



## Supplement to "Robust linear least squares regression"

Jean-Yves Audibert, Olivier Catoni

### ► To cite this version:

Jean-Yves Audibert, Olivier Catoni. Supplement to "Robust linear least squares regression". Annals of Statistics, Institute of Mathematical Statistics, 2011, 39 (5), 19 p. <10.1214/11-AOS918SUPP>. <hal-00624459>

**HAL Id: hal-00624459**

**<https://hal.archives-ouvertes.fr/hal-00624459>**

Submitted on 17 Sep 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# SUPPLEMENT TO “ROBUST LINEAR LEAST SQUARES REGRESSION”

BY JEAN-YVES AUDIBERT<sup>\*,†</sup>, AND OLIVIER CATONI<sup>‡,§</sup>

This supplementary material provides the proofs of Theorems 2.1, 2.2 and 3.1 of the article “Robust linear least squares regression”.

## CONTENTS

1	Proofs of Theorems 2.1 and 2.2 . . . . .	1
1.1	Proof of Theorem 2.1 . . . . .	8
1.2	Proof of Theorem 2.2 . . . . .	9
2	Proof of Theorem 3.1 . . . . .	11
	References . . . . .	19

**1. Proofs of Theorems 2.1 and 2.2.** To shorten the formulae, we will write  $X$  for  $\varphi(X)$ , which is equivalent to considering without loss of generality that the input space is  $\mathbb{R}^d$  and that the functions  $\varphi_1, \dots, \varphi_d$  are the coordinate functions. Therefore, the function  $f_\theta$  maps an input  $x$  to  $\langle \theta, x \rangle$ . With a slight abuse of notation,  $R(\theta)$  will denote the risk of this prediction function.

Let us first assume that the matrix  $Q_\lambda = Q + \lambda I$  is positive definite. This indeed does not restrict the generality of our study, even in the case when  $\lambda = 0$ , as we will discuss later (Remark 1.1).

Consider the change of coordinates

$$\bar{X} = Q_\lambda^{-1/2} X.$$

Let us introduce

$$\bar{R}(\theta) = \mathbb{E}[(\langle \theta, \bar{X} \rangle - Y)^2],$$

---

\*Université Paris-Est, LIGM, Imagine, 6 avenue Blaise Pascal, 77455 Marne-la-Vallée, France, E-mail: [audibert@imagine.enpc.fr](mailto:audibert@imagine.enpc.fr)

†CNRS/École Normale Supérieure/INRIA, LIENS, Sierra – UMR 8548, 23 avenue d’Italie, 75214 Paris cedex 13, France.

‡École Normale Supérieure, CNRS – UMR 8553, Département de Mathématiques et Applications, 45 rue d’Ulm, 75230 Paris cedex 05, France, E-mail: [olivier.catoni@ens.fr](mailto:olivier.catoni@ens.fr)

§ INRIA Paris-Rocquencourt - CLASSIC team.

*AMS 2000 subject classifications:* 62J05, 62J07

*Keywords and phrases:* Linear regression, Generalization error, Shrinkage, PAC-Bayesian theorems, Risk bounds, Robust statistics, Resistant estimators, Gibbs posterior distributions, Randomized estimators, Statistical learning theory

so that

$$\overline{R}(Q_\lambda^{1/2}\theta) = R(\theta) = \mathbb{E}[(\langle \theta, X \rangle - Y)^2].$$

Let

$$\overline{\Theta} = \{Q_\lambda^{1/2}\theta; \theta \in \Theta\}.$$

Consider

$$(1.1) \quad r(\theta) = \frac{1}{n} \sum_{i=1}^n (\langle \theta, X_i \rangle - Y_i)^2,$$

$$(1.2) \quad \overline{r}(\theta) = \frac{1}{n} \sum_{i=1}^n (\langle \theta, \overline{X}_i \rangle - Y_i)^2,$$

$$(1.3) \quad \theta_0 = \arg \min_{\theta \in \overline{\Theta}} \overline{R}(\theta) + \lambda \|Q_\lambda^{-1/2}\theta\|^2,$$

$$(1.4) \quad \hat{\theta} \in \arg \min_{\theta \in \Theta} r(\theta) + \lambda \|\theta\|^2,$$

$$(1.5) \quad \theta_1 = Q_\lambda^{1/2}\hat{\theta} \in \arg \min_{\theta \in \overline{\Theta}} \overline{r}(\theta) + \lambda \|Q_\lambda^{-1/2}\theta\|^2.$$

For  $\alpha > 0$ , let us introduce the notation

$$W_i(\theta) = \alpha \left\{ (\langle \theta, \overline{X}_i \rangle - Y_i)^2 - (\langle \theta_0, \overline{X}_i \rangle - Y_i)^2 \right\},$$

$$W(\theta) = \alpha \left\{ (\langle \theta, \overline{X} \rangle - Y)^2 - (\langle \theta_0, \overline{X} \rangle - Y)^2 \right\}.$$

For any  $\theta_2 \in \mathbb{R}^d$  and  $\beta > 0$ , let us consider the Gaussian distribution centered at  $\theta_2$

$$\rho_{\theta_2}(d\theta) = \left(\frac{\beta}{2\pi}\right)^{d/2} \exp\left(-\frac{\beta}{2}\|\theta - \theta_2\|^2\right) d\theta.$$

LEMMA 1.1. *For any  $\eta > 0$  and  $\alpha > 0$ , with probability at least  $1 - \exp(-\eta)$ , for any  $\theta_2 \in \mathbb{R}^d$ ,*

$$\begin{aligned} & -n \int \log \left\{ 1 - \mathbb{E}[W(\theta)] + \mathbb{E}[W(\theta)^2]/2 \right\} \rho_{\theta_2}(d\theta) \\ & \leq -\sum_{i=1}^n \int \log \left\{ 1 - W_i(\theta) + W_i(\theta)^2/2 \right\} \rho_{\theta_2}(d\theta) + \mathcal{K}(\rho_{\theta_2}, \rho_{\theta_0}) + \eta, \end{aligned}$$

where  $\mathcal{K}(\rho_{\theta_2}, \rho_{\theta_0})$  is the Kullback-Leibler divergence function :

$$\mathcal{K}(\rho_{\theta_2}, \rho_{\theta_0}) = \int \log \left[ \frac{d\rho_{\theta_2}}{d\rho_{\theta_0}}(\theta) \right] \rho_{\theta_2}(d\theta).$$

PROOF. Since

$$\mathbb{E} \left( \int \rho_{\theta_0}(d\theta) \prod_{i=1}^n \frac{1 - W_i(\theta) + W_i(\theta)^2/2}{1 - \mathbb{E}[W(\theta)] + \mathbb{E}[W(\theta)^2]/2} \right) \leq 1,$$

with probability at least  $1 - \exp(-\eta)$

$$\log \left( \int \rho_{\theta_0}(d\theta) \prod_{i=1}^n \frac{1 - W_i(\theta) + W_i(\theta)^2/2}{1 - \mathbb{E}[W(\theta)] + \mathbb{E}[W(\theta)^2]/2} \right) \leq \eta.$$

We conclude the proof using the convex inequality (see [2], [3, Proposition 1.4.2] or [1, page 159])

$$\log \left( \int \rho_{\theta_0}(d\theta) \exp[h(\theta)] \right) \geq \int \rho_{\theta_2}(d\theta) h(\theta) - \mathcal{K}(\rho_{\theta_2}, \rho_{\theta_0}).$$

