

# Sparse Exponential Weighting and Langevin Monte-Carlo

Alexandre Tsybakov,  
joint work with Arnak Dalalyan

Laboratoire de Statistique, CREST  
and  
Laboratoire de Probabilités et Modèles Aléatoires,  
Université Paris 6

Cumberland Lodge, April 1, 2009

# Nonparametric regression model (fixed design)

Assume that we observe the pairs  $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^d \times \mathbb{R}$  where

$$Y_i = f(X_i) + \xi_i, \quad i = 1, \dots, n.$$

- Regression function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is unknown
- Errors  $\xi_i$  are independent Gaussian  $\mathcal{N}(0, \sigma^2)$  random variables.
- $X_i \in \mathbb{R}^d$  are arbitrary fixed (non-random) points.

We want to estimate  $f$  based on the data  $(X_1, Y_1), \dots, (X_n, Y_n)$ .

## Approximating function, dictionary

We assume that there exists a function  $f_\lambda(x)$  (known as a function of  $\lambda$  and  $x$ ) such that

$$f \approx f_\lambda$$

for some  $\lambda = (\lambda_1, \dots, \lambda_M)$ .

Possibly  $M \gg n$

## Example: linear approximation, dictionary

Let  $f_1, \dots, f_M$  be a finite **dictionary of functions**,  $f_j : \mathbb{R}^d \rightarrow \mathbb{R}$ .  
We approximate the regression function  $f$  by linear combination

$$f_\lambda(x) = \sum_{j=1}^M \lambda_j f_j(x) \quad \text{with weights} \quad \lambda = (\lambda_1, \dots, \lambda_M).$$

We believe that

$$f(x) \approx \sum_{j=1}^M \lambda_j f_j(x)$$

for some  $\lambda = (\lambda_1, \dots, \lambda_M)$ .

# Scenarios for linear approximation

(LinReg) Exact equality: there exists  $\lambda^* \in \mathbb{R}^M$  such that

$$f = f_{\lambda^*} = \sum_{j=1}^M \lambda_j^* f_j$$

(**linear regression**, with possibly  $M \gg n$  parameters);

## Scenarios for linear approximation

(LinReg) Exact equality: there exists  $\lambda^* \in \mathbb{R}^M$  such that

$$f = f_{\lambda^*} = \sum_{j=1}^M \lambda_j^* f_j$$

(**linear regression**, with possibly  $M \gg n$  parameters);

(NPRReg)  $f_1, \dots, f_M$  are the first  $M$  functions of a basis (usually orthonormal) and  $M \leq n$ , there exists  $\lambda^*$  such that  $f - f_{\lambda^*}$  is small: **nonparametric estimation of regression**;

## Scenarios for linear approximation

(LinReg) Exact equality: there exists  $\lambda^* \in \mathbb{R}^M$  such that

$$f = f_{\lambda^*} = \sum_{j=1}^M \lambda_j^* f_j$$

(**linear regression**, with possibly  $M \gg n$  parameters);

(NPRReg)  $f_1, \dots, f_M$  are the first  $M$  functions of a basis (usually orthonormal) and  $M \leq n$ , there exists  $\lambda^*$  such that  $f - f_{\lambda^*}$  is small: **nonparametric estimation of regression**;

(Agg) **aggregation of arbitrary estimators**: in this case  $f_1, \dots, f_M$  are preliminary estimators of  $f$  based on a training sample independent of the observations  $(X_1, Y_1), \dots, (X_n, Y_n)$ ;

Weak learning, additive models etc.

## Example: nonlinear approximation

We consider the generalized linear (or single-index) model

$$f_\lambda(x) = G(\lambda^T x)$$

where  $G : \mathbb{R} \rightarrow \mathbb{R}$  is a known (or unknown) function.

Thus, we believe that

$$f(x) \approx G(\lambda^T x)$$

for some  $\lambda = (\lambda_1, \dots, \lambda_M)$ , possibly with  $M \gg n$ .



