

Optimal Aggregation of Affine Estimators

Arnak Dalalyan, École des Ponts ParisTech
Joseph Salmon, Duke University

Georgia Tech, September 2011

Introduction

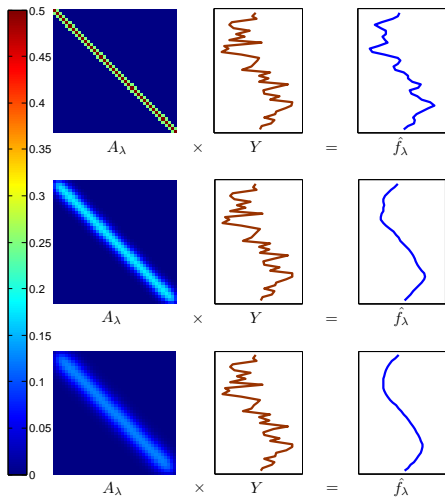
Motivations

- ▶ Theoretical : sharp oracle inequalities (high dimension, sparsity), Adaptation in the regression model
- ▶ Applications : image processing, genetics, inverse problems (derivative estimation, deconvolution with a known kernel, tomography), etc.

Underlying Heuristic

- ▶ Aggregating/mixing estimators can be better than selecting only one estimator

Doing as good as the best filter (1D example)



$Y \in \mathbb{R}^n$: noisy signal

\hat{f}_λ : estimated signal

A_λ : convolution/filter/kernel
matrix indexed by some
smoothing parameter
(bandwidth) λ in a family Λ
 \mathcal{F}_Λ : family of estimators

Doing as good as the best filter (2D example)

FIGURE: Evolution of the PSNR/MSE with respect to the bandwidth λ

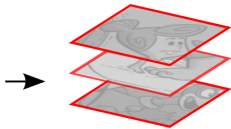
Doing as good as the best dictionary approximation

Patch-based methods are State-of-the-Art for denoising images

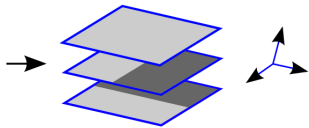
- ▶ Dabov et al. [07] (Wavelet),
- ▶ Mairal et al. [09] (Dictionary learning),
- ▶ Deledalle et al. [11] (PCA)



Patch extraction



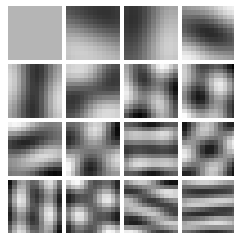
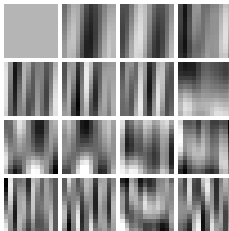
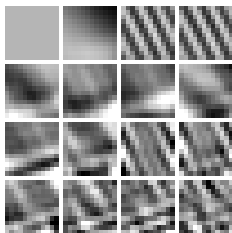
Patch stack



Patch denoising

Doing as good as the best dictionary approximation

Image denoising with patches



Dictionary (Dictionaries?)

Estimate an image/patch f by $\hat{f}_\lambda = f_\lambda = \sum_{j=1}^M \lambda_j \varphi_j$, for some dictionary/frame/orthonormal basis $\{\varphi_j, j = 1, \dots, M\}$

$\mathcal{F}_\Lambda = \text{Span}(\varphi_1, \dots, \varphi_M)$ and the $\lambda = (\lambda_1, \dots, \lambda_M)$ are the coefficients

Penalization Methods

Assume $\hat{f}_\lambda = f_\lambda = \sum_{j=1}^M \lambda_j \varphi_j$, for some features $\varphi_j \in \mathbb{R}^n$, $\Lambda = \mathbb{R}^M$

$$\hat{r}_\lambda = \|Y - \hat{f}_\lambda\|_n^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}_{\lambda,i})^2 : \text{empirical quadratic risk}$$

Penalization Methods

$$\hat{f}^{\text{Pen}} = f_{\hat{\lambda}}, \quad \text{where} \quad \hat{\lambda} = \arg \min_{\lambda \in \mathbb{R}^M} \left(\underbrace{\hat{r}_\lambda}_{\text{data-fitting}} + \underbrace{\text{Pen}(\lambda)}_{\text{regularization}} \right)$$

- $\text{Pen}(\lambda) = \beta \|\lambda\|_2^2$: Ridge **Tikhonov** [43]
- $\text{Pen}(\lambda) = \beta \|\lambda\|_0$: AIC, BIC, ... **Akaike** [74], **Schwarz** [78]
- $\text{Pen}(\lambda) = \beta \|\lambda\|_1$: LASSO **Tibshirani** [96]

Rem 1 : β smoothing parameter

Rem 2 : possible blocks/mixture versions (eg. Elastic Net)

Rem 3 : one usually uses only one estimate in the end : $f_{\hat{\lambda}}$

Mixing classical filtering and dictionary learning : known results

- ▶ Y : noisy vector/patch of pixels intensities, f the true one.
- ▶ Classical filtering : estimate f by AY , A convolution matrix.
 - Sharp oracle inequality for mixing estimators of the form AY with A projection matrix (Countable family) **Leung and Barron [06]**
- ▶ Dictionary learning : estimate f combining features b that are essentially independent of Y .
 - Sharp oracle inequality for mixing estimators built on an independent sample **Dalalyan and Tsybakov [07,08]**
- ▶ Goal : extending those results to aggregate estimates of the form $AY + b$ with A and b independent of Y .

