

Nonparametric Regression by Projection on Non-compactly Supported Bases

Florian Dussap



53^e Journées de Statistique
Juin 2022

Regression model with random design

Let $A \subset \mathbb{R}^p$, we observe $n \geq 1$ r.v. $(\mathbf{X}_i, Y_i) \in A \times \mathbb{R}$ given by:

$$Y_i = b(\mathbf{X}_i) + \varepsilon_i,$$

where:

- (\mathbf{X}_i) are i.i.d. with unknown distribution μ .
- (ε_i) are i.i.d. with zero mean and known variance σ^2 .
- (\mathbf{X}_i) and (ε_i) are independent.

Our goal is to estimate the regression function $b: A \rightarrow \mathbb{R}$. To quantify the error of an estimator, we consider two norms:

$$\|t\|_n^2 := \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i)^2, \quad \|t\|_\mu^2 := \int_A t(\mathbf{x})^2 d\mu(\mathbf{x}).$$

The error relative to the norm $\|\cdot\|_\mu$ can be viewed as a prediction error:

$$\forall \hat{b} \text{ estimator, } \|b - \hat{b}\|_\mu^2 = \mathbb{E}_{\mathbf{X} \sim \mu} \left[(b(\mathbf{X}) - \hat{b}(\mathbf{X}))^2 \mid \mathbf{X}_1, \dots, \mathbf{X}_n \right].$$

Assumptions

- 1 We assume that $\mu \ll \nu$ for a fixed measure ν , and that $\frac{d\mu}{d\nu}$ is bounded on A . Hence, we have $L^2(A, \mu) \subset L^2(A, \nu)$.

- 2 If A is compact, we assume that:

$$\forall \mathbf{x} \in A, \quad \frac{d\mu}{d\nu}(\mathbf{x}) \geq f_0 > 0.$$

Hence, the norms $\|\cdot\|_\mu$ and $\|\cdot\|_\nu$ are equivalent, and we have $L^2(A, \mu) = L^2(A, \nu)$.

- 3 We assume that $b \in L^{2r}(A, \mu)$ for some $r \in (1, +\infty]$. We consider $r' \in [1, +\infty)$ such that $\frac{1}{r} + \frac{1}{r'} = 1$.
- 4 We assume that $A = A_1 \times \cdots \times A_p$ and that $\nu = \nu_1 \otimes \cdots \otimes \nu_p$.

Projection estimator

Let $(\varphi_k^i)_{k \in \mathbb{N}}$ be an orthonormal basis of $L^2(A_i, \nu_i)$. We construct a basis of $L^2(A, \nu)$ by tensorization. For all $\mathbf{k} = (k_1, \dots, k_p) \in \mathbb{N}^p$ we define:

$$\varphi_{\mathbf{k}}(\mathbf{x}) := (\varphi_{k_1}^1 \otimes \dots \otimes \varphi_{k_p}^p)(\mathbf{x}) := \varphi_{k_1}^1(x_1) \times \dots \times \varphi_{k_p}^p(x_p).$$

For $\mathbf{m} \in \mathbb{N}_+^p$, we consider the model:

$$S_{\mathbf{m}} := \text{Span}(\varphi_{\mathbf{k}} : \forall i, 0 \leq k_i < m_i), \quad D_{\mathbf{m}} := \dim(S_{\mathbf{m}}) = m_1 \cdots m_p,$$

and we estimate b by a least squares minimization on $S_{\mathbf{m}}$:

$$\hat{b}_{\mathbf{m}} := \arg \min_{t \in S_{\mathbf{m}}} \frac{1}{n} \sum_{i=1}^n [Y_i - t(\mathbf{X}_i)]^2.$$

Example

- 1 For $A = [-\pi, \pi]$ and $\nu = \text{Leb}$, we choose the trigonometric basis.
- 2 For $A = \mathbb{R}$ and $\nu = \text{Leb}$, we choose $\varphi_k(x) = c_k H_k(x) e^{-x^2/2}$ with H_k the k -th Hermite polynomial.

