

# Joint estimation of the conditional mean and the conditional variance in high-dimensions

Joint work with M. Hebiri, K. Meziani, J. Salmon

SAPS IX, Le Mans, FRANCE



Arnak S. Dalalyan

ENSAE / CREST / GENES

# I. Problem presentation

---

## Continuous-time autoregression

**Observations** :  $(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \mathbb{R}$  obeying

$$\dot{y}_t = \mathbf{b}^*(\mathbf{x}_t) + \mathbf{s}^*(\mathbf{x}_t) \dot{W}_t, \quad t \in \mathcal{T} = [0, T],$$

where  $\mathbf{b}^* : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\mathbf{s}^* : \mathbb{R}^d \rightarrow \mathbb{R}_+$  such that

$$\text{Drift} : \mathbf{b}^*(\mathbf{x}_t) = \mathbf{E}[\dot{y}_t | \mathbf{x}_t].$$

$$\text{Volatility} : \mathbf{s}^{*2}(\mathbf{x}_t) = \mathbf{Var}[\dot{y}_t | \mathbf{x}_t].$$

$(W_t)_{t \geq 0}$  is a standard Wiener process (w.r.t a filtration  $(\mathcal{F}_t)_{t \geq 0}$ ).

The goal is to estimate the function  $\mathbf{b}^*$ .

## Discrete time (auto-)regression

**Observations** : finite collection  $(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \mathbb{R}$  obeying

$$y_t = \mathbf{b}^*(\mathbf{x}_t) + \mathbf{s}^*(\mathbf{x}_t) \xi_t, \quad t \in \mathcal{T} = \{1, \dots, T\},$$

where  $\mathbf{b}^* : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\mathbf{s}^* : \mathbb{R}^d \rightarrow \mathbb{R}_+$  such that

Conditional mean :  $\mathbf{E}[y_t | \mathbf{x}_t] = \mathbf{b}^*(\mathbf{x}_t)$ .

Conditional variance :  $\mathbf{Var}[y_t | \mathbf{x}_t] = \mathbf{s}^{*2}(\mathbf{x}_t)$ .

Therefore,  $\xi_t$ 's are such that  $\mathbf{E}[\xi_t | \mathbf{x}_t] = 0$  and  $\mathbf{Var}[\xi_t | \mathbf{x}_t] = 1$ .

They are often assumed Gaussian  $\mathcal{N}(0, 1)$  for simplicity.

The goal is to estimate the functions  $\mathbf{b}^*$  and  $\mathbf{s}^*$ .

## Sparsity Assumption

- In these settings, estimating  $\mathbf{b}^*$  and  $\mathbf{s}^*$  under no further assumption is an ill-posed problem.
- Sparsity scenario :  $\mathbf{b}^*$  and  $\mathbf{s}^*$  belong to some low dimensional spaces.

### Example : Homoscedastic regression

$$\forall \mathbf{x}, \quad \mathbf{b}^*(\mathbf{x}) = \sum_{j=1}^p f_j(\mathbf{x})\beta_j^* = [f_1(\mathbf{x}), \dots, f_p(\mathbf{x})]\beta^*, \quad \text{and} \quad \mathbf{s}^*(\mathbf{x}) \equiv \sigma^*$$

↪ Dictionary  $\{f_1, \dots, f_p\}$  of functions from  $\mathbb{R}^d$  to  $\mathbb{R}$

↪ Unknown vector  $(\beta^*, \sigma^*) \in \mathbb{R}^p \times \mathbb{R}$ , sparse vector  $\beta^*$

↪ Sparsity index :  $p^* = |\beta^*|_0 := \sum_{j=1}^p \mathbb{1}(\beta_j^* \neq 0)$  with  $p^* \ll p$

## Main assumptions on $b^*$ and $s^*$

Re-parametrize by the inverse of the conditional volatility  $s^*$

$$r^*(\mathbf{x}) = \frac{1}{s^*(\mathbf{x})} \quad \text{and} \quad f^*(\mathbf{x}) = \frac{b^*(\mathbf{x})}{s^*(\mathbf{x})}.$$

### Group Sparsity Assumption

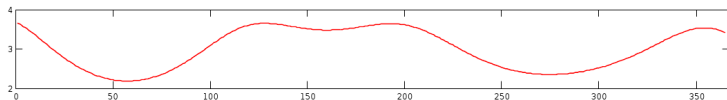
For  $p$  given functions  $f_1, \dots, f_p$  mapping  $\mathbb{R}^d$  into  $\mathbb{R}$ , there is a vector  $\phi^* \in \mathbb{R}^p$  such that  $f^*(\mathbf{x}) = \sum_{j=1}^p \phi_j^* f_j(\mathbf{x})$ . Furthermore, for a given partition  $G_1, \dots, G_K$  of  $\{1, \dots, p\}$ , the vector  $\phi^*$  is group-sparse that is  $\text{Card}(\{k : |\phi_{G_k}^*|_2 \neq 0\}) \ll K$ .