Notation and model

Gaussian Heteroscedastic Model

$$Y_i = f_i + \sigma_i \varepsilon_i, \quad i = 1, \dots, n \quad (\star)$$

$$\varepsilon_i \text{ i.i.d } \mathcal{N}(0, 1) \quad \text{and} \quad \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \quad (\Sigma \text{ known})$$

- Rem 1 : $f_i = f(x_i)$, $(x_i)_{i=1, \dots, n}$ fixed design (cf. pixels)
- Rem 2 : $\Sigma = \sigma^2 I_n$, homoscedastic model

Goal : estimate f by \hat{f} , with a small (quadratic) risk

$$r = \mathbb{E} \left(\|f - \hat{f}\|_n^2 \right) = \mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n (f_i - \hat{f}_i)^2 \right)$$

Rem : link with inverse problems with known operator [Cavalier \[08\]](#)

Link inverse problems / heteroscedatic

T : **known** operator on a Hilbert space $(\mathcal{H}, \langle \cdot | \cdot \rangle_{\mathcal{H}})$ Y : random process on \mathcal{H} , for a $h \in \mathcal{H}$

$$Y = Th + \varepsilon \xi \quad \Longleftrightarrow \quad Y(g) = \langle Th | g \rangle_{\mathcal{H}} + \varepsilon \xi(g), \quad \forall g \in \mathcal{H},$$

T^* : Hermitian adjoint of T ; when $T^* T$ is compact, using the SVD

$$T\phi_k = b_k\psi_k, \quad T^*\psi_k = b_k\phi_k, \quad k \in \mathbb{N},$$

b_k : singular values, $\{\phi_k\}$: orthonormal basis of \mathcal{H} ,

$\{\psi_k\}$: orthonormal basis of $\text{Im}(T) \subset \mathcal{H}$. Model could be written

$$Y(\psi_k) = \langle h | \phi_k \rangle_{\mathcal{H}} b_k + \varepsilon \xi(\psi_k), \quad k \in \mathbb{N}.$$

If $b_k \neq 0$ the model is equivalent to (\star) , with $f_i = \langle h | \phi_i \rangle_{\mathcal{H}}$ and $\sigma_i = \varepsilon b_i^{-1}$

Aggregation of Estimators and Oracle Inequalities

Family of « pre-estimators » : $\mathcal{F}_\Lambda = \{\hat{f}_\lambda \in \mathbb{R}^n, \lambda \in \Lambda\}, \Lambda \subset \mathbb{R}^M$

Goal : providing a **non asymptotic** bound on the risk of an estimator \hat{f}_{aggr} build upon \mathcal{F}_Λ

Oracle Inequality / Aggregation Nemirovski [00]

$$\mathbb{E}\|\hat{f}_{aggr} - f\|_n^2 \leq C_n \inf_{\lambda \in \Lambda} \mathbb{E}\|\hat{f}_\lambda - f\|_n^2 + R_{n,\Lambda}$$

- An **Oracle** is any \hat{f}_{λ^*} s.t. $\lambda^* \in \arg \min_{\lambda \in \mathcal{F}_\Lambda} \mathbb{E}\|\hat{f}_\lambda - f\|_n^2$
- $C_n \geq 1$. When $C_n = 1$: the inequality is said **Sharp**
- $R_{n,\Lambda} \xrightarrow{n \rightarrow \infty} 0$: price to pay for not knowing the Oracle, depends on the complexity of Λ and on the noise intensity

Rem 1 : \hat{f}_{aggr} might not be in \mathcal{F}_Λ

Rem 2 : Optimality (lower bound) for some sets Λ Tsybakov [03]

EWA : classical point of view

Family of « pre-estimators » : $\mathcal{F}_\Lambda = \{\hat{f}_\lambda \in \mathbb{R}^n, \lambda \in \Lambda\}, \Lambda \subset \mathbb{R}^M$

EWA/Gibbs Measure

$$\hat{\pi}^{\text{EWA}}(d\lambda) \propto \exp(-n\hat{r}_\lambda/\beta)\pi(d\lambda)$$

- ▶ $\hat{\pi}^{\text{EWA}}$: posterior over Λ
- ▶ π : prior over Λ
- ▶ β : smoothing parameter/temperature
- ▶ \hat{r}_λ : unbiased risk estimate $\mathbb{E}(\hat{r}_\lambda) = \mathbb{E}\|\hat{f}_\lambda - f\|_n^2 = r_\lambda$

Posterior expectation :
$$\hat{f}^{\text{EWA}} = \int_\Lambda \hat{f}_\lambda \hat{\pi}^{\text{EWA}}(d\lambda)$$

