

# NEW PAC-BAYESIAN BOUNDS FOR $k$ -MEANS ALGORITHMS

OLIVIER CATONI

Olivier.Catoni@ensae.fr

<http://ocatoni.perso.math.cnrs.fr/>

CREST, CNRS — UMR 9194

Université Paris-Saclay

Three-day meeting of statisticians in Paris

I.H.P., July 18-20 2022

*Joint work with Gautier Appert*

# The $k$ -means criterion

Consider a random variable  $X \in H$ , where  $H$  is a separable Hilbert space and the quantization problem

$$\inf_{c \in H^k} \mathbb{P}_X \left( \min_{j \in \llbracket 1, k \rrbracket} \|X - c_j\|^2 \right).$$

Given a sample  $\bar{X} = (X_1, \dots, X_n)$  made of  $n$  independent copies of  $X$ , we want an estimator  $\widehat{c}(\bar{X}) \in H^k$  such that

$$\mathbb{P}_X \left( \min_{j \in \llbracket 1, k \rrbracket} \|X - \widehat{c}_j\|^2 \right).$$

is small. Since it is a r. v. we can bound either its mean or its deviations. The aim of this talk is to present a series of ideas that lead to new bounds and new estimators.

## A min-linear criterion

The first thing we propose is to rewrite the criterion as a min-linear problem.

$$\begin{aligned}\min_{j \in \llbracket 1, k \rrbracket} \|X - c_j\|^2 &= \min_{j \in \llbracket 1, k \rrbracket} \|X\|^2 + \|c_j\|^2 - 2\langle X, c_j \rangle \\ &= \min_{j \in \llbracket 1, k \rrbracket} \|W_1\|^2/4 + \langle \theta_j, W \rangle,\end{aligned}$$

where  $\theta_j = (c_j, \gamma^{-1}\|c_j\|^2)$ ,  $W = (-2X, \gamma) \in H \times \mathbb{R}$  and  $W_1 = -2X \in H$  is the first component of  $W$ .

## A more general loss function

More generally we will consider a loss function

$$f(\theta, w) \in \mathbb{R}, \quad \theta \in H^m, w \in H,$$

such that

$$|f(\theta', w) - f(\theta, w)| \leq (a + b\|w\|^{\alpha_2}) \max_{j \in \llbracket 1, k \rrbracket} \left| \sum_{\ell=1}^m A_{j,\ell} \langle \theta'_\ell - \theta_\ell, w \rangle \right|^{\alpha_1},$$
$$\alpha_1 \in ]0, 1], \alpha_2 \in \mathbb{R}_+, A \in \mathbb{R}^{k \times m}.$$

The minimization problem

$$\inf_{\theta \in \Theta} \mathbb{P}_W[f(\theta, W)]$$

covers in particular the case

$$\inf_{c \in H^k} \mathbb{P}_X \left( \mu(X) \min_{j \in \llbracket 1, k \rrbracket} \|X - c_j\|^{2\alpha_1} \right), \quad \alpha_1 \in ]0, 1], \quad (1)$$

$$\text{where } \mu(X) \in [0, a' + b'\|X\|^{\alpha_2}], \quad a', b', \alpha_2 \in \mathbb{R}_+. \quad (2)$$

# Structured $k$ -means

We use  $k$  centers depending on  $T$  parameters.

$$\inf_{\xi \in \Xi} \mathbb{P}_X \left( \mu(X) \min_{j \in \llbracket 1, k \rrbracket} \left\| X - \sum_{t=1}^T B_{j,t} \xi_t \right\|^{2\alpha_1} \right)$$

The  $(\theta, W)$  parametrization is given by  $m = 2T + T(T-1)/2$ ,  
 $W = (-2X, \gamma) \in H \times \mathbb{R}$ , and for  $1 \leq t \leq T$ ,  $1 \leq s < T$ ,

$$A_{j,t} = B_{j,t}, \quad \theta_t = (\xi_t, 0),$$

$$A_{j,T+t} = B_{j,t}^2, \quad \theta_{T+j} = \left( 0, \gamma^{-1} \|\xi_t\|^2 \right),$$

$$A_{j,2T+t(t-1)/2+s-1} = 2B_{j,t}B_{j,s}, \quad \theta_{2T+t(t-1)/2+s-1} = \left( 0, \gamma^{-1} \langle \xi_t, \xi_s \rangle \right)$$

With this choice of coordinates

$$\left\| X - \sum_{t=1}^T B_{j,t} \xi_t \right\|^{2\alpha_1} = \left| \|W_1\|^2/4 + \left\langle \sum_{\ell=1}^m A_{j,\ell} \theta_\ell, W \right\rangle \right|^{\alpha_1}$$