This estimator can be computed using matrix calculus. Let $(\phi_1, \dots, \phi_{D_m})$ be an orthonormal basis of S_m for the inner product $\langle \cdot, \cdot \rangle_\nu$, we have:

$$\begin{aligned} \hat{b}_m &= \sum_{j=1}^{D_m} \hat{a}_j^{(m)} \phi_j, & \hat{\mathbf{a}}^{(m)} &:= \arg \min_{\mathbf{a} \in \mathbb{R}^{D_m}} \left\| \mathbf{Y} - \hat{\Phi}_m \mathbf{a} \right\|_{\mathbb{R}^n}^2 \\ & & &= \hat{\mathbf{G}}_m^{-1} \hat{\Phi}_m^* \mathbf{Y}, \end{aligned}$$

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

$$\hat{\mathbf{G}}_m := \left[\langle \phi_j, \phi_k \rangle_n \right]_{j,k} \in \mathbb{R}^{D_m \times D_m}, \quad \hat{\Phi}_m := \left[\phi_j(\mathbf{X}_i) \right]_{i,j} \in \mathbb{R}^{n \times D_m}.$$

In the following, we also consider the expectation of $\hat{\mathbf{G}}_m$:

$$\mathbf{G}_m := \mathbb{E}[\hat{\mathbf{G}}_m] = \left[\langle \phi_j, \phi_k \rangle_\mu \right]_{j,k} \in \mathbb{R}^{D_m \times D_m}.$$

Basic bound on the empirical risk

We recall the classical bias-variance decomposition of the empirical risk.

Proposition

If $\hat{\mathbf{G}}_m$ is invertible, then we have:

$$\begin{aligned}\mathbb{E}_{\mathbf{X}} \left[\|b - \hat{b}_m\|_n^2 \right] &:= \mathbb{E} \left[\|b - \hat{b}_m\|_n^2 \mid \mathbf{X}_1, \dots, \mathbf{X}_n \right] \\ &= \inf_{t \in S_m} \|b - t\|_n^2 + \sigma^2 \frac{D_m}{n}.\end{aligned}$$

If $\hat{\mathbf{G}}_m$ is invertible a.s., then we have:

$$\mathbb{E} \|b - \hat{b}_m\|_n^2 \leq \inf_{t \in S_m} \|b - t\|_\mu^2 + \sigma^2 \frac{D_m}{n}.$$

From the empirical norm to the design norm

We introduce the event:

$$\Omega_{\mathbf{m}}(\delta) := \left\{ \sup_{t \in S_{\mathbf{m}} \setminus \{0\}} \frac{\|t\|_{\mu}^2}{\|t\|_n^2} \leq \frac{1}{1 - \delta} \right\}, \quad \delta \in (0, 1).$$

Using matrix concentration inequalities from [Tropp, 2012], the following bound holds.

Lemma

For all $\delta \in (0, 1)$ and all $\mathbf{m} \in \mathbb{N}_+^p$, we have:

$$\mathbb{P}[\Omega_{\mathbf{m}}(\delta)^c] \leq D_{\mathbf{m}} \exp\left(-h(\delta) \frac{n}{L(\mathbf{m}) \|\mathbf{G}_{\mathbf{m}}^{-1}\|_{\text{op}}}\right),$$

where $h(\delta) := (1 - \delta) \log(1 - \delta) + \delta$, and where:

$$L(\mathbf{m}) := \left\| \sum_{k \leq m-1} \varphi_k^2 \right\|_{\infty} = \sup_{t \in S_{\mathbf{m}} \setminus \{0\}} \frac{\|t\|_{\infty}^2}{\|t\|_{\nu}^2}.$$

Remarks on the lemma

- For the trigonometric basis, we have $L(m) \leq m$.
- For the Hermite basis, we have $L(m) \leq C\sqrt{m}$.
- If A is compact, then we have $\|\mathbf{G}_m^{-1}\|_{\text{op}} \leq 1/f_0$.
- If $A = \mathbb{R}$ and $(\varphi_k)_{k \in \mathbb{N}}$ is the Hermite basis, then we have $\|\mathbf{G}_m^{-1}\|_{\text{op}} \geq C(\mu)\sqrt{m}$ [Comte and Genon-Catalot, 2020].

