

# Penalized QMLEs of time series regressions

Christian Francq, Sébastien Laurent and Julie Schnaitmann

ECODEP CONFERENCE

September 30-October 1, 2024, IHP

## Framework

We consider a **linear regression model** applied to the **components of a time series**, with the aim to

- identify the **constant conditional beta** coefficients;
- identify the **zero conditional betas**.

To address the **non-identifiability** of the parameters when a conditional beta is constant, and to reduce the **numerical complexity** of the estimation procedure

- we employ a **lasso-type\*** **multi-step<sup>†</sup>** QMLE.

---

\*this penalized estimator simplifies the model by shrinking the estimates to their simplest form when the beta is constant.

<sup>†</sup>which first captures the dynamics of the regressors before estimating the dynamics of the betas.

# Outline

- 1 Time series model for a beta coefficient
  - The autoregressive conditional beta (ACB) model
  - Unpenalized multi-step QMLE
  - Lasso multi-step QMLE
- 2 Asymptotic properties of Lasso multi-step QMLEs
  - Conditions for (in)consistency
  - Adaptive version
  - Penalized estimator for detecting relevant betas
- 3 Numerical study
  - Optimization algorithm
  - Monte-Carlo Simulations
  - Financial application

## Regressions with time-varying betas

Let  $(y_t, \mathbf{x}_t^\top)$  be a time series of  $1 + p$  random variables.

Under **stationarity** and existence of second-order moments, it makes sense to predict  $y_t$  at a long horizon by a **time-constant mean or regression**:

$$Ey_1 \quad \text{or} \quad Ey_1 + \boldsymbol{\beta}^\top (\mathbf{x}_t - E\mathbf{x}_1).$$

Short-term predictions can incorporate the **information**  $\mathcal{F}_t$  available at time  $t$  (given by  $\{y_u, \mathbf{x}_u; u \leq t\}$  and possibly some vector  $\mathbf{z}_t$  of exogenous variables): given  $\mathcal{F}_{t-1}$ ,  $y_t$  is better predicted by

$$E_{t-1}y_t := E(y_t \mid \mathcal{F}_{t-1}) \quad \text{or} \quad E_{t-1}y_t + \boldsymbol{\beta}_t^\top (\mathbf{x}_t - E_{t-1}\mathbf{x}_t)$$

with

$$\boldsymbol{\beta}_t = \text{Cov}_{t-1}(y_t, \mathbf{x}_t) \text{Var}_{t-1}^{-1} \mathbf{x}_t.$$

## Regressions with time-varying betas

Let  $(y_t, \mathbf{x}_t^\top)$  be a time series of  $1 + p$  random variables.

Under **stationarity** and existence of second-order moments, it makes sense to predict  $y_t$  at a long horizon by a **time-constant mean or regression**:

$$Ey_1 \quad \text{or} \quad Ey_1 + \boldsymbol{\beta}^\top (\mathbf{x}_t - E\mathbf{x}_1).$$

Short-term predictions can incorporate the **information**  $\mathcal{F}_t$  available at time  $t$  (given by  $\{y_u, \mathbf{x}_u; u \leq t\}$  and possibly some vector  $\mathbf{z}_t$  of exogenous variables): given  $\mathcal{F}_{t-1}$ ,  $y_t$  is better predicted by

$$E_{t-1}y_t := E(y_t \mid \mathcal{F}_{t-1}) \quad \text{or} \quad E_{t-1}y_t + \boldsymbol{\beta}_t^\top (\mathbf{x}_t - E_{t-1}\mathbf{x}_t)$$

with

$$\boldsymbol{\beta}_t = \text{Cov}_{t-1}(y_t, \mathbf{x}_t) \text{Var}_{t-1}^{-1} \mathbf{x}_t.$$