## Bound the excess risk

Consider a non random reference value of the parameter  $\theta^\star$  and work on the excess risk

$$h_0(\theta, w) = f(\theta, w) - f(\theta^\star, w).$$

Consider a statistical sample  $\overline{W} = (W_1, \dots, W_n)$  made of  $n$  independent copies of  $W \in H$ . From a bound in expectation

$$\mathbb{P}_{\overline{W}} \left[ \sup_{\theta \in \Theta} \left( (\mathbb{P}_W h_0)(\theta) - B(\theta, \overline{W}) \right) \right] \leq 0 \quad (3)$$

and the  $\epsilon$ -minimizer

$$B(\widehat{\theta}, \overline{W}) \leq \inf_{\theta \in \Theta} B(\theta, \overline{W}) + \epsilon,$$

we get

$$\mathbb{P}_{\overline{W}} \left[ \mathbb{P}_W [f(\widehat{\theta}, W)] \right] \leq \inf_{\theta^\star \in \Theta} \left[ \mathbb{P}_W [f(\theta^\star, W)] + \mathbb{P}_{\overline{W}} [B(\theta^\star, \overline{W})] \right] + \epsilon.$$

when the choice of  $\theta^\star$  is optimal.

# Deviation bounds

From a deviation bound

$$\mathbb{P}_{\overline{W}} \left[ \sup_{\theta \in \Theta} \left( (\mathbb{P}_W h_0)(\theta) - B(\theta, \overline{W}) - \log(\delta^{-1}) \right) \leq 0 \right] \geq 1 - \delta. \quad (4)$$

and the  $\epsilon$ -minimizer

$$B(\widehat{\theta}, \overline{W}) \leq \inf_{\theta \in \Theta} B(\theta, \overline{W}) + \epsilon,$$

we get

$$\mathbb{P}_{\overline{W}} \left[ \mathbb{P}_W [f(\widehat{\theta}, W)] \leq \inf_{\theta^* \in \Theta} \mathbb{P}_W [f(\theta^*, W)] + B(\theta^*, \overline{W}) + \epsilon \right] \geq 1 - \delta.$$

when  $\theta^*$  is optimal.



## Deduce everything from exponential moments

We will deduce both bounds in expectation and deviation bounds from bounds on exponential moments of the form

$$\mathbb{P}_{\bar{W}}\left\{\exp\left[\lambda\left(\sup_{\theta\in\Theta}(\mathbb{P}_W h_0)(\theta) - B(\theta, \bar{W})\right)\right]\right\} \leq 1, \quad (5)$$

where  $\lambda$  is a positive real exponent. This implies (3) and

$$\mathbb{P}_{\bar{W}}\left[\sup_{\theta\in\Theta}\left((\mathbb{P}_W h_0)(\theta) - B(\theta, \bar{W}) - \log(\delta^{-1})/\lambda\right) \leq 0\right] \geq 1 - \delta. \quad (6)$$

# Lemma on the expectation of the supremum of Gaussian random variables

Let  $\alpha$  be some positive real exponent, and let  $\epsilon_j$ ,  $1 \leq j \leq k$  be  $k$  centered Gaussian random variables with variances  $\sigma_j^2$ ,  $1 \leq j \leq k$ . We do not assume independence nor the fact that the vector  $(\epsilon_1, \dots, \epsilon_k)$  has a joint Gaussian distribution. Let  $m_j \in \mathbb{R}$ ,  $1 \leq j \leq k$  be mean parameters. Assume that  $k \geq \exp(\alpha - 1)$ .

$$\mathbb{P}_{\epsilon_1, \dots, \epsilon_k} \left( \max_{j \in \llbracket 1, k \rrbracket} |m_j + \epsilon_j|^\alpha \right) \leq \left( \sqrt{2 \log(2k)} \max_{j \in \llbracket 1, k \rrbracket} \sigma_j + \max_{j \in \llbracket 1, k \rrbracket} |m_j| \right)^\alpha.$$

# Gaussian perturbations

- Assume w.l.o.g. that  $H = \ell^2$ .
- Let  $\rho_{\theta' | \theta} = \bigotimes_{\ell=1}^m \left( \bigotimes_{i \in \mathbb{N}} \mathcal{N}(\theta_{\ell, i}, \sigma_{\ell}^2 / \beta) \right) : (\mathbb{R}^{\mathbb{N}})^m \rightarrow \mathcal{M}_+^1((\mathbb{R}^{\mathbb{N}})^m)$ .
- Let  $\langle \theta, w \rangle = \begin{cases} \lim_{s \rightarrow +\infty} \sum_{i=0}^s \theta_i w_i, & \text{when } \overline{\lim}_{s \rightarrow +\infty} \sum_{i=0}^s \theta_i w_i = \underline{\lim}_{s \rightarrow +\infty} \sum_{i=0}^s \theta_i w_i \in \mathbb{R}, \\ 0, & \text{otherwise} \end{cases}$

be a non bilinear but measurable extension of the scalar product from  $\ell^2$  to  $\mathbb{R}^{\mathbb{N}}$ .