□

Let us compute some useful quantities

$$(1.6) \quad \mathcal{K}(\rho_{\theta_2}, \rho_{\theta_0}) = \frac{\beta}{2} \|\theta_2 - \theta_0\|^2,$$

$$\int \rho_{\theta_2}(d\theta) [W(\theta)] = \alpha \int \rho_{\theta_2}(d\theta) \langle \theta - \theta_2, \bar{X} \rangle^2 + W(\theta_2)$$

$$(1.7) \quad = W(\theta_2) + \frac{\alpha \|\bar{X}\|^2}{\beta},$$

$$(1.8) \quad \int \rho_{\theta_2}(d\theta) \langle \theta - \theta_2, \bar{X} \rangle^4 = \frac{3 \|\bar{X}\|^4}{\beta^2},$$

$$(1.9) \quad \begin{aligned} \int \rho_{\theta_2}(d\theta) [W(\theta)^2] &= \alpha^2 \int \rho_{\theta_2}(d\theta) \langle \theta - \theta_0, \bar{X} \rangle^2 (\langle \theta + \theta_0, \bar{X} \rangle - 2Y)^2 \\ &= \alpha^2 \int \rho_{\theta_2}(d\theta) \left[ \langle \theta - \theta_2 + \theta_2 - \theta_0, \bar{X} \rangle (\langle \theta - \theta_2 + \theta_2 + \theta_0, \bar{X} \rangle - 2Y) \right]^2 \\ &= \int \rho_{\theta_2}(d\theta) \left[ \alpha \langle \theta - \theta_2, \bar{X} \rangle^2 + 2\alpha \langle \theta - \theta_2, \bar{X} \rangle (\langle \theta_2, \bar{X} \rangle - Y) + W(\theta_2) \right]^2 \\ &= \int \rho_{\theta_2}(d\theta) \left[ \alpha^2 \langle \theta - \theta_2, \bar{X} \rangle^4 + 4\alpha^2 \langle \theta - \theta_2, \bar{X} \rangle^2 (\langle \theta_2, \bar{X} \rangle - Y)^2 + W(\theta_2)^2 \right. \\ &\quad \left. + 2\alpha \langle \theta - \theta_2, \bar{X} \rangle^2 W(\theta_2) \right] \\ &= \frac{3\alpha^2 \|\bar{X}\|^4}{\beta^2} + \frac{2\alpha \|\bar{X}\|^2}{\beta} \left[ 2\alpha (\langle \theta_2, \bar{X} \rangle - Y)^2 + W(\theta_2) \right] + W(\theta_2)^2. \end{aligned}$$

Using the fact that

$$2\alpha (\langle \theta_2, \bar{X} \rangle - Y)^2 + W(\theta_2) = 2\alpha (\langle \theta_0, \bar{X} \rangle - Y)^2 + 3W(\theta_2),$$

and that for any real numbers  $a$  and  $b$ ,  $6ab \leq 9a^2 + b^2$ , we get

LEMMA 1.2.

(1.10)

$$\int \rho_{\theta_2}(d\theta) [W(\theta)] = W(\theta_2) + \frac{\alpha \|\bar{X}\|^2}{\beta},$$

$$\int \rho_{\theta_2}(d\theta) [W(\theta)^2] = W(\theta_2)^2 + \frac{2\alpha \|\bar{X}\|^2}{\beta} \left[ 2\alpha (\langle \theta_0, \bar{X} \rangle - Y)^2 + 3W(\theta_2) \right]$$

(1.11)

$$+ \frac{3\alpha^2 \|\bar{X}\|^4}{\beta^2}$$

$$(1.12) \quad \leq 10W(\theta_2)^2 + \frac{4\alpha^2 \|\bar{X}\|^2}{\beta} (\langle \theta_0, \bar{X} \rangle - Y)^2 + \frac{4\alpha^2 \|\bar{X}\|^4}{\beta^2},$$

and the same holds true when  $W$  is replaced with  $W_i$  and  $(\bar{X}, Y)$  with  $(\bar{X}_i, Y_i)$ .

Another important thing to realize is that

$$(1.13) \quad \begin{aligned} \mathbb{E}[\|\bar{X}\|^2] &= \mathbb{E}[\text{Tr}(\bar{X}\bar{X}^T)] &&= \mathbb{E}[\text{Tr}(Q_\lambda^{-1/2} X X^T Q_\lambda^{-1/2})] \\ &= \mathbb{E}[\text{Tr}(Q_\lambda^{-1} X X^T)] &&= \text{Tr}[Q_\lambda^{-1} \mathbb{E}(X X^T)] \\ &= \text{Tr}(Q_\lambda^{-1}(Q_\lambda - \lambda I)) &&= d - \lambda \text{Tr}(Q_\lambda^{-1}) = D. \end{aligned}$$

We can weaken Lemma 1.1 (page 2) noticing that for any real number  $x$ ,

$$\begin{aligned} x - \frac{x^2}{2} &\leq -\log\left(1 - x + \frac{x^2}{2}\right) = \log\left(\frac{1 + x + x^2/2}{1 + x^4/4}\right) \\ &\leq \log\left(1 + x + \frac{x^2}{2}\right) \leq x + \frac{x^2}{2}. \end{aligned}$$

We obtain with probability at least  $1 - \exp(-\eta)$

$$\begin{aligned} &n\mathbb{E}[W(\theta_2)] + \frac{n\alpha}{\beta} \mathbb{E}[\|\bar{X}\|^2] - 5n\mathbb{E}[W(\theta_2)^2] \\ &\quad - \mathbb{E}\left\{ \frac{2n\alpha^2 \|\bar{X}\|^2}{\beta} (\langle \theta_0, \bar{X} \rangle - Y)^2 + \frac{2n\alpha^2 \|\bar{X}\|^4}{\beta^2} \right\} \\ &\quad \leq \sum_{i=1}^n \left\{ W_i(\theta_2) + 5W_i(\theta_2)^2 \right. \\ &\quad \left. + \frac{\alpha \|\bar{X}_i\|^2}{\beta} + \frac{2\alpha^2 \|\bar{X}_i\|^2}{\beta} (\langle \theta_0, \bar{X}_i \rangle - Y)^2 + \frac{2\alpha^2 \|\bar{X}_i\|^4}{\beta^2} \right\} \end{aligned}$$

$$+ \frac{\beta}{2} \|\theta_2 - \theta_0\|^2 + \eta.$$

Noticing that for any real numbers  $a$  and  $b$ ,  $4ab \leq a^2 + 4b^2$ , we can then bound

$$\begin{aligned} \alpha^{-2} W(\theta_2)^2 &= \langle \theta_2 - \theta_0, \bar{X} \rangle^2 (\langle \theta_2 + \theta_0, \bar{X} \rangle - 2Y)^2 \\ &= \langle \theta_2 - \theta_0, \bar{X} \rangle^2 \left[ \langle \theta_2 - \theta_0, \bar{X} \rangle + 2(\langle \theta_0, \bar{X} \rangle - Y) \right]^2 \\ &= \langle \theta_2 - \theta_0, \bar{X} \rangle^4 + 4\langle \theta_2 - \theta_0, \bar{X} \rangle^3 (\langle \theta_0, \bar{X} \rangle - Y) \\ &\quad + 4\langle \theta_2 - \theta_0, \bar{X} \rangle^2 (\langle \theta_0, \bar{X} \rangle - Y)^2 \\ &\leq 2\langle \theta_2 - \theta_0, \bar{X} \rangle^4 + 8\langle \theta_2 - \theta_0, \bar{X} \rangle^2 (\langle \theta_0, \bar{X} \rangle - Y)^2. \end{aligned}$$