# Sparsity of a vector

The number of non-zero coordinates of  $\lambda$ :

$$M(\lambda) = \sum_{j=1}^M \mathbb{I}_{\{\lambda_j \neq 0\}}$$

The value  $M(\lambda)$  characterizes the **sparsity** of vector  $\lambda \in \mathbb{R}^M$ : the smaller  $M(\lambda)$ , the “sparser”  $\lambda$ .

# Sparsity of the model

## Intuitive formulation of sparsity assumption:

$$f(x) \approx f_\lambda \quad (\text{"}f\text{ is well approximated by }f_\lambda\text{"})$$

where the vector  $\lambda = (\lambda_1, \dots, \lambda_M)$  is sparse:

$$M(\lambda) \ll M.$$

# Sparsity and dimension reduction

Let  $\hat{\lambda}_{\text{OLS}}$  be the ordinary least squares (OLS) estimator. Let  $f_\lambda$  be **linear** approximation. Elementary result:

$$\mathbb{E} \|f_{\hat{\lambda}_{\text{OLS}}} - f\|_n^2 \leq \|f - f_\lambda\|_n^2 + \frac{\sigma^2 M}{n}$$

for any  $\lambda \in \mathbb{R}^M$  where  $\|\cdot\|_n$  is the empirical norm:

$$\|f\|_n = \sqrt{\frac{1}{n} \sum_{i=1}^n f^2(X_i)}.$$

# Sparsity and dimension reduction

For any  $\lambda \in \mathbb{R}^M$  the “oracular” OLS that acts only on the relevant  $M(\lambda)$  coordinates satisfies

$$\mathbb{E} \|\hat{f}_{\hat{\lambda}_{\text{OLS}}}^{\text{oracle}} - f\|_n^2 \leq \|f - f_\lambda\|_n^2 + \frac{\sigma^2 M(\lambda)}{n}.$$

This is only an OLS oracle, not an estimator. The set of relevant coordinates should be known.

## Sparsity oracle inequalities

Do there exist true estimators with similar behavior? Basic idea: Choose some suitable data-driven weights  $\hat{\lambda} = (\hat{\lambda}_1, \dots, \hat{\lambda}_M)$  and estimate  $f$  by

$$\hat{f}(x) = f_{\hat{\lambda}}(x) = \sum_{j=1}^M \hat{\lambda}_j f_j(x).$$

- What to do when the approximation is non-linear (ex.  $G(\lambda^T x)$ )? Should we also plug in an estimator  $\hat{\lambda}$ ?
- Can we find  $\hat{\lambda}$  such that  $\tilde{f} = f_{\hat{\lambda}}$  or  $\tilde{f}$  defined in differently satisfies

$$\mathbb{E} \|\tilde{f} - f\|_n^2 \lesssim \|f - f_{\lambda}\|_n^2 + \frac{\sigma^2 M(\lambda)}{n}, \quad \forall \lambda?$$

## Sparsity oracle inequalities (SOI)

Realizable task: Construct an estimator  $\tilde{f}$  satisfying a **sparsity oracle inequality (SOI)**

$$\mathbb{E} \|\tilde{f} - f\|_n^2 \leq \inf_{\lambda \in \mathbb{R}^M} \left\{ C \|f - f_\lambda\|_n^2 + C' \frac{M(\lambda) (" \log M'' )}{n} \right\}$$

with some constants  $C \geq 1$ ,  $C' > 0$  and an inevitable extra  $" \log M''$  in the variance term.

$C = 1 \Rightarrow$  **sharp SOI**.

## “Ideal” requirements for SOI

We would like to construct an estimator  $\tilde{f}$  such that it satisfies:

- SOI with leading constant 1 (sharp SOI);
- this holds under no assumptions on the approximation function; in the linear case, under no assumptions on the dictionary  $f_1, \dots, f_M$ ;
- the estimator is computationally feasible