- The linear operator  $\rho$  operates on suitably integrable functions of  $\theta$  and  $w$  according to the rule

$$(\rho f)(\theta, w) = \rho_{\theta' | \theta}(f(\theta', w)).$$

# PAC-Bayesian lemma

Consider the increasing function  $g(t) = \frac{2}{t^2} [\exp(t) - 1 - t]$ ,  $t \in \mathbb{R}$  defined by continuity at  $t = 0$ , where  $g(0) = 1$ . For any measurable bounded real valued function  $h(w)$ ,  $w \in H$ , such that  $\sup_{w \in H} |h(w)| \leq \eta$ , for any positive exponent  $\lambda$ ,

$$\mathbb{P}_{\bar{W}} \left\{ \exp \left[ n\lambda (\bar{\mathbb{P}}_W - \mathbb{P}_W)h - n \frac{\lambda^2}{2} g(2\eta\lambda) \mathbb{P}_W [(h - \mathbb{P}_W h)^2] \right] \right\} \leq 1.$$

If  $h(\theta, w) \in \mathbb{R}$ ,  $\theta \in H^m$ ,  $w \in H$  depends also on  $\theta$  and if  $\sup_{\theta \in H^m, w \in H} |h(\theta, w)| \leq \eta$ ,

$$\mathbb{P}_{\bar{W}} \left\{ \exp \left[ \sup_{\rho \in \mathcal{M}_+^1(\Theta)} n\lambda (\bar{\mathbb{P}}_W - \mathbb{P}_W) \rho h - n \frac{\lambda^2}{2} g(2\eta\lambda) \rho \mathbb{P}_W [(h - \mathbb{P}_W h)^2] - \mathcal{K}(\rho, \pi) \right] \right\} \leq 1.$$

# Multi-scale decomposition of the excess risk

The decomposition

$$h_0 = \left( I - \rho + \sum_{q=1}^p (\rho^{2^{q-1}} - \rho^{2^q}) + \rho^{2^p} \right) h_0, \quad (7)$$

can be written as

$$h_0 = h_{p+1} + \sum_{q=0}^p (h_q - h_{q+1}), \quad (8)$$

where

$$h_0(\theta, w) = f(\theta, w) - f(\theta^\star, w), \quad \theta \in H^m, w \in H,$$

( with implicit dependence on  $\theta^\star$  ),

$$h_q = \rho^{2^{q-1}} h_0 = \rho(2^{-q+1}\beta)h_0, \quad 1 \leq q \leq p+1.$$

# Thresholds

We will also need a sequence of truncation operators  $T_q$  using threshold levels  $\eta_q > 0$ . They are defined as

$$\begin{aligned}T_q(z) &= \min\{\eta_q, z\}, \quad z \in \mathbb{R}_+, \\T_q(w) &= T_q(\|w\|) \frac{w}{\|w\|}, \quad w \in H, \\(T_q f)(\theta, w) &= f(\theta, T_q(w)), \quad \theta \in H^m, w \in H,\end{aligned}$$

and the operator  $T_q f$  acting on functions is linear ( although the truncation acting on vectors or positive real numbers is not ). Moreover,  $T_q$  commutes with  $\rho$  :

$$T_q \rho = \rho T_q.$$

We also have the composition rules

$$\mathbb{P}_W \rho = \rho \mathbb{P}_W, \quad \text{and} \quad \mathbb{P}_W T_q = \mathbb{P}_{T_q(W)}.$$

At stage  $p + 1$ , we will need the non linear threshold operator  $T_{p+1}$  defined as

$$(T_{p+1} h)(\theta, w) = \min\left\{\eta_{p+1}, \max\{-\eta_{p+1}, h(\theta, w)\}\right\},$$

## Bounding the expected excess risk

We will prove a bound that compares the expected excess risk  $\mathbb{P}_W h_0$  with the possibly truncated empirical excess risk  $\overline{\mathbb{P}}_W T_0 h_0$ . We decompose the risk into

$$\mathbb{P}_W h_0 = \mathbb{P}_W (\mathbf{I} - T_0) h_0 + (\mathbb{P}_W - \overline{\mathbb{P}}_W) T_0 h_0 + \overline{\mathbb{P}}_W T_0 h_0,$$

leading to

$$\mathbb{P}_W h_0 = \mathbb{P}_W (\mathbf{I} - T_0) h_0 + (\mathbb{P}_W - \overline{\mathbb{P}}_W) T_0 \left( h_{p+1} + \sum_{q=0}^p (h_q - h_{q+1}) \right) + \overline{\mathbb{P}}_W T_0 h_0. \quad (9)$$

We will analyze each term of this decomposition separately.