Rem 1 : -if $\beta \rightarrow 0$, $\hat{f}^{\text{EWA}} \rightarrow \hat{f}_{\lambda^*}$ with $\lambda^* = \arg \min_{\lambda \in \Lambda} \hat{r}_\lambda$
-if $\beta \rightarrow \infty$, $\hat{f}^{\text{EWA}} \rightarrow \int_\Lambda \hat{f}_\lambda \pi(d\lambda)$

Rem 2 : unbiased risk estimates \hat{r}_λ by Stein's Lemma Stein [81]

EWA : Penalty point of view

- ▶ Extension : enlarge the parameter space and adapt the penalty
- ▶ Parameter space : $\mathcal{P}_\Lambda = \{p : \text{probability over } \Lambda\}$
- ▶ Extended penalty : $\hat{f}^{\text{Pen}} = \int_\Lambda \hat{f}_\lambda \hat{\pi}^{\text{Pen}}(d\lambda)$ with

$$\hat{\pi}^{\text{Pen}} = \arg \min_{p \in \mathcal{P}_\Lambda} \left(\int_\Lambda \hat{r}_\lambda p(d\lambda) + \int_\Lambda \text{Pen}(\lambda) p(d\lambda) \right)$$

EWA/Kullback-Leibler penalty

$$\text{EWA : } \begin{cases} \hat{\pi}^{\text{EWA}} &= \arg \min_{p \in \mathcal{P}_\Lambda} \left(\int_\Lambda \hat{r}_\lambda p(d\lambda) + \frac{\beta}{n} \mathcal{K}(p, \pi) \right) \\ \hat{f}^{\text{EWA}} &= \int_\Lambda \hat{f}_\lambda \hat{\pi}^{\text{EWA}}(d\lambda) \end{cases}$$

- ▶ π prior over Λ ; β smoothing parameter (aka « temperature »)
- ▶ $\mathcal{K}(p, \pi)$: KL-divergence between probabilities $p, \pi \in \mathcal{P}_\Lambda$,

$$\mathcal{K}(p, \pi) = \begin{cases} \int_\Lambda \log \left(\frac{dp}{d\pi}(\lambda) \right) p(d\lambda) & \text{if } p \ll \pi, \\ +\infty & \text{otherwise.} \end{cases}$$

EWA : Penalty point of view

- ▶ Extension : enlarge the parameter space and adapt the penalty
- ▶ Parameter space : $\mathcal{P}_\Lambda = \{p : \text{probability over } \Lambda\}$
- ▶ Extended penalty : $\hat{f}^{\text{Pen}} = \int_\Lambda \hat{f}_\lambda \hat{\pi}^{\text{Pen}}(d\lambda)$ with

$$\hat{\pi}^{\text{Pen}} = \arg \min_{p \in \mathcal{P}_\Lambda} \left(\int_\Lambda \hat{r}_\lambda p(d\lambda) + \int_\Lambda \text{Pen}(\lambda) p(d\lambda) \right)$$

EWA/Kullback-Leibler penalty

$$\text{EWA : } \begin{cases} \hat{\pi}^{\text{EWA}} &= \arg \min_{p \in \mathcal{P}_\Lambda} \left(\int_\Lambda \hat{r}_\lambda p(d\lambda) + \frac{\beta}{n} \mathcal{K}(p, \pi) \right) \\ \hat{f}^{\text{EWA}} &= \int_\Lambda \hat{f}_\lambda \hat{\pi}^{\text{EWA}}(d\lambda) \end{cases}$$

- ▶ π prior over Λ ; β smoothing parameter (aka « temperature »)
- ▶ $\mathcal{K}(p, \pi)$: KL-divergence between probabilities $p, \pi \in \mathcal{P}_\Lambda$,

$$\mathcal{K}(p, \pi) = \begin{cases} \int_\Lambda \log \left(\frac{dp}{d\pi}(\lambda) \right) p(d\lambda) & \text{if } p \ll \pi, \\ +\infty & \text{otherwise.} \end{cases}$$

Affine estimators

Affine estimators

$$\hat{f}_\lambda = A_\lambda Y + b_\lambda$$

- ▶ A_λ : $n \times n$ matrix; b_λ : deterministic vector in \mathbb{R}^n
- ▶ A_λ , b_λ : independent of Y
- ▶ Λ : possibly non-countable

Constant case : $A_\lambda = 0$, $\hat{f}_\lambda = b_\lambda$

$\{\varphi_1, \dots, \varphi_M\}$ is a finite « dictionary » of features

- ▶ $\mathcal{F}_\Lambda = \{\varphi_1, \dots, \varphi_M\}$ finite family
- ▶ $\mathcal{F}_\Lambda = \text{conv}(\varphi_1, \dots, \varphi_M)$ convex combinations
- ▶ $\mathcal{F}_\Lambda = \text{Span}(\varphi_1, \dots, \varphi_M)$ linear combinations
- ▶ $\mathcal{F}_\Lambda = \text{Span}_S(\varphi_1, \dots, \varphi_M)$ S -sparse combinations

Lower bounds : Tsybakov [03], Bunea et al. [07] , Lounici [07]

Linear case : $\hat{f}_\lambda = A_\lambda Y$ ($b_\lambda = 0$)