## Which time series model for a conditional beta ?

Blasques, Francq and Laurent (2024) (BFL) propose the autoregressive conditional beta (ACB) model

$$\left\{ \begin{array}{l} x_{it} = \mu_{0i} + \varepsilon_{it}, \quad \varepsilon_{it} = g_{it}\eta_{it}, \\ g_{i,t+1}^2 = \omega_{0i} + \alpha_{0i}\varepsilon_{it}^2 + \beta_{0i}g_{it}^2, \\ y_t = \beta_{1t}x_{1t} + \dots + \beta_{pt}x_{pt} + v_t, \quad v_t = g_t\eta_t, \\ \beta_{i,t+1} = \varpi_{0i} + \xi_{0i}\frac{v_t x_{it}}{\mu_{0i}^2 + g_{it}^2} + c_{0i}\beta_{it} + \gamma_{01i}z_{1t} + \dots + \gamma_{0qi}z_{qt}, \\ g_{t+1}^2 = \omega_0 + \alpha_0 v_t^2 + \beta_0 g_t^2, \end{array} \right.$$

with obvious notations ( $\mathbf{x}_t = (x_{1t}, \dots, x_{pt})^\top$  are the regressors and  $\mathbf{z}_t = (z_{1t}, \dots, z_{qt})^\top$  is a vector of exogenous variables).

► More about ACB

## Unidentifiability of the constant betas

The beta  $\beta_{it}$  is constant if (and only if, under some regularity conditions)

$$\xi_{0i} = \gamma_{01i} = \dots = \gamma_{0qi} = 0.$$

Note that, when this relation holds the parameter  $c_{0i}$  is not well defined because the model remains the same for all values of  $(\varpi_{i0}, c_{0i})$  such that  $\varpi_{i0}/(1 - c_{0i})$  is fixed:

$$\beta_i = \varpi_{0i} + 0 \times \frac{v_t x_{it}}{\mu_{0i}^2 + g_{it}^2} + c_{0i} \beta_i$$

## First-step estimator of the regressor GARCH models

The observations are  $(y_t, \mathbf{x}_t, \mathbf{z}_t)$  for  $t = 1, \dots, n$ . The unknown parameter is

$$\varphi_0 = (\boldsymbol{\theta}_0^\top, \boldsymbol{\nu}_0^\top)^\top, \quad \boldsymbol{\theta}_0 = \underbrace{(\boldsymbol{\theta}_0^{(1)\top}, \dots, \boldsymbol{\theta}_0^{(p)\top})^\top}_{\text{GARCH parameters of the regressors}}.$$

The regressor GARCH(1,1) models can be [estimated in parallel](#) by the standard QMLE  $\hat{\boldsymbol{\theta}}^{(i)}$  involving  $(x_{i1}, \dots, x_{in})$ .

Let

$$\hat{\boldsymbol{\theta}} = \left( \hat{\boldsymbol{\theta}}^{(1)\top}, \dots, \hat{\boldsymbol{\theta}}^{(p)\top} \right)^\top.$$



## The remaining parameters

Let

$$\boldsymbol{\vartheta}_0^{(0)} = (\omega_0, \alpha_0, \beta_0)^\top$$

be the GARCH(1,1) parameters of the regression error term. Let

$$\boldsymbol{\vartheta}_0^{(i)} = (\varpi_{0i}, \xi_{0i}, c_{0i}, \gamma_{01i}, \dots, \gamma_{0qi})^\top$$

the parameters that are specific to  $\beta_{it}$ , for  $i \in \{1, \dots, p\}$ . The vector of the remaining parameters is thus

$$\boldsymbol{\vartheta}_0 = \left( \boldsymbol{\vartheta}_0^{(0)\top}, \boldsymbol{\vartheta}_0^{(1)\top}, \dots, \boldsymbol{\vartheta}_0^{(p)\top} \right)^\top \in \mathbb{R}^{d_2}.$$

Denote by  $\boldsymbol{\vartheta}$  be a generic element of the parameter space  $\Theta_{\boldsymbol{\vartheta}} \subset (0, \infty) \times [0, \infty)^2 \times \mathbb{R}^{p(3+q)}$  and let the generic parameter

$$\boldsymbol{\varphi} = (\boldsymbol{\theta}^\top, \boldsymbol{\vartheta}^\top)^\top.$$

