

# PAC-Bayes bounds for the Gram matrix and least squares regression

Olivier Catoni  
CREST – EXCESS,  
CNRS UMR 9194  
`Olivier.Catoni@ensae.fr`

*SMILE seminar,*

PARISTECH,  
ENS, 45 RUE D'ULM, 75230 PARIS CEDEX,

*May 4, 2015*

- ① Estimate the Gram matrix

$$G \stackrel{\text{def}}{=} \int x x^\top dP(x),$$

or equivalently the quadratic form

$$N(\theta) = \theta^\top G \theta = \int \langle x, \theta \rangle^2 dP(x)$$

for all  $\theta \in \mathbb{S}_d$ , the sphere of  $\mathbb{R}^d$ , from an i.i.d. sample  $(X_i)_{i=1}^n \sim P^{\otimes n}$ .

- ② Estimate

$$\arg \min_{\theta \in \mathbb{R}^d} \int (y - \langle \theta, x \rangle)^2 dP(x, y).$$

from an i.i.d. sample  $(X_i, Y_i)_{i=1}^n \sim P^{\otimes n}$ .

- ③ We will assume that

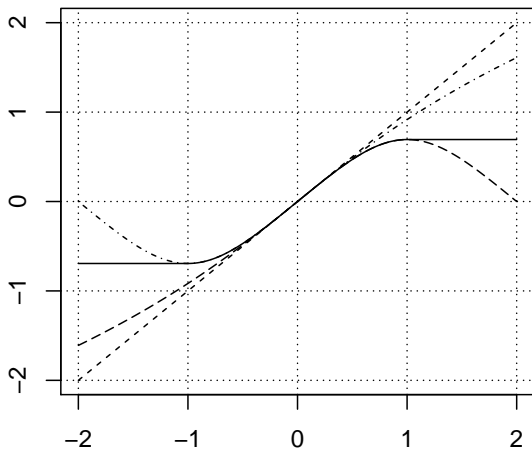
$$\int \|x\|^2 dP(x) = \mathbf{Tr}(G) < +\infty.$$

# Robust Gram matrix estimate

Let us introduce the **influence function**

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

## Robust Gram matrix estimate



$z \mapsto \psi(z)$ , compared with  $z \mapsto z$   
 $z \mapsto \log(1 + z + z^2/2)$ , and  $z \mapsto -\log(1 - z + z^2/2)$

## Robust Gram matrix estimate

It is symmetric, non decreasing, bounded and satisfies for any  $z \in \mathbb{R}$ ,

$$-\log(1 - z + z^2/2) \leq \psi(z) \leq \log(1 + z + z^2/2),$$

$$-\log(2) \leq \psi(z) \leq \log(2).$$

Let  $\bar{P} = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  and

$$r_\lambda(\theta) = \lambda^{-1} \int \psi \left[ \lambda (\langle \theta, x \rangle^2 - 1) \right] d\bar{P}(x),$$

where  $\lambda > 0$  will be chosen later. Let us consider the estimator

$$\hat{N}(\theta) = \hat{\alpha}(\theta)^{-2}, \quad \text{where} \quad \hat{\alpha}(\theta) = \sup \{ \alpha \in \mathbb{R}_+ : r_\lambda(\alpha\theta) \leq 0 \}.$$

This makes sense since

$$\lim_{\lambda \rightarrow 0} r_\lambda(\alpha\theta) = \alpha^2 \int \langle \theta, x \rangle^2 d\bar{P}(x) - 1$$

## Proposition

Let us assume that  $\kappa = \sup_{\theta \neq 0} \frac{\int \langle \theta, x \rangle^4 dP(x)}{\left( \int \langle \theta, x \rangle^2 dP(x) \right)^2} < \infty$  and

$$\text{let us put } \lambda = \sqrt{\frac{2}{(\kappa - 1)n} [\log(\epsilon^{-1}) + 0.73 d]}.$$