**THEOREM 1.3.** *Let us put*

$$\hat{D} = \frac{1}{n} \sum_{i=1}^n \|\bar{X}_i\|^2 \quad (\text{let us remind that } D = \mathbb{E}[\|\bar{X}\|^2] \text{ from (1.13)}),$$

$$B_1 = 2\mathbb{E} \left[ \|\bar{X}\|^2 (\langle \theta_0, \bar{X} \rangle - Y)^2 \right],$$

$$\hat{B}_1 = \frac{2}{n} \sum_{i=1}^n \left[ \|\bar{X}_i\|^2 (\langle \theta_0, \bar{X}_i \rangle - Y_i)^2 \right],$$

$$B_2 = 2\mathbb{E} \left[ \|\bar{X}\|^4 \right],$$

$$\hat{B}_2 = \frac{2}{n} \sum_{i=1}^n \|\bar{X}_i\|^4,$$

$$B_3 = 40 \sup \left\{ \mathbb{E} \left[ \langle u, \bar{X} \rangle^2 (\langle \theta_0, \bar{X} \rangle - Y)^2 \right] : u \in \mathbb{R}^d, \|u\| = 1 \right\},$$

$$\hat{B}_3 = \sup \left\{ \frac{40}{n} \sum_{i=1}^n \langle u, \bar{X}_i \rangle^2 (\langle \theta_0, \bar{X}_i \rangle - Y_i)^2 : u \in \mathbb{R}^d, \|u\| = 1 \right\},$$

$$B_4 = 10 \sup \left\{ \mathbb{E} \left[ \langle u, \bar{X} \rangle^4 \right] : u \in \mathbb{R}^d, \|u\| = 1 \right\},$$

$$\hat{B}_4 = \sup \left\{ \frac{10}{n} \sum_{i=1}^n \langle u, \bar{X}_i \rangle^4 : u \in \mathbb{R}^d, \|u\| = 1 \right\}.$$

With probability at least  $1 - \exp(-\eta)$ , for any  $\theta_2 \in \mathbb{R}^d$ ,

$$\begin{aligned} n\mathbb{E}[W(\theta_2)] &- \left[ n\alpha^2(B_3 + \hat{B}_3) + \frac{\beta}{2} \right] \|\theta_2 - \theta_0\|^2 \\ &- n\alpha^2(B_4 + \hat{B}_4) \|\theta_2 - \theta_0\|^4 \end{aligned}$$

$$\leq \sum_{i=1}^n W_i(\theta_2) + \frac{n\alpha}{\beta}(\widehat{D} - D) + \frac{n\alpha^2}{\beta}(B_1 + \widehat{B}_1) + \frac{n\alpha^2}{\beta^2}(B_2 + \widehat{B}_2) + \eta.$$

Let us now assume that  $\theta_2 \in \overline{\Theta}$  and let us use the fact that  $\overline{\Theta}$  is a convex set and that  $\theta_0 = \arg \min_{\theta \in \overline{\Theta}} \overline{R}(\theta) + \lambda \|Q_\lambda^{-1/2} \theta\|^2$ . Introduce  $\theta_* = \arg \min_{\theta \in \mathbb{R}^d} \overline{R}(\theta) + \lambda \|Q_\lambda^{-1/2} \theta\|^2$ . As we have

$$\overline{R}(\theta) + \lambda \|Q_\lambda^{-1/2} \theta\|^2 = \|\theta - \theta_*\|^2 + \overline{R}(\theta_*) + \lambda \|Q_\lambda^{-1/2} \theta_*\|^2,$$

the vector  $\theta_0$  is uniquely defined as the projection of  $\theta_*$  on  $\overline{\Theta}$  for the Euclidean distance, and for any  $\theta_2 \in \overline{\Theta}$

$$\begin{aligned} (1.14) \quad & \alpha^{-1} \mathbb{E}[W(\theta_2)] + \lambda \|Q_\lambda^{-1/2} \theta_2\|^2 - \lambda \|Q_\lambda^{-1/2} \theta_0\|^2 \\ &= \overline{R}(\theta_2) - \overline{R}(\theta_0) + \lambda \|Q_\lambda^{-1/2} \theta_2\|^2 - \lambda \|Q_\lambda^{-1/2} \theta_0\|^2 \\ &= \|\theta_2 - \theta_*\|^2 - \|\theta_0 - \theta_*\|^2 \\ &= \|\theta_2 - \theta_0\|^2 + 2\langle \theta_2 - \theta_0, \theta_0 - \theta_* \rangle \geq \|\theta_2 - \theta_0\|^2. \end{aligned}$$

This and the inequality

$$\alpha^{-1} \sum_{i=1}^n W_i(\theta_1) + n\lambda \|Q_\lambda^{-1/2} \theta_1\|^2 - n\lambda \|Q_\lambda^{-1/2} \theta_0\|^2 \leq 0$$

leads to the following result.

**THEOREM 1.4.** *With probability at least  $1 - \exp(-\eta)$ ,*

$$\begin{aligned} R(\hat{\theta}) + \lambda \|\hat{\theta}\|^2 - \inf_{\theta \in \overline{\Theta}} [R(\theta) + \lambda \|\theta\|^2] \\ = \alpha^{-1} \mathbb{E}[W(\theta_1)] + \lambda \|Q_\lambda^{-1/2} \theta_1\|^2 - \lambda \|Q_\lambda^{-1/2} \theta_0\|^2 \end{aligned}$$

*is not greater than the smallest positive non degenerate root of the following polynomial equation as soon as it has one*

$$\begin{aligned} & \left\{ 1 - \left[ \alpha(B_3 + \widehat{B}_3) + \frac{\beta}{2n\alpha} \right] \right\} x - \alpha(B_4 + \widehat{B}_4)x^2 \\ &= \frac{1}{\beta} \max(\widehat{D} - D, 0) + \frac{\alpha}{\beta}(B_1 + \widehat{B}_1) + \frac{\alpha}{\beta^2}(B_2 + \widehat{B}_2) + \frac{\eta}{n\alpha}. \end{aligned}$$

**PROOF.** Let us remark first that when the polynomial appearing in the theorem has two distinct roots, they are of the same sign, due to the sign of its constant coefficient. Let  $\widehat{\Omega}$  be the event of probability at least  $1 - \exp(-\eta)$

described in Theorem 1.3 (page 5). For any realization of this event for which the polynomial described in Theorem 1.4 does not have two distinct positive roots, the statement of Theorem 1.4 is void, and therefore fulfilled. Let us consider now the case when the polynomial in question has two distinct positive roots  $x_1 < x_2$ . Consider in this case the random (trivially nonempty) closed convex set

$$\widehat{\Theta} = \left\{ \theta \in \Theta : R(\theta) + \lambda \|\theta\|^2 \leq \inf_{\theta' \in \Theta} [R(\theta') + \lambda \|\theta'\|^2] + \frac{x_1 + x_2}{2} \right\}.$$

Let  $\theta_3 \in \arg \min_{\theta \in \widehat{\Theta}} r(\theta) + \lambda \|\theta\|^2$  and  $\theta_4 \in \arg \min_{\theta \in \Theta} r(\theta) + \lambda \|\theta\|^2$ . We see from Theorem 1.3 that

$$(1.15) \quad R(\theta_3) + \lambda \|\theta_3\|^2 < R(\theta_0) + \lambda \|\theta_0\|^2 + \frac{x_1 + x_2}{2},$$

because it cannot be larger from the construction of  $\widehat{\Theta}$ . On the other hand, since  $\widehat{\Theta} \subset \Theta$ , the line segment  $[\theta_3, \theta_4]$  is such that  $[\theta_3, \theta_4] \cap \widehat{\Theta} \subset \arg \min_{\theta \in \widehat{\Theta}} r(\theta) + \lambda \|\theta\|^2$ . We can therefore apply equation (1.15) to any point of  $[\theta_3, \theta_4] \cap \widehat{\Theta}$ , which proves that  $[\theta_3, \theta_4] \cap \widehat{\Theta}$  is an open subset of  $[\theta_3, \theta_4]$ . But it is also a closed subset by construction, and therefore, as it is non empty and  $[\theta_3, \theta_4]$  is connected, it proves that  $[\theta_3, \theta_4] \cap \widehat{\Theta} = [\theta_3, \theta_4]$ , and thus that  $\theta_4 \in \widehat{\Theta}$ . This can be applied to any choice of  $\theta_3 \in \arg \min_{\theta \in \widehat{\Theta}} r(\theta) + \lambda \|\theta\|^2$  and  $\theta_4 \in \arg \min_{\theta \in \Theta} r(\theta) + \lambda \|\theta\|^2$ , proving that  $\arg \min_{\theta \in \Theta} r(\theta) + \lambda \|\theta\|^2 \subset \arg \min_{\theta \in \widehat{\Theta}} r(\theta) + \lambda \|\theta\|^2$  and therefore that any  $\theta_4 \in \arg \min_{\theta \in \Theta} r(\theta) + \lambda \|\theta\|^2$  is such that

$$R(\theta_4) + \lambda \|\theta_4\|^2 \leq \inf_{\theta \in \Theta} [R(\theta) + \lambda \|\theta\|^2] + x_1.$$

because the values between  $x_1$  and  $x_2$  are excluded by Theorem 1.3.  $\square$

