norm is competitive with the variance term $\mathbf{p}_{m,m^\circ} = \text{tr}(\mathbb{V}_{m,m^\circ})$.

We say that m° is a “good” choice if there is no significant bias b_{m,m° for any $m > m^\circ$. This condition can be quantified as “bias-variance trade-off”:

$$\|b_{m,m^\circ}\|^2 \leq \beta^2 \mathbf{p}_{m,m^\circ}, \quad m > m^\circ \quad (4)$$

for a given parameter β .

Define the *oracle* m^* as the minimal m° with the property (4):

$$m^* \stackrel{\text{def}}{=} \min \left\{ m^\circ : \max_{m > m^\circ} \{ \|b_{m,m^\circ}\|^2 - \beta^2 \mathbf{p}_{m,m^\circ} \} \leq 0 \right\}.$$

Let the noise distribution be known. A particular example is the case of Gaussian errors $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$. Then the distribution of the stochastic component ξ_{m,m° is known as well. Introduce for each pair $m > m^\circ$ from \mathcal{M} a *tail function* $z_{m,m^\circ}(t)$ of the argument t such that

$$\mathbb{P}\left(\|\xi_{m,m^\circ}\| > z_{m,m^\circ}(t)\right) = e^{-t}. \quad (5)$$

Here we assume that the distribution of $\|\xi_{m,m^\circ}\|$ is continuous and the value $z_{m,m^\circ}(t)$ is well defined. Otherwise one has to define $z_{m,m^\circ}(t)$ as a smallest value providing the prescribing error probability e^{-t} .

For checking the propagation condition, we need a uniform in $m > m^\circ$ version of the probability bound (5). Let

$$\mathcal{M}^+(m^\circ) \stackrel{\text{def}}{=} \{m \in \mathcal{M} : m > m^\circ\}.$$

Given \mathbf{x} , by $q_{m^\circ} = q_{m^\circ}(\mathbf{x})$ denote the corresponding multiplicity correction:

$$\mathbb{P}\left(\bigcup_{m \in \mathcal{M}^+(m^\circ)} \{\|\xi_{m,m^\circ}\| \geq z_{m,m^\circ}(\mathbf{x} + q_{m^\circ})\}\right) = e^{-\mathbf{x}}.$$

A simple way of computing the multiplicity correction q_{m° is based on the Bonferroni bound:

$$q_{m^\circ} = \log(\#\mathcal{M}^+(m^\circ)).$$

However, it is well known that the Bonferroni bound is very conservative and leads to very large correction q_{m° , especially if the random vectors ξ_{m,m° are strongly correlated.

As the joint distribution of the ξ_{m,m° 's is precisely known, define the correction $q_{m^\circ} = q_{m^\circ}(\mathbf{x})$ just by condition

$$\mathbb{P}\left(\bigcup_{m \in \mathcal{M}^+(m^\circ)} \{\|\xi_{m,m^\circ}\| \geq z_{m,m^\circ}(\mathbf{x} + q_{m^\circ})\}\right) = e^{-x}.$$

Finally we define the critical values \mathfrak{z}_{m,m° by one more correction for the bias:

$$\mathfrak{z}_{m,m^\circ} \stackrel{\text{def}}{=} z_{m,m^\circ}(\mathbf{x} + q_{m^\circ}) + \beta \sqrt{\mathbf{p}_{m,m^\circ}} \quad (6)$$

for $\mathbf{p}_{m,m^\circ} = \text{tr}(\mathbb{V}_{m,m^\circ})$.

In practice $x = 3$ and $\beta = 0$ provide a reasonable choice.

Define the selector \hat{m} by the “smallest accepted” (SmA) rule. Namely, with \mathfrak{J}_{m,m° from (6), the acceptance rule reads as follows:

$$\{m^\circ \text{ is accepted}\} \Leftrightarrow \left\{ \max_{m \in \mathcal{M}^+(m^\circ)} \{\mathbb{T}_{m,m^\circ} - \mathfrak{J}_{m,m^\circ}\} \leq 0 \right\}.$$

The SmA rule is

$$\begin{aligned} \hat{m} &\stackrel{\text{def}}{=} \text{“smallest accepted”} \\ &= \min \left\{ m^\circ : \max_{m \in \mathcal{M}^+(m^\circ)} \{\mathbb{T}_{m,m^\circ} - \mathfrak{J}_{m,m^\circ}\} \leq 0 \right\}. \end{aligned}$$

Our study mainly focuses on the behavior of the selector \hat{m} . The performance of the resulting estimator $\hat{\phi} = \tilde{\phi}_{\hat{m}}$ is a kind of corollary from statements about the selected model \hat{m} . The ideal solution would be $\hat{m} \equiv m^*$, then the adaptive estimator $\hat{\phi}$ coincides with the oracle estimate $\tilde{\phi}_{m^*}$.