For any  $\epsilon > 0$ , any  $n$  such that

$$n > \left( 20\sqrt{\kappa d} + \left( \frac{5}{2} + \frac{1}{2(\kappa - 1)} \right) \sqrt{2(\kappa - 1) [\log(\epsilon^{-1}) + 0.73 d]} \right)^2,$$

with probability at least  $1 - 2\epsilon$ , for any  $\theta \in \mathbb{R}^d$ ,

$$\left| \frac{N(\theta)}{\widehat{N}(\theta)} - 1 \right| \leq \frac{\mu}{1 - 2\mu},$$

where  $\mu = \sqrt{\frac{2(\kappa - 1)}{n} [\log(\epsilon^{-1}) + 0.73 d]} + 6.81 \sqrt{\frac{2\kappa d}{n}}$ .

## Obtaining a true quadratic form

Let us assume that with probability  $1 - 2\epsilon$ ,

$$B_-(\theta) \leq N(\theta) \leq B_+(\theta), \quad \theta \in \mathbb{R}^d.$$

We can take for example

$$\hat{N}(\theta) \left( 1 - \frac{\mu}{1 - 2\mu} \right) \leq N(\theta) \leq \hat{N}(\theta) \left( 1 + \frac{\mu}{1 - 2\mu} \right)$$

Let  $\Theta \subset \mathbb{R}^d$  be any finite set (e.g. a  $\delta$ -net of  $\mathbb{S}_d$ ). Let

$$\begin{aligned} \hat{G} = \sum_{\theta \in \Theta} \xi(\theta) \theta \theta^\top, \text{ where } \xi \in \arg \min \frac{1}{2} \sum_{(\theta_1, \theta_2) \in \Theta^2} \xi(\theta_1) \xi(\theta_2) \langle \theta_1, \theta_2 \rangle^2 \\ - \sum_{\theta \in \Theta} \xi(\theta) \frac{B_-(\theta) + B_+(\theta)}{2} + |\xi(\theta)| \frac{B_+(\theta) - B_-(\theta)}{2}. \end{aligned}$$

## Obtaining a true quadratic form

Then with probability at least  $1 - 2\epsilon$ ,

$$\|\widehat{G}\|_F^2 \stackrel{\text{def}}{=} \mathbf{Tr}(\widehat{G}\widehat{G}^\top) \leq \mathbf{Tr}(GG^\top) \leq \mathbf{Tr}(G)^2,$$

so that

$$\theta_2^\top \widehat{G}\theta_2 - \theta_1^\top \widehat{G}\theta_1 \leq \|\theta_1 + \theta_2\| \mathbf{Tr}(G) \|\theta_2 - \theta_1\|, \quad \theta_1, \theta_2 \in \mathbb{R}^d,$$

and

$$B_-(\theta) \leq \theta^\top \widehat{G}\theta \leq B_+(\theta), \quad \theta \in \Theta.$$



# Proof

$\widehat{G}$  minimizes

$$\sup_{\xi_+ \geq 0, \xi_- \geq 0} \frac{1}{2} \|\widehat{G}\|_F^2 + \sum_{\theta \in \Theta} \xi_+(\theta) [B_-(\theta) - \theta^\top \widehat{G} \theta] + \xi_-(\theta) [\theta^\top \widehat{G} \theta - B_+(\theta)]$$

and there is no duality gap, so you can minimize in  $\widehat{G}$  first and then as  $\widehat{\xi}_+(\theta)\widehat{\xi}_-(\theta) = 0$ , you can assume that  $\xi_-(\theta)\xi_+(\theta) = 0$  and introduce  $\xi(\theta) = \xi_+(\theta) - \xi_-(\theta)$ , so that  $|\xi(\theta)| = \xi_+(\theta) + \xi_-(\theta)$ .

