

Simultaneous adaptation for several criteria using an extended Lepskii principle

G. Blanchard

Université Paris-Saclay

Symposium on Inverse Problems: From experimental data to models and back
Potsdam/Griebnitzsee 19-21/09/2022

Based on joint work with: P. Mathé (Weierstrass Institute, Berlin), N. Mücke (TU Braunschweig)

université
PARIS-SACLAY

Setting: linear regression in Hilbert space

We consider the observation model

$$Y_i = \langle f_o, X_i \rangle + \xi_i,$$

where

- ▶ X_i takes its values in a Hilbert space \mathcal{H} , with $\|X_i\| \leq 1$ a.s.;
- ▶ ξ_i is a random variable with $\mathbb{E}[\xi_i | X_i] = 0$, $\mathbb{E}[\xi_i^2 | X_i] \leq \sigma^2$, $|\xi_i| \leq M$ a.s.;
- ▶ $(X_i, \xi_i)_{1 \leq i \leq n}$ are i.i.d. (the distribution of X is **not known**.)

The goal is to estimate f_o (in a sense to be specified) from the data.

Note that if $\dim(\mathcal{H}) = \infty$, this is essentially a non-parametric model.

Why this model?

- ▶ Hilbert-space valued variables appear in standard models of **Functional Data Analysis**, where the observed data are modeled (idealized) as function-valued.
- ▶ Such models also appear in **reproducing kernel Hilbert space (RKHS) methods** in machine learning:
 - ▶ assume observations X_i take values in some space \mathcal{X}
 - ▶ let $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ be a “feature mapping” in a Hilbert space \mathcal{H} , and $\tilde{X} = \Phi(X)$, then one considers the model

$$Y_i = \langle f_0, \tilde{X}_i \rangle + \zeta_i = \tilde{f}_0(X_i) + \zeta_i,$$

where $\tilde{f} \in \tilde{\mathcal{H}} := \{x \mapsto \langle f, \Phi(x) \rangle; f \in \mathcal{H}\}$ is a nonparametric model of functions (nonlinear in x !).

- ▶ Usually all computations don't require explicit knowledge of Φ but only access to the **kernel** $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$.

Why this model (II) - inverse learning

Of interest is also the **inverse learning** problem:

- ▶ X_i takes value in \mathcal{X} ;
- ▶ if A is a (**known**) linear operator from a Hilbert \mathcal{G} to a real function space on \mathcal{X} ;
- ▶ inverse regression learning model:

$$Y_i = (Ag_*)(X_i) + \xi_i,$$

- ▶ where A is a Carleman operator
(i.e. evaluation functionals $f \mapsto (Af)(x)$ are continuous for all x),
- ▶ In this case the goal is to recover $g_* \in \mathcal{G}$.

Why this model (III) - inverse learning, continued

- ▶ inverse regression learning model: $Y_i = \underbrace{(Ag_*)}_{f_o}(X_i) + \zeta_i.$

- ▶ If A is a Carleman operator $\mathcal{G} \rightarrow \mathcal{F}_{\text{mes.}}(\mathcal{X}, \mathbb{R})$, then of evaluation functionals

$$\text{For all } f \in \mathcal{G}, x \in \mathcal{X} : \quad (Af)(x) = \langle F_x, f \rangle_{\mathcal{G}} \text{ for some } F_x \in \mathcal{G}$$

- ▶ Then $\mathcal{H} := \text{Im}(A)$ can be equipped with a RKHS structure with kernel

$$k(x, x') := \langle F_x, F_{x'} \rangle_{\mathcal{G}}.$$

- ▶ Furthermore, A is then a **partial isometry** $\mathcal{H} \rightarrow \mathcal{G}$.

- ▶ Therefore, if $\hat{f} \in \mathcal{H}$ is an estimate of $f_o = Ag_*$ and if we assume $g_* \in \text{Ker}(A)^\perp$:

$$\text{Put } \hat{g} := A^{-1}\hat{f}, \text{ then } \|\hat{g} - g_*\|_{\mathcal{G}} = \left\| A^{-1}(\hat{f} - f_o) \right\|_{\mathcal{G}} = \left\| \hat{f} - f_o \right\|_{\mathcal{H}}$$

- ▶ Here the RKHS \mathcal{H} is entirely determined by A . Mathematically speaking, we are back in the RKHS learning scenario, but **the convergence in \mathcal{H} -norm** is of major importance.