The decomposition

$$\mathbb{T}_{m,m^\circ} = \|\tilde{\phi}_m - \tilde{\phi}_{m^\circ}\| = \|b_{m,m^\circ} + \xi_{m,m^\circ}\| \leq \|b_{m,m^\circ}\| + \|\xi_{m,m^\circ}\|,$$

and the bounds

$$\mathbb{P}\left(\bigcup_{m \in \mathcal{M}^+(m^\circ)} \{\|\xi_{m,m^\circ}\| \geq z_{m,m^\circ}(\mathbf{x} + q_{m^\circ})\}\right) = e^{-x},$$

$$\|b_{m,m^\circ}\|^2 \leq \beta^2 p_{m,m^\circ}, \quad m > m^\circ$$

automatically ensures the desired *propagation property*.

Theorem. Any good model m° will be accepted with probability at least $1 - e^{-x}$:

$$\mathbb{P}(m^* \text{ is rejected}) \leq e^{-x}.$$

Corollary. By definition, the oracle m^* is also a “good” choice, thus accepted.

The oracle m^* is also a “good” choice, this yields

$$\mathbb{P}(m^* \text{ is rejected}) \leq e^{-x}.$$

Therefore, the selector \hat{m} typically takes its value in $\mathcal{M}^-(m^*)$, where

$$\mathcal{M}^-(m^*) = \{m \in \mathcal{M} : m < m^*\}$$

is the set of all models in \mathcal{M} smaller than m^* .

Zone of insensitivity: a subset \mathcal{M}° of $\mathcal{M}^-(m^*)$ of possible \hat{m} -values.

The definition of m^* implies that there is a significant bias for each $m \in \mathcal{M}^-(m^*)$.

Intuition: Zone of insensitivity is composed of m -values for which the bias is significant but not very large.

Theorem. For any subset $\mathcal{M}^c \subseteq \mathcal{M}^-(m^*)$ s.t.

$$\|b_{m^*,m}\| > \mathfrak{z}_{m^*,m} + z_{m^*,m}(\mathbf{x}_s), \quad m \in \mathcal{M}^c,$$

for $\mathbf{x}_s \stackrel{\text{def}}{=} \mathbf{x} + \log(|\mathcal{M}^c|)$ with $|\mathcal{M}^c|$ being the cardinality of \mathcal{M}^c , it holds

$$\mathbb{P}(\hat{m} \in \mathcal{M}^c) \leq e^{-x}.$$

The SmA estimator $\hat{\phi} = \tilde{\phi}_{\hat{m}}$ satisfies the following bound:

$$\mathbb{P}\left(\|\hat{\phi} - \tilde{\phi}_{m^*}\| > \bar{\mathfrak{z}}_{m^*}\right) \leq 2e^{-x},$$

where $\bar{\mathfrak{z}}_{m^*}$ is defined with $\mathcal{M}^\circ \stackrel{\text{def}}{=} \mathcal{M}^-(m^*) \setminus \mathcal{M}^c$ as

$$\bar{\mathfrak{z}}_{m^*} \stackrel{\text{def}}{=} \max_{m \in \mathcal{M}^\circ} \mathfrak{z}_{m^*,m}.$$

1 Introduction

2 SmA procedure for known noise variance

3 Bootstrap tuning

4 Numerical results

Consider the bootstrap estimates $\tilde{\phi}_m = W\tilde{\theta}_m^b$ in the form

$$\tilde{\phi}_m^b = W(\Psi_m \Psi_m^\top)^{-1} \Psi_m W^b \mathbf{Y} = \mathcal{K}_m W^b \mathbf{Y}.$$

Here $W^b \mathbf{Y}$ means the vector with entries $w_i^b Y_i$ for $i \leq n$, where w_i^b are i.i.d. bootstrap weights with $E^b w_i^b = \text{Var}(w_i^b) = 1$. We are interested if the distribution of the differences

$$\tilde{\phi}_m^b - \tilde{\phi}_{m^\circ}^b = \mathcal{K}_{m,m^\circ} W^b \mathbf{Y}, \quad m > m^\circ$$

mimics their real world counterparts.