## Proposition

Assume that  $B_{+/-}(\theta) = \|\theta\| B_{+/-}(\|\theta\|^{-1}\theta)$  and are continuous functions. With probability at least  $1 - 2\epsilon$ , for any  $\eta > 0$ , there is  $\delta > 0$ , such that if  $\Theta$  is a  $\delta$ -net on the sphere  $\mathbb{S}_d$ , then

$$B_-(\theta) - \eta\|\theta\|^2 \leq \theta^\top \widehat{G}\theta \leq B_+(\theta) + \eta\|\theta\|^2, \quad \theta \in \mathbb{R}^d.$$

# The empirical Gram matrix estimator

Here we consider

$$\overline{G} = \int x x^\top d\overline{P}(x) \text{ and } \overline{N}(\theta) = \theta^\top \overline{G} \theta.$$

Remark that, with probability  $1 - \epsilon$ , for any  $\theta \in \mathbb{S}_d$ ,  $\widehat{N}(\theta) > 0$ , and therefore

$$\int \psi \left[ \lambda (\langle \theta, x \rangle^2 \widehat{N}(\theta)^{-1} - 1) \right] d\overline{P}(x) = 0,$$

so that

$$\frac{\overline{N}(\theta)}{\widehat{N}(\theta)} - 1 = \lambda^{-1} \int g \left[ \lambda (\langle \theta, x \rangle^2 \overline{N}(\theta)^{-1} - 1) \right] d\overline{P}(x),$$

where  $g(z) = z - \psi(z) \leq z_+^3/3$ .

# The empirical Gram matrix estimator

## Proposition

With probability at least  $1 - \epsilon$ ,

$$\begin{aligned} -\frac{\lambda^2}{3} \int \left(1 - \langle \theta, x \rangle^2 \widehat{N}(\theta)^{-1}\right)_+^3 d\bar{\mathbb{P}}(x) &\leq \frac{\overline{N}(\theta)}{\widehat{N}(\theta)} - 1 \\ &\leq \frac{\lambda^2}{3} \int \left(\langle \theta, x \rangle^2 \widehat{N}(\theta)^{-1} - 1\right)_+^3 d\bar{\mathbb{P}}(x), \end{aligned}$$

where

$$\lambda^2 = \frac{2[\log(\epsilon^{-1}) + 0.73d]}{(\kappa - 1)n}.$$

## The empirical Gram matrix estimator

Consider  $R = \max_{i=1, \dots, n} \|G^{-1/2} X_i\|$ . With probability at least  $1 - 2\epsilon$ , for any  $\theta \in \mathbb{R}^d$ ,

$$-\delta - \gamma_+ \leq \frac{N(\theta)}{\bar{N}(\theta)} - 1 \leq \frac{\delta + \gamma_-}{1 - \gamma_-},$$

where

$$\delta = \frac{\mu}{1 - 2\mu}, \quad \gamma_- = \frac{2[\log(\epsilon^{-1}) + 0.73d]}{3(\kappa - 1)n}, \quad \text{and } \gamma_+ = \gamma_- R^4(1 + \delta^2).$$

## The empirical Gram matrix estimator

Moreover, when, for some exponent  $p \in ]0, 1]$ ,

$$\mathbb{E} \left\{ \exp \left[ \frac{\alpha}{2} \left( \|G^{-1/2} X\|^{2p} - d^p - \eta^p \right) \right] \right\} \leq 1,$$

with probability at least  $1 - \epsilon$ ,

$$R^4 \leq \left( \frac{2 \log(n/\epsilon)}{\alpha} + d^p + \eta^p \right)^{2/p}.$$

When  $X \in \mathbb{R}^d$  is a Gaussian vector, for any  $\alpha \in ]0, 1[$ , with probability at least  $1 - \epsilon$ ,

$$R^4 \leq \left( \frac{2 \log(n/\epsilon) + d \log[(1 - \alpha)^{-1}]}{\alpha} \right)^2.$$