The actual convergence speed of the least squares estimator  $\widehat{\theta}$  on  $\Theta$  will depend on the speed of convergence of the “empirical bounds”  $\widehat{B}_k$  towards their expectations. We can rephrase the previous theorem in the following more practical way:

**THEOREM 1.5.** *Let  $\eta_0, \eta_1, \dots, \eta_5$  be positive real numbers. With probability at least*

$$1 - \mathbb{P}(\widehat{D} > D + \eta_0) - \sum_{k=1}^4 \mathbb{P}(\widehat{B}_k - B_k > \eta_k) - \exp(-\eta_5),$$

$R(\widehat{\theta}) + \lambda \|\widehat{\theta}\|^2 - \inf_{\theta \in \Theta} [R(\theta) + \lambda \|\theta\|^2]$  is smaller than the smallest non degenerate positive root of



$$(1.16) \quad \left\{ 1 - \left[ \alpha(2B_3 + \eta_3) + \frac{\beta}{2n\alpha} \right] \right\} x - \alpha(2B_4 + \eta_4)x^2 \\ = \frac{\eta_0}{\beta} + \frac{\alpha}{\beta}(2B_1 + \eta_1) + \frac{\alpha}{\beta^2}(2B_2 + \eta_2) + \frac{\eta_5}{n\alpha},$$

where we can optimize the values of  $\alpha > 0$  and  $\beta > 0$ , since this equation has non random coefficients. For example, taking for simplicity

$$\alpha = \frac{1}{8B_3 + 4\eta_3}, \\ \beta = \frac{n\alpha}{2},$$

we obtain

$$x - \frac{2B_4 + \eta_4}{4B_3 + 2\eta_3}x^2 = \frac{16\eta_0(2B_3 + \eta_3)}{n} + \frac{8B_1 + 4\eta_1}{n} \\ + \frac{32(2B_3 + \eta_3)(2B_2 + \eta_2)}{n^2} + \frac{8\eta_5(2B_3 + \eta_3)}{n}.$$

1.1. *Proof of Theorem 2.1.* Let us now deduce Theorem 2.1 from Theorem 1.5. Let us first remark that with probability at least  $1 - \varepsilon/2$

$$\widehat{D} \leq D + \sqrt{\frac{B_2}{\varepsilon n}},$$

because the variance of  $\widehat{D}$  is less than  $\frac{B_2}{2n}$ . For a given  $\varepsilon > 0$ , let us take  $\eta_0 = \sqrt{\frac{B_2}{\varepsilon n}}$ ,  $\eta_1 = B_1$ ,  $\eta_2 = B_2$ ,  $\eta_3 = B_3$  and  $\eta_4 = B_4$ . We get that  $R_\lambda(\widehat{\theta}) - \inf_{\theta \in \Theta} R_\lambda(\theta)$  is smaller than the smallest positive non degenerate root of

$$x - \frac{B_4}{2B_3}x^2 = \frac{48B_3}{n} \sqrt{\frac{B_2}{n\varepsilon}} + \frac{12B_1}{n} + \frac{288B_2B_3}{n^2} + \frac{24 \log(3/\varepsilon)B_3}{n},$$

with probability at least

$$1 - \frac{5\varepsilon}{6} - \sum_{k=1}^4 \mathbb{P}(\widehat{B}_k > B_k + \eta_k).$$

According to the weak law of large numbers, there is  $n_\varepsilon$  such that for any  $n \geq n_\varepsilon$ ,

$$\sum_{k=1}^4 \mathbb{P}(\widehat{B}_k > B_k + \eta_k) \leq \varepsilon/6.$$

Thus, increasing  $n_\varepsilon$  and the constants to absorb the second order terms, we see that for some  $n_\varepsilon$  and any  $n \geq n_\varepsilon$ , with probability at least  $1 - \varepsilon$ , the excess risk is less than the smallest positive root of

$$x - \frac{B_4}{2B_3}x^2 = \frac{13B_1}{n} + \frac{24 \log(3/\varepsilon)B_3}{n}.$$

Now, as soon as  $ac < 1/4$ , the smallest positive root of  $x - ax^2 = c$  is  $\frac{2c}{1 + \sqrt{1 - 4ac}}$ . This means that for  $n$  large enough, with probability at least  $1 - \varepsilon$ ,

$$R_\lambda(\hat{\theta}) - \inf_{\theta} R_\lambda(\theta) \leq \frac{15B_1}{n} + \frac{25 \log(3/\varepsilon)B_3}{n},$$

which is precisely the statement of Theorem 2.1, up to some change of notation.

*1.2. Proof of Theorem 2.2.* Let us now weaken Theorem 1.4 in order to make a more explicit non asymptotic result and obtain Theorem 2.2. From now on, we will assume that  $\lambda = 0$ . We start by giving bounds on the quantity defined in Theorem 1.3 in terms of

$$B = \sup_{f \in \text{span}\{\varphi_1, \dots, \varphi_d\} - \{0\}} \|f\|_\infty^2 / \mathbb{E}[f(X)]^2.$$

Since we have

$$\|\bar{X}\|^2 = \|Q_\lambda^{-1/2} X\|^2 \leq dB,$$

we get

$$\hat{d} = \frac{1}{n} \sum_{i=1}^n \|\bar{X}_i\|^2 \leq dB,$$

$$B_1 = 2\mathbb{E} \left[ \|\bar{X}\|^2 (\langle \theta_0, \bar{X} \rangle - Y)^2 \right] \leq 2dB R(f^*),$$

$$\hat{B}_1 = \frac{2}{n} \sum_{i=1}^n \left[ \|\bar{X}_i\|^2 (\langle \theta_0, \bar{X}_i \rangle - Y_i)^2 \right] \leq 2dB r(f^*),$$

$$B_2 = 2\mathbb{E} \left[ \|\bar{X}\|^4 \right] \leq 2d^2 B^2,$$

$$\hat{B}_2 = \frac{2}{n} \sum_{i=1}^n \|\bar{X}_i\|^4 \leq 2d^2 B^2,$$

$$B_3 = 40 \sup \left\{ \mathbb{E} \left[ \langle u, \bar{X} \rangle^2 (\langle \theta_0, \bar{X} \rangle - Y)^2 \right] : u \in \mathbb{R}^d, \|u\| = 1 \right\} \leq 40B R(f^*),$$

$$\hat{B}_3 = \sup \left\{ \frac{40}{n} \sum_{i=1}^n \langle u, \bar{X}_i \rangle^2 (\langle \theta_0, \bar{X}_i \rangle - Y_i)^2 : u \in \mathbb{R}^d, \|u\| = 1 \right\} \leq 40B r(f^*),$$

$$B_4 = 10 \sup \left\{ \mathbb{E} \left[ \langle u, \bar{X} \rangle^4 \right] : u \in \mathbb{R}^d, \|u\| = 1 \right\} \leq 10B^2,$$

$$\widehat{B}_4 = \sup \left\{ \frac{10}{n} \sum_{i=1}^n \langle u, \bar{X}_i \rangle^4 : u \in \mathbb{R}^d, \|u\| = 1 \right\} \leq 10B^2.$$