## Penalized techniques (BIC, Lasso)

Penalize the residual sum of squares directly by  $M(\lambda)$  (BIC criterion, Schwarz (1978), Foster and George (1994)):

$$\hat{\lambda}^{BIC} = \arg \min_{\lambda \in \mathbb{R}^M} \left\{ \|\mathbf{y} - \mathbf{f}_\lambda\|_n^2 + \gamma \frac{M(\lambda) \log M}{n} \right\},$$

where  $\gamma > 0$  and

$$\|\mathbf{y} - \mathbf{f}_\lambda\|_n^2 \triangleq \frac{1}{n} \sum_{i=1}^n \left( Y_i - \mathbf{f}_\lambda(X_i) \right)^2, \quad \mathbf{y} = (Y_1, \dots, Y_n).$$

Remarks:

- If the matrix  $X = (f_j(X_i))_{i,j}$  has orthonormal columns, BIC is equivalent to hard thresholding of the components of  $X^T \mathbf{y} / n$  at the level  $\sqrt{\gamma(\log M)/n}$ .
- Non-convex, discontinuous minimization problem.



## Sparsity oracle inequality for BIC (linear approximation)

**Theorem.** [Bunea/ T/ Wegkamp (2004)]: if  $\gamma > K_0\sigma^2$  for an absolute constant  $K_0$ , and **with no assumption on the dictionary**  $f_1, \dots, f_M$ , the BIC estimator satisfies, with probability close to 1,

$$\|f_{\hat{\lambda}^{BIC}} - f\|_n^2 \leq (1+\varepsilon) \inf_{\lambda \in \mathbb{R}^M} \left\{ \|f - f_\lambda\|_n^2 + C(\varepsilon) \frac{M(\lambda) \log M}{n} \right\}, \quad \forall \varepsilon > 0.$$

Remarks:

- the BIC is realizable only for small  $M$  (say,  $M \leq 20$ ),
- the leading constant is **not** 1,
- $C(\varepsilon) \sim 1/\varepsilon$ .
- no result for non-linear approximation

# LASSO

Second penalization technique: LASSO [Frank and Friedman (1993, Bridge regression), Tibshirani (1996), Chen and Donoho (1998, basis pursuit)]. Penalize the residual sum of squares not by  $M(\lambda)$ , as in the BIC, but by the  $\ell_1$ -norm of  $\lambda$ :

$$\hat{\lambda}^L = \arg \min_{\lambda \in \mathbb{R}^M} \{ \|\mathbf{y} - f_\lambda\|_n^2 + 2r|\lambda|_1 \},$$

where  $|\lambda|_1 = \sum_{j=1}^M |\lambda_j|$ ,  $r > 0$  a tuning constant. A sensible choice:

$$r \sim \sqrt{\frac{\log M}{n}}.$$

- If the matrix  $X = (f_j(X_i))_{i,j}$  has orthonormal columns, LASSO is equivalent to soft thresholding of the components of  $X^T \mathbf{y}/n$  at the level  $r$ .

Advantages of the LASSO: computationally simple, selects the sparsity pattern [Bühlmann and Meinshausen (2004), Zhao and Yu (2006)], ...

Disadvantages of the LASSO:

- SOI for the LASSO holds under strong assumptions on the dictionary involving minimal “restricted eigenvalues”. Moreover, the assumptions depend on the (unknown) number  $s$  of non-zero components of the oracle vector, or eventually on the upper bound on this number. Such assumptions are unavoidable: Candès and Plan (2008).
- The leading constant in SOI is **not** 1.
- How to deal with non-linear approximations?

Same problems with other  $\ell_1$  penalized techniques (Dantzig selector, modifications of the Lasso).

## Exponential weighting

Estimate  $f(x)$  by

$$\tilde{f}^{EW}(x) = \int_{\mathbb{R}^M} f_{\lambda}(x) S_n(d\lambda)$$

where the probability measure  $S_n$  is given by

$$S_n(d\lambda) = \frac{\exp \left\{ -n \|\mathbf{y} - f_{\lambda}\|_n^2 / \beta \right\} \pi(d\lambda)}{\int_{\mathbb{R}^M} \exp \left\{ -n \|\mathbf{y} - f_w\|_n^2 / \beta \right\} \pi(dw)}$$

with some  $\beta > 0$  and some prior measure  $\pi$ .