Ordinary Least Squares

$\{\mathcal{S}_\lambda : \lambda \in \Lambda\}$ family of subspaces of \mathbb{R}^n A_λ : orthogonal projectors over \mathcal{S}_λ Leung and Barron [06], Alquier and Lounici [10], Rigollet and Tsybakov [11,11']

Diagonal Matrices : $A_\lambda = \text{diag}(a_1, \dots, a_n)$

- ▶ Ordered projections : $a_k = \mathbb{1}_{(k \leq \lambda)}$ for λ integer, ie. $\Lambda = \{1, \dots, n\}$
- ▶ Pinsker's Filter : $a_k = (1 - \frac{k^\alpha}{w})_+$, with $x_+ = \max(x, 0)$ and $w, \alpha > 0$, i.e., $\Lambda = (\mathbb{R}_+^*)^2$
- ▶ ...

Linear case : $\hat{f}_\lambda = A_\lambda Y$ ($b_\lambda = 0$)

Ordinary Least Squares

$\{\mathcal{S}_\lambda : \lambda \in \Lambda\}$ family of subspaces of \mathbb{R}^n A_λ : orthogonal projectors over \mathcal{S}_λ Leung and Barron [06], Alquier and Lounici [10], Rigollet and Tsybakov [11,11']

Diagonal Matrices : $A_\lambda = \text{diag}(a_1, \dots, a_n)$

- ▶ Ordered projections : $a_k = \mathbb{1}_{(k \leq \lambda)}$ for λ integer, ie. $\Lambda = \{1, \dots, n\}$
- ▶ Pinsker's Filter : $a_k = (1 - \frac{k^\alpha}{w})_+$, with $x_+ = \max(x, 0)$ and $w, \alpha > 0$, i.e., $\Lambda = (\mathbb{R}_+^*)^2$
- ▶ ...

Affine estimators and risk estimation

Stein Unbiased Risk Estimate (Gaussian Noise) Stein [81]

SURE : If \hat{f} is almost everywhere differentiable in Y and $\partial_{Y_i} \hat{f}_i$ is integrable, then

$$\hat{r} = \| \mathbf{Y} - \hat{f} \|_n^2 + \frac{2}{n} \sum_{i=1}^n \sigma_i^2 \partial_{Y_i} \hat{f}_i - \frac{1}{n} \sum_{i=1}^n \sigma_i^2$$

is an unbiased risk estimate $\mathbb{E}(\hat{r}) = r$

SURE, Affine case : $\hat{f}_\lambda = A_\lambda Y + b_\lambda$

$$\hat{r}_\lambda = \| \mathbf{Y} - \hat{f}_\lambda \|_n^2 + \frac{2}{n} \text{Tr}(\Sigma A_\lambda) - \frac{1}{n} \text{Tr}(\Sigma)$$

is an unbiased risk estimate $\mathbb{E}(\|f - \hat{f}_\lambda\|_n^2) = r_\lambda$ where $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$

Pré-estimateurs affines et estimation du risque

Formule de Stein [81] (bruit gaussien)

Si \hat{f} est un estimateur différentiable presque partout en Y et que $\partial_{Y_i} \hat{f}_i$ est intégrable alors

$$\hat{r}_n = \| \mathbf{Y} - \hat{f} \|_n^2 + \frac{2}{n} \sum_{i=1}^n \partial_{Y_i} \hat{f}_i \sigma_i^2 - \frac{1}{n} \sum_{i=1}^n \sigma_i^2$$

est un estimateur sans biais du risque $\mathbb{E}(\hat{r}_n) = r$

Cas affine : $\hat{f}_\lambda = A_\lambda \mathbf{Y} + b_\lambda$

Conclusion :
$$\hat{r}_\lambda = \| \mathbf{Y} - \hat{f}_\lambda \|_n^2 + \frac{2}{n} \text{Tr}(\Sigma A_\lambda) - \frac{1}{n} \text{Tr}(\Sigma)$$

est un estimateur sans biais du risque $r_\lambda = \mathbb{E}(\|f - \hat{f}_\lambda\|_n^2)$

Rappel : $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$

Orthonormal family

Family of « pre-estimators » : $\mathcal{F}_\Lambda = \{f_\lambda \in \mathbb{R}^n, \lambda \in \Lambda\}, \Lambda = \mathbb{R}^M$

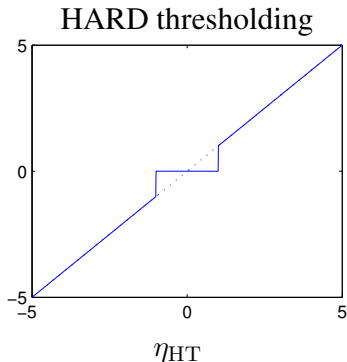
Assume $f_\lambda = \sum_{j=1}^n \lambda_j \varphi_j$, for some features $\varphi_j \in \mathbb{R}^n$

As before if $f_{\hat{\lambda}} = \arg \min_{\lambda \in \mathbb{R}^M} \left(\|Y - f_\lambda\|_n^2 + \beta \text{Pen}(\lambda) \right)$

- For $\text{Pen}(\lambda) = \beta \|\lambda\|_0$:

$$f_{\hat{\lambda}} = \sum_{j=1}^M \eta_{\text{HT}}(\langle \varphi_j | Y \rangle) \varphi_j$$

where $\eta_{\text{HT}}(x) = x \cdot \mathbb{1}(\sqrt{\beta} < |x|)$.