## Second-step estimator of the beta dynamics

We estimate  $\vartheta_0$  by

$$\hat{\vartheta} = \arg \min_{\vartheta \in \Theta_{\vartheta}} \tilde{O}(\hat{\theta}, \vartheta), \quad \tilde{O}_n(\varphi) = \frac{1}{n} \sum_{t=2}^n \tilde{\ell}_t(\varphi),$$

where, with standard notation,

$$\tilde{\ell}_t(\varphi) = \frac{\tilde{v}_t^2(\varphi)}{\tilde{g}_t^2(\varphi)} + \log \tilde{g}_t^2(\varphi), \quad \tilde{v}_t(\varphi) = y_t - \sum_{i=1}^p \tilde{\beta}_{it}(\varphi) x_{it},$$

$$\tilde{g}_t^2(\varphi) = \omega + \alpha \tilde{v}_{t-1}^2(\varphi) + \beta \tilde{g}_{t-1}^2(\varphi),$$

$$\tilde{\beta}_{it}(\varphi) = \varpi_i + \xi_i \frac{\tilde{v}_{t-1}(\varphi) x_{i,t-1}}{\mu_i^2 + \tilde{g}_{i,t-1}^2(\theta)} + c_i \tilde{\beta}_{i,t-1}(\varphi) + \sum_{j=1}^q \gamma_{ji} z_{j,t-1}.$$

## Motivation for penalized estimators

BFL showed the consistency and asymptotic normality (CAN) of

$$\hat{\varphi} = (\hat{\theta}^T, \hat{\vartheta}^T)^T$$

under the assumption  $\xi_{0i} \neq 0$  for  $i = 1, \dots, p$  (and regularity conditions).

We define and study Lasso-type estimators which present the advantages:

- of being CAN under a more general framework than the multi-step QMLE;
- to lead to sparsity of the parameters, and thus to more parsimonious models;
- of being variable selection consistent (impossibility to use BIC for, say  $p \geq 5$  which corresponds to more than 32000 models).

## The penalized components

We want to penalize non-zero estimated values of the parameters  $\xi_i, c_i, \gamma_{1i}, \dots, \gamma_{qi}$ , for  $i = 1, \dots, p$ , to solve the non-identifiability problem by favoring the solution

$$\widehat{\boldsymbol{\vartheta}}_n^{(i)} = (\widehat{\omega}_i, 0, 0, 0, \dots, 0)^\top$$

if  $\beta_{it}$  is constant.

Let  $S = \{5, \dots, q + 6, q + 8, \dots, d_2\}$  be the set of the components that we want to shrink.

## Partially penalized estimator for constant beta detection

Thus consider the penalized QMLE

$$\hat{\boldsymbol{\vartheta}}_n = \arg \min_{\boldsymbol{\vartheta} \in \Theta_{\boldsymbol{\vartheta}}} \tilde{Q}_n(\hat{\boldsymbol{\theta}}, \boldsymbol{\vartheta}), \quad \tilde{Q}_n(\hat{\boldsymbol{\theta}}, \boldsymbol{\vartheta}) = \tilde{O}_n(\hat{\boldsymbol{\theta}}, \boldsymbol{\vartheta}) + \lambda_n p(\boldsymbol{\vartheta}),$$

where  $\lambda_n \geq 0$  and  $p(\boldsymbol{\vartheta}) = \sum_{i \in S} |\vartheta_i|$ .

The Lasso multi-step QMLE  $\hat{\boldsymbol{\varphi}}_n = \left( \hat{\boldsymbol{\theta}}^\top, \hat{\boldsymbol{\vartheta}}_n^\top \right)^\top$  encourages sparsity of specific components of  $\hat{\boldsymbol{\vartheta}}_n$  in order to allow constant  $\beta_{it}(\hat{\boldsymbol{\varphi}}_n)$ 's.

## Consistency to a biased model

At the price of some bias, the Lasso estimator solves the lack of identifiability when  $\lambda_n \rightarrow \lambda_0 > 0$  and

**A**( $\lambda_0$ ):

$$Q_{\lambda_0}(\boldsymbol{\vartheta}) = E\ell_1(\boldsymbol{\theta}_0, \boldsymbol{\vartheta}) + \lambda_0 p(\boldsymbol{\vartheta})$$

admits a minimum over  $\Theta_{\boldsymbol{\vartheta}}$  at some unique point  $\boldsymbol{\vartheta}^*$ .

Assume standard regularity conditions and **A**( $\lambda_0$ ) for some  $\lambda_0 > 0$ .  
If  $\lambda_n \rightarrow \lambda_0$ , then  $\hat{\boldsymbol{\vartheta}}_n \rightarrow \boldsymbol{\vartheta}^*$  in probability as  $n \rightarrow \infty$ .