# Exponential weighting

- For the linear approximation:  $\tilde{f}^{EW} = f_{\hat{\lambda}^{EW}}$  where

$$\hat{\lambda}_j^{EW} = \int_{\mathbb{R}^M} \lambda_j S_n(d\lambda), \quad j = 1, \dots, M,$$

- Bayesian estimator if  $\beta = 2\sigma^2$ , but we need a larger  $\beta$ .
- Non-discrete  $\pi$ : Computational issues?

## A PAC-Bayesian bound

Lemma [Dalalyan and T., 2007]

The estimator with exponential weights  $\tilde{f}^{EW}$  defined with  $\beta \geq 4\sigma^2$  and any prior  $\pi$  satisfies:

$$\mathbb{E} \|\tilde{f}^{EW} - f\|_n^2 \leq \inf_P \left\{ \int \|f_\lambda - f\|_n^2 P(d\lambda) + \frac{\beta \mathcal{K}(P, \pi)}{n} \right\}$$

where the infimum is taken over all probability measures  $P$  on  $\mathbb{R}^M$  and  $\mathcal{K}(P, \pi)$  denotes the Kullback-Leibler divergence between  $P$  and  $\pi$ .

## Sparsity prior

Choose a specific prior measure  $\pi$  with Lebesgue density  $q$ :

$$q(\lambda) = \prod_{j=1}^M \tau^{-1} q_0(\lambda_j/\tau), \quad \forall \lambda \in \mathbb{R}^M,$$

where  $q_0$  is the Student  $t_3$  density,

$$q_0(t) \sim |t|^{-4}, \quad \text{for large } |t|$$

and  $\tau \sim (Mn)^{-1/2}$ . We will call this prior the **sparsity prior**. The resulting estimator  $\tilde{f}^{EW}$  is called the **Sparse Exponential Weighting (SEW)** estimator.

## SOI for the SEW estimator: Linear approximation case

### Theorem [Dalalyan and T., 2007]

Let  $\max_{1 \leq j \leq M} \|f_j\|_n \leq c_0 < \infty$ . Let  $f_\lambda$  be linear in  $\lambda$ . Then for  $\beta \geq 4\sigma^2$  the estimator  $f_{\hat{\lambda}^{EW}}$  with the **sparsity prior** satisfies:

$$\mathbb{E} \|f_{\hat{\lambda}^{EW}} - f\|_n^2 \leq \inf_{\lambda \in \mathbb{R}^M} \left\{ \|f_\lambda - f\|_n^2 + \frac{CM(\lambda)}{n} \log \left( 1 + \frac{|\lambda|_1 \sqrt{Mn}}{M(\lambda)} \right) \right\}$$

where  $|\lambda|_1$  is the  $\ell_1$ -norm of  $\lambda$ .

- **No assumption on the dictionary.**
- **Leading constant 1.**
- $\ell_1$ -norm of  $\lambda$ , but under the log.
- Fast computation for at least  $M \sim 10^3$ .



## SOI for the SEW estimator: generalized linear models

Assume now:

$$f_\lambda(x) = G(x^T \lambda).$$

Then we have the same result as for the linear approximation case provided that

$$\sup_{\lambda} \text{Spec} \left\{ \frac{1}{n} \sum_{i=1}^n G''(X_i^T \lambda) X_i X_i^T \right\} \leq c_0 < \infty.$$

## SEW estimator: discussion

- *SEW is not a penalized estimator.*

$$\hat{\lambda}_j^{EW} = \int_{\mathbb{R}^M} \lambda_j S_n(d\lambda) = \int_{\mathbb{R}^M} \lambda_j g_n(\lambda) d\lambda, \quad j = 1, \dots, M,$$

with posterior density  $g_n(\lambda) = S_n(d\lambda)/d\lambda$ :

$$g_n(\lambda) \propto \exp \left\{ -n \|\mathbf{y} - \mathbf{f}_\lambda\|_n^2 / \beta - C \sum_{j=1}^M \log(1 + \lambda_j^2 / \tau) \right\}$$