# The empirical Gram matrix estimator

Polynomial moment assumptions: Define

$$\check{\gamma}_+ = \gamma_-(1 + \delta)^3 \left[ \underbrace{\mathbb{E}(\|G^{-1/2}X\|^6)}_{\geq d^3} + \left( \frac{\mathbb{E}(\|G^{-1/2}X\|^{12})}{n\epsilon} \right)^{1/2} \right].$$

With probability at least  $1 - 3\epsilon$ , for any  $\theta \in \mathbb{S}_d$ ,

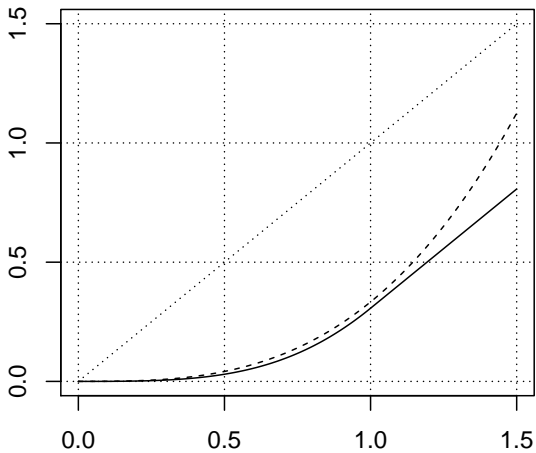
$$-\frac{\delta + \check{\gamma}_+}{1 + \check{\gamma}_+} \leq \frac{N(\theta)}{\bar{N}(\theta)} - 1 \leq \frac{\delta + \gamma_-}{1 - \gamma_-}$$

**Proof.**

Remark that  $\langle \theta, X_i \rangle^2 \leq N(\theta) \|G^{-1/2}X_i\|^2$ , so that

$$\int \left( \langle \theta, x \rangle^2 \hat{N}(\theta)^{-1} - 1 \right)_+^3 d\bar{\mathbb{P}}(x) \leq (1 + \delta)^3 \int \|G^{-1/2}X\|^6 d\bar{\mathbb{P}}(x),$$

and conclude using the Bienaymé Chebyshev inequality.  $\square$



$z \mapsto g(z)$ , compared with  $z \mapsto z^3/3$



# Least squares regression with a random design

Solve

$$\inf_{\theta \in \mathbb{R}^d} \int (y - \langle \theta, x \rangle)^2 dP(x, y), \text{ where } P \in \mathcal{M}_+^1(\mathbb{R}^d \times \mathbb{R}),$$

knowing an i.i.d. sample  $(X_i, Y_i)_{i=1, \dots, n}$ ,  $X_i \in \mathbb{R}^d$ ,  $Y_i \in \mathbb{R}$ .

Introduce the risk function

$$R(\theta) = \int (y - \langle \theta, x \rangle)^2 dP(x, y), \quad \theta \in \mathbb{R}^d,$$

and the homogeneous quadratic form

$$\begin{aligned} N(\theta, \gamma) &= \int (\gamma y - \langle \theta, x \rangle)^2 dP(x, y), \\ &= \begin{pmatrix} \theta \\ \gamma \end{pmatrix}^\top \underbrace{\left[ \int \begin{pmatrix} x \\ -y \end{pmatrix} \begin{pmatrix} x \\ -y \end{pmatrix}^\top dP(x, y) \right]}_{\stackrel{\text{def}}{=} G} \begin{pmatrix} \theta \\ \gamma \end{pmatrix}, \quad \theta \in \mathbb{R}^d, \gamma \in \mathbb{R}. \end{aligned}$$

## Least squares regression with a random design

Let  $\theta_* \in \arg \min_{\theta \in \mathbb{R}^d} R(\theta) = \arg \min_{\theta \in \mathbb{R}^d} N(\theta, 1)$ .