### Low dimensional volatility assumption

For  $q$  given functions  $r_1, \dots, r_q$  mapping  $\mathbb{R}^d$  into  $\mathbb{R}_+$ , there is a vector  $\alpha^* \in \mathbb{R}^q$  such that  $r^*(\mathbf{x}) = \sum_{\ell=1}^q \alpha_\ell^* r_\ell(\mathbf{x})$  for almost every  $\mathbf{x} \in \mathbb{R}^d$ .

## Motivations for these assumptions

- ▶ **Group sparsity assumption** is relevant in a sparse additive model, that is when
  - ↪  $\mathbf{f}^*(\mathbf{x}) = \psi_1(x_1) + \dots + \psi_d(x_d)$  s.t.  $\psi_j \equiv 0$  for most  $j$ ,
  - ↪ projection on basis  $\psi_j(x_j) \approx \sum_{\ell=1}^{K_j} \phi_{\ell,j}^* f_{\ell}(x_j)$ ,
  - ↪ group sparsity of  $\phi = (\phi_{\ell,j})$ .
- ▶ **Low dimensionality of the volatility** occurs, for instance when the noise is block-wise homoscedastic or periodic.



## II. Relation to previous work

---



# Homoscedastic regression

The model is

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \sigma^* \boldsymbol{\xi}$$

Observations :  $\mathbf{Y} = [y_1, \dots, y_T]^\top \in \mathbb{R}^T$

Noise :  $\boldsymbol{\xi} = [\xi_1, \dots, \xi_T]^\top \in \mathbb{R}^T$

Design Matrix :  $\mathbf{X} = x_{t,j}$  with  $x_{t,j} = [f_j(\mathbf{x}_t)] \in \mathbb{R}$

Coefficients :  $\boldsymbol{\beta}^* = [\beta_1^*, \dots, \beta_p^*]^\top \in \mathbb{R}^p$

Standard deviation :  $\mathbf{s}^*(\mathbf{x}_t) \equiv \sigma^* \in \mathbb{R}_*^+$

Recall that the sparsity assumption postulates that  $|\boldsymbol{\beta}^*|_0 = p^* \ll p$ .

## Most popular methods : Lasso and Dantzig selector

- ◀ The LASSO of Tibshirani (1996) is defined as

$$\hat{\beta}^{\text{Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{\|\mathbf{Y} - \mathbf{X}\beta\|_2^2}{2\sigma^{*2}} + \lambda \sum_{j=1}^p |\mathbf{x}_j|_2 |\beta_j| \right)$$

- ◀ The Dantzig selector of Candès and Tao (2007) is

$$\hat{\beta}^{\text{DS}} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{j=1}^p |\mathbf{x}_j|_2 |\beta_j| : \max_{j=1, \dots, p}, \frac{|\mathbf{x}_j^\top (\mathbf{Y} - \mathbf{X}\beta)|}{|\mathbf{x}_j|_2} \leq \lambda \right\}$$

For a tuning parameter satisfying  $\lambda \propto 1/\sigma^*$ , sharp oracle inequalities are available, e.g., Bickel *et al.* (2009).

$$\mathbf{E} \left( \frac{1}{T} \|\mathbf{X}(\hat{\beta}^\bullet - \beta)\|_2^2 \right) \leq C \frac{p^* \log(p)}{T}.$$

To correctly tune the parameter  $\lambda$ , the knowledge of  $\sigma^*$  is necessary.

## Joint estimation of $\beta^*$ and $\sigma^*$

- ◀ Scaled Lasso, **Städler *et al.* (2010)**,

$$(\hat{\beta}^{\text{ScL}}, \hat{\sigma}^{\text{ScL}}) = \arg \min_{\beta, \sigma} \left( T \log(\sigma) + \frac{\|\mathbf{Y} - \mathbf{X}\beta\|_2^2}{2\sigma^2} + \frac{\lambda}{\sigma} \sum_{j=1}^p \|\mathbf{x}_j\|_2 |\beta_j| \right).$$

This can be recast in a convex problem (do  $\rho := \frac{1}{\sigma}$  and  $\phi := \frac{\beta}{\sigma}$ ):

$$\arg \min_{\phi, \rho} \left( -T \log(\rho) + \frac{\|\rho \mathbf{Y} - \mathbf{X}\phi\|_2^2}{2} + \lambda \sum_{j=1}^p \|\mathbf{x}_j\|_2 |\phi_j| \right).$$