Let us put

$$a_0 = \frac{2dB + 4dB\alpha[R(f^*) + r(f^*)] + \eta}{\alpha n} + \frac{16B^2d^2}{\alpha n^2},$$

$$a_1 = 3/4 - 40\alpha B[R(f^*) + r(f^*)],$$

and

$$a_2 = 20\alpha B^2.$$

Theorem 1.4 applied with  $\beta = n\alpha/2$  implies that with probability at least  $1 - \eta$  the excess risk  $R(\hat{f}^{(\text{erm})}) - R(f^*)$  is upper bounded by the smallest positive root of  $a_1x - a_2x^2 = a_0$  as soon as  $a_1^2 > 4a_0a_2$ . In particular, setting  $\varepsilon = \exp(-\eta)$  when (1.17) holds, we have

$$R(\hat{f}^{(\text{erm})}) - R(f^*) \leq \frac{2a_0}{a_1 + \sqrt{a_1^2 - 4a_0a_2}} \leq \frac{2a_0}{a_1}.$$

We conclude that

**THEOREM 1.6.** *For any  $\alpha > 0$  and  $\varepsilon > 0$ , with probability at least  $1 - \varepsilon$ , if the inequality*

$$(1.17) \quad 80 \left( \frac{(2 + 4\alpha[R(f^*) + r(f^*)])Bd + \log(\varepsilon^{-1})}{n} + \left( \frac{4Bd}{n} \right)^2 \right) < \left( \frac{3}{4B} - 40\alpha[R(f^*) + r(f^*)] \right)^2$$

*holds, then we have*

$$(1.18) \quad R(\hat{f}^{(\text{erm})}) - R(f^*) \leq \mathcal{J} \left( \frac{(2 + 4\alpha[R(f^*) + r(f^*)])Bd + \log(\varepsilon^{-1})}{n} + \left( \frac{4Bd}{n} \right)^2 \right),$$

*where  $\mathcal{J} = 8/(3\alpha - 160\alpha^2 B[R(f^*) + r(f^*)])$*

Now, the Bienaymé-Chebyshev inequality implies

$$\mathbb{P}(r(f^*) - R(f^*) \geq t) \leq \frac{\mathbb{E}(r(f^*) - R(f^*))^2}{t^2} \leq \mathbb{E}[Y - f^*(X)]^4 / nt^2.$$

Under the finite moment assumption of Theorem 2.2, we obtain that for any  $\varepsilon \geq 1/n$ , with probability at least  $1 - \varepsilon$ ,

$$r(f^*) < R(f^*) + \sqrt{\mathbb{E}[Y - f^*(X)]^4}.$$

From Theorem 1.6 and a union bound, by taking

$$\alpha = \left(80B[2R(f^*) + \sqrt{\mathbb{E}[Y - f^*(X)]^4}\right)^{-1},$$

we get that with probability  $1 - 2\varepsilon$ ,

$$R(\hat{f}^{\text{(erm)}}) - R(f^*) \leq \mathcal{J}_1 B \left( \frac{3Bd + \log(\varepsilon^{-1})}{n} + \left(\frac{4Bd}{n}\right)^2 \right),$$

with  $\mathcal{J}_1 = 640 \left(2R(f^*) + \sqrt{\mathbb{E}\{[Y - f^*(X)]^4\}}\right)$ . This concludes the proof of Theorem 2.2.

**REMARK 1.1.** *Let us indicate now how to handle the case when  $Q$  is degenerate. Let us consider the linear subspace  $S$  of  $\mathbb{R}^d$  spanned by the eigenvectors of  $Q$  corresponding to positive eigenvalues. Then almost surely  $\text{Span}\{X_i, i = 1, \dots, n\} \subset S$ . Indeed for any  $\theta$  in the kernel of  $Q$ ,  $\mathbb{E}(\langle \theta, X \rangle^2) = 0$  implies that  $\langle \theta, X \rangle = 0$  almost surely, and considering a basis of the kernel, we see that  $X \in S$  almost surely,  $S$  being orthogonal to the kernel of  $Q$ . Thus we can restrict the problem to  $S$ , as soon as we choose*

$$\hat{\theta} \in \text{span}\{X_1, \dots, X_n\} \cap \arg \min_{\theta} \sum_{i=1}^n (\langle \theta, X_i \rangle - Y_i)^2,$$

or equivalently with the notation  $\mathbf{X} = (\varphi_j(X_i))_{1 \leq i \leq n, 1 \leq j \leq d}$  and  $\mathbf{Y} = [Y_j]_{j=1}^n$ ,

$$\hat{\theta} \in \text{im } \mathbf{X}^T \cap \arg \min_{\theta} \|\mathbf{X}\theta - \mathbf{Y}\|^2$$

*This proves that the results of this section apply to this special choice of the empirical least squares estimator. Since we have  $\mathbb{R}^d = \ker \mathbf{X} \oplus \text{im } \mathbf{X}^T$ , this choice is unique. Finally, we also have that inequality (2.3) of the paper still holds by replacing  $d$  by  $\text{rank}(Q)$ .*

**2. Proof of Theorem 3.1.** We use the same notations as in Section 1. We write  $X$  for  $\varphi(X)$ , therefore, the function  $f_{\theta}$  maps an input  $x$  to  $\langle \theta, x \rangle$ . We consider the change of coordinates

$$\bar{X} = Q_{\lambda}^{-1/2} X.$$

Thus, from (1.13), we have  $\mathbb{E}[\|\overline{X}\|^2] = D$ . We will use

$$\overline{R}(\theta) = \mathbb{E}[(\langle \theta, \overline{X} \rangle - Y)^2],$$

so that  $\overline{R}(Q_\lambda^{1/2}\theta) = \mathbb{E}[(\langle \theta, X \rangle - Y)^2] = R(f_\theta)$ . Let

$$\overline{\Theta} = \{Q_\lambda^{1/2}\theta; \theta \in \Theta\},$$

and consider

$$\theta_0 = \arg \min_{\theta \in \overline{\Theta}} \left\{ \overline{R}(\theta) + \lambda \|Q_\lambda^{-1/2}\theta\|^2 \right\}.$$

With these notations,

$$\begin{aligned} \tilde{\theta} &= Q_\lambda^{-1/2}\theta_0, \\ \sigma &= \sqrt{\mathbb{E}[(\langle \theta_0, \overline{X} \rangle - Y)^2]}, \\ \chi &= \sup_{u \in \mathbb{R}^d} \frac{\mathbb{E}(\langle u, \overline{X} \rangle^4)^{1/2}}{\mathbb{E}(\langle u, \overline{X} \rangle^2)}, \\ \kappa &= \frac{\mathbb{E}(\|\overline{X}\|^4)^{1/2}}{\mathbb{E}(\|\overline{X}\|^2)} = \frac{\mathbb{E}(\|\overline{X}\|^4)^{1/2}}{D}, \\ \kappa' &= \frac{\mathbb{E}[(\langle \theta_0, \overline{X} \rangle - Y)^4]^{1/2}}{\sigma^2}, \\ \text{and } T &= \|\overline{\Theta}\| = \max_{\theta, \theta' \in \overline{\Theta}} \|\theta - \theta'\|. \end{aligned}$$

For  $\alpha > 0$ , we introduce

$$\begin{aligned} J_i(\theta) &= \langle \theta, \overline{X}_i \rangle - Y_i, & J(\theta) &= \langle \theta, \overline{X} \rangle - Y \\ \overline{L}_i(\theta) &= \alpha(\langle \theta, \overline{X}_i \rangle - Y_i)^2, & \overline{L}(\theta) &= \alpha(\langle \theta, \overline{X} \rangle - Y)^2 \\ W_i(\theta) &= \overline{L}_i(\theta) - \overline{L}_i(\theta_0), & W(\theta) &= \overline{L}(\theta) - \overline{L}(\theta_0), \end{aligned}$$

and

$$r'(\theta, \theta') = \lambda \left( \|Q_\lambda^{-1/2}\theta\|^2 - \|Q_\lambda^{-1/2}\theta'\|^2 \right) + \frac{1}{n\alpha} \sum_{i=1}^n \psi(\overline{L}(\theta) - \overline{L}(\theta')).$$