Inverse regression vs inverse “learning”

- ▶ Bissantz, Hohage, Munk and Ruymgaart (2007) propose a very general analysis of general regularization methods for **statistical** inverse problems.
- ▶ Their model includes applications to the inverse regression model where the **design distribution** (X -marginal) is assumed to be known (the exact integral operator is used to construct the estimator).
- ▶ A proper characteristic of inverse **“learning”** is the absence of information a priori on the X -marginal – it has to be “learned” as well.

Two notions of risk

We will consider two notions of error (risk) for a candidate estimate \hat{f} of f_0 :

- ▶ Squared prediction error:

$$\mathcal{E}(\hat{f}) := \mathbb{E} \left[\left(\langle \hat{f}, X \rangle - Y \right)^2 \right].$$

- ▶ The associated (excess error) risk is

$$\mathcal{E}(\hat{f}) - \mathcal{E}(f_0) = \mathbb{E} \left[\left(\langle \hat{f} - f_0, X \rangle \right)^2 \right] = \left\| \hat{f} - f_0 \right\|_{2, X}^2,$$

- ▶ Reconstruction error risk (especially relevant for **inverse learning**):

$$\left\| \hat{f} - f_0 \right\|_{\mathcal{H}}^2.$$

The goal is to find a suitable estimator \hat{f} of f_0 from the data having “optimal” convergence properties with respect to these two risks.

Finite-dimensional case

- ▶ The final dimensional case: $\mathcal{X} = \mathbb{R}^p$, f_{\circ} now denoted β_{\circ}
- ▶ In usual matrix form:

$$Y = X\beta_{\circ} + \xi.$$

- ▶ X_j^T form the lines of the (n, p) design matrix X
 - ▶ $Y = (Y_1, \dots, Y_n)^T$
 - ▶ $\xi = (\xi_1, \dots, \xi_n)^T$
- ▶ “Reconstruction” risk corresponds to $\|\beta_{\circ} - \hat{\beta}\|^2$.
 - ▶ Prediction risk corresponds to

$$\mathbb{E} \left[\langle \beta_{\circ} - \hat{\beta}, X \rangle^2 \right] = \|\Sigma^{1/2}(\beta_{\circ} - \hat{\beta})\|^2,$$

where $\Sigma := \mathbb{E}[XX^T]$.

- ▶ In Hilbert space, same relation with $\Sigma := \mathbb{E}[X \otimes X^*]$.

Convergence of OLS in finite dimension

- ▶ The “ordinary” least squares (OLS) solution:

$$\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

- ▶ We want to understand the behavior of $\hat{\beta}_{OLS}$, when the data size n grows large. Will we be close to the truth β_o ?
- ▶ Recall

$$\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \underbrace{\left(\frac{1}{n} \mathbf{X}^T \mathbf{X}\right)^{-1}}_{:=\hat{\Sigma}} \underbrace{\left(\frac{1}{n} \mathbf{X}^T \mathbf{Y}\right)}_{:=\hat{\gamma}} = \hat{\Sigma}^{-1} \hat{\gamma},$$

- ▶ Observe by a vectorial LLN, as $n \rightarrow \infty$:

$$\hat{\Sigma} := \frac{1}{n} \mathbf{X}^T \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \underbrace{X_i X_i^T}_{=: Z_i^T} \longrightarrow \mathbb{E}[X_1 X_1^T] =: \Sigma;$$

$$\hat{\gamma} := \frac{1}{n} \mathbf{X}^T \mathbf{Y} = \frac{1}{n} \sum_{i=1}^n \underbrace{X_i Y_i}_{=: Z_i} \longrightarrow \mathbb{E}[X_1 Y_1] = \Sigma \beta_o =: \gamma;$$

- ▶ Hence $\hat{\beta} = \hat{\Sigma}^{-1} \hat{\gamma} \rightarrow \Sigma^{-1} \gamma = \beta_o$. (Assuming Σ invertible.)