The identity $E^b W^b = I_n$ yields $E^b \tilde{\phi}_{m,m^\circ}^b = \tilde{\phi}_{m,m^\circ}$, and the natural idea would be to use the difference

$$\tilde{\phi}_{m,m^\circ}^b - \tilde{\phi}_{m,m^\circ} = \mathcal{K}_{m,m^\circ} (W^b - I_n) \mathbf{Y}$$

as a proxy for the stochastic component $\xi_{m,m^\circ} = \mathcal{K}_{m,m^\circ} \varepsilon$.

Unfortunately, this can only be justified if the bias component b_{m,m° of $\tilde{\phi}_{m,m^\circ}^b$ is small relative to its stochastic variance p_{m,m° . But this is exactly what we would like to test!

To avoid this problem we apply a presmoothing which removes a pilot prediction of the regression function from the data. This presmoothing requires some minimal smoothness of the regression function, and this condition seems to be unavoidable if no information about noise is given: otherwise one cannot separate between signal and noise.

Suppose that a linear predictor $\tilde{f}_0 = \Pi Y$ is given where Π is a sub-projector in the space \mathbb{R}^n . In most of cases one can take $\Pi = \Psi_{m^\dagger}^\top (\Psi_{m^\dagger} \Psi_{m^\dagger}^\top)^{-1} \Psi_{m^\dagger}$ where m^\dagger is a large index, e.g. the largest index M in our collection.

Idea: computes the residuals $\check{Y} = Y - \Pi Y$ and uses them in place of the original data. This allows to remove the bias while keeping the noise variance only slightly changed.

For each bootstrap realization $w^b = (w_i^b)$, we apply the procedure to the data vector $\mathcal{W}^b \check{Y}$ with entries $\check{Y}_i w_i^b$ for $i \leq n$. The bootstrap stochastic components ξ_{m, m°^b are defined as

$$\xi_{m, m^\circ}^b \stackrel{\text{def}}{=} \mathcal{K}_{m, m^\circ} \mathcal{E}^b \check{Y}, \quad m > m^\circ,$$

where $\mathcal{E}^b = \mathcal{W}^b - I_n$ is the diagonal matrix of bootstrap errors $\varepsilon_i^b = w_i^b - 1 \sim \mathcal{N}(0, 1)$.

The bootstrap quantiles $z_{m,m^\circ}^b(t)$ are given by the analog of (5):

$$\mathbb{P}^b\left(\|\boldsymbol{\xi}_{m,m^\circ}^b\| > z_{m,m^\circ}^b(t)\right) = e^{-t}.$$

The multiplicity correction $q_{m^\circ}^b = q_{m^\circ}^b(\mathbf{x})$ is specified by the condition

$$\mathbb{P}^b\left(\bigcup_{m \in \mathcal{M}^+(m^\circ)} \{\|\boldsymbol{\xi}_{m,m^\circ}^b\| \geq z_{m,m^\circ}^b(\mathbf{x} + q_{m^\circ}^b)\}\right) = e^{-x}.$$

Finally, the bootstrap critical values are fixed by the analog of (6):

$$\mathfrak{z}_{m,m^\circ}^b \stackrel{\text{def}}{=} z_{m,m^\circ}^b(\mathbf{x} + q_{m^\circ}^b) + \beta \sqrt{\mathbb{P}_{m,m^\circ}^b}$$

for $\mathbb{P}_{m,m^\circ}^b = \mathbb{E}^b \|\boldsymbol{\xi}_{m,m^\circ}^b\|^2$. Remind that all these quantities are data-driven and depend upon the original data. Now we apply the SmA procedure with such defined critical values

$\mathfrak{z}_{m,m^\circ}^b$.