Bound on the prediction risk

Let us consider the collection:

$$\mathcal{M}_{n,\alpha} := \left\{ \mathbf{m} \in \mathbb{N}_+^p \mid L(\mathbf{m})(\|\mathbf{G}_{\mathbf{m}}^{-1}\|_{\text{op}} \vee 1) \leq \alpha \frac{n}{\log n} \right\}.$$

If $\mathbf{m} \in \mathcal{M}_{n,\alpha}$, then we have $\mathbb{P}[\Omega_{\mathbf{m}}(\delta)^c] \leq D_{\mathbf{m}} n^{-\alpha} \leq n^{-\alpha+1}$.

Theorem

For all $\alpha \in (0, \frac{1}{2r'+1})$ and for all $\mathbf{m} \in \mathcal{M}_{n,\alpha}$ we have:

$$\mathbb{E} \|b - \hat{b}_{\mathbf{m}}\|_{\mu}^2 \leq C_n(\alpha, r') \inf_{t \in S_{\mathbf{m}}} \|b - t\|_{\mu}^2 + C'(\alpha, r') \sigma^2 \frac{D_{\mathbf{m}}}{n} + o\left(\frac{1}{n}\right).$$

A model selection result in a fixed design setting

Let $\widehat{\mathcal{M}}_n$ a model collection that may depend on the (\mathbf{X}_i) , and let:

$$\hat{\mathbf{m}} := \arg \min_{\mathbf{m} \in \widehat{\mathcal{M}}_n} \left(-\|\hat{\mathbf{b}}_{\mathbf{m}}\|_n^2 + \text{pen}(\mathbf{m}) \right), \quad \text{pen}(\mathbf{m}) := (1 + \theta) \sigma^2 \frac{D_{\mathbf{m}}}{n}.$$

Theorem ([Baraud, 2000])

If $\mathbb{E}|\varepsilon_1|^q$ is finite for some $q > 4$, then the following upper bound holds:

$$\mathbb{E}_{\mathbf{X}} \|b - \hat{\mathbf{b}}_{\hat{\mathbf{m}}}\|_n^2 \leq C(\theta) \inf_{\mathbf{m} \in \widehat{\mathcal{M}}_n} \left(\inf_{t \in S_{\mathbf{m}}} \|b - t\|_n^2 + \sigma^2 \frac{D_{\mathbf{m}}}{n} \right) + \sigma^2 \frac{\Sigma_n(\theta, q)}{n},$$

with $\Sigma_n(\theta, q) := C'(\theta, q) \frac{\mathbb{E}|\varepsilon_1|^q}{\sigma^q} \sum_{\mathbf{m} \in \widehat{\mathcal{M}}_n} D_{\mathbf{m}}^{-\left(\frac{q}{2}-2\right)}$.

Oracle bound for the empirical risk

We choose the model collection:

$$\widehat{\mathcal{M}}_{n,\beta} := \left\{ \mathbf{m} \in \mathbb{N}_+^p \mid L(\mathbf{m}) (\|\widehat{\mathbf{G}}_{\mathbf{m}}^{-1}\|_{\text{op}} \vee 1) \leq \beta \frac{n}{\log n} \right\}.$$