From OLS to Hilbert-space regression

- ▶ For ordinary linear regression with $\mathcal{X} = \mathbb{R}^p$ (fixed $p, n \rightarrow \infty$):
 - ▶ LLN implies $\hat{\beta}_{OLS} (= \hat{\Sigma}^{-1} \hat{\gamma}) \rightarrow \beta_o (= \Sigma^{-1} \gamma)$;
 - ▶ CLT+Delta Method imply asymptotic normality and convergence in $\mathcal{O}(n^{-\frac{1}{2}})$.
- ▶ How to generalize to $\mathcal{X} = \mathcal{H}$?
- ▶ **Main issue:** $\Sigma = \mathbb{E}[X \otimes X^*]$ does not have a continuous inverse.
(\rightarrow ill-posed problem)

From OLS to Hilbert-space regression

- ▶ For ordinary linear regression with $\mathcal{X} = \mathbb{R}^p$ (fixed $p, n \rightarrow \infty$):
 - ▶ LLN implies $\hat{\beta}_{OLS}(= \hat{\Sigma}^{-1}\hat{\gamma}) \rightarrow \beta_o(= \Sigma^{-1}\gamma)$;
 - ▶ CLT+Delta Method imply asymptotic normality and convergence in $\mathcal{O}(n^{-\frac{1}{2}})$.
- ▶ How to generalize to $\mathcal{X} = \mathcal{H}$?
- ▶ **Main issue:** $\Sigma = \mathbb{E}[X \otimes X^*]$ does not have a continuous inverse.
(\rightarrow ill-posed problem)
- ▶ Need to consider a suitable approximation $\zeta(\hat{\Sigma})$ of Σ^{-1} (**regularization**), where

$$\hat{\Sigma} := \frac{1}{n} \sum_{i=1}^m X_i \otimes X_i^*$$

is the empirical second moment operator.

Regularization methods

- ▶ Main idea: replace $\widehat{\Sigma}^{-1}$ by an approximate inverse, such as
- ▶ Ridge regression/Tikhonov:

$$\widehat{f}_{\text{Ridge}(\lambda)} = (\widehat{\Sigma} + \lambda I_p)^{-1} \widehat{\gamma}$$

- ▶ PCA projection/spectral cut-off: restrict $\widehat{\Sigma}$ on its k first eigenvectors

$$\widehat{f}_{\text{PCA}(k)} = (\widehat{\Sigma})_{|k}^{-1} \widehat{\gamma}$$

- ▶ Gradient descent/Landweber Iteration/ L^2 boosting:

$$\begin{aligned} \widehat{f}_{\text{LW}(k)} &= \widehat{f}_{\text{LW}(k-1)} + (\widehat{\gamma} - \widehat{\Sigma} \widehat{f}_{\text{LW}(k-1)}) \\ &= \sum_{i=0}^k (I - \widehat{\Sigma})^i \widehat{\gamma}, \end{aligned}$$

(assuming $\|\widehat{\Sigma}\|_{op} \leq 1$).

General form spectral linearization

Bauer, Rosasco, Pereverzev 2007

- ▶ **General form** regularization method:

$$\hat{f}_\lambda = \zeta_\lambda(\hat{\Sigma})\hat{\gamma}$$

for some well-chosen function $\zeta_\lambda : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ acting on the spectrum and “approximating” the function $x \mapsto x^{-1}$.

- ▶ $\lambda > 0$: regularization parameter; $\lambda \rightarrow 0 \Leftrightarrow$ less regularization
- ▶ Notation of (autoadjoint) functional calculus, i.e.

$$\hat{\Sigma} = \mathbf{Q}^T \text{diag}(\mu_1, \mu_2, \dots) \mathbf{Q} \Rightarrow \zeta(\hat{\Sigma}) := \mathbf{Q}^T \text{diag}(\zeta(\mu_1), \zeta(\mu_2), \dots) \mathbf{Q}$$

- ▶ Examples (revisited):

- ▶ **Tikhonov**: $\zeta_\lambda(t) = (t + \lambda)^{-1}$
- ▶ **Spectral cut-off**: $\zeta_\lambda(t) = t^{-1} \mathbf{1}\{t \geq \lambda\}$
- ▶ **Landweber iteration**: $\zeta_k(t) = \sum_{i=0}^k (1 - t)^i$.