Let  $\bar{\theta} = Q_\lambda^{1/2}\hat{\theta} \in \overline{\Theta}$ . We have

$$(2.1) \quad -r'(\theta_0, \bar{\theta}) = r'(\bar{\theta}, \theta_0) \leq \max_{\theta_1 \in \overline{\Theta}} r'(\bar{\theta}, \theta_1) \leq \gamma + \max_{\theta_1 \in \overline{\Theta}} r'(\theta_0, \theta_1),$$

where the quantity  $\gamma = \max_{\theta_1 \in \bar{\Theta}} r'(\bar{\theta}, \theta_1) - \inf_{\theta \in \bar{\Theta}} \max_{\theta_1 \in \bar{\Theta}} r'(\theta, \theta_1)$  can be made arbitrary small by a proper choice of the estimator. Using an upper bound  $r'(\theta_0, \theta_1)$  that holds uniformly in  $\theta_1$ , we will control both left and right hand sides of (2.1).

To achieve this, we will upper bound

$$(2.2) \quad r'(\theta_0, \theta_1) = \lambda \left( \|Q_\lambda^{-1/2} \theta_0\|^2 - \|Q_\lambda^{-1/2} \theta_1\|^2 \right) + \frac{1}{n\alpha} \sum_{i=1}^n \psi[-W_i(\theta_1)]$$

by the expectation of a distribution depending on  $\theta_1$  of a *quantity that does not depend on  $\theta_1$* , and then use the PAC-Bayesian argument to control this expectation uniformly in  $\theta_1$ . The distribution depending on  $\theta_1$  should therefore be taken such that for any  $\theta_1 \in \bar{\Theta}$ , its Kullback-Leibler divergence with respect to some fixed distribution is small (at least when  $\theta_1$  is close to  $\theta_0$ ).

Let us start with the following result.

LEMMA 2.1. *Let  $f, g : \mathbb{R} \rightarrow \mathbb{R}$  be two Lebesgue measurable functions such that  $f(x) \leq g(x)$ ,  $x \in \mathbb{R}$ . Let us assume that there exists  $h \in \mathbb{R}$  such that  $x \mapsto g(x) + hx^2/2$  is convex. Then for any probability distribution  $\mu$  on the real line,*

$$f\left(\int x \mu(dx)\right) \leq \int g(x) \mu(dx) + \min\left\{\sup f - \inf f, \frac{h}{2} \text{Var}(\mu)\right\}.$$

PROOF. Let us put  $x_0 = \int x \mu(dx)$ . The function

$$x \mapsto g(x) + \frac{h}{2}(x - x_0)^2$$

is convex. Thus, by Jensen's inequality

$$f(x_0) \leq g(x_0) \leq \int \mu(dx) \left[ g(x) + \frac{h}{2}(x - x_0)^2 \right] = \int g(x) \mu(dx) + \frac{h}{2} \text{Var}(\mu).$$

On the other hand

$$\begin{aligned} f(x_0) &\leq \sup f \leq \sup f + \int [g(x) - \inf f] \mu(dx) \\ &= \int g(x) \mu(dx) + \sup f - \inf f. \end{aligned}$$

The lemma is a combination of these two inequalities.  $\square$

The above lemma will be used with  $f = g = \psi$ , where  $\psi$  is the increasing influence function

$$\psi(x) = \begin{cases} -\log(2), & x \leq -1, \\ \log(1 + x + x^2/2), & -1 \leq x \leq 0, \\ -\log(1 - x + x^2/2), & 0 \leq x \leq 1, \\ \log(2), & x \geq 1. \end{cases}$$

Since we have for any  $x \in \mathbb{R}$

$$-\log\left(1 - x + \frac{x^2}{2}\right) = \log\left(\frac{1 + x + \frac{x^2}{2}}{1 + \frac{x^2}{4}}\right) < \log\left(1 + x + \frac{x^2}{2}\right),$$

the function  $\psi$  satisfies for any  $x \in \mathbb{R}^*$

$$-\log\left(1 - x + \frac{x^2}{2}\right) < \psi(x) < \log\left(1 + x + \frac{x^2}{2}\right).$$

Moreover

$$\psi'(x) = \frac{1-x}{1-x+\frac{x^2}{2}}, \quad \psi''(x) = \frac{x(x-2)}{2(1-x+\frac{x^2}{2})^2} \geq -2, \quad 0 \leq x \leq 1,$$

showing (by symmetry) that the function  $x \mapsto \psi(x) + 2x^2$  is convex on the real line.

For any  $\theta' \in \mathbb{R}^d$  and  $\beta > 0$ , we consider the Gaussian distribution with mean  $\theta'$  and covariance  $\beta^{-1}I$ :

$$\rho_{\theta'}(d\theta) = \left(\frac{\beta}{2\pi}\right)^{d/2} \exp\left(-\frac{\beta}{2}\|\theta - \theta'\|^2\right) d\theta.$$

From Lemmas 1.2 and 2.1 (with  $\mu$  the distribution of  $-W_i(\theta) + \frac{\alpha\|\bar{X}_i\|^2}{\beta}$  when  $\theta$  is drawn from  $\rho_{\theta_1}$  and for a fixed pair  $(X_i, Y_i)$ ), we can see that

$$\begin{aligned} \psi[-W_i(\theta_1)] &= \psi\left\{\int \rho_{\theta_1}(d\theta) \left[-W_i(\theta) + \frac{\alpha\|\bar{X}_i\|^2}{\beta}\right]\right\} \\ &\leq \int \rho_{\theta_1}(d\theta) \psi\left[-W_i(\theta) + \frac{\alpha\|\bar{X}_i\|^2}{\beta}\right] \\ &\quad + \min\left\{\log(4), \text{Var}_{\rho_{\theta_1}}[\bar{L}_i(\theta)]\right\}. \end{aligned}$$

Let us compute

$$\frac{1}{\alpha^2} \text{Var}_{\rho_{\theta_1}}(\bar{L}_i(\theta)) = \text{Var}_{\rho_{\theta_1}}[J_i^2(\theta) - J_i^2(\theta_1)]$$

$$\begin{aligned}
&= \int \rho_{\theta_1}(d\theta) [J_i^2(\theta) - J_i^2(\theta_1)]^2 - \frac{\|\bar{X}_i\|^4}{\beta^2} \\
&= \int \rho_{\theta_1}(d\theta) [\langle \theta - \theta_1, \bar{X}_i \rangle^2 + 2\langle \theta - \theta_1, \bar{X}_i \rangle J_i(\theta_1)]^2 - \frac{\|\bar{X}_i\|^4}{\beta^2} \\
(2.3) \quad &= \frac{2\|\bar{X}_i\|^4}{\beta^2} + \frac{4\bar{L}_i(\theta_1)\|\bar{X}_i\|^2}{\alpha\beta}.
\end{aligned}$$

Let  $\xi \in (0, 1)$ , and let us remark that

$$\bar{L}_i(\theta_1) \leq \frac{\bar{L}_i(\theta)}{\xi} + \frac{\alpha \langle \theta - \theta_1, \bar{X}_i \rangle^2}{1 - \xi}.$$