- **Design Regularity** is measured by the value δ_Ψ

$$\delta_\Psi \stackrel{\text{def}}{=} \max_{i=1, \dots, n} \|S^{-1/2} \Psi_i\|_{\sigma_i}, \quad \text{where} \quad S \stackrel{\text{def}}{=} \sum_{i=1}^n \Psi_i \Psi_i^\top \sigma_i^2; \quad (7)$$

Obviously

$$\sum_{i=1}^n \|S^{-1/2} \Psi_i\|^2 \sigma_i^2 = \text{tr} \left(\sum_{i=1}^n S^{-2} \Psi_i \Psi_i^\top \sigma_i^2 \right) = \text{tr} I_p = p,$$

and therefore in typical situations the value δ_Ψ is of order $\sqrt{p/n}$.

- **Presmoothing bias** for a projector Π is described by the vector

$$\mathbf{B} = \Sigma^{-1/2} (\mathbf{f}^* - \Pi \mathbf{f}^*).$$

We will use the sup-norm $\|\mathbf{B}\|_\infty = \max_i |b_i|$ and the squared ℓ_2 -norm $\|\mathbf{B}\|^2 = \sum_i b_i^2$ to measure the bias after presmoothing.

- Stochastic noise after presmoothing** is described via the covariance matrix $\text{Var}(\check{\varepsilon})$ of the smoothed noise $\check{\varepsilon} = \Sigma^{-1/2}(\varepsilon - \Pi\varepsilon)$. Namely, this matrix is assumed to be sufficiently close to the unit matrix I_n , in particular, its diagonal elements should be close to one. This is measured by the operator norm of $\text{Var}(\check{\varepsilon}) - I_n$ and by deviations of the individual variances $\mathbb{E}\check{\varepsilon}_i^2$ from one:

$$\delta_1 \stackrel{\text{def}}{=} \|\text{Var}(\check{\varepsilon}) - I_n\|_{\text{op}},$$

$$\delta_\varepsilon \stackrel{\text{def}}{=} \max_i |\mathbb{E}\check{\varepsilon}_i^2 - 1|.$$

In particular, in the case of homogeneous errors $\Sigma = \sigma^2 I_n$ and the smoothing operator Π as a p -dimensional projector, it holds

$$\text{Var}(\check{\varepsilon}) = (I_n - \Pi)^2 = I_n - \Pi \leq I_n,$$

$$\delta_1 = \|\text{Var}(\check{\varepsilon}) - I_n\|_{\text{op}} = \|\Pi\|_{\text{op}} = 1,$$

$$\delta_\varepsilon = \max_i |\mathbb{E}\check{\varepsilon}_i^2 - 1| = \max_i |\Pi_{ii}|.$$

- Regularity of the smoothing operator** Π is required in Theorems 2, 3, and 4. This condition will be expressed via the norm of the rows \mathcal{Y}_i^\top of the matrix $\mathcal{Y} \stackrel{\text{def}}{=} \Sigma^{-1/2} \Pi \Sigma^{1/2}$ fulfill

$$\|\mathcal{Y}_i^\top\| \leq \delta_\Psi, \quad i = 1, \dots, n. \quad (8)$$

This condition is in fact very close to the design regularity condition (7). To see this, consider the case of a homogeneous noise with $\Sigma = \sigma^2 I_n$ and $\Pi = \Psi^\top (\Psi \Psi^\top)^{-1} \Psi$. Then $\mathcal{Y} = \Pi$ and (7) implies

$$\|\mathcal{Y}_i^\top\| = \|\Psi^\top (\Psi \Psi^\top)^{-1} \Psi_i\| = \|(\Psi \Psi^\top)^{-1/2} \Psi_i\| \leq \delta_\Psi.$$

In general one can expect that (8) is fulfilled with some other constant which however, is of the same magnitude as δ_Ψ . For simplicity we use the same letter.