Assumptions on regularization function

Standard assumptions on the regularization family $\zeta_\lambda : [0, 1] \rightarrow \mathbb{R}$ are:

(i) There exists a constant $D < \infty$ such that

$$\sup_{0 < \lambda \leq 1} \sup_{0 < t \leq 1} |t\zeta_\lambda(t)| \leq D,$$

(ii) There exists a constant $E < \infty$ such that

$$\sup_{0 < \lambda \leq 1} \sup_{0 < t \leq 1} \lambda |\zeta_\lambda(t)| \leq E,$$

(iii) **Qualification:** for **residual** $r_\lambda(t) := 1 - t\zeta_\lambda(t)$,

$$\forall \lambda \leq 1: \quad \sup_{0 < t \leq 1} |r_\lambda(t)| t^\nu \leq \gamma_\nu \lambda^\nu,$$

holds for $\nu = 0$ and $\nu = q > 0$.

Structural Assumptions (I)

- ▶ Denote $(\mu_i)_{i \geq 1}$ the sequence of positive eigenvalues of Σ in nonincreasing order.
- ▶ **Assumptions on spectrum decay:** for $s \in (0, 1); \alpha > 0$:

$$\mathbf{IP}^<(s, \alpha) : \mu_i \leq \alpha i^{-\frac{1}{s}}$$

Structural Assumptions (I)

- ▶ Denote $(\mu_i)_{i \geq 1}$ the sequence of positive eigenvalues of Σ in nonincreasing order.
- ▶ **Assumptions on spectrum decay:** for $s \in (0, 1); \alpha > 0$:

$$\mathbf{IP}^<(s, \alpha) : \mu_i \leq \alpha i^{-\frac{1}{s}}$$

- ▶ This implies quantitative estimates of the “effective dimension”

$$\mathcal{N}(\lambda) := \text{Tr}((\Sigma + \lambda)^{-1} \Sigma) \lesssim \lambda^{-s}.$$

Structural Assumptions (II)

- ▶ Denote $(\mu_i)_{i \geq 1}$ the sequence of positive eigenvalues of Σ in nonincreasing order.
- ▶ **Source condition** for the signal: for $r > 0$, define

$$\mathbf{SC}(r, R) : f_{\circ} = \Sigma^r h_{\circ} \text{ for some } h_{\circ} \text{ with } \|h_{\circ}\| \leq R,$$

or equivalently, as a **Sobolev-type regularity**

$$\mathbf{SC}(r, R) : f_{\circ} \in \left\{ f \in \mathcal{H} : \sum_{i \geq 1} \mu_i^{-2r} f_i^2 \leq R^2 \right\},$$

where f_i are the coefficients of h in the eigenbasis of Σ .

- ▶ Under $(\mathbf{SC})(r, R)$ it is assumed that the **qualification** q of the regularization method satisfies $q \geq r + \frac{1}{2}$.

A general upper bound risk estimate

Theorem

Assume the source condition **(SC)** (r, R) holds.

If λ is such that $\lambda \gtrsim (\mathcal{N}(\lambda) \vee \log(\eta)^2) / n$, then with probability at least $1 - \eta$, it holds:

$$\begin{aligned} & \left\| (\Sigma + \lambda)^{1/2} (f_{\circ} - \hat{f}_{\lambda}) \right\|_{\mathcal{H}} \\ & \lesssim \log(\eta)^2 \left(R\lambda^{r+\frac{1}{2}} + \sigma \sqrt{\frac{\mathcal{N}(\lambda)}{n}} + \frac{1}{n\sqrt{\lambda}} + \mathcal{O}(n^{-\frac{1}{2}}) \right). \end{aligned}$$

A general upper bound risk estimate

Theorem

Assume the source condition **(SC)** (r, R) holds.

If λ is such that $\lambda \gtrsim (\mathcal{N}(\lambda) \vee \log(\eta)^2) / n$, then with probability at least $1 - \eta$, it holds:

$$\begin{aligned} \left\| (\Sigma + \lambda)^{1/2} (f_{\circ} - \hat{f}_{\lambda}) \right\|_{\mathcal{H}} \\ \lesssim \log(\eta)^2 \left(R\lambda^{r+\frac{1}{2}} + \sigma \sqrt{\frac{\mathcal{N}(\lambda)}{n}} + \frac{1}{n\sqrt{\lambda}} + \mathcal{O}(n^{-\frac{1}{2}}) \right). \end{aligned}$$