- ◀ Square-Root Lasso **Antoniadis (2010)**, **Belloni *et al.* (2011)**, **Sun & Zhang (2012)**, **Gautier & Tsybakov (2011)**,

$$\hat{\beta}^{\text{SqR-Lasso}} = \arg \min_{\beta} \left( \|\mathbf{Y} - \mathbf{X}\beta\|_2 + \lambda \sum_{j=1}^p \|\mathbf{x}_{:,j}\|_2 |\beta_j| \right)$$

$$\hat{\sigma}^{\text{SqR-Lasso}} = T^{-1/2} \|\mathbf{Y} - \mathbf{X}\hat{\beta}^{\text{SqR-Lasso}}\|_2.$$

- ◀ Scaled DS version proposed by **Dalalyan & Chen (2012)**: sharp analysis and computational advantages.

### III. Main results

---

# Estimation for heteroscedastic regression

## Continuous time

**Observations** :  $(\mathbf{x}_t, y_t)_{t \leq T}$  obeying  $\dot{y}_t = \mathbf{b}^*(\mathbf{x}_t) + \mathbf{s}^*(\mathbf{x}_t) \dot{W}_t$ , thus

$$\frac{\dot{y}_t}{\mathbf{s}^*(\mathbf{x}_t)} = \left( \sum_{j=1}^p \phi_j^* f_j(\mathbf{x}_t) \right) + \dot{W}_t = \mathbf{X}(t)\phi^* + \dot{W}_t, \quad t \in \mathcal{T} = [0, T].$$

### Scaled Heteroscedastic Lasso :

$$\hat{\phi}^{\text{ScHeL}} = \arg \min_{\phi \in \mathbb{R}^p} \left\{ \frac{1}{2} \int_0^T \left( \frac{\dot{y}_t}{\mathbf{s}^*(\mathbf{x}_t)} - \mathbf{X}(t)\phi \right)^2 dt + \lambda \sum_{k=1}^K |\mathbf{X}_{G_k} \phi_{G_k}|_2 \right\},$$

$$\text{with } |\mathbf{X}_G \phi_G|_2^2 = \int_0^T \left( \sum_{j \in G} \mathbf{X}_j(t) \phi_j \right)^2 dt = \int_0^T \left( \sum_{j \in G} f_j(\mathbf{x}_t) \phi_j \right)^2 dt.$$

**Computation** : the estimator  $\hat{\phi}^{\text{ScHeL}}$  can be efficiently computed even for very large dimensions  $p$  by solving a second-order cone program.

# Estimation for heteroscedastic regression

## Discrete time

**Observations** :  $(\mathbf{x}_t, y_t)_{t=1, \dots, T}$  obeying

$$y_t = \mathbf{b}^*(\mathbf{x}_t) + \mathbf{s}^*(\mathbf{x}_t) \xi_t = \mathbf{r}^*(\mathbf{x}_t)^{-1} (\mathbf{f}^*(\mathbf{x}_t) + \xi_t).$$

Under our assumptions

$$\mathbf{f}^*(\mathbf{x}_t) = \sum_{j=1}^p \phi_j^* \mathbf{f}_j(\mathbf{x}_t) = \mathbf{X}(t) \phi^*, \quad \mathbf{r}^*(\mathbf{x}_t) = \sum_{\ell=1}^q \alpha_\ell^* r_\ell(\mathbf{x}_t) = \mathbf{R}(t) \alpha^*.$$

Thus,  $[\mathbf{f}^*(\mathbf{x}_1), \dots, \mathbf{f}^*(\mathbf{x}_T)]^\top = \mathbf{X} \phi^*$  and  $[\mathbf{r}^*(\mathbf{x}_1), \dots, \mathbf{r}^*(\mathbf{x}_T)]^\top = \mathbf{R} \alpha^*$ . This leads to

$$\mathbf{D}_Y \mathbf{R} \alpha^* = \mathbf{X} \phi^* + \xi, \quad \mathbf{D}_Y = \text{diag}(y_t; t = 1, \dots, T).$$

**Scaled Heteroscedastic Lasso** :  $(\hat{\phi}^{\text{ScHeL}}, \hat{\alpha}^{\text{ScHeL}})$  solution to

$$\min_{\phi \in \mathbb{R}^p} \left\{ \underbrace{- \sum_{t=1}^T \log(\mathbf{R}(t) \alpha) + \frac{1}{2} \|\mathbf{D}_Y \mathbf{R} \alpha - \mathbf{X} \phi\|_2^2}_{\text{Gaussian log-likelihood}} + \underbrace{\lambda \sum_{k=1}^K \|\mathbf{X}_{G_k} \phi_{G_k}\|_2}_{\text{sparsity promoting penalty}} \right\}.$$