## Bounding $A_{q,0} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)T_0(h_q - h_{q+1})$

Let us deal first with the case when  $q \geq 1$ . Decompose further  $A_{q,0}$  into

$$A_{q,0} = A_{q,1} + A_{q,2},$$

where

$$A_{q,1} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)T_0(\mathbf{I} - T_q)(h_q - h_{q+1})$$

$$A_{q,2} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)T_q(h_q - h_{q+1}).$$

As  $h_q = \rho_q h_0$ ,

$$h_q - h_{q+1} = \rho_q(\mathbf{I} - \rho_q)h_0 = \rho_q(\mathbf{I} - \rho_q)f,$$

since  $h_0(\theta, w) = f(\theta, w) - f(\theta^*, w)$ , where  $f(\theta^*, w)$  does not depend on  $\theta$ . Thus

$$A_{q,2} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)\rho_q(\mathbf{I} - \rho_q)T_q f.$$



Moreover

$$\begin{aligned} |(I - \rho_q)T_q f| &= |\rho(\beta_q)_{\theta' | \theta} [f(\theta, T_q(w)) - f(\theta', T_q(w))]| \\ &\leq \rho(\beta_q)_{\theta' | \theta} [|f(\theta, T_q(w)) - f(\theta', T_q(w))|], \end{aligned}$$

so that

$$\begin{aligned} |(I - \rho_q)T_q f| &\leq (a + bT_q(\|w\|)^{\alpha_2}) \rho(\beta_q)_{\theta' | \theta} \left( \max_{j \in \llbracket 1, k \rrbracket} \left| \langle T_q(w), \sum_{\ell=1}^m A_{j,\ell} (\theta_\ell - \theta'_\ell) \rangle \right|^{\alpha_1} \right) \\ &\leq (a + bT_q(\|w\|)^{\alpha_2}) (\sqrt{2 \log(2k)} / \beta \sigma_\star T_q(\|w\|))^{\alpha_1}, \end{aligned}$$

where

$$\sigma_\star = \max_{j \in \llbracket 1, k \rrbracket} \sqrt{\sum_{\ell=1}^m A_{j,\ell}^2 \sigma_\ell^2}.$$

Introduce the bounds

$$B_{q,0}(w) = B_{q,1}(\|w\|)$$

where

$$B_{q,1}(t) = (a + bt^{\alpha_2})(2 \log(2k)/\beta_q)^{\alpha_1/2} (\sigma_{\star} t)^{\alpha_1}, \quad t \in \mathbb{R}_+.$$

We get

$$|(I - \rho_q)T_q f| \leq T_q B_{q,0}.$$

Let us put

$$K_q(\theta) = \sum_{\ell=1}^m \frac{\beta_q}{2\sigma_{\ell}^2} \|\theta_{\ell} - \tilde{\theta}_{\ell}\|^2 = \mathcal{K}(\rho(\beta_q)_{\theta' | \theta}, \rho(\beta_q)_{\theta' | \theta = \tilde{\theta}}), \quad \theta \in H^m,$$

where  $\tilde{\theta} \in H^m$  is a non random reference that may or may not be equal to  $\theta^{\star}$ .

For any  $\lambda_q > 0$ ,

$$\mathbb{P}_{\overline{W}} \left\{ \exp \left[ \sup_{\theta \in \Theta} n\lambda A_{q,2} - n \frac{\lambda^2}{2} g(2\lambda B_{q,1}(\eta_q)) \mathbb{P}_W(B_{q,0}(W)^2) - K_q(\theta) \right] \right\} \leq 1.$$

Let us now bound

$$A_{q,1} = (\mathbb{P}_W - \overline{\mathbb{P}}_W)T_0(\mathbf{I} - T_q)\rho_q(\mathbf{I} - \rho_q)f.$$

We can write

$$\begin{aligned} |A_{q,1}| &\leq (\mathbb{P}_W + \overline{\mathbb{P}}_W)T_0(\mathbf{I} + T_q) \left[ \mathbf{1}(\|W\| \geq \eta_q)\rho_q |(\mathbf{I} - \rho_q)f| \right] \\ &\leq 2(\mathbb{P}_W + \overline{\mathbb{P}}_W)T_0 \left[ \mathbf{1}(\|W\| \geq \eta_q)B_{q,0} \right] \\ &\leq 2(\mathbb{P}_W + \overline{\mathbb{P}}_W) \left[ T_0 B_{q,0}^2 / B_{q,1}(\eta_q) \right]. \end{aligned}$$