We get

$$\begin{aligned}
&\min \left\{ \log(4), \text{Var}_{\rho_{\theta_1}} [\bar{L}_i(\theta)] \right\} \\
&= \min \left\{ \log(4), \frac{4\alpha \|\bar{X}_i\|^2 \bar{L}_i(\theta_1)}{\beta} + \frac{2\alpha^2 \|\bar{X}_i\|^4}{\beta^2} \right\} \\
&\leq \int \rho_{\theta_1}(d\theta) \min \left\{ \log(4), \right. \\
&\quad \left. \frac{4\alpha \|\bar{X}_i\|^2 \bar{L}_i(\theta)}{\beta\xi} + \frac{2\alpha^2 \|\bar{X}_i\|^4}{\beta^2} + \frac{4\alpha^2 \|\bar{X}_i\|^2 \langle \theta - \theta_1, \bar{X}_i \rangle^2}{\beta(1 - \xi)} \right\} \\
&\leq \int \rho_{\theta_1}(d\theta) \min \left\{ \log(4), \frac{4\alpha \|\bar{X}_i\|^2 \bar{L}_i(\theta)}{\beta\xi} + \frac{2\alpha^2 \|\bar{X}_i\|^4}{\beta^2} \right\} \\
&\quad + \min \left\{ \log(4), \frac{4\alpha^2 \|\bar{X}_i\|^4}{\beta^2(1 - \xi)} \right\}.
\end{aligned}$$

Let us now put  $a = \frac{3}{\log(4)} < 2.17$ ,  $b = a + a^2 \log(4) < 8.7$  and let us remark that

$$\begin{aligned}
&\min \{ \log(4), x \} + \min \{ \log(4), y \} \\
&\leq \log[1 + a \min \{ \log(4), x \}] + \log(1 + ay) \\
&\leq \log(1 + ax + by), \quad x, y \in \mathbb{R}_+.
\end{aligned}$$

Thus

$$\begin{aligned}
&\min \left\{ \log(4), \text{Var}_{\rho_{\theta_1}} [\bar{L}_i(\theta)] \right\} \\
&\leq \int \rho_{\theta_1}(d\theta) \log \left[ 1 + \frac{4a\alpha \|\bar{X}_i\|^2 \bar{L}_i(\theta)}{\beta\xi} + \frac{2\alpha^2 \|\bar{X}_i\|^4}{\beta^2} \left( a + \frac{2b}{1 - \xi} \right) \right].
\end{aligned}$$



We can then remark that

$$\begin{aligned} \psi(x) + \log(1 + y) &= \log[\exp[\psi(x)] + y \exp[\psi(x)]] \\ &\leq \log[\exp[\psi(x)] + 2y] \leq \log\left(1 + x + \frac{x^2}{2} + 2y\right), \quad x \in \mathbb{R}, y \in \mathbb{R}_+. \end{aligned}$$

Thus, putting  $c_0 = a + \frac{2b}{1 - \xi}$ , we get

$$(2.4) \quad \psi[-W_i(\theta_1)] \leq \int \rho_{\theta_1}(d\theta) \log[A_i(\theta)],$$

with

$$\begin{aligned} A_i(\theta) &= 1 - W_i(\theta) + \frac{\alpha \|\bar{X}_i\|^2}{\beta} + \frac{1}{2} \left( -W_i(\theta) + \frac{\alpha \|\bar{X}_i\|^2}{\beta} \right)^2 \\ &\quad + \frac{8a\alpha \|\bar{X}_i\|^2 \bar{L}_i(\theta)}{\beta \xi} + \frac{4c_0 \alpha^2 \|\bar{X}_i\|^4}{\beta^2}. \end{aligned}$$

Similarly, we define  $A(\theta)$  by replacing  $(X_i, Y_i)$  by  $(X, Y)$ . Since we have

$$\mathbb{E} \left[ \exp \left( \sum_{i=1}^n \log[A_i(\theta)] - n \log[\mathbb{E}A(\theta)] \right) \right] = 1,$$

from the usual PAC-Bayesian argument, we have with probability at least  $1 - \varepsilon$ , for any  $\theta_1 \in \mathbb{R}^d$ ,

$$\begin{aligned} \int \rho_{\theta_1}(d\theta) \left( \sum_{i=1}^n \log[A_i(\theta)] \right) - n \int \rho_{\theta_1}(d\theta) \log[A(\theta)] &\leq K(\rho_{\theta_1}, \rho_{\theta_0}) + \log(\varepsilon^{-1}) \\ &\leq \frac{\beta \|\theta_1 - \theta_0\|^2}{2} + \log(\varepsilon^{-1}). \end{aligned}$$

From (2.2) and (2.4), with probability at least  $1 - \varepsilon$ , for any  $\theta_1 \in \mathbb{R}^d$ , we get

$$\begin{aligned} r'(\theta_0, \theta_1) &\leq \frac{1}{\alpha} \log \left\{ 1 + \mathbb{E} \left[ \int \rho_{\theta_1}(d\theta) \left( -W(\theta) + \frac{\alpha \|\bar{X}\|^2}{\beta} \right. \right. \right. \\ &\quad \left. \left. + \frac{1}{2} \left( -W(\theta) + \frac{\alpha \|\bar{X}\|^2}{\beta} \right)^2 + \frac{8a\alpha \|\bar{X}\|^2 \bar{L}(\theta)}{\beta \xi} + \frac{4c_0 \alpha^2 \|\bar{X}\|^4}{\beta^2} \right) \right] \right\} \\ &\quad + \frac{\beta \|\theta_1 - \theta_0\|^2}{2n\alpha} + \frac{\log(\varepsilon^{-1})}{n\alpha} + \lambda \left( \|Q_\lambda^{-1/2} \theta_0\|^2 - \|Q_\lambda^{-1/2} \theta_1\|^2 \right). \end{aligned}$$

Moreover from (2.3) and  $\frac{\alpha\|\bar{X}\|^2}{\beta} = -\bar{L}(\theta_1) + \int \rho_{\theta_1}(d\theta)\bar{L}(\theta)$ , we deduce that

$$\begin{aligned} \int \rho_{\theta_1}(d\theta) \left( -W(\theta) + \frac{\alpha\|\bar{X}\|^2}{\beta} \right)^2 &= \text{Var}_{\rho_{\theta_1}}[\bar{L}(\theta)] + W(\theta_1)^2 \\ &= W(\theta_1)^2 + \frac{4\alpha\bar{L}(\theta_1)\|\bar{X}\|^2}{\beta} + \frac{2\alpha^2\|\bar{X}\|^4}{\beta^2}. \end{aligned}$$

PROPOSITION 2.2. *With probability at least  $1 - \varepsilon$ , for any  $\theta_1 \in \mathbb{R}^d$ ,*

$$\begin{aligned} r'(\theta_0, \theta_1) &\leq \frac{1}{\alpha} \log \left\{ 1 + \mathbb{E} \left[ -W(\theta_1) + \frac{W(\theta_1)^2}{2} + \frac{(2 + 8a/\xi)\alpha\|\bar{X}\|^2\bar{L}(\theta_1)}{\beta} \right. \right. \\ &\quad \left. \left. + \frac{(1 + 8a/\xi + 4c_0)\alpha^2\|\bar{X}\|^4}{\beta^2} \right] \right\} + \frac{\beta\|\theta_1 - \theta_0\|^2}{2n\alpha} + \frac{\log(\varepsilon^{-1})}{n\alpha} \\ &\quad + \lambda \left( \|Q_\lambda^{-1/2}\theta_0\|^2 - \|Q_\lambda^{-1/2}\theta_1\|^2 \right) \\ &\leq \mathbb{E} \left[ J(\theta_0)^2 - J(\theta_1)^2 + \frac{1}{2\alpha} W(\theta_1)^2 + \frac{(2 + 8a/\xi)\|\bar{X}\|^2\bar{L}(\theta_1)}{\beta} \right. \\ &\quad \left. + \frac{(1 + 8a/\xi + 4c_0)\alpha\|\bar{X}\|^4}{\beta^2} \right] + \frac{\beta\|\theta_1 - \theta_0\|^2}{2n\alpha} + \frac{\log(\varepsilon^{-1})}{n\alpha} \\ &\quad + \lambda \left( \|Q_\lambda^{-1/2}\theta_0\|^2 - \|Q_\lambda^{-1/2}\theta_1\|^2 \right). \end{aligned}$$