Let $\mathbf{Y} = \mathbf{f}^* + \varepsilon \sim \mathcal{N}(\mathbf{f}^*, \Sigma)$ for $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$.

$$\delta_\Psi \stackrel{\text{def}}{=} \max_{i=1, \dots, n} \|S^{-1/2} \Psi_i\| \sigma_i, \quad \text{where} \quad S \stackrel{\text{def}}{=} \sum_{i=1}^n \Psi_i \Psi_i^\top \sigma_i^2;$$

$$\delta_\varepsilon = \max_i |\mathbb{E} \check{\varepsilon}_i^2 - 1|, \quad \mathbf{B} = \Sigma^{-1/2} (\mathbf{f}^* - \Pi \mathbf{f}^*).$$

Consider

$$\mathbb{Q} = \mathcal{L}(\boldsymbol{\xi}_{m, m^\circ}, m, m^\circ \in \mathcal{M}), \quad \mathbb{Q}^b = \mathcal{L}^b(\boldsymbol{\xi}_{m, m^\circ}^b, m, m^\circ \in \mathcal{M}).$$

Theorem 2. *It holds on a random set $\Omega_2(\mathbf{x})$ with $\mathbb{P}(\Omega_2(\mathbf{x})) \geq 1 - 3e^{-x}$:*

$$\|\mathbb{Q} - \mathbb{Q}^b\|_{\text{TV}} \leq \frac{1}{2} \Delta_2(\mathbf{x}),$$

$$\Delta_2(\mathbf{x}) \stackrel{\text{def}}{=} 2\sqrt{\delta_\Psi^2 p \mathbf{x}_n} + \sqrt{\delta_\varepsilon^2 p} + \sqrt{\|\mathbf{B}\|_\infty^4 p} + 4\delta_\Psi^2 \|\mathbf{B}\| (1 + \sqrt{\mathbf{x}}).$$

where $\mathbf{x}_n = \mathbf{x} + \log(n)$.

Theorem 3. (Bootstrap validity) *Assume the conditions of Theorem 2, and let the rows Υ_i^\top of the matrix $\Upsilon \stackrel{\text{def}}{=} \Sigma^{-1/2} \Pi \Sigma^{1/2}$ fulfill (8). Then for each $m^\circ \in \mathcal{M}$*

$$\mathbb{P} \left(\max_{m > m^\circ} \left\{ \|\boldsymbol{\xi}_{m, m^\circ}\| - z_{m, m^\circ}^b(\mathbf{x} + q_{m^\circ}^b) \right\} \geq 0 \right) \leq 6e^{-\mathbf{x}} + \sqrt{p} \Delta_0(\mathbf{x}),$$

where with $\mathbf{x}_n = \mathbf{x} + \log(n)$ and $\mathbf{x}_p = \mathbf{x} + \log(2p)$

$$\Delta_0(\mathbf{x}) \stackrel{\text{def}}{=} \|\mathbf{B}\|_\infty^2 + \delta_\Psi^2 \|\mathbf{B}\| \sqrt{2\mathbf{x}} + 2\delta_\Psi \mathbf{x}_n + \delta_\Psi^2 \mathbf{x}_n + 2\delta_\Psi \sqrt{\mathbf{x}_p} + 2\delta_\Psi^2 \mathbf{x}_p.$$

The SmA procedure also involves the values \mathbf{p}_{m,m° which are unknown and depend on the noise ε . The next result shows the bootstrap counterparts \mathbf{p}_{m,m°^b can be well used in place of \mathbf{p}_{m,m° .

Theorem 4. *Assume the conditions of Theorem 2. Then it holds on a set $\Omega_1(\mathbf{x})$ with $\mathbb{P}(\Omega_1(\mathbf{x})) \geq 1 - 3e^{-x}$ for all pairs $m < m^\circ \in \mathcal{M}$*

$$\left| \frac{\mathbf{p}_{m,m^\circ}^b}{\mathbf{p}_{m,m^\circ}} - 1 \right| \leq \Delta_p,$$

$$\Delta_p \stackrel{\text{def}}{=} \|\mathbf{B}\|_\infty^2 + 4\mathbf{x}_{\mathcal{M}}^{1/2} \delta_n^2 \|\mathbf{B}\| + 4\mathbf{x}_{\mathcal{M}}^{1/2} \delta_n + 4\mathbf{x}_{\mathcal{M}} \delta_n^2 + \delta_\varepsilon,$$

where $\mathbf{p}_{m,m^\circ}^b = \mathbb{E}^b \|\boldsymbol{\xi}_{m,m^\circ}^b\|^2$, $\mathbf{p}_{m,m^\circ} = \mathbb{E} \|\boldsymbol{\xi}_{m,m^\circ}\|^2$, and $\mathbf{x}_{\mathcal{M}} = \mathbf{x} + 2 \log(|\mathcal{M}|)$.