Theorem

If $\mathbb{E}|\varepsilon_1|^q$ is finite for some $q > 6$, then there exists a constant $\alpha_{\beta,r'} > 0$ such that for all $\alpha \in (0, \alpha_{\beta,r'})$, we have:

$$\mathbb{E} \|b - \hat{b}_{\hat{\mathbf{m}}}\|_n^2 \leq C(\theta) \inf_{\mathbf{m} \in \mathcal{M}_{n,\alpha}} \left(\inf_{t \in S_{\mathbf{m}}} \|b - t\|_{\mu}^2 + \sigma^2 \frac{D_{\mathbf{m}}}{n} \right) + \sigma^2 \frac{\Sigma(\theta, q)}{n} + o\left(\frac{1}{n}\right),$$

$$\text{with } \Sigma(\theta, q) := C'(\theta, q) \frac{\mathbb{E}|\varepsilon_1|^q}{\sigma^q} \sum_{\mathbf{m} \in \mathbb{N}_+^p} D_{\mathbf{m}}^{-(\frac{q}{2}-2)}.$$

Oracle bound for the prediction risk

Theorem

If A is compact:

If $\mathbb{E}|\varepsilon_1|^q$ is finite for some $q > 6$, then there exists $\beta^* > 0$ such that for all $\beta \in (0, \beta^*)$, there exists $\alpha_{\beta, r'} > 0$ such that for all $\alpha \in (0, \alpha_{\beta, r'})$, we have:

$$\begin{aligned} \mathbb{E} \|b - \hat{b}_{\hat{m}}\|_{\mu}^2 &\leq C(\theta, \beta, r) \inf_{m \in \mathcal{M}_{n, \alpha}} \left(\inf_{t \in S_m} \|b - t\|_{\mu}^2 + \sigma^2 \frac{D_m}{n} \right) \\ &\quad + C'(\beta, r) \sigma^2 \frac{\Sigma(\theta, q)}{n} + o\left(\frac{1}{n}\right), \end{aligned}$$

with:

$$\begin{aligned} \mathcal{M}_{n, \alpha} &:= \left\{ m \in \mathbb{N}_+^p \mid L(m) \left(\|\mathbf{G}_m^{-1}\|_{\text{op}} \vee 1 \right) \leq \alpha \frac{n}{\log n} \right\}, \\ \widehat{\mathcal{M}}_{n, \beta} &:= \left\{ m \in \mathbb{N}_+^p \mid L(m) \left(\|\widehat{\mathbf{G}}_m^{-1}\|_{\text{op}} \vee 1 \right) \leq \beta \frac{n}{\log n} \right\}. \end{aligned}$$

Oracle bound for the prediction risk

Theorem

If A is *not* compact:

If $\mathbb{E}|\varepsilon_1|^q$ is finite for some $q > 6$, then there exists $\beta^* > 0$ such that for all $\beta \in (0, \beta^*)$, there exists $\alpha_{\beta, r'} > 0$ such that for all $\alpha \in (0, \alpha_{\beta, r'})$, we have:

$$\begin{aligned} \mathbb{E}\|b - \hat{b}_{\hat{m}}\|_{\mu}^2 &\leq C(\theta, \beta, r) \inf_{m \in \mathcal{M}_{n, \alpha}} \left(\inf_{t \in S_m} \|b - t\|_{\mu}^2 + \sigma^2 \frac{D_m}{n} \right) \\ &\quad + C'(\beta, r) \sigma^2 \frac{\Sigma(\theta, q)}{n} + o\left(\frac{1}{n}\right), \end{aligned}$$

with:

$$\begin{aligned} \mathcal{M}_{n, \alpha} &:= \left\{ m \in \mathbb{N}_+^p \mid L(m) \left(\|\mathbf{G}_m^{-1}\|_{\text{op}}^2 \vee 1 \right) \leq \alpha \frac{n}{\log n} \right\}, \\ \widehat{\mathcal{M}}_{n, \beta} &:= \left\{ m \in \mathbb{N}_+^p \mid L(m) \left(\|\widehat{\mathbf{G}}_m^{-1}\|_{\text{op}}^2 \vee 1 \right) \leq \beta \frac{n}{\log n} \right\}. \end{aligned}$$

Conclusion and perspective

- We obtain bounds for the empirical risk from the results for fixed design regression.
- To obtain a bound on the prediction risk, we need to study the minimum eigenvalue of a random matrix. We do so by using concentration inequalities of [Gittens and Tropp, 2011] and [Tropp, 2012].
- We improve the results of [Baraud, 2002] and [Comte and Genon-Catalot, 2020].
- I think that these results can be extended to more general approximation spaces $(\mathcal{S}_m)_{m \in \mathcal{M}_n}$, that are not constructed from an orthonormal basis.