## Comments on $\mathbf{A}(\lambda_0)$

- Note that  $\mathbf{A}(\lambda_0)$  would be satisfied for all  $\lambda_0 \geq 0$  if the function  $\vartheta \mapsto E\ell_1(\boldsymbol{\theta}_0, \vartheta)$  were strictly convex.
- It can be shown that  $\mathbf{A}(\lambda_0)$  generally holds true, at least when  $\lambda_0 > 0$  is sufficiently small.
- Note however that one can not take  $\lambda_0 = 0$ , that is  $\mathbf{A}(0)$  does not hold true when a beta is constant.
- Note also that, in general  $\vartheta^* \neq \vartheta_0$ . Thus the penalization introduces an asymptotic bias, but it can be shown that the bias is small when  $\lambda_0 > 0$  is small.

► Illustration

## FOC satisfied by $\vartheta^*$

Let  $\partial p(\vartheta^*)$  be the set of the subgradients of  $p$  on  $\Theta_\vartheta$  at  $\vartheta^* = (\vartheta_1^*, \dots, \vartheta_{d_2}^*)$ . When the limit  $\vartheta^*$  belongs to  $\overset{\circ}{\Theta}_\vartheta$ , the interior of  $\Theta_\vartheta$ , it must satisfy **the subgradient first-order condition**

$$0 \in \partial^\circ Q_{\lambda_0}(\vartheta^*) := \left\{ \frac{\partial E\ell_1(\theta_0, \vartheta^*)}{\partial \vartheta} \right\} + \lambda_0 \partial p(\vartheta^*).$$

When  $E\ell_1(\theta_0, \cdot)$  is convex,  $\partial^\circ Q_{\lambda_0}(\vartheta)$  is the subdifferential of  $Q_{\lambda_0}(\vartheta)$  (the FOC is necessary and sufficient).

More generally, Clarke (1975) showed that, since the functions  $E\ell_1(\theta_0, \cdot)$  and  $p(\cdot)$  are locally Lipschitz,  $\partial^\circ Q_{\lambda_0}(\vartheta)$  is the set of the **generalized gradients** of  $Q_{\lambda_0}(\vartheta)$ , which contains its subgradients (the **FOC is necessary**).



## Sparsity of $\vartheta^*$

The set  $\partial p(\vartheta)$  consists of the vectors  $\mathbf{u} = (u_1, \dots, u_{d_2})^\top$  where:

- for  $i \in \bar{S}$ ,  $u_i = 0$ ;
- for  $i \in S$ ,
  - $u_i = \text{sign}(\vartheta_i)$  when  $\vartheta_i \neq 0$ ,
  - $u_i \in [-1, 1]$  when  $\vartheta_i = 0$ .

It follows that if  $\vartheta^* \in \mathring{\Theta}_{\vartheta}$ , for all  $i \in S$ ,

- $\left| \frac{\partial E\ell_1(\boldsymbol{\theta}_0, \vartheta^*)}{\partial \vartheta_i} \right| \leq \lambda_0$  if  $\vartheta_i^* = 0$ ;
- $\frac{\partial E\ell_1(\boldsymbol{\theta}_0, \vartheta^*)}{\partial \vartheta_i} = -\lambda_0 \text{sign}(\vartheta_i^*)$  if  $\vartheta_i^* \neq 0$ .

► Estimator's FOC

## Upper bound for the penalty term

If  $\lambda_n > \bar{\lambda}$  for some sufficiently large  $\bar{\lambda}$ , then the penalized estimator  $\widehat{\boldsymbol{\vartheta}}_n$  is equal to the constrained QMLE  $\widehat{\boldsymbol{\varphi}}_n^c = (\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\vartheta}}_n^c)$  such that

$$\widehat{\boldsymbol{\vartheta}}_n^c = \arg \min_{\boldsymbol{\vartheta} \in \Theta_{\mathcal{J}}^c} \widetilde{O}_n(\widehat{\boldsymbol{\theta}}, \boldsymbol{\vartheta}),$$

where  $\Theta_{\mathcal{J}}^c$  denotes the set of the parameters  $\boldsymbol{\vartheta} \in \Theta_{\mathcal{J}}$  with  $i$ -th element  $\vartheta_i = 0$  for all  $i \in S$ .