Assume that we have some estimator  $\widehat{N}(\theta, \gamma)$ , such that with probability at least  $1 - 2\epsilon$ , for any  $(\theta, \gamma) \in \mathbb{R}^{d+1}$ ,

$$\left| \frac{N(\theta, \gamma)}{\widehat{N}(\theta, \gamma)} - 1 \right| \leq \eta.$$

Consider  $\widehat{\theta} \in \arg \min_{\theta \in \mathbb{R}^d} \widehat{N}(\theta, 1)$ . With probability at least  $1 - 2\epsilon$ ,

$$\begin{aligned} (1 - \eta) \widehat{N}(\widehat{\theta} - \theta_*, 0) &\leq N(\widehat{\theta} - \theta_*, 0) = R(\widehat{\theta}) - R(\theta_*) \\ &\leq \frac{\eta^2}{1 - \eta} \widehat{N}(\widehat{\theta}, 1) \leq \frac{\eta^2}{1 - \eta} \widehat{N}(\theta_*, 1) \leq \frac{\eta^2}{(1 - \eta)^2} R(\theta_*). \end{aligned}$$

## Least squares regression with a random design

Let us remark that in any case,

$$R(\theta_*) \leq \int y^2 dP(y),$$

and that in the case when  $Y_i = f(X_i) + W_i$ , where  $W_i$  is independent of  $X_i$ , centered and  $\mathbb{E}(W_i) = \sigma^2$ , then

$$R(\theta_*) \leq \int f(x)^2 dP(x) + \sigma^2 \leq \|f\|_\infty^2 + \sigma^2,$$

whereas when  $Y_i = \langle \theta_*, X_i \rangle + W_i$ , then  $R(\theta_*) = \sigma^2$ . Remind that when  $\widehat{N}(\theta, \gamma)$  is the robust estimator described above,  $\eta$  can be taken arbitrarily close to

$$\frac{2\mu}{1-2\mu}, \text{ where } \mu = \sqrt{\frac{2(\kappa-1)}{n} [\log(\epsilon^{-1}) + 0.73 d]} + 6.81 \sqrt{\frac{2\kappa d}{n}}.$$

## Proof

Remark that  $\begin{pmatrix} \theta_* \\ 1 \end{pmatrix}^\top G \begin{pmatrix} \xi \\ 0 \end{pmatrix} = 0$ , for any  $\xi \in \mathbb{R}^d$ .

$$\begin{aligned} R(\hat{\theta}) - R(\theta_*) &= \begin{pmatrix} \hat{\theta} - \theta_* \\ 0 \end{pmatrix}^\top G \begin{pmatrix} \hat{\theta} - \theta_* \\ 0 \end{pmatrix} \\ &= \sup \left\{ \rho^{-1} \begin{pmatrix} \hat{\theta} - \theta_* \\ 0 \end{pmatrix}^\top G \begin{pmatrix} \xi \\ 0 \end{pmatrix}, \quad \xi \in \mathbb{R}^d, \begin{pmatrix} \xi \\ 0 \end{pmatrix}^\top G \begin{pmatrix} \xi \\ 0 \end{pmatrix} = \rho^2 \right\}^2 \\ &= \sup \left\{ \rho^{-1} \begin{pmatrix} \hat{\theta} \\ 1 \end{pmatrix}^\top G \begin{pmatrix} \xi \\ 0 \end{pmatrix}, \quad \xi \in \mathbb{R}^d, N(\xi, 0) = \rho^2 \right\}^2 \\ &= \sup \left\{ \frac{1}{4\rho} [N(\hat{\theta} + \xi, 1) - N(\hat{\theta} - \xi, 1)], \quad \xi \in \mathbb{R}^d, N(\xi, 0) = \rho^2 \right\}^2 \end{aligned}$$