## Finite sample risk bounds for the ScHeL

**Theorem** Consider either continuous-time model or discrete-time with sub-Gaussian errors  $\xi$ . Let  $K^*$  (resp.  $p^*$ ) be the number of relevant groups (resp. coordinates of  $\phi^*$ ). Let  $\varepsilon \in (0, 1)$  be a tolerance level and set

$$\lambda_k = 4(\sqrt{|G_k|} + \sqrt{\log(K/\varepsilon)}).$$

Under some assumptions, with probability at least  $1 - 2\varepsilon$ ,

$$\|\mathbf{X}(\hat{\phi} - \phi^*)\|_2 \leq D_{T,\varepsilon}^{3/2} \sqrt{q \log(2q/\varepsilon)} + D_{T,\varepsilon} \sqrt{p^* + K^* \log(K/\varepsilon)}.$$

$$\|\mathbf{R}(\hat{\alpha} - \alpha^*)\|_2 \leq D_{T,\varepsilon}^{3/2} \sqrt{q \log(2q/\varepsilon)} + D_{T,\varepsilon} \sqrt{p^* + K^* \log(K/\varepsilon)},$$

where  $D_{T,\varepsilon} \propto (\max_t |f^*(\mathbf{x}_t)| + \log(2T/\varepsilon))$ .

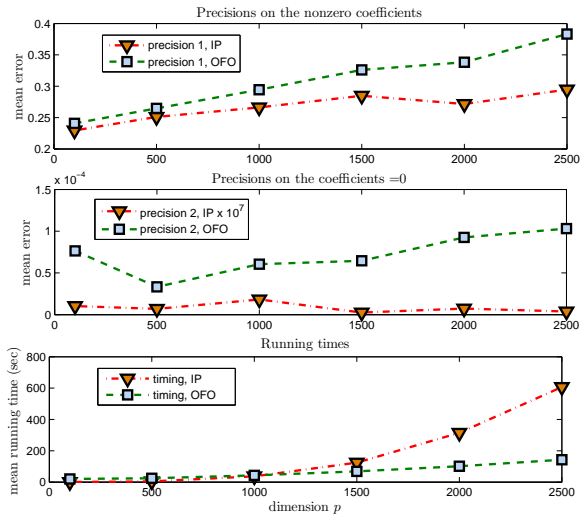
## Summary

New procedures named ScHeL and ScHeDs :

- ◀ Suitable for fitting the heteroscedastic regression model.
- ◀ Simultaneous estimation of the mean and the variance functions.
- ◀ Takes into account group sparsity.
- ◀ Implemented using two different solvers :
  - ↪ primal-dual interior point method (highly accurate),
  - ↪ optimal first-order method (moderately accurate but with cheap iterations).
- ◀ Competitive with state-of-the art algorithms
  - ↪ applicable in a much more general framework.










Manuscript, proofs and codes available on request.





**FIGURE:** Comparing interior point (IP) vs. optimal first-order (OFO) method. Top & middle : MSE of  $\hat{\beta}^{\text{ScHeDs}}$ . Bottom : running times.

# References I

-  A. Antoniadis, [Comments on :  \$\ell\_1\$ -penalization for mixture regression models](#), TEST **19** (2010), no. 2, 257–258. MR 2677723
-  A. Belloni, V. Chernozhukov, and L. Wang, [Square-root Lasso : Pivotal recovery of sparse signals via conic programming](#), Biometrika **98** (2011), no. 4, 791–806.
-  P. J. Bickel, Y. Ritov, and A. B. Tsybakov, [Simultaneous analysis of Lasso and Dantzig selector](#), Ann. Statist. **37** (2009), no. 4, 1705–1732.
-  E. J. Candès and T. Tao, [The Dantzig selector : statistical estimation when  \$p\$  is much larger than  \$n\$](#) , Ann. Statist. **35** (2007), no. 6, 2392–2404.
-  Arnak S. Dalalyan and Yin Chen, [Fused sparsity and robust estimation for linear models with unknown variance](#), Advances in Neural Information Processing Systems 25 : NIPS, 2012, pp. 1268–1276.
-  E. Gautier and A. Tsybakov, [High-dimensional instrumental variables regression and confidence sets](#), September 2011.
-  N. Städler, P. Bühlmann, and Sara s van de Geer,  [\$\ell\_1\$ -penalization for mixture regression models](#), TEST **19** (2010), no. 2, 209–256.
-  T. Sun and C.-H. Zhang, [Scaled sparse linear regression](#), Biometrika **99** (2012), no. 4, 879–898.
-  R. Tibshirani, [Regression shrinkage and selection via the Lasso](#), J. Roy. Statist. Soc. Ser. B **58** (1996), no. 1, 267–288.