Maximizer of this density (the MAP estimator):

$$\hat{\lambda}^{MAP} = \arg \min_{\lambda \in \mathbb{R}^M} \left\{ \|\mathbf{y} - \mathbf{f}_\lambda\|_n^2 + \frac{\gamma}{n} \sum_{j=1}^M \log(1 + \lambda_j^2 / \tau) \right\} \neq \hat{\lambda}^{EW}.$$

## SEW estimator: discussion

- *Precursors of SEW for the “diagonal” sequence model.*

Rivoirard (2004): minimax Bayes priors with heavy tails,  
Johnstone and Silverman (2005): “quasi-Cauchy” prior.

## Exponential weights: models with i.i.d. data

- An i.i.d. sample  $Z_1, \dots, Z_n$  from the distribution of an abstract random variable  $Z \in \mathcal{Z}$ .
- $Q(Z, f_\lambda)$  a given real-valued loss (prediction loss).

Define the probability measure  $S_n$  on  $\mathbb{R}^M$  by

$$S_n(d\lambda) = \frac{\exp \left\{ - \sum_{i=1}^n Q(Z_i, f_\lambda) / \beta \right\} \pi(d\lambda)}{\int_{\mathbb{R}^M} \exp \left\{ - \sum_{i=1}^n Q(Z_i, f_w) / \beta \right\} \pi(dw)}$$

with some  $\beta > 0$  and some prior measure  $\pi$ . Generalization of the previous definition: we replace

$$n \|\mathbf{y} - f_\lambda\|_n^2 \rightsquigarrow \sum_{i=1}^n Q(Z_i, f_\lambda).$$

## Mirror averaging

- Cumulative exponential weights (**mirror averaging**):

$$\hat{\lambda}_j^{MA} = \int_{\mathbb{R}^M} \lambda_j S(d\lambda), \quad j = 1, \dots, M, \quad \text{with } S = \frac{1}{n} \sum_{i=1}^n S_i$$

cf. Juditsky/Rigollet/T (2005) [even more general method: Juditsky/Nazin/T/Vayatis (2005)]. In a particular case we get the “progressive mixture method” of Catoni and Yang.

- Choose a prior measure  $\pi$  supported on a convex compact  $\Lambda \subset \mathbb{R}^M$  (e.g., on an  $\ell_1$  ball).

### Assumption JRT (2005).

The mapping  $\lambda \mapsto Q(Z, f_\lambda)$  is convex for all  $Z$  and there exists  $\beta > 0$  such that the function

$$\lambda \mapsto \mathbb{E} \exp \left( \frac{Q(Z, f_{\lambda'}) - Q(Z, f_\lambda)}{\beta} \right)$$

is concave on a convex compact set  $\Lambda \subset \mathbb{R}^M$  for all  $\lambda' \in \Lambda$ .

Roughly: “strong convexity on the average”.

## Example: Gaussian regression, squared loss

- Gaussian regression with random design :

$$Z = (X, Y), \quad X \in \mathbb{R}^d, \quad Y \in \mathbb{R} \quad \text{such that}$$

$$Y = f(X) + \xi,$$

$$\xi|X \sim \mathcal{N}(0, \sigma^2), \quad X \sim P_X, \quad \|f\|_\infty \leq L.$$

- Assumption on the dictionary  $\|f_j\|_\infty \leq L, j = 1, \dots, M.$

- The loss function

$$Q(Z, f_\lambda) = (Y - f_\lambda(X))^2 \quad \text{where} \quad f_\lambda = \sum_{j=1}^M \lambda_j f_j.$$

- Then  $A(\lambda) = \mathbb{E} Q(Z, f_\lambda) = \|f_\lambda - f\|_X^2 + \sigma^2, \quad \|f\|_X^2 \triangleq \int f^2 dP_X.$
- $\beta \geq 2\sigma^2 + 8L^2.$

## Example: density estimation with $L_2$ loss

- $Z = X \in \mathbb{R}^d$  with density  $f$ , such that  $\|f\|_\infty \leq L$ .
- Assumption on the dictionary:  $f_1, \dots, f_M$  are probability densities such that  $\|f_j\|_\infty \leq L$ .
- The loss function:

$$Q(X, f_\lambda) = \|f_\lambda\|^2 - 2f_\lambda(X) \quad \text{where} \quad \|f\|^2 = \int f^2(x) dx.$$

- The associated risk:

$$A(\lambda) = \mathbb{E} Q(X, f_\lambda) = \|f - f_\lambda\|^2 - \|f\|^2.$$

- $\beta \geq 12L$ .