Orthonormal family $n = M$

Family of « pre-estimators » : $\mathcal{F}_\Lambda = \{f_\lambda \in \mathbb{R}^n, \lambda \in \mathbb{R}^n\}, \Lambda = \mathbb{R}^n$

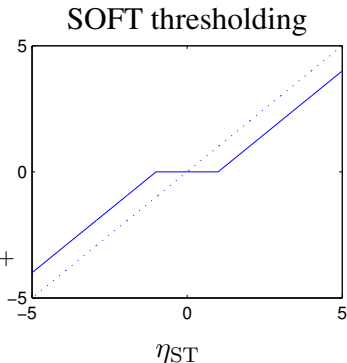
Assume $f_\lambda = \sum_{j=1}^n \lambda_j \varphi_j$, for some features $\varphi_j \in \mathbb{R}^n$

As before if $f_{\hat{\lambda}} = \arg \min_{\lambda \in \mathbb{R}^n} (\|Y - f_\lambda\|_n^2 + \beta \text{Pen}(\lambda))$

- For $\text{Pen}(\lambda) = \beta \|\lambda\|_1$:

$$f_{\hat{\lambda}} = \sum_{j=1}^n \eta_{\text{ST}}(\langle \varphi_j | Y \rangle) \varphi_j$$

where $\eta_{\text{ST}}(x) = \text{sign}(x) \cdot (|x| - \beta)_+$



Orthonormal family

Family of « pre-estimators » : $\mathcal{F}_\Lambda = \{\text{Proj}_\lambda Y, \lambda \in \Lambda\}, \Lambda = \{0, 1\}^n$

Proj_λ : projection on the φ_i associated to the “support” vector λ

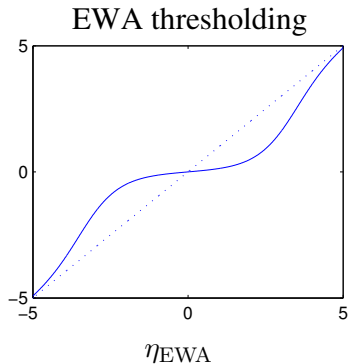
$$\hat{\pi}^{\text{EWA}}(\lambda) \propto \exp(-n\hat{r}_\lambda/\beta)\pi(\lambda)$$

$$\text{then } f_{\hat{\lambda}} = \sum_{j=1}^n \eta_{\text{EWA}}(\langle \varphi_j | Y \rangle) \varphi_j$$

- ▶ π : prior over Λ s.t
 $\pi(m) \propto c^{-\|\lambda\|_0}$ for any
subspace m in
- ▶ $\hat{r}_\lambda = \|Y - f_\lambda\|_n^2 + \frac{2\sigma^2\|\lambda\|_0}{n} - \sigma^2$

$$\eta_{\text{EWA}}(x) = \frac{x}{1 + ce^{-2\sigma^2/\beta} e^{-x^2/\beta}}$$

Giraud [08]



Main theorem conditions

$$\hat{f}_\lambda = A_\lambda Y + b_\lambda$$

Condition C₁

- ▶ Matrices A_λ : orthogonal projections ($A_\lambda^2 = A_\lambda^\top = A_\lambda$)
- ▶ Vectors b_λ : $A_\lambda b_\lambda = 0$

Example : A_λ projectors on subspaces Leung and Barron [06]

Condition C₂

- ▶ Matrices A_λ : symmetric, positive semi-definite
- ▶ $A_\lambda A_{\lambda'} = A_{\lambda'} A_\lambda, \forall \lambda, \lambda' \in \Lambda$ and $A_\lambda \Sigma = \Sigma A_\lambda, \forall \lambda \in \Lambda$
- ▶ Vectors b_λ : $A_{\lambda'} b_\lambda = 0, \forall \lambda, \lambda' \in \Lambda$

Example : two-blocks James-Stein shrinking estimators Leung [04]

Main theorem conditions

$$\hat{f}_\lambda = A_\lambda Y + b_\lambda$$

Condition C₁

- ▶ Matrices A_λ : orthogonal projections ($A_\lambda^2 = A_\lambda^\top = A_\lambda$)
- ▶ Vectors b_λ : $A_\lambda b_\lambda = 0$

Example : A_λ projectors on subspaces [Leung and Barron \[06\]](#)

Condition C₂

- ▶ Matrices A_λ : symmetric, positive semi-definite
- ▶ $A_\lambda A_{\lambda'} = A_{\lambda'} A_\lambda, \forall \lambda, \lambda' \in \Lambda$ and $A_\lambda \Sigma = \Sigma A_\lambda, \forall \lambda \in \Lambda$
- ▶ Vectors b_λ : $A_{\lambda'} b_\lambda = 0, \forall \lambda, \lambda' \in \Lambda$