The above results immediately imply all the oracle bounds for probabilistic loss with the obvious correction of the error terms.

- 1 Introduction
- 2 SmA procedure for known noise variance
- 3 Bootstrap tuning
- 4 Numerical results**

Consider a regression problem for an unknown univariate function on $[0, 1]$ with unknown inhomogeneous noise. The aim is to compare the bootstrap-calibrated procedure with the Sma procedure for the known noise and with the oracle estimator. We also check the sensitivity of the method to the choice of the presmoothing parameter m^\dagger .

We use a uniform design on $[0, 1]$ and the Fourier basis $\{\psi_j(x)\}_{j=1}^\infty$ for approximation of the regression function f which is modelled in the form

$$f(x) = c_1\psi_1(x) + \dots + c_p\psi_p(x),$$

where the $(c_j)_{1 \leq j \leq p}$ are chosen randomly: with γ_j i.i.d. standard normal

$$c_j = \begin{cases} \gamma_j, & 1 \leq j \leq 10, \\ \gamma_j/(j-10)^2, & 11 \leq j \leq 200. \end{cases}$$

The noise intensity grows from low to high as x increases to one. We use

$n_{\text{sim-bs}} = n_{\text{sim-theo}} = n_{\text{sim-calib}} = 1000$ samples for computing the bootstrap marginal quantiles and the theoretical quantiles and for checking the calibration condition. The maximal model dimension is $M = 34$ and we also choose $m^\dagger = M$. The calibration is run with $x = 2$ and $\beta = 1$.

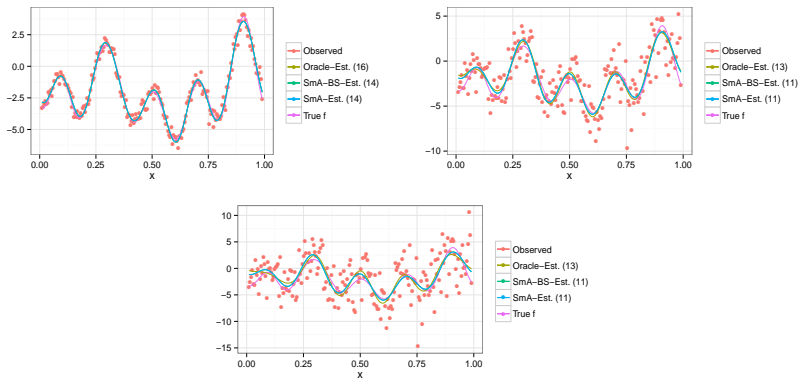


Figure : True functions and observed values plotted with oracle estimator, the known-variance SmA-Estimator (SmA-Est.) and the Bootstrap-SmA-Estimator (SmA-BS-Est.) for 3 different functions with different noise structure going from low noise to high noise. The numbers in parentheses indicate the chosen model dimension.

The oracles are respectively $m^* = 12$ for $n = 100, 200$ and $m^* = 10$ for $n = 50$.

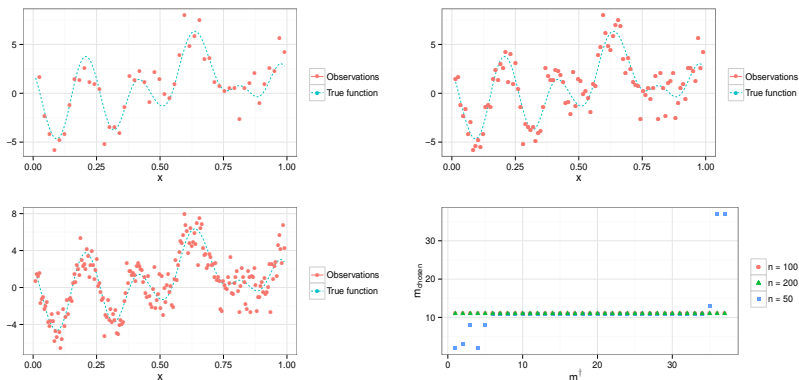


Figure : The first three plots show an exemplary function with $n = 50, 100, 200$ observations. The right plot shows the \hat{m} chosen by the Bootstrap-SmA-Method as a function of the calibration dimension m^\dagger and the number of observations.