## Bounding $A_{0,0}$

In the case when  $q = 0$ ,

$$A_{0,0} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)T_0(\mathbf{I} - \rho)f,$$

so that

$$|A_{0,0}| \leq (\mathbb{P}_W + \bar{\mathbb{P}}_W)T_0|(\mathbf{I} - \rho)f| \leq (\mathbb{P}_W + \bar{\mathbb{P}}_W)T_0B_{1,0}.$$

Let us now come to

$$A_{-1,0} = \mathbb{P}_W(\mathbf{I} - T_0)h_0.$$

Remark that

$$|h_0(\theta, w)| \leq B_{-1,0}(\theta, w),$$

where

$$B_{-1,0}(\theta, w) = B_{-1,1}(\theta, \|w\|),$$

where

$$B_{-1,1}(\theta, t) = (a + bt^{\alpha_2}) \left( t \max_{j \in \llbracket 1, k \rrbracket} \left\| \sum_{\ell=1}^m A_{j,\ell} (\theta_\ell - \theta_\ell^\star) \right\| \right)^{\alpha_1}$$

is not decreasing in  $t \in \mathbb{R}_+$ . Therefore,

$$\begin{aligned} |A_{-1,0}| &= |\mathbb{P}_W(\mathbf{I} - T_0)h_0| \leq \mathbb{P}_W(|(\mathbf{I} - T_0)h_0|) \\ &= \mathbb{P}_W(|(\mathbf{I} - T_0)[\mathbf{1}(\|w\| \geq \eta_0)h_0]|) \\ &\leq \mathbb{P}_W((\mathbf{I} + T_0)[\mathbf{1}(\|w\| \geq \eta_0)|h_0|]) \\ &\leq \mathbb{P}_W((\mathbf{I} + T_0)[\mathbf{1}(\|w\| \geq \eta_0)B_{-1,0}]) \\ &\leq 2\mathbb{P}_W([\mathbf{1}(\|w\| \geq \eta_0)B_{-1,0}]) \\ &\leq 2\mathbb{P}_W\left(B_{-1,0}(\theta, W)^2 / B_{-1,1}(\theta, \eta_0)\right). \end{aligned}$$

Finally, let us bound

$$A_{p+1,3} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)T_0h_{p+1} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)\rho_{p+1}T_0h_0.$$

Introduce the non linear threshold operator  $T_{p+1}$  defined by the equation

$$(T_{p+1}h)(\theta, w) = \min\left\{\eta_{p+1}, \max\{-\eta_{p+1}, h(\theta, w)\}\right\},$$

where  $\eta_{p+1} > 0$  is some threshold level. We can decompose  $A_{p+1,3}$  into

$$A_{p+1,3} = A_{p+1,4} + A_{p+1,5},$$

where

$$A_{p+1,4} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)\rho_{p+1}T_{p+1}T_0h_0$$

and

$$A_{p+1,5} = (\mathbb{P}_W - \bar{\mathbb{P}}_W)\rho_{p+1}(\mathbf{I} - T_{p+1})T_0h_0.$$

Remark first that

$$\begin{aligned}
|A_{p+1,5}| &\leq (\mathbb{P}_W + \bar{\mathbb{P}}_W)\rho_{p+1}|(\mathbf{I} - T_{p+1})T_0h_0| \\
&= (\mathbb{P}_W + \bar{\mathbb{P}}_W)\rho_{p+1}[ (|T_0h_0| - \eta_{p+1})_+ ] \\
&\leq (\mathbb{P}_W + \bar{\mathbb{P}}_W)\rho_{p+1}[\mathbb{1}(|T_0h_0| \geq \eta_{p+1})|T_0h_0|] \\
&\leq \frac{1}{\eta_{p+1}}(\mathbb{P}_W + \bar{\mathbb{P}}_W)\rho_{p+1}(|T_0h_0|^2).
\end{aligned}$$