## PAC-Bayesian bound for mirror averaging

Define the average risk:  $A(\lambda) = \mathbb{E}Q(Z, f_\lambda)$ .

**Lemma (PAC-Bayesian bound).**

Let  $f_{\hat{\lambda}^{MA}}$  be a mirror averaging estimator defined with  $\beta$  satisfying Assumption JRT and any prior  $\pi$  supported on a convex compact set  $\Lambda$ . Then

$$\mathbb{E} A(\hat{\lambda}^{MA}) \leq \inf_P \left\{ \int A(\lambda) P(d\lambda) + \frac{\beta \mathcal{K}(P, \pi)}{n+1} \right\}$$

where the infimum is taken over all probability measures  $P$  on  $\Lambda$  and  $\mathcal{K}(P, \pi)$  is the Kullback-Leibler divergence between  $P$  and  $\pi$ .

Proof follows the scheme of Juditsky, Rigollet and T. (2005), cf. Rigollet and Zhao (2006), Audibert (2006), Lounici (2007).

## SOI for Mirror Averaging

### Theorem

Assume that  $\sup_{|\lambda|_1 \leq 2R} \text{Spec}\{\nabla^2 A(\lambda)\} < \infty$  for some  $R > 0$ . Let  $f_{\hat{\lambda}^{MA}}$  be a mirror averaging estimator satisfying assumptions of the PAC lemma, with the **sparsity prior**  $\pi$  truncated to  $\{\lambda : |\lambda|_1 \leq 2R\}$  and  $\tau \sim 1/\sqrt{M(n \vee M)}$ . Then

$$\mathbb{E} A(\hat{\lambda}^{MA}) \leq \inf_{|\lambda|_1 \leq R} \left\{ A(\lambda) + \frac{CR^2 M(\lambda)}{n} \log \left( \frac{C'R \sqrt{M(n \vee M)}}{M(\lambda)} \right) \right\}.$$

- No restrictive assumption on the dictionary.
- Leading constant 1.

## Comparison with SOI for the LASSO

The LASSO type estimators

$$\hat{\lambda} = \arg \min_{\lambda \in \mathbb{R}^M} \left\{ \frac{1}{n} \sum_{i=1}^n Q(Z_i, f_\lambda) + r \sum_{j=1}^M |\lambda_j| \right\} .$$

van de Geer (2007,2008), Koltchinskii (2007,2008):

$$\mathbb{E} A(\hat{\lambda}) \leq \inf_{|\lambda|_1 \leq R} \left( \boxed{3} A(\lambda) + \frac{CR^2 M(\lambda) \log M}{\boxed{\kappa} n} \right)$$

where  $\kappa$  is a “Restricted Eigenvalue”, can be very small.

## Modified SEW estimators

Take the modified sparsity prior

$$q(\lambda) \propto \left( \prod_{j=1}^M \frac{e^{-\omega(\alpha\lambda_j)}}{(1 + \lambda_j/\tau)^2} \right) \mathbf{1}\{|\lambda|_1 \leq R\}$$

where  $\omega(\cdot)$  is Huber's function

$$\omega(t) = \begin{cases} t^2, & \text{if } |t| \leq 1, \\ 2|t|, & \text{if } |t| > 1, \end{cases}$$

$\alpha$  and  $\tau$  are small (ex.:  $\alpha \sim M^{-1}$ ,  $\tau \sim n^{-1/2}$ ),  $R$  is large ( $R \sim M$ ).