Example : two-blocks James-Stein shrinking estimators [Leung \[04\]](#)

Main Theorem

PAC (EAC) - Bayesian Bound

If \mathbf{C}_1 or \mathbf{C}_2 is satisfied, then for any prior π , \hat{f}^{EWA} satisfies :

$$\mathbb{E}(\|\hat{f}^{\text{EWA}} - f\|_n^2) \leq \inf_{p \in \mathcal{P}_\Lambda} \left(\int_\Lambda \mathbb{E} \|\hat{f}_\lambda - f\|_n^2 p(d\lambda) + \frac{\beta}{n} \mathcal{K}(p, \pi) \right)$$

$$\text{where } \beta \geq 4 \max_{i=1, \dots, n} \sigma_i^2 \text{ under } \mathbf{C}_1$$

$$\beta \geq 8 \max_{i=1, \dots, n} \sigma_i^2 \text{ under } \mathbf{C}_2$$

with $\mathcal{K}(p, \pi)$ the KL divergence between p and π

Corollary : finite case

Oracle Inequality : $\Lambda = \llbracket 1, M \rrbracket$, π uniform

If \mathbf{C}_1 or \mathbf{C}_2 is satisfied, and if π is uniform on $\llbracket 1, M \rrbracket$, then

$$\mathbb{E}(\|\hat{f}^{\text{EWA}} - f\|_n^2) \leq \inf_{\lambda \in \llbracket 1, M \rrbracket} \left(\mathbb{E} \|\hat{f}_\lambda - f\|_n^2 \right) + \frac{\beta \log(M)}{n}$$

where $\beta \geq 4 \max_{i=1, \dots, n} \sigma_i^2$ under \mathbf{C}_1

$\beta \geq 8 \max_{i=1, \dots, n} \sigma_i^2$ under \mathbf{C}_2

- ▶ If $b_\lambda = 0$: extends result by **Leung and Barron [06]**
- ▶ If $A_\lambda = 0$, $\Sigma = \sigma I_n$: optimal inequality **Tsybakov [03]** and no selector can achieve a rate faster than $\sqrt{\frac{\log(M)}{n}}$ for some function f and dictionary $\mathcal{F}_\Lambda = \{f_1, \dots, f_M\}$! **Juditsky et al. [08]**, **Rigollet and Tsybakov [11]**

Corollary : Sparse Oracle Inequality

Sparse scenario : $\Lambda = \mathbb{R}^M$ and $\exists \lambda^*$ s.t. $\hat{f}_{\lambda^*} \approx f$. Let π be a sparsifying (heavy-tailed) prior and $\tau > 0$ a scale parameter.

$$\text{e.g. } \pi(d\lambda) \propto \prod_{j=1}^M 1/(1 + |\lambda_j/\tau|^2)^2 d\lambda.$$

Oracle Inequality

With such a π , under \mathbf{C}_1 or \mathbf{C}_2 and if $\exists \mathcal{M} \in \mathbb{R}^{M \times M}$ s.t :

$$r_{\lambda} - r_{\lambda'} - \nabla r_{\lambda'}^{\top}(\lambda - \lambda') \leq (\lambda - \lambda')^{\top} \mathcal{M}(\lambda - \lambda'), \quad \forall \lambda, \lambda' \in \Lambda.$$

$$\mathbb{E}(\|\hat{f}^{\text{EWA}} - f\|_n^2) \leq \inf_{\lambda \in \mathbb{R}^M} \left\{ \mathbb{E}\|\hat{f}_{\lambda} - f\|_n^2 + \frac{4\beta}{n\tau} \|\lambda\|_1 + \text{Tr}(\mathcal{M})\tau^2 \right\}$$

where $\beta \geq 4 \max_{i=1,\dots,n} \sigma_i^2$ under \mathbf{C}_1 (and $8 \max_{i=1,\dots,n} \sigma_i^2$ under \mathbf{C}_2)

Rem : for a dictionary $\varphi_1, \dots, \varphi_M$, \mathcal{M} could be $G = \langle \varphi_i, \varphi_j \rangle_n$ (Gramm Matrix)

No assumption on the dictionary Dalalyan and Tsybakov [07]

Minimax point of view ($\Sigma = \sigma^2 I_n$)

$\theta_k(f) = \langle f | \varphi_k \rangle_n$: Discrete Fourier coefficients

$\mathcal{D}f$: Discrete Fourier Transform of f

Sobolev Ellipsoid : $\mathcal{E}(\alpha, R) = \{f \in \mathbb{R}^n : \sum_{k=1}^n k^{2\alpha} \theta_k(f)^2 \leq R\}$

Pinsker's Theorem : linear estimates are minimax on ellipsoids

$$\begin{aligned} \inf_{\hat{f}} \sup_{f \in \mathcal{E}(\alpha, R)} \mathbb{E}(\|\hat{f} - f\|_n^2) &\sim \inf_A \sup_{f \in \mathcal{F}(\alpha, R)} \mathbb{E}(\|AY - f\|_n^2) \\ &\sim \inf_{w > 0} \sup_{f \in \mathcal{E}(\alpha, R)} \mathbb{E}(\|A_{\alpha, w} Y - f\|_n^2) \end{aligned}$$

the inf is taken among all the possible estimators \hat{f} and

$A_{\alpha, w} = \mathcal{D}^\top \text{diag}((1 - k^\alpha/w)_+; k = 1, \dots, n) \mathcal{D}$: Pinsker's Filter