Define  $f^\star(\theta, w) = f(\theta^\star, w)$  depending on  $w$  only. Remark that

$$\begin{aligned}
\rho_{p+1}[|T_0h_0|^2] &= T_0\rho_{p+1}[(f - f^\star)^2] \\
&\leq T_0\rho_{p+1}\left[ (a + b\|w\|^{\alpha_2})^2 \left( \max_{j \in \llbracket 1, k \rrbracket} \left| \langle w, \sum_{\ell=1}^m A_{j,\ell}(\theta_\ell - \theta_\ell^\star) \rangle \right| \right)^{2\alpha_1} \right] \\
&\leq T_0(a + b\|w\|^{\alpha_2})^2 \left( (2 \log(2k) / \beta_{p+1})^{1/2} \sigma_\star \|w\| + \max_{j \in \llbracket 1, k \rrbracket} \left| \langle w, \sum_{\ell=1}^m A_{j,\ell}(\theta_\ell - \theta_\ell^\star) \rangle \right| \right)^{2\alpha_1} \\
&\leq T_0 B_{p+1,0},
\end{aligned}$$

where

$$B_{p+1,0}(\theta, w) = (a + b\|w\|^{\alpha_2})^2 \|w\|^{2\alpha_1} \\ \times \left( (2 \log(2k)/\beta_{p+1})^{1/2} \sigma_{\star} + \max_{j \in \llbracket 1, k \rrbracket} \left\| \sum_{\ell=1}^m A_{j,\ell} (\theta_{\ell} - \theta_{\ell}^{\star}) \right\| \right)^{2\alpha_1}.$$

Therefore

$$|A_{p+1,5}| \leq (\mathbb{P}_W + \overline{\mathbb{P}}_W) (T_0 B_{p+1,0} / \eta_{p+1}).$$



## Low frequency component

The last term to bound is  $A_{p+1,4}$ . We get for any  $\lambda_{p+1} > 0$

$$\mathbb{P}_{\bar{W}} \left\{ \exp \left[ \sup_{\theta \in \Theta} n \lambda_{p+1} A_{p+1,4} - n \frac{\lambda_{p+1}^2}{2} g(2\lambda_{p+1} \eta_{p+1}) \mathbb{P}_W \rho_{p+1}((T_0 h_0)^2) - K_{p+1}(\theta) \right] \right\} \leq 1,$$

so that

$$\mathbb{P}_{\bar{W}} \left\{ \exp \left[ \sup_{\theta \in \Theta} \left( n \lambda_{p+1} A_{p+1,4} - n \frac{\lambda_{p+1}^2}{2} g(2\lambda_{p+1} \eta_{p+1}) \mathbb{P}_W (T_0 B_{p+1,0}) - K_{p+1}(\theta) \right) \right] \right\} \leq 1.$$

# Summary

We have written

$$\mathbb{P}_W h_0 = \overline{\mathbb{P}}_W T_0 h_0 + A_{-1,0} + A_{0,0} + A_{p+1,4} + A_{p+1,5} + \sum_{q=1}^p (A_{q,1} + A_{q,2})$$

and provided bounds for each  $A_{q,\ell}$ , either almost sure bounds or exponential moment bounds.

Based on the bounds

$$B_{q,0}(w) = B_{q,1}(\|w\|)$$

where  $B_{q,1}(t) = (a + bt^{\alpha_2})(2 \log(2k)/\beta_q)^{\alpha_1/2}(\sigma_\star t)^{\alpha_1}$ ,  $t \in \mathbb{R}_+$

$$B_{q,0}(w) = \xi(\|w\|)\mathfrak{S}_3^{\alpha_1}\beta_q^{-\alpha_1/2},$$

$$\begin{aligned} B_{-1,0}(\theta, w) &= (a + b\|w\|^{\alpha_2}) \left( \|w\| \max_{j \in \llbracket 1, k \rrbracket} \left\| \sum_{\ell=1}^m A_{j,\ell}(\theta_\ell - \theta_\ell^\star) \right\| \right)^{\alpha_1} \\ &= \xi(\|w\|)\mathfrak{S}_2^{\alpha_1}, \end{aligned}$$

$$\begin{aligned} B_{p+1,0}(\theta, w) &= (a + b\|w\|^{\alpha_2})^2 \|w\|^{2\alpha_1} \\ &\quad \times \left( (2 \log(2k)/\beta_{p+1})^{1/2} \sigma_\star + \max_{j \in \llbracket 1, k \rrbracket} \left\| \sum_{\ell=1}^m A_{j,\ell}(\theta_\ell - \theta_\ell^\star) \right\| \right)^{2\alpha_1}, \\ &= \xi(\|w\|)^2 \left( \mathfrak{S}_3 \beta_{p+1}^{-1/2} + \mathfrak{S}_2 \right)^{2\alpha_1} \end{aligned}$$

$$\text{and } K_q(\theta) = \sum_{\ell=1}^m \frac{\beta_q}{2\sigma_\ell^2} \|\theta_\ell - \tilde{\theta}_\ell\|^2,$$