This gives rise to estimates in both norms of interest since

$$\left\| f_{\circ} - \hat{f}_{\lambda} \right\|_{\mathcal{H}} \leq \lambda^{-\frac{1}{2}} \left\| (\Sigma + \lambda)^{1/2} (f_{\circ} - \hat{f}_{\lambda}) \right\|_{\mathcal{H}},$$

and

$$\left\| f_{\circ}^* - \hat{f}_{\lambda}^* \right\|_{L^2(P_X)} = \left\| \Sigma^{\frac{1}{2}} (f_{\circ} - \hat{f}_{\lambda}) \right\|_{\mathcal{H}} \leq \left\| (\Sigma + \lambda)^{1/2} (f_{\circ} - \hat{f}_{\lambda}) \right\|_{\mathcal{H}}.$$

Upper bound on rates

Optimizing the obtained bound over λ (i.e. balancing the main terms) one obtains

Theorem

Assume r, R, s, α are fixed positive constants and assume \mathbb{P}_{XY} satisfies $(\mathbf{IP}^<)(s, \alpha)$, $(\mathbf{SC})(r, R)$ and $\|X\| \leq 1, \|Y\| \leq M, \text{Var}[Y|X]_\infty \leq \sigma^2$ a.s. Define

$$\hat{\beta}_n = \zeta_{\lambda_n}(\hat{\Sigma})\hat{\gamma},$$

using a regularization family (ζ_λ) satisfying the standard assumptions with qualification $q \geq r + \frac{1}{2}$, and the parameter choice rule

$$\lambda_n = (R^2 \sigma^2 / n)^{-\frac{1}{2r+1+s}}.$$

Then it holds for any $p \geq 1$:

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{E}^{\otimes n} \left(\left\| f_\circ - \hat{f}_{\lambda_n} \right\|^p \right)^{1/p} / R \left(\frac{\sigma^2}{R^2 n} \right)^{\frac{r}{2r+1+s}} &\leq C_\blacktriangle; \\ \limsup_{n \rightarrow \infty} \mathbb{E}^{\otimes n} \left(\left\| f_\circ^* - \hat{f}_{\lambda_n} \right\|_{2,X}^p \right)^{1/p} / R \left(\frac{\sigma^2}{R^2 n} \right)^{\frac{r+1/2}{2r+1+s}} &\leq C_\blacktriangle. \end{aligned}$$

Extensions to nonlinear operators

- ▶ Extensions possible to **nonlinear** inverse problems
- ▶ Need stronger assumptions (output space of A has RKHS structure, A is Lipschitz continuous)
- ▶ For Tikhonov regularization: see Abishake R, Blanchard, Mathé 2020.

Towards adaptivity: existing approaches

- ▶ Cross-validation (or hold-out) will yield a tuning of the parameter which is **adaptive in the prediction risk**, it is based on a unbiased estimate of the risk (**URE**) principle.
- ▶ Standard Lepski's principle parameter selection can be applied for any fixed norm (provided a good estimate of the "variance" term $\sigma \sqrt{\mathcal{N}(\lambda)/n}$ is available)
- ▶ Despite the **existence** of a regularization parameter λ being optimal for both norms, there is no guarantee that **any** (close to) optimal parameter for prediction risk (eg. selected by cross-validation) will be close to optimal in reconstruction risk, or vice-versa.
- ▶ We want to construct a **simultaneously (for both norms) adaptive** data-driven parameter selection.

Standard Lepskii's principle

We consider the following “deterministic” assumption to highlight the construction.

Assumption

Let $\Lambda \subset \mathbb{R}_+$ be a finite set of candidate regularization parameters,

$$\Lambda := \{\lambda_j, \lambda_0 > \lambda_1 > \dots > \lambda_m = \lambda_{\min} > 0\},$$

The (known) family of elements of \mathcal{H} , $(f_\lambda)_{\lambda \in \Lambda}$, satisfies for any $\lambda \in \Lambda$:

$$\|f_0 - f_\lambda\|_{\mathcal{H}} \leq C(\mathcal{A}(\lambda) + \mathcal{S}(\lambda)),$$

where

- ▶ the function $\lambda \in \Lambda \mapsto \mathcal{A}(\lambda) \in \mathbb{R}_+$ is **non-decreasing** with $\mathcal{A}(0) = 0$ and possibly **unknown**;
- ▶ the function $\lambda \in \Lambda \mapsto \sqrt{\lambda} \mathcal{S}(\lambda) \in \mathbb{R}_+$ is **non-increasing** and **known**.