## Computation of SEW estimators

Consider the linear regression scenario:

$$\mathbf{y} = X\lambda + \xi.$$

$X$  is a  $n \times M$  deterministic design matrix,  $\lambda \in \mathbb{R}^M$  is an unknown vector and  $\xi \in \mathbb{R}^n$  is a Gaussian vector with i.i.d. components, with variances  $\sigma^2$ . The SEW estimator

$$\hat{\lambda}^{EW} \triangleq \int_{\mathbb{R}^M} \mathbf{u} g(\mathbf{u}) d\mathbf{u}$$

where the posterior density

$$g(\mathbf{u}) \propto \exp(-V(\mathbf{u}))$$

$$V(\mathbf{u}) = \beta^{-1} \|\mathbf{y} - X\mathbf{u}\|^2 + 2 \sum_{j=1}^M \log(\tau^2 + u_j^2).$$

## Langevin Monte Carlo

**Remark:** the posterior density  $g(\cdot)$  is the invariant density of the Langevin diffusion

$$\mathbf{L}_t = -\nabla V(\mathbf{L}_t) dt + \sqrt{2} d\mathbf{W}_t, \quad \mathbf{L}_0 = 0, \quad t > 0.$$

Here  $\mathbf{W}_t$  is the  $M$ -dimensional Brownian motion.

Let now  $\eta_1, \eta_2, \dots$  be i.i.d. standard normal random vectors. Set

$$\bar{\mathbf{L}}_0 = 0, \quad \bar{\mathbf{L}}_{k+1} = \bar{\mathbf{L}}_k - h\nabla V(\bar{\mathbf{L}}_k) + \sqrt{2h} \eta_k, \quad k = 0, 1, \dots$$

Then

$$\frac{1}{\lfloor Th^{-1} \rfloor} \sum_{k=1}^{\lfloor Th^{-1} \rfloor} \bar{\mathbf{L}}_k \approx \frac{1}{T} \int_0^T \mathbf{L}_t dt \xrightarrow[T \rightarrow \infty]{a.s.} \int_{\mathbb{R}^M} \mathbf{u} g(\mathbf{u}) d\mathbf{u} = \hat{\lambda}^{EW}.$$

## Simulations

### Example 1: Compressive sensing

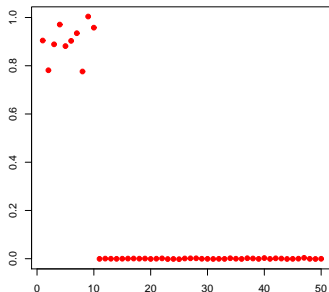
The entries of the matrix  $X$  are i.i.d. Rademacher random variables independent of the noise  $\xi$ .

$$\lambda_j = \mathbf{1}\{j \leq S\} \quad \text{and} \quad \sigma^2 = \frac{S}{9n}.$$

We apply the SEW estimator using Langevin Monte-Carlo with

$$\tau = 4\sigma/\sqrt{M}, \quad \beta = 4\sigma^2, \quad h = 0.0001.$$

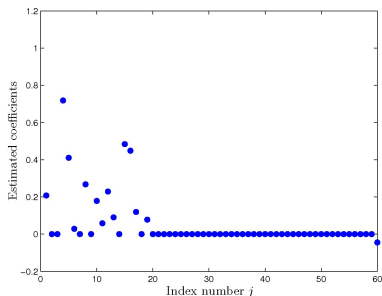
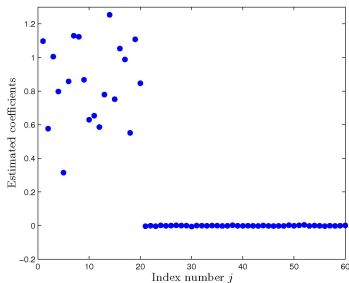
## Simulations



**Figure:** Typical result for Example 1 with  $n = 200$ ,  $M = 500$ ,  $S = 10$ ,  $h = 10^{-4}$ ,  $T = 5$ . The estimates of first 50 coefficients are plotted. In this example, we have  $\frac{1}{n} \|X(\hat{\lambda} - \lambda)\|^2 = 0.0021$ . The time of computation of the estimator was about 30 seconds.