Rem : $\lambda = (\alpha, w)$ and $\Lambda = (\mathbb{R}_+^*)^2$

Corollary : Adaptation

EWA on Pinsker filters : $\hat{f}_\lambda = \hat{f}_{\alpha,w} = \mathcal{D}^\top A_{\alpha,w} \mathcal{D} Y$ (\mathcal{D} : DCT),
with $A_{\alpha,w} = \text{diag}((1 - \frac{k^\alpha}{w})_+, k = 1, \dots, n)$

Choose the prior π over $\Lambda = (\mathbb{R}_+^*)^2$:

- ▶ Draw α according to an exponential distribution with parameter 1
- ▶ Knowing α , draw w according to the density

$$w \rightarrow \frac{2n_\sigma^{-\alpha/(2\alpha+1)}}{(1+n_\sigma^{-\alpha/(2\alpha+1)}w)^3} \text{ with } n_\sigma = n/\sigma^2$$

Performance

- ▶ Theoretical : adaptive in the exact minimax sense on Sobolev ellipsoids
- ▶ Practical : performance as good as other classical adaptive methods such as SURE/ Soft Thresholding [Donoho and Johnstone \[95\]](#) , Block James-Stein [Cai \[99\]](#) , empirical risk minimization [Cavalier et al. \[02\]](#)

Extension to non symmetric matrices : SEWA, Symmetrized EWA

1. For every λ , compute the risk estimate
$$\hat{r}_\lambda^{\text{unb}} = \| \mathbf{Y} - \hat{f}_\lambda \|^2_n + \frac{2}{n} \text{Tr}(\Sigma A_\lambda) - \frac{1}{n} \text{Tr}[\Sigma].$$
2. Define the prob. distribution $\hat{\pi}^{\text{EWA}}(d\lambda) = \theta(\lambda)\pi(d\lambda)$ with $\theta(\lambda) \propto \exp(-n\hat{r}_\lambda^{\text{unb}}/\beta)$.
3. For every λ , build the symmetrized linear smoothers
$$\tilde{f}_\lambda = (A_\lambda + A_\lambda^\top - A_\lambda^\top A_\lambda) \mathbf{Y}.$$
4. Average out the symmetrized smoothers w.r.t. posterior

$$\hat{f}_{\text{SEWA}} = \int_{\Lambda} \tilde{f}_\lambda \hat{\pi}^{\text{EWA}}(d\lambda)$$

$$\text{Rem : } \hat{f}^{\text{EWA}} = \int_{\Lambda} \hat{f}_\lambda \hat{\pi}^{\text{EWA}}(d\lambda)$$

Main Theorem (II)

Condition C_3

- ▶ Matrices $A_\lambda : \text{Tr}(\Sigma A_\lambda) \leq \text{Tr}(\Sigma A_{\lambda'}^\top A_\lambda), \forall \lambda \in \Lambda$
- ▶ Vectors $b_\lambda : b_\lambda = 0, \forall \lambda, \lambda' \in \Lambda$

PAC (EAC) - Bayesian Bound

If the matrices A_λ satisfies condition C_3 , then for any prior π , \hat{f}_{SEWA} satisfies :

$$\mathbb{E}(\|\hat{f}_{\text{SEWA}} - f\|_n^2) \leq \inf_{p \in \mathcal{P}_\Lambda} \left(\int_\Lambda \mathbb{E} \|\hat{f}_\lambda - f\|_n^2 p(d\lambda) + \frac{\beta}{n} \mathcal{K}(p, \pi) \right)$$

$$\text{where } \beta \geq 4 \max_{i=1, \dots, n} \sigma_i^2$$

with $\mathcal{K}(p, \pi)$ the KL divergence between p and π

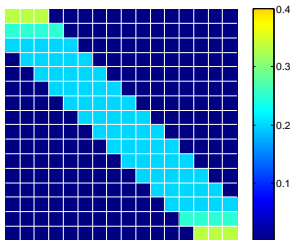
Condition C_3

$$\text{Tr}(\Sigma A_\lambda) \leq \text{Tr}(\Sigma A_\lambda^\top A_\lambda), \forall \lambda \in \Lambda$$

1. Orth. projection $A_\lambda^\top = A_\lambda = A_\lambda^2$: $\text{Tr}(\Sigma A_\lambda) = \text{Tr}(\Sigma A_\lambda^\top A_\lambda)$
2. If Σ diagonal, and $A_{ii} \leq \sum_{j=1}^n A_{ji}^2$. $A_{ii} = 0$ e.g. k -NN filter in which the weight of the observation Y_i is replaced by 0
3. If $\Sigma = \sigma^2 I_n$ and
 - ▶ all the non-zero elements of each row are equal
 - ▶ each row sums up to some $c \geq 1$.

then $\text{Tr}(A_\lambda) = \text{Tr}(A_\lambda^\top A_\lambda)$.

e.g. : Nadaraya-Watson estimators with rectangular kernel and nearest neighbor filters.