## Proof

$$\begin{aligned} &\leq \sup \left\{ \frac{1}{4\rho} \underbrace{\left[ \widehat{N}(\widehat{\theta} + \xi, 1) - \widehat{N}(\widehat{\theta} - \xi, 1) \right]}_{=0} \right. \\ &\quad \left. + \frac{\eta}{4\rho} \underbrace{\left[ \widehat{N}(\widehat{\theta} + \xi, 1) + \widehat{N}(\widehat{\theta} - \xi, 1) \right]}_{=2\widehat{N}(\widehat{\theta}, 1) + 2\widehat{N}(\xi, 0)}, \quad \xi \in \mathbb{R}^d, N(\xi, 0) = \rho^2 \right\}^2 \\ &\leq \sup \left\{ \frac{\eta}{2\rho} \left[ \widehat{N}(\widehat{\theta}, 1) + (1 - \eta)^{-1} N(\xi, 0) \right], N(\xi, 0) = \rho^2 \right\}^2 \\ &= \left[ \frac{\eta}{2} \left( \frac{\widehat{N}(\widehat{\theta}, 1)}{\rho} + \frac{\rho}{1 - \eta} \right) \right]^2. \end{aligned}$$

Taking  $\rho = \widehat{N}(\widehat{\theta}, 1)^{1/2}(1 - \eta)^{1/2}$ , we get

$$R(\widehat{\theta}) - R(\theta_*) \leq \frac{\eta^2}{1 - \eta} \widehat{N}(\widehat{\theta}, 1).$$

## Estimation of parameters

The kurtosis coefficient  $\kappa$  can be deduced from a separate analysis of the distribution of  $X$  and of the distribution of the "noise"  $Y - \langle \theta_*, X \rangle$ . Indeed

$$\begin{aligned} \kappa &= \sup_{\theta \in \mathbb{R}^d, \gamma \in \mathbb{R}} \frac{\mathbb{E}[(\gamma Y - \langle \theta, X \rangle)^4]}{\mathbb{E}[(\gamma Y - \langle \theta, X \rangle)^2]^2} \\ &\leq 4 \max \left\{ \frac{\mathbb{E}[(Y - \langle \theta_*, X \rangle)^4]}{\mathbb{E}[(Y - \langle \theta_*, X \rangle)^2]^2}, \sup_{\theta \in \mathbb{R}^d} \frac{\mathbb{E}(\langle \theta, X \rangle^4)}{\mathbb{E}(\langle \theta, X \rangle^2)^2} \right\}. \end{aligned}$$

## Estimation of parameters

To analyse the empirical risk minimizer  $\hat{\theta}$ , we need a bound on the empirical Gram estimator  $\bar{N}(\theta, \gamma)$ . Assume that, for some exponents  $p, q \in ]0, 1]$  and some constants  $\alpha > 0$  and  $\beta > 0$ ,

$$\mathbb{E} \left\{ \exp \left[ \frac{\alpha}{2} \left( \|G^{-1/2} X\|^{2p} - d^p - \eta^p \right) \right] \right\} \leq 1,$$
$$\mathbb{E} \left\{ \exp \left[ \frac{\beta}{2} \left( \frac{(Y - \langle \theta_*, X \rangle)^{2q}}{\mathbb{E}[(Y - \langle \theta_*, X \rangle)^2]^q} - 1 - \gamma^q \right) \right] \right\} \leq 1.$$

Remark that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left( \frac{(\gamma Y_i - \langle \theta, X_i \rangle)^2}{\hat{N}(\theta, \gamma)} - 1 \right)_+^3 \\ \leq \frac{\bar{N}(\theta, \gamma)}{\hat{N}(\theta, \gamma)} \max_{i=1, \dots, n} \frac{(\gamma Y_i - \langle \theta, X_i \rangle)^4}{N(\theta, \gamma)^2} (1 + \delta)^2. \end{aligned}$$