## Example 1: Compressive sensing

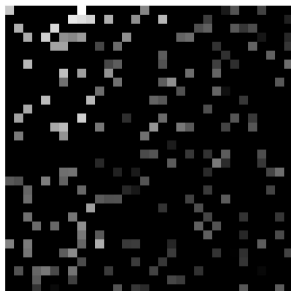
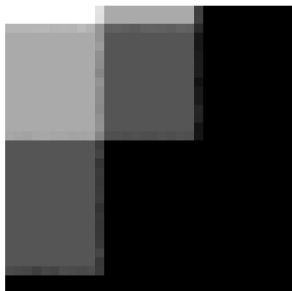


Typical outcome for  $n = 200$ ,  $M = 500$  and  $S = 20$ .

## Image denoising: A simple example

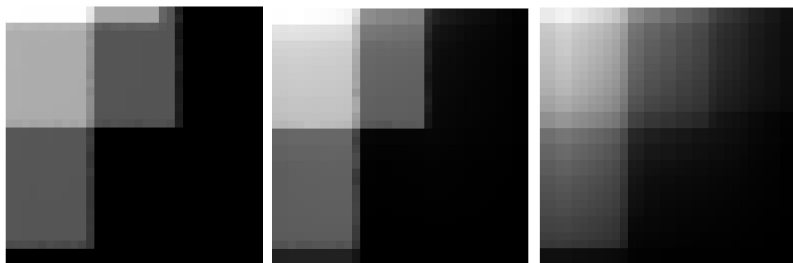
- Input:  $n, k$  positive integers and  $\sigma > 0$ .
- We generate  $n$  vectors  $U_i$  of  $\mathbb{R}^2$  uniformly distributed in  $[0, 1]^2$ .
- Covariates  $\phi_j(u) = \mathbf{1}_{\{[0, j_1/k] \times [0, j_2/k]\}}(u)$ .
- Errors: we generate a centered Gaussian vector  $\xi$  with covariance matrix  $\sigma^2 I$ .
- Response:  $Y_i = (\phi_1(U_i), \dots, \phi_{k^2}(U_i))^T \lambda + \xi_i$  where  $\lambda = [\mathbf{1}_{\{j \in \{10, 100, 200\}\}}]'$ .
- Tuning parameters: the same rule as previously.

## Image denoising



The original image and its sampled noisy version.

## Image denoising



Estimated images from observations with noise magnitudes 0.1, 0.5 and 1.

## Image denoising

$\sigma$	$n = 100$			$n = 200$		
	EWA	Lasso	Ideal LG	EWA	Lasso	Ideal LG
2	0.210 (0.072) $T = 1$	0.759 (0.562)	0.330 (0.145)	0.187 (0.048) $T = 1$	0.661 (0.503)	0.203 (0.086)
4	0.420 (0.222) $T = 1$	2.323 (1.257)	0.938 (0.631)	0.278 (0.132) $T = 1$	2.230 (1.137)	0.571 (0.324)

BICKEL, P.J., RITOV, Y. and TSYBAKOV, A.B. (2007) Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, to appear.

BUNEA, F., TSYBAKOV, A.B. and WEGKAMP, M.H. (2007) Aggregation for Gaussian regression. *Annals of Statistics*, v.35, 1674-1697.

BUNEA, F., TSYBAKOV, A.B. and WEGKAMP, M.H. (2007) Sparsity oracle inequalities for the Lasso. *Electronic Journal of Statistics*, v.1, 169-194.

DALALYAN, A. and TSYBAKOV, A.B. (2007) Aggregation by exponential weighting and sharp oracle inequalities. *COLT-2007*, 97-111.

DALALYAN, A. and TSYBAKOV, A.B. (2008) Aggregation by exponential weighting, sharp PAC-Bayesian bounds and sparsity. *Machine Learning*, v.72, 39-61.

JUDITSKY, A., RIGOLLET, P. and TSYBAKOV, A.B. Learning by mirror averaging. *Annals of Statistics*, to appear.