References



Baraud, Y. (2000).
Model selection for regression on a fixed design.
Probability Theory and Related Fields, 117(4):467–493.



Baraud, Y. (2002).
Model selection for regression on a random design.
ESAIM: Probability and Statistics, 6:127–146.



Cohen, A., Davenport, M. A., and Leviatan, D. (2013).
On the Stability and Accuracy of Least Squares Approximations.
Foundations of Computational Mathematics, 13(5):819–834.



Comte, F. and Genon-Catalot, V. (2020).
Regression function estimation as a partly inverse problem.
Annals of the Institute of Statistical Mathematics, 72(4):1023–1054.



Gittens, A. and Tropp, J. A. (2011).
Tail bounds for all eigenvalues of a sum of random matrices.
arXiv:1104.4513 [math].



Tropp, J. A. (2012).
User-Friendly Tail Bounds for Sums of Random Matrices.
Foundations of Computational Mathematics, 12(4):389–434.

Sketch of the proof of the lemma

The proof is inspired by [Cohen et al., 2013]. Let $(\phi_1, \dots, \phi_{D_m})$ be an orthonormal of S_m for the inner product $\langle \cdot, \cdot \rangle_\mu$, and let \mathbf{H}_m be their Gram matrix relative to the empirical inner product, that is:

$$\mathbf{H}_m := [\langle \phi_j, \phi_k \rangle_n]_{j,k} \in \mathbb{R}^{D_m \times D_m}.$$

Then, we have:

$$\sup_{t \in S_m \setminus \{0\}} \frac{\|t\|_\mu^2}{\|t\|_n^2} = \|\mathbf{H}_m^{-1}\|_{\text{op}} = \frac{1}{\lambda_{\min}(\mathbf{H}_m)}.$$

Hence, we can rewrite the event as:

$$\Omega_m(\delta)^c = \{\lambda_{\min}(\mathbf{H}_m) < 1 - \delta\} = \{\lambda_{\min}(\mathbf{H}_m) < (1 - \delta)\lambda_{\min}(\mathbb{E}[\mathbf{H}_m])\},$$

since $\mathbb{E}[\mathbf{H}_m] = \mathbf{Id}_{D_m}$.

We conclude using the following concentration inequality.

Theorem ([Gittens and Tropp, 2011], [Tropp, 2012])

Let $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ be independent random self-adjoint positive semi-definite matrices with dimension d , such that $\sup_k \lambda_{\max}(\mathbf{Z}_k) \leq R$ a.s. If we define:

$$\mu_{\min} := \lambda_{\min} \left(\sum_{k=1}^n \mathbb{E}[\mathbf{Z}_k] \right),$$

then we have:

$$\mathbb{P} \left[\lambda_{\min} \left(\sum_{k=1}^n \mathbf{Z}_k \right) \leq (1 - \delta) \mu_{\min} \right] \leq d \times \left(\frac{e^{-\delta}}{(1 - \delta)^{(1-\delta)}} \right)^{\mu_{\min}/R},$$
$$\mathbb{P} \left[\lambda_{\min} \left(\sum_{k=1}^n \mathbf{Z}_k \right) \geq (1 + \delta) \mu_{\min} \right] \leq \left(\frac{e^{\delta}}{(1 + \delta)^{(1+\delta)}} \right)^{\mu_{\min}/R}.$$