Standard Lepskii's principle (II)

Set

$$\hat{\lambda} := \max \left\{ \lambda \in \Lambda : \|f_\lambda - f_{\lambda'}\|_{\mathcal{H}} \leq 4\mathcal{C}\mathcal{S}(\lambda'), \forall \lambda' \in \Lambda, \text{ s.t. } \lambda' \leq \lambda \right\},$$

Theorem

Under the assumptions made previously, if

$$\lambda_* := \max\{\lambda \in \Lambda : \mathcal{A}(\lambda) \leq \mathcal{S}(\lambda)\},$$

▶ *then it holds:*

$$\|f_\circ - f_{\hat{\lambda}}\|_{\mathcal{H}} \lesssim \mathcal{S}(\lambda_*);$$

▶ *Assuming it holds $\mathcal{S}(\lambda_k) \leq \mathcal{C}_S \mathcal{S}(\lambda_{k-1})$ for $k = 1, \dots, m$, then:*

$$\|f_\circ - f_{\hat{\lambda}}\|_{\mathcal{H}} \lesssim \min_{\lambda \in \Lambda} (\mathcal{A}(\lambda) + \mathcal{S}(\lambda)).$$

Generalized Lepskii's principle

We consider the following “deterministic” assumption to highlight the construction.

Assumption

Let $\Lambda \subset \mathbb{R}_+$ be a finite set of candidate regularization parameters,

$$\Lambda := \{\lambda_j, \lambda_0 > \lambda_1 > \dots > \lambda_m = \lambda_{\min} > 0\},$$

The (known) family of elements of \mathcal{H} , $(f_\lambda)_{\lambda \in \Lambda}$, satisfies for any $\lambda \in \Lambda$:

$$\|(\Sigma + \lambda)^{1/2}(f_0 - f_\lambda)\|_{\mathcal{H}} \leq C\sqrt{\lambda}(\mathcal{A}(\lambda) + \mathcal{S}(\lambda)),$$

where

- ▶ the function $\lambda \in \Lambda \mapsto \mathcal{A}(\lambda) \in \mathbb{R}_+$ is **non-decreasing** with $\mathcal{A}(0) = 0$ and possibly **unknown**;
- ▶ the function $\lambda \in \Lambda \mapsto \sqrt{\lambda}\mathcal{S}(\lambda) \in \mathbb{R}_+$ is **non-increasing** and **known**.

Generalized Lepskii's principle (II)

► Set

$$\mathcal{M}(\Lambda) := \left\{ \lambda \in \Lambda : \left\| (\Sigma + \lambda')^{1/2} (f_\lambda - f_{\lambda'}) \right\|_{\mathcal{H}} \leq 4C\sqrt{\lambda'}\mathcal{S}(\lambda'), \right. \\ \left. \forall \lambda' \in \Lambda, \text{ s.t. } \lambda' \leq \lambda \right\}.$$

► The balancing parameter is given as

$$\hat{\lambda} := \max \mathcal{M}(\Lambda) ;$$

(this quantity is always well-defined since $\lambda_{\min} \in \mathcal{M}(\Lambda)$.)

Generalized Lepskii's principle: bound

Theorem

Under the assumptions made previously, if

$$\lambda_* := \max\{\lambda \in \Lambda : \mathcal{A}(\lambda) \leq \mathcal{S}(\lambda)\},$$

and $\hat{\lambda}$ is the parameter choice defined previously, then:

► It holds

$$\left\| (\Sigma + \lambda_*)^{\frac{1}{2}} (f_o - f_{\hat{\lambda}}) \right\|_{\mathcal{H}} \lesssim \sqrt{\lambda_*} \mathcal{S}(\lambda_*);$$