# What we proved

We proved that

$$\mathbb{P}_{\overline{W}} \left\{ \exp \left[ \sup_{\theta \in \Theta} n\lambda A_{q,2} - n \frac{\lambda^2}{2} g(2\lambda B_{q,1}(\eta_q)) \mathbb{P}_W(B_{q,0}(W)^2) - K_q(\theta) \right] \right\} \leq 1,$$

$$\mathbb{P}_{\overline{W}} \left\{ \exp \left[ \sup_{\theta \in \Theta} \left( n\lambda A_{p+1,4} - n \frac{\lambda^2}{2} g(2\lambda \eta_{p+1}) \mathbb{P}_W(T_0 B_{p+1,0}) - K_{p+1}(\theta) \right) \right] \right\} \leq 1.$$

We also bounded the remaining terms by

$$\begin{aligned} |A_{q,1}| &\leq 2(\mathbb{P}_W + \overline{\mathbb{P}}_W) T_0 [\mathbf{1}(\|W\| \geq \eta_q) B_{q,0}] \\ &\leq 2(\mathbb{P}_W + \overline{\mathbb{P}}_W) [T_0 B_{q,0}(W)^2 / B_{q,1}(\eta_q)] \end{aligned}$$

$$|A_{0,0}| \leq (\mathbb{P}_W + \overline{\mathbb{P}}_W) T_0 B_{1,0}.$$

$$|A_{-1,0}| \leq 2\mathbb{P}_W \left( B_{-1,0}(\theta, W)^2 / B_{-1,1}(\theta, \eta_0) \right),$$

$$|A_{p+1,5}| \leq (\mathbb{P}_W + \overline{\mathbb{P}}_W) (T_0 B_{p+1,0} / \eta_{p+1})$$

## Bound in expectation

Assume that  $n \geq 2\mathcal{S}_1$  and set  $\eta_0 = +\infty$ . We get

$$\mathbb{P}_{\bar{W}} \left\{ \sup_{\theta \in \Theta} \mathbb{P}_W h_0 - \bar{\mathbb{P}}_W h_0 - \gamma \right\} \leq 0, \text{ where}$$

$$\begin{aligned} \gamma = 2 & \left[ \left( \frac{g(1)}{2} + 8 \right)^{1/2} \mathbb{P}_W (\xi(\|W\|)^2)^{1/2} + \mathbb{P}_W (\xi(\|W\|)) \right] \\ & \times \mathcal{S}_3^{\alpha_1} \left( \frac{\log(2n/\mathcal{S}_1)}{\log(2)} \right)^{\alpha_1} \left( \frac{\mathcal{S}_1}{n} \right)^{\alpha_1/2} \\ & + 2 \left[ \left( \frac{g(1)}{2} + 4 \right) \mathbb{P}_W (\xi(\|W\|)^2) (\mathcal{S}_3 + \mathcal{S}_2)^{2\alpha_1} \frac{\mathcal{S}_1}{n} \right]^{1/2} \end{aligned}$$

$$\text{and } \xi(t) = (a + bt^{\alpha_2})t^{\alpha_1}, \quad \mathcal{S}_1 = \sup_{\theta \in \Theta} \sum_{\ell=1}^m \frac{\|\theta_\ell - \tilde{\theta}_\ell\|^2}{2\sigma_\ell^2}$$

$$\mathcal{S}_2 = \sup_{\theta \in \Theta} \max_{j \in \llbracket 1, k \rrbracket} \left\| \sum_{\ell=1}^m A_{j,\ell} (\theta_\ell - \theta_\ell^\star) \right\|$$

$$\text{and } \mathcal{S}_3 = (2 \log(2k))^{1/2} \max_{j \in \llbracket 1, k \rrbracket} \sqrt{\sum_{\ell=1}^m A_{j,\ell}^2 \sigma_\ell^2}.$$

# Proposition

With the above definitions, consider an  $\epsilon$ -minimizer of the empirical risk  $\widehat{\theta}(\overline{W}) \in \Theta$ , that is an estimator satisfying  $\mathbb{P}_{\overline{W}}$  almost surely

$$\overline{\mathbb{P}}(f(\widehat{\theta}, W)) \leq \inf_{\theta \in \Theta} \overline{\mathbb{P}}(f(\theta, W)) + \epsilon.$$

Its mean excess risk satisfies

$$\mathbb{P}_{\overline{W}}[\mathbb{P}(f(\widehat{\theta}, W))] \leq \inf_{\theta \in \Theta} \mathbb{P}(f(\theta, W)) + \gamma + \epsilon.$$