Using the triangular inequality and Cauchy-Schwarz's inequality, we get

$$\begin{aligned} \frac{1}{\alpha^2} \mathbb{E}[W(\theta_1)^2] &= \mathbb{E} \left\{ [\langle \theta_1 - \theta_0, \bar{X} \rangle^2 + 2\langle \theta_1 - \theta_0, \bar{X} \rangle J(\theta_0)]^2 \right\} \\ (2.5) \quad &\leq \left\{ \mathbb{E}[\langle \theta_1 - \theta_0, \bar{X} \rangle^4]^{1/2} + 2\mathbb{E}[\langle \theta_1 - \theta_0, \bar{X} \rangle^4]^{1/4} \mathbb{E}[J(\theta_0)^4]^{1/4} \right\}^2 \\ &\leq \left\{ \chi \|\theta_1 - \theta_0\|^2 \mathbb{E} \left[ \left\langle \frac{\theta_1 - \theta_0}{\|\theta_1 - \theta_0\|}, \bar{X} \right\rangle^2 \right] \right. \\ &\quad \left. + 2\|\theta_1 - \theta_0\| \sigma \sqrt{\kappa' \chi} \sqrt{\mathbb{E} \left[ \left\langle \frac{\theta_1 - \theta_0}{\|\theta_1 - \theta_0\|}, \bar{X} \right\rangle^2 \right]} \right\}^2 \\ &\leq \frac{\chi q_{\max}}{q_{\max} + \lambda} \|\theta_1 - \theta_0\|^2 \left\{ \|\theta_1 - \theta_0\| \sqrt{\frac{\chi q_{\max}}{q_{\max} + \lambda}} + 2\sigma \sqrt{\kappa'} \right\}^2, \end{aligned}$$

and

$$\frac{1}{\alpha} \mathbb{E}[\|\bar{X}\|^2 \bar{L}(\theta_1)] = \mathbb{E} \left\{ [\|\bar{X}\| \langle \theta_1 - \theta_0, \bar{X} \rangle + \|\bar{X}\| J(\theta_0)]^2 \right\}$$

$$\begin{aligned}
(2.6) \quad &\leq \mathbb{E}[\|\bar{X}\|^4]^{1/2} \left\{ \mathbb{E}[\langle \theta_1 - \theta_0, \bar{X} \rangle^4]^{1/4} + \mathbb{E}[J(\theta_0)^4]^{1/4} \right\}^2 \\
&\leq \kappa D \left\{ \|\theta_1 - \theta_0\| \sqrt{\frac{\chi q_{\max}}{q_{\max} + \lambda}} + 2\sigma\sqrt{\kappa'} \right\}^2.
\end{aligned}$$

Let us put

$$\begin{aligned}
\tilde{R}(\theta) &= \bar{R}(\theta) + \lambda \|Q_\lambda^{-1/2}\theta\|^2, \\
c_1 &= 4(2 + 8a/\xi), \\
c_2 &= 4(1 + 8a/\xi + 4c_0), \\
\delta &= \frac{c_1 \kappa \kappa' D \sigma^2}{n} + \frac{2\chi \left( \frac{\log(\varepsilon^{-1})}{n} + \frac{c_2 \kappa^2 D^2}{n^2} \right) \left( 2\sqrt{\kappa'}\sigma + \|\bar{\Theta}\|\sqrt{\chi} \right)^2}{1 - \frac{4c_1 \kappa \chi D}{n}}.
\end{aligned}$$

We have proved the following result.

**PROPOSITION 2.3.** *With probability at least  $1 - \varepsilon$ , for any  $\theta_1 \in \mathbb{R}^d$ ,*

$$\begin{aligned}
r'(\theta_0, \theta_1) &\leq \tilde{R}(\theta_0) - \tilde{R}(\theta_1) + \frac{\alpha}{2} \chi \|\theta_1 - \theta_0\|^2 [2\sqrt{\kappa'}\sigma + \|\theta_1 - \theta_0\|\sqrt{\chi}]^2 \\
&\quad + \frac{c_1 \alpha}{4\beta} \kappa D [\sqrt{\kappa'}\sigma + \|\theta_1 - \theta_0\|\sqrt{\chi}]^2 + \frac{c_2 \alpha \kappa^2 D^2}{4\beta^2} \\
&\quad \quad \quad + \frac{\beta \|\theta_1 - \theta_0\|^2}{2n\alpha} + \frac{\log(\varepsilon^{-1})}{n\alpha}.
\end{aligned}$$

Let us assume from now on that  $\theta_1 \in \bar{\Theta}$ , our convex bounded parameter set. In this case, as seen in (1.14), we have  $\|\theta_0 - \theta_1\|^2 \leq \tilde{R}(\theta_1) - \tilde{R}(\theta_0)$ . We can also use the fact that

$$[\sqrt{\kappa'}\sigma + \|\theta_1 - \theta_0\|\sqrt{\chi}]^2 \leq 2\kappa'\sigma^2 + 2\chi\|\theta_1 - \theta_0\|^2.$$

We deduce from these remarks that with probability at least  $1 - \varepsilon$ ,

$$\begin{aligned}
r'(\theta_0, \theta_1) &\leq \left\{ -1 + \frac{\alpha\chi}{2} [2\sqrt{\kappa'}\sigma + \|\bar{\Theta}\|\sqrt{\chi}]^2 + \frac{\beta}{2n\alpha} + \frac{c_1 \alpha \kappa D \chi}{2\beta} \right\} [\tilde{R}(\theta_1) - \tilde{R}(\theta_0)] \\
&\quad + \frac{c_1 \alpha \kappa D \kappa' \sigma^2}{2\beta} + \frac{c_2 \alpha \kappa^2 D^2}{4\beta^2} + \frac{\log(\varepsilon^{-1})}{n\alpha}.
\end{aligned}$$

Let us assume that  $n > 4c_1 \kappa \chi D$  and let us choose

$$\beta = \frac{n\alpha}{2},$$

$$\alpha = \frac{1}{2\chi[2\sqrt{\kappa'}\sigma + \|\bar{\Theta}\|\sqrt{\chi}]^2} \left(1 - \frac{4c_1\kappa\chi D}{n}\right),$$

to get

$$r'(\theta_0, \theta_1) \leq -\frac{\tilde{R}(\theta_1) - \tilde{R}(\theta_0)}{2} + \delta.$$

Plugging this into (2.1), we get

$$\frac{\tilde{R}(\bar{\theta}) - \tilde{R}(\theta_0)}{2} - \delta \leq r'(\bar{\theta}, \theta_0) \leq \max_{\theta_1 \in \bar{\Theta}} \left(\frac{\tilde{R}(\theta_0) - \tilde{R}(\theta_1)}{2}\right) + \gamma + \delta = \gamma + \delta,$$

hence

$$\tilde{R}(\bar{\theta}) - \tilde{R}(\theta_0) \leq 2\gamma + 4\delta.$$

Computing the numerical values of the constants when  $\xi = 0.8$  gives  $c_1 < 95$  and  $c_2 < 1511$ .

### References.

- [1] O. Catoni. *Statistical Learning Theory and Stochastic Optimization, Lectures on Probability Theory and Statistics, École d'Été de Probabilités de Saint-Flour XXXI – 2001*, volume 1851 of *Lecture Notes in Mathematics*. Springer, 2004. Pages 1–269.
- [2] M.D. Donsker and S.R.S. Varadhan. Asymptotic evaluation of certain Markov process expectations for large time, I. *Communications on Pure and Applied Mathematics*, 28(1):1–47, 1975.
- [3] P. Dupuis and R.S. Ellis. *A weak convergence approach to the theory of large deviations*. Wiley-Interscience, 1997.