## Estimation of parameters

Remark also that

$$\begin{aligned} \frac{(\gamma Y_i + \langle \theta, X_i \rangle)^4}{N(\theta, \gamma)^2} &\leq 4 \max \left\{ \frac{(Y_i - \langle \theta_*, X_i \rangle)^4}{(\int (y - \langle \theta_*, x \rangle)^2 dP(x, y))^2}, \|G^{-1/2} X_i\|^4 \right\} \\ &\leq 4 \max \left\{ \left( d^p + \eta^p + 2\alpha^{-1} \log(n/\epsilon) \right)^{2/p}, \left( 1 + \gamma^q + 2\beta^{-1} \log(n/\epsilon) \right)^{2/q} \right\} \end{aligned}$$

Consider

$$\begin{aligned} \tilde{\gamma}_+ &= 4\gamma_- (1 + \delta)^2 \\ &\times \max \left\{ \left( d^p + \eta^p + 2\alpha^{-1} \log(n/\epsilon) \right)^{2/p}, \left( 1 + \gamma^q + 2\beta^{-1} \log(n/\epsilon) \right)^{2/q} \right\}, \end{aligned}$$

With probability at least  $1 - 4\epsilon$ , for any  $(\theta, \gamma) \in \mathbb{R}^{d+1}$ ,

$$-\delta - \tilde{\gamma}_+ \leq \frac{N(\theta, \gamma)}{\bar{N}(\theta, \gamma)} - 1 \leq \frac{\delta + \gamma_-}{1 - \gamma_-}.$$



## Estimation of parameters

**Polynomial moment assumptions.** Define  $G_d = \mathbb{E}(XX^\top)$ , so

that  $\theta^\top G_d \theta = \begin{pmatrix} \theta \\ 0 \end{pmatrix}^\top G \begin{pmatrix} \theta \\ 0 \end{pmatrix}$  and consider

$$\check{\gamma}_+ = 8\gamma_-(1+\delta)^3 \max \left\{ \mathbb{E}(\|G_d^{-1/2}X\|^6) + \left( \frac{\mathbb{E}(\|G_d^{-1/2}X\|^{12})}{n\epsilon} \right)^{1/2}, \right. \\ \left. \frac{\mathbb{E}[(Y - \langle \theta_*, X \rangle)^6]}{\mathbb{E}[(Y - \langle \theta_*, X \rangle)^2]^3} + \left( \frac{\mathbb{E}[(Y - \langle \theta_*, X \rangle)^{12}]}{n\epsilon \mathbb{E}[(Y - \langle \theta_*, X \rangle)^2]^6} \right)^{1/2} \right\}$$

With probability at least  $1 - 4\epsilon$ , for any  $(\theta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ ,

$$-\frac{\delta + \check{\gamma}_+}{1 + \check{\gamma}_+} \leq \frac{N(\theta, \gamma)}{\bar{N}(\theta, \gamma)} - 1 \leq \frac{\delta + \gamma_-}{1 - \gamma_-}.$$

## Dimension free bounds (*Ilaria Giulini*)

With probability  $1 - 2\epsilon$ , for any  $\theta \in \mathbb{S}_d$ , and some robust estimator  $\widehat{N}$ ,

$$\left| \frac{\widehat{N}(\theta) \vee \sigma}{N(\theta) \vee \sigma} - 1 \right| \leq \frac{\mu(N(\theta) \vee \sigma)}{1 - 4\mu(N(\theta) \vee \sigma)},$$

where, for  $n \leq 10^{20}$ ,

$$\mu(t) = \sqrt{\frac{2.13(\kappa - 1)}{n} \left[ \log(\epsilon^{-1}) + 4.64 + \frac{0.73 \times \mathbf{Tr}(G)}{t} \right]} + \sqrt{\frac{118\kappa \mathbf{Tr}(G)}{nt}},$$

and  $\sigma$  is such that  $8\mu(\sigma) \leq 1$ . Let us recall that

$$\mathbf{Tr}(G) = \int \|x\|^2 dP(x) = \sum_{i=1}^d N(\theta_i), \text{ for any orthonormal basis } (\theta_i, 1 \leq i \leq n).$$