Conclusion

Contributions

- ▶ Sharp oracle inequalities for some affine estimators
- ▶ Adaptive results with respect to the signal smoothness
- ▶ Reasonable experimental performance
- ▶ New estimator : SEWA, symmetrized version of the EWA

Details : available on-line

- ▶ COLT 2011 (EWA)
- ▶ ALT 2011 (SEWA)
- ▶ Long version arXiv
- ▶ Code

Et Voilà !

Conclusion

Contributions

- ▶ Sharp oracle inequalities for some affine estimators
- ▶ Adaptive results with respect to the signal smoothness
- ▶ Reasonable experimental performance
- ▶ New estimator : SEWA, symmetrized version of the EWA

Details : available on-line

- ▶ COLT 2011 (EWA)
- ▶ ALT 2011 (SEWA)
- ▶ Long version arXiv
- ▶ Code

Et Voilà !

References I

► H. Akaike.

A new look at the statistical model identification.

IEEE Trans. Automatic Control, AC-19 :716–723, 1974.

System identification and time-series analysis.

► P. Alquier and K. Lounici.

Pac-bayesian bounds for sparse regression estimation with exponential weights.

Electron. J. Statist., 5 :127–145, 2010.

► F. Bunea, A. B. Tsybakov, and M. H. Wegkamp.

Aggregation for Gaussian regression.

Ann. Statist., 35(4) :1674–1697, 2007.

► T. T. Cai.

Adaptive wavelet estimation : a block thresholding and oracle inequality approach.

Ann. Statist., 27(3) :898–924, 1999.

References II

- ▶ L. Cavalier.
Nonparametric statistical inverse problems.
Inverse Problems, 24(3) :034004, 19, 2008.
- ▶ L. Cavalier, G. K. Golubev, D. Picard, and A. B. Tsybakov.
Oracle inequalities for inverse problems.
Ann. Statist., 30(3) :843–874, 2002.
- ▶ K. Dabov, A. Foi, V. Katkovnik, and K. O. Egiazarian.
Image denoising by sparse 3-D transform-domain collaborative filtering.
IEEE Trans. Image Process., 16(8) :2080–2095, 2007.
- ▶ D. L. Donoho and I. M. Johnstone.
Adapting to unknown smoothness via wavelet shrinkage.
J. Amer. Statist. Assoc., 90(432) :1200–1224, 1995.
- ▶ C.-A. Deledalle, J. Salmon, and A. S. Dalalyan.
Image denoising with patch based PCA : local versus global.
In *BMVC*, 2011.

References III

- ▶ A. S. Dalalyan and A. B. Tsybakov.
Aggregation by exponential weighting, sharp oracle inequalities and sparsity.
In COLT, pages 97–111, 2007.
- ▶ A. S. Dalalyan and A. B. Tsybakov.
Aggregation by exponential weighting, sharp pac-bayesian bounds and sparsity.
Mach. Learn., 72(1-2) :39–61, 2008.
- ▶ Ch. Giraud.
Mixing least-squares estimators when the variance is unknown.
Bernoulli, 14(4) :1089–1107, 2008.
- ▶ A. B. Juditsky, Ph. Rigollet, and A. B. Tsybakov.
Learning by mirror averaging.
Ann. Statist., 36(5) :2183–2206, 2008.
- ▶ G. Leung and A. R. Barron.
Information theory and mixing least-squares regressions.
IEEE Trans. Inf. Theory, 52(8) :3396–3410, 2006.

References IV

- ▶ G. Leung.
Information Theory and Mixing Least Squares Regression.
PhD thesis, Yale University, 2004.
- ▶ K. Lounici.
Generalized mirror averaging and D -convex aggregation.
Math. Methods Statist., 16(3) :246–259, 2007.
- ▶ A. S. Nemirovski.
Topics in non-parametric statistics, volume 1738 of *Lecture Notes in Math.*
Springer, Berlin, 2000.
- ▶ Ph. Rigollet and A. B. Tsybakov.
Sparse estimation by exponential weighting.
Ann. Statist.
- ▶ Ph. Rigollet and A. B. Tsybakov.
Exponential screening and optimal rates of sparse estimation.
Ann. Statist., 39(2) :731–471, 2011.

References V

- ▶ G. Schwarz.
Estimating the dimension of a model.
Ann. Statist., 6(2) :461–464, 1978.
- ▶ C. M. Stein.
Estimation of the mean of a multivariate normal distribution.
Ann. Statist., 9(6) :1135–1151, 1981.
- ▶ R. Tibshirani.
Regression shrinkage and selection via the lasso.
J. Roy. Statist. Soc. Ser. B, 58(1) :267–288, 1996.
- ▶ A. N. Tikhonov.
On the stability of inverse problems.
C. R. (Doklady) Acad. Sci. URSS (N.S.), 39 :176–179, 1943.
- ▶ A. B. Tsybakov.
Optimal rates of aggregation.
In *COLT*, pages 303–313, 2003.