► Assuming it holds $\mathcal{S}(\lambda_k) \leq C_S \mathcal{S}(\lambda_{k-1})$ for $k = 1, \dots, m$, then:

$$\begin{aligned} \|f_o - f_{\hat{\lambda}}\|_{\mathcal{H}} &\lesssim \min_{\lambda \in \Lambda} (\mathcal{A}(\lambda) + \mathcal{S}(\lambda)); \\ \left\| \Sigma^{\frac{1}{2}} (f_o - f_{\hat{\lambda}}) \right\|_{\mathcal{H}} &\lesssim \min_{\lambda \in \Lambda} \sqrt{\lambda} (\mathcal{A}(\lambda) + \mathcal{S}(\lambda)). \end{aligned}$$

Applying Lepskii's principle

Looking at the main error bound obtained earlier, with high probability the assumption

$$\left\| (\Sigma + \lambda)^{1/2} (f_{\circ} - f_{\lambda}) \right\|_{\mathcal{H}} \leq C\sqrt{\lambda}(\mathcal{A}(\lambda) + \mathcal{S}(\lambda))$$

is satisfied with

$$\begin{aligned}\mathcal{A}(\lambda) &:= \left(R\lambda^r + \mathcal{O}(n^{-\frac{1}{2}}) \right), \\ \mathcal{S}(\lambda) &:= \frac{\sigma\sqrt{\mathcal{N}(\lambda)} + \mathcal{O}(1)}{\sqrt{\lambda n}}.\end{aligned}$$

Remaining issues:

- ▶ Σ is not known;
- ▶ $\mathcal{N}(\lambda) = \text{Tr}((\Sigma + \lambda)^{-1}\Sigma)$ is not known;
- ▶ the noise variance σ^2 might not be known (issue ignored for now).

Replacing $\Sigma, \mathcal{N}(\lambda)$ by empirical quantities

Proposition

If λ is such that $\lambda \gtrsim (\mathcal{N}(\lambda) \vee \log(\eta)^2) / n$, then with probability at least $1 - \eta$, it holds:

$$\left\| (\Sigma + \lambda)^{\frac{1}{2}} (\widehat{\Sigma} + \lambda)^{-\frac{1}{2}} \right\| \lesssim 1 + \log(\eta^{-1}).$$

Proposition

If $\lambda \gtrsim n^{-1}$, it holds with probability at least $1 - \eta$, for $\widehat{\mathcal{N}}(\lambda) := \text{Tr}(\widehat{\Sigma}(\widehat{\Sigma} + \lambda)^{-1})$:

$$\max \left(\frac{\mathcal{N}(\lambda) \vee 1}{\widehat{\mathcal{N}}(\lambda) \vee 1}, \frac{\widehat{\mathcal{N}}(\lambda) \vee 1}{\mathcal{N}(\lambda) \vee 1} \right) \lesssim (1 + \log \eta^{-1})^2.$$

Fully empirical procedure (σ, M known)

- ▶ Put $L := 2 \log(8 \log n / (\eta \log q))$ and let

$$\hat{\Lambda} := \left\{ \lambda_i = q^{-i}, i \in \mathbb{N}, \text{ s.t. } \lambda_i \geq 100(\hat{\mathcal{N}}(\lambda) \vee L^2/n) \right\}.$$

- ▶ Define the parameter choice

$$\hat{\lambda} = \max \left\{ \lambda \in \hat{\Lambda} : \forall \lambda' \in \hat{\Lambda}, \text{ s.t. } \lambda' \leq \lambda : \right. \\ \left. \left\| (\hat{\Sigma} + \lambda')^{\frac{1}{2}} (\hat{f}_\lambda - \hat{f}_{\lambda'}) \right\| \leq cL \sqrt{\lambda' \hat{S}(\lambda')} \right\},$$

where

$$\hat{S}(\lambda) := \frac{\sigma \sqrt{2(\hat{\mathcal{N}}(\lambda) \vee 1) + M/5}}{\sqrt{\lambda n}}.$$

Result for the empirical selection procedure

Theorem

Assume the source condition **(SC)** (r, R) holds.