## Proof of main result

Replacing  $\mathbb{R}^d$  by  $\mathbf{Im}(G)$ , we can assume that  $\mathbf{Ker}(G) = \{0\}$ , without loss of generality. For any  $\lambda > 0$ , any  $x, \theta \in \mathbb{R}^d$ ,

$$\begin{aligned} \psi\left\{\lambda\left[\langle\theta, x\rangle^2-1\right]\right\} &\leq \int \log \left\{1+\lambda\left[\langle\theta', x\rangle^2-1-\frac{\|x\|^2}{\beta}\right]\right. \\ &+ \left.\frac{\lambda^2}{2}\left[\langle\theta', x\rangle^2-1-\frac{\|x\|^2}{\beta}\right]^2+\frac{c\lambda^2\|x\|^2}{\beta}\left(\langle\theta', x\rangle^2+\frac{\|x\|^2}{2\beta}\right)\right\} d\pi_{\theta}\left(\theta'\right), \end{aligned}$$

where  $\pi_{\theta} = \mathcal{N}\left(\theta, \beta^{-1} G^{-1}\right)$  is a Gaussian perturbation of the parameter  $\theta$  and

$$c = \frac{15}{8 \log (2)(\sqrt{2}-1)} \exp \left(\frac{1+2 \sqrt{2}}{2}\right) \leq 44.3.$$

## Proof of main result

Let  $\nu \in \mathcal{M}_+^1(\Theta)$  be a prior probability measure on the parameter space  $\Theta$ . Under suitable assumptions, with probability at least  $1 - \epsilon$  for any posterior probability measure  $\rho \in \mathcal{M}_+^1(\Theta)$  such that  $\mathcal{K}(\rho, \nu) < +\infty$ ,

$$\iint \log[1 + f(x, \theta)] d\rho(\theta) d\bar{\mathbb{P}}(x) \leq \iint f(x, \theta) d\mathbb{P}(x) d\rho(\theta) + \frac{\mathcal{K}(\rho, \nu) - \log(\epsilon)}{n}.$$

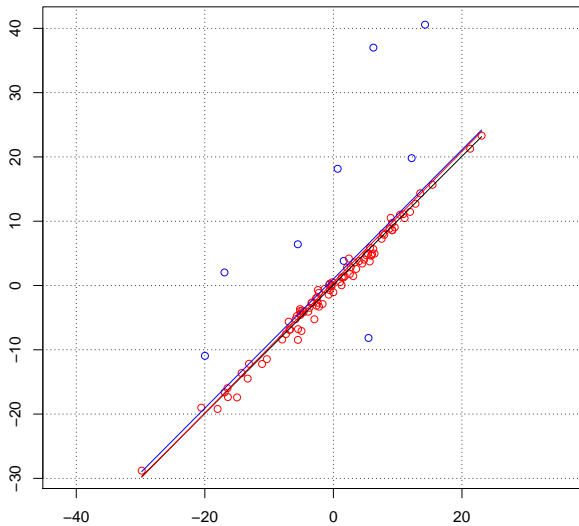
## Proof of main result

With probability at least  $1 - 2\epsilon$ , for any  $\theta \in \mathbb{R}^d$ ,

$$r_\lambda(\theta) \leq \xi[N(\theta) - 1]^2 + (1 + \mu)[N(\theta) - 1] + \mu,$$

$$r_\lambda(\theta) \geq -\xi[N(\theta) - 1]^2 + (1 - \mu)[N(\theta) - 1] - \mu,$$

where  $\xi = \frac{\kappa\lambda}{2}$ . We end the proof by showing that when  $N(\theta) > 0$ , with probability  $1 - 2\epsilon$ ,  $r_\lambda(\hat{\alpha}\theta) = 0$  and solving the resulting inequalities in  $N(\hat{\alpha}\theta) = \frac{N(\theta)}{\hat{N}(\theta)}$ .



n1 = 10, n2 = 90