Figure 3 again demonstrates the dependence of the ratios on m^\dagger . It is remarkable that the ratio is varying very slowly above $m^* = 12$.

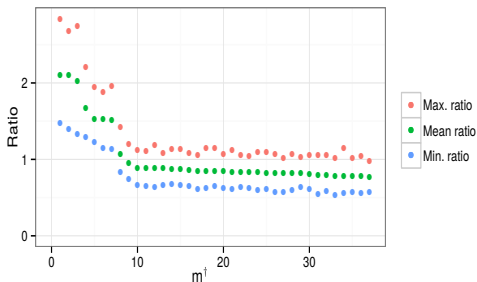


Figure : Maximal, minimal and mean ratio of the bootstrap and theoretical tail functions at $x = 2$, $|\hat{\mathfrak{J}}_{m_1, m_2}^b / \hat{\mathfrak{J}}_{m_1, m_2}|^2$ as a function of m^\dagger .

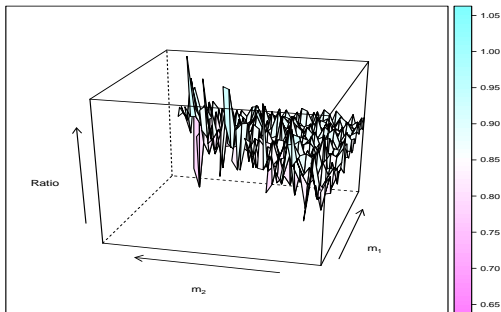


Figure : Ratio of quantiles $|\hat{\beta}_{m_1, m_2}^b / \hat{\beta}_{m_1, m_2}|^2$ for $m^\dagger = 20$ and $n = 200$ with the data and true function as in Fig. 2.

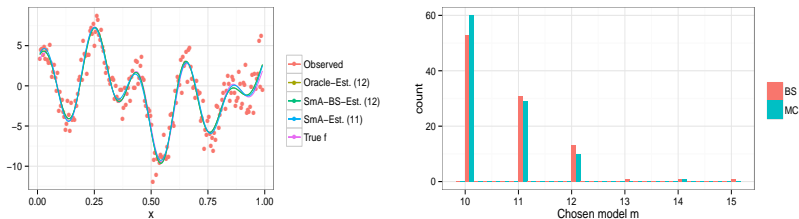


Figure : In the left plot, the true function and observed values are plotted for one realization together with the oracle estimator, the known-variance SmA-Estimator (SmA-Est.) and the Bootstrap-SmA-Estimator (SmA-BS-Est.). The numbers in parentheses indicate the chosen model dimension. In the right plot, histograms for the selected model are given for the bootstrap (BS) and the known-variance method (MC) for repeated observations of the same underlying function with a simulation size $n_{\text{hist}} = 100$.

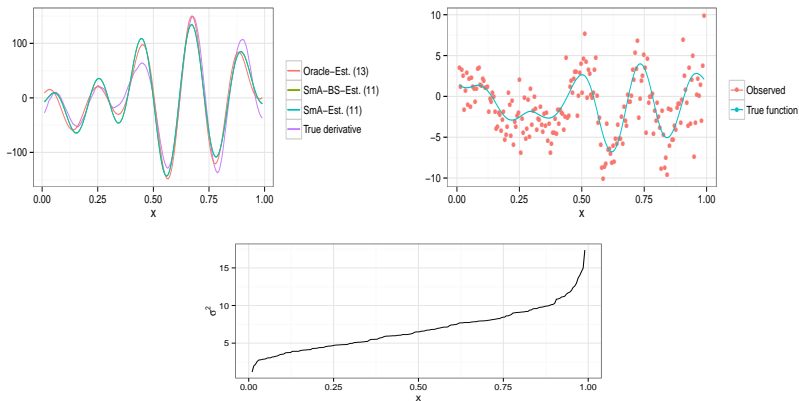









Figure : The upper left plot shows the true derivative, the oracle estimator, the known-variance SmA-Estimator (SmA-Est.) and the Bootstrap-SmA-Estimator (SmA-BS-Est.). The upper right plot shows the true function and the observations and in the lower plot one can find the standard deviation of the errors.

- A unified fully adaptive procedure for ordered model selection.
- Sharp oracle bounds
- Impact of the bias in the size of bootstrap confidence sets.