Then for the generalized-Lepski parameter choice $\hat{\lambda}$, with probability at least $1 - \eta$:

$$\left\| (\Sigma + \lambda)^{\frac{1}{2}} (\hat{f}_{\hat{\lambda}} - f_{\circ}) \right\| \lesssim L^3 \min_{\lambda \in [\lambda_{\min}, 1]} \left(R\lambda^{r+\frac{1}{2}} + \sigma \sqrt{\frac{\mathcal{N}(\lambda)}{n}} + \frac{1}{n\sqrt{\lambda}} + \mathcal{O}(n^{-\frac{1}{2}}) \right).$$

where

$$\lambda_{\min} = \min \left\{ \lambda \in [0, 1] : \lambda \gtrsim (\mathcal{N}(\lambda) \vee L^2/n) \right\}.$$

Conclusion: as a direct byproduct we get the same rates (up to $\log \log n$ factor) as the optimal choice of λ in the original bound, for **both norms of interest**.

Estimating the unknown noise variance σ^2 ?

- ▶ Observe that in general, there is no identifiability in the model

$$y_i = f(x_i) + \sigma \xi_i,$$

if the function f can be “arbitrary”.

- ▶ There is a hope when we assumed that f has some regularity (here: linearity)

▶ Idea:

- ▶ Take λ small so that the “bias” $\mathcal{A}(\lambda)$ is expected to be much lower than the “variance” $\mathcal{S}(\lambda)$ (e.g., close to $\hat{\lambda}_{\min}$).
 - ▶ Split the sample into two subsamples giving rise to $\hat{f}_\lambda^{(1)}, \hat{f}_\lambda^{(2)}$.
 - ▶ The hope is that by considering $\left\| \hat{f}_\lambda^{(1)} - \hat{f}_\lambda^{(2)} \right\|^2$ in a suitable norm, we cancel the bias and observe twice the “variance”.
- ▶ Need somewhat precise concentration (upper and lower) for this quantity.

Estimation of the variance σ^2

- ▶ Assume we have two independent sample of the same size n , giving rise to estimators $\hat{f}_\lambda^{(1)}, \hat{f}_\lambda^{(2)}$ (using the same regularization parameter $\lambda > 0$).
- ▶ Consider the statistic

$$\begin{aligned}\Delta^2 &:= \left\| \frac{1}{2} (\hat{\Sigma}^{(1)} + \hat{\Sigma}^{(2)} + \lambda)^{\frac{1}{2}} (\hat{f}_\lambda^{(1)} - \hat{f}_\lambda^{(2)}) \right\|_{\mathcal{H}}^2 \\ &= \frac{1}{2n} \sum_{i=1}^{2n} (\hat{f}_\lambda^{(1)} - \hat{f}_\lambda^{(2)})^2(x_i) + \lambda \left\| \hat{f}_\lambda^{(1)} - \hat{f}_\lambda^{(2)} \right\|_{\mathcal{H}}^2,\end{aligned}$$

and

$$\hat{\sigma}^2 := \frac{\Delta^2}{\sum_{i,j=1}^2 \|A_{ij}\|_{HS}^2},$$

where $A_{ij} = (\hat{\Sigma}^{(i)} + \lambda)^{\frac{1}{2}} \zeta_\lambda(\hat{\Sigma}^{(j)}) (\hat{\Sigma}^{(j)})^{\frac{1}{2}}$.

Estimation of the variance σ^2

Theorem

If $\lambda \geq \hat{\lambda}_{min}$, where

$$\hat{\lambda}_{min} = \min \left\{ \lambda > 0 : \lambda \geq 100(\widehat{\mathcal{N}}(\lambda) \vee \log(\eta^{-1})/2) \right\},$$

then with probability at least $1 - \eta$, it holds

$$\hat{\sigma}^2 \in \left[\sigma^2 \pm \left(\lambda \sigma^2 + F(\lambda) \log(\eta^{-1}) \right) \right],$$

with $F(\lambda) \rightarrow 0$ as $\lambda \rightarrow 0$.

Conclusion: the estimator $\hat{\sigma}^2$ is consistent, and can be used as a proxy for σ^2 in the procedure, with the same conclusions (up to changes in numerical constants, and for n big enough).

Thank you for your attention

G. Blanchard, P. Mathé, N. Mücke. Lepskii Principle in Supervised Learning. ArXiv 1902.05404

G. Blanchard, N. Mücke. Optimal Rates For Regularization Of Statistical Inverse Learning Problems. Foundations of Computational Mathematics , 2017.

Abhishake Rastogi, G. Blanchard, P. Mathé. Convergence analysis of Tikhonov regularization for non-linear statistical inverse learning problems. Electronic Journal of Statistics 14 (2): 2798-2841, 2020.