# Deviation bounds

Introduce the increasing function

$$\tilde{g}(t) = \frac{1}{t} [\exp(t) - 1].$$

Remark that with probability at least  $1 - \delta$

$$\bar{\mathbb{P}}_W(\xi(\|T_0(W)\|)^2) \leq \tilde{g}(\lambda_0 \xi(\eta_0)^2) \mathbb{P}(\xi(\|W\|)^2) + \frac{\log(\delta^{-1})}{n\lambda_0}$$

Take

$$\lambda_0 = \xi(\eta_0)^{-2},$$

and choose  $\eta_0$  such that

$$\frac{\log(2/\delta)}{n} \xi(\eta_0)^2 = \tilde{g}(1) \mathbb{P}(\xi(W)^2),$$

so that with probability at least  $1 - \delta/2$

$$\bar{\mathbb{P}}_W(\xi(\|T_0(W)\|)^2) \leq 2\tilde{g}(1) \mathbb{P}(\xi(\|W\|)^2).$$

## Choosing thresholds

Choose as previously  $\eta_q$  and  $\eta_{p+1}$  such that

$$B_{q,1}(\eta_q) = \frac{1}{2\lambda_q} \text{ and } \eta_{p+1} = \frac{1}{2\lambda_{p+1}}.$$

We get with probability at least  $1 - \delta/2$

$$|A_{q,1}| \leq 4(1 + 2\tilde{g}(1))\lambda_q \mathbb{P}_W [T_0 B_{q,0}(W)^2],$$

$$|A_{0,0}| \leq (1 + \sqrt{2\tilde{g}(1)}) \mathbb{P}_W (B_{1,0}^2)^{1/2},$$

$$|A_{-1,0}| \leq 2\mathbb{P}_W (B_{-1,0}^2 / B_{-1,1}(\eta_0)),$$

$$|A_{p+1,5}| \leq 2(1 + 2\tilde{g}(1))\lambda_{p+1} \mathbb{P}_W (T_0 B_{p+1,0}).$$



# PAC-Bayesian inequality

We also have

$$\mathbb{P}_{\bar{W}} \left\{ \exp \left[ n\lambda \sup_{\theta \in \Theta} \sum_{q=1}^p A_{q,2} - \frac{\lambda_q}{2} g(1) \mathbb{P}_W(B_{q,0}^2) - \frac{\beta_q \mathcal{S}_1}{n\lambda_q} \right. \right. \\ \left. \left. + A_{p+1,4} - \frac{\lambda_{p+1}}{2} g(1) \mathbb{P}_W(B_{p+1,0}) - \frac{\beta_{p+1} \mathcal{S}_1}{n\lambda_{p+1}} \right] \right\} \leq 1,$$

where

$$\frac{1}{\lambda} = \frac{1}{\lambda_{p+1}} + \sum_{q=1}^p \frac{1}{\lambda_q}.$$

# Proposition

Assume that  $\mathbb{P}[\xi(\|W\|)^2]^{1/2} \leq B$ , where  $B$  is known and that  $n \geq 2\mathcal{S}_1$ . Consider the threshold  $\eta_0$  such that

$$\frac{\log(2/\delta)}{n} \xi(\eta_0)^2 = (e - 1)B.$$

Let  $\hat{\theta} \in \Theta$ , be such that

$$\bar{\mathbb{P}}_W(f(\hat{\theta}, T_0 W)) \leq \inf_{\theta \in \Theta} \bar{\mathbb{P}}(f(\theta, T_0 W)) + \epsilon.$$

With probability at least  $1 - \delta$ , its excess risk is such that

$$\mathbb{P}_W(f(\hat{\theta}, W)) \leq \inf_{\theta \in \Theta} \mathbb{P}_W(f(\theta, W)) + B\gamma + \epsilon$$

$$\begin{aligned} \text{where } \gamma &= 12p^{\alpha_1} \mathcal{S}_3^{\alpha_1} (\mathcal{S}_1/n)^{\alpha_1/2} + 7(\mathcal{S}_3 + \mathcal{S}_2/p)^{\alpha_1} p^{-(1-\alpha_1)} (\mathcal{S}_1/n)^{1/2} \\ &+ [(9(2^{\alpha_1/2} - 1)^{-1} + 7)\mathcal{S}_3^{\alpha_1} p^{\alpha_1} + 12\mathcal{S}_2^{\alpha_1}] (\log(2/\delta)/n)^{1/2} \\ &= \mathcal{O}\left(\log(n/\mathcal{S}_1)^{\alpha_1} (\mathcal{S}_1/n)^{\alpha_1/2}\right), \text{ with } p = \lceil \log(n/\mathcal{S}_1)/\log(2) \rceil. \end{aligned}$$