In progress:

- Model selection for **unordered case** like anisotropic classes. Theory applies but has to be extended;
- Active set selection. Theory applies but the problem of algorithmic efficient implementation.
- Large p s.t. p^2/n large; Use of sparse or complexity **penalty**.
- Extension to other problems like Hidden Markov Chain modeling;
- Multiscale change-point detection.

-  Arlot, S. (2009).
Model selection by resampling penalization.
Electron. J. Statist., 3:557–624.
-  Arlot, S. and Bach, F. R. (2009).
Data-driven calibration of linear estimators with minimal penalties.
In Bengio, Y., Schuurmans, D., Lafferty, J., Williams, C., and Culotta, A., editors, *Advances in Neural Information Processing Systems 22*, pages 46–54. Curran Associates, Inc.
-  Barron, A., Birgé, L., and Massart, P. (1999).
Risk bounds for model selection via penalization.
Probab. Theory Related Fields, 113(3):301–413.
-  Birgé, L. (2001).
An alternative point of view on Lepski's method, volume Volume 36 of *Lecture Notes–Monograph Series*, pages 113–133.
Institute of Mathematical Statistics, Beachwood, OH.
-  Birgé, L. and Massart, P. (2007).
Minimal penalties for gaussian model selection.
Probability Theory and Related Fields, 138(1-2):33–73.
-  Cavalier, L., Golubev, G. K., Picard, D., and Tsybakov, A. B. (2002).
Oracle inequalities for inverse problems.
Ann. Statist., 30(3):843–874.
-  Cavalier, L. and Golubev, Y. (2006).
Risk hull method and regularization by projections of ill-posed inverse problems.
Ann. Statist., 34(4):1653–1677.



Chernozhukov, V., Chetverikov, D., and Kato, K. (2014).
Anti-concentration and honest, adaptive confidence bands.
Ann. Statist., 42(5):1787–1818.



Dalalyan, A. S. and Salmon, J. (2012).
Sharp oracle inequalities for aggregation of affine estimators.
Ann. Statist., 40(4):2327–2355.



Gach, F., Nickl, R., and Spokoiny, V. (2013).
Spatially Adaptive Density Estimation by Localised Haar Projections.
Annales de l'Institut Henri Poincaré - Probability and Statistics, 49(3):900–914.
DOI: 10.1214/12-AIHP485; arXiv:1111.2807.



Goldenshluger, A. (2009).
A universal procedure for aggregating estimators.
Ann. Statist., 37(1):542–568.



Kneip, A. (1994).
Ordered linear smoothers.
Ann. Statist., 22(2):835–866.



Lepski, O. V. (1990).
A problem of adaptive estimation in Gaussian white noise.
Teor. Veroyatnost. i Primenen., 35(3):459–470.



Lepski, O. V. (1991).
Asymptotically minimax adaptive estimation. I. Upper bounds. Optimally adaptive estimates.
Teor. Veroyatnost. i Primenen., 36(4):645–659.



Lepski, O. V. (1992).

Asymptotically minimax adaptive estimation. II. Schemes without optimal adaptation. Adaptive estimates.
Teor. Veroyatnost. i Primenen., 37(3):468–481.



Lepski, O. V., Mammen, E., and Spokoiny, V. G. (1997).

Optimal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selectors.
The Annals of Statistics, 25(3):929–947.



Lepski, O. V. and Spokoiny, V. G. (1997).

Optimal pointwise adaptive methods in nonparametric estimation.
Ann. Statist., 25(6):2512–2546.



Massart, P. (2007).

Concentration inequalities and model selection.
Number 1896 in Ecole d'Été de Probabilités de Saint-Flour. Springer.



Spokoiny, V. and Vial, C. (2009).

Parameter tuning in pointwise adaptation using a propagation approach.
Ann. Statist., 37:2783–2807.



Spokoiny, V., Wang, W., and Härdle, W. (2013).

Local quantile regression (with rejoinder).
J. of Statistical Planning and Inference, 143(7):1109–1129.
ArXiv:1208.5384.



Tsybakov, A. (2000).

On the best rate of adaptive estimation in some inverse problems.
Comptes Rendus de l'Academie des Sciences - Series I - Mathematics, 330(9):835 – 840.