The FOC entails that

$$\bar{\lambda} \geq \lambda^i := \max_{j \in S} \left| \frac{1}{n} \sum_{t=2}^n \frac{\partial}{\partial \vartheta_j} \tilde{\ell}_t(\widehat{\boldsymbol{\varphi}}_n^c) \right|.$$

## Another bound for $\bar{\lambda}$

Let  $\Theta_{\vartheta_*}^{j-}$  (resp.  $\Theta_{\vartheta_*}^{j+}$ ) denotes the set of the parameters  $\vartheta \in \Theta_{\vartheta}$  with  $j$ -th element  $\vartheta_j \leq 0$  (resp.  $\vartheta_j \geq 0$ ) and the other component of  $\vartheta$  are that of  $\vartheta_*$ . Assume  $\hat{\vartheta}_n^c \in \mathring{\Theta}_{\vartheta}$ .

One can take

$$\bar{\lambda} = \lambda^s := \max_{j \in S} \max \left\{ \sup_{\vartheta \in \Theta_{\hat{\vartheta}_n^c}^{j-}} \frac{1}{n} \sum_{t=2}^n \frac{\partial}{\partial \vartheta_j} \tilde{\ell}_t(\hat{\theta}, \vartheta), \right. \\ \left. - \inf_{\vartheta \in \Theta_{\hat{\vartheta}_n^c}^{j+}} \frac{1}{n} \sum_{t=2}^n \frac{\partial}{\partial \vartheta_j} \tilde{\ell}_t(\hat{\theta}, \vartheta) \right\}.$$

If  $\vartheta \mapsto \tilde{O}_n(\hat{\theta}, \vartheta)$  is a strictly convex, then  $\bar{\lambda} = \lambda^i = \lambda^s$ . ▶ case  $\lambda^i < \lambda^s$

## Adaptive estimator for constant beta detection

Let the data-driven weights

$$\widehat{\delta}_i = \frac{1}{|\widehat{\vartheta}_{ni}|} 1_{\widehat{\vartheta}_{ni} \neq 0} + \infty 1_{\widehat{\vartheta}_{ni} = 0} \text{ for } i \in S,$$

and the adaptive penalized QMLE

$$\widehat{\vartheta}_n^a = \arg \min_{\vartheta \in \Theta_\vartheta} \widetilde{O}_n(\widehat{\theta}, \vartheta) + \lambda_n^a \sum_{i \in S} \widehat{\delta}_i |\vartheta_i|.$$

Let  $\mathcal{A} = \mathcal{A}(\vartheta_0)$  be the subset of the **active** (and shrunk) components of the model of parameter  $\vartheta_0$ , *i.e.* the set of indices  $i \in S$  such that  $\vartheta_{0i} \neq 0$ . Let  $\mathcal{I} = S \cap \overline{\mathcal{A}(\vartheta_0)}$  be the subset of the **inactive** components. Let  $\mathbf{A}$  be the  $d_6 \times d_2$  selector matrix which selects the active (or not shrunk) components of  $\vartheta_0$ .

## Asymptotics of the adaptive penalized multi-step QMLE

Let the previous assumptions and  $\mathcal{A}(\vartheta_0) = \mathcal{A}(\vartheta^*)$  for some  $\lambda_0 > 0$ . If  $\lambda_n \rightarrow \lambda_0$ , if there exists  $n_0$  such that  $\lambda_n^a > 0$  for all  $n \geq n_0$ , and  $\sqrt{n}\lambda_n^a \rightarrow \lambda_0^a \geq 0$ , then the components of  $\widehat{\vartheta}_n^a$  whose indices belong to  $\mathcal{I}$  are zero with probability tending to 1, and

$$\sqrt{n}\mathbf{A} \left( \widehat{\vartheta}_n^a - \vartheta_0 \right) \xrightarrow{d} \arg \min_{\mathbf{u} \in \mathbb{R}^{d_6}} V(\mathbf{u}),$$

where, for  $\mathbf{u} = (u_1, \dots, u_{d_6})^\top$ ,

$$V(\mathbf{u}) = \mathbf{u}^\top \mathbf{W}_2 - \mathbf{u}^\top \mathbf{J}_{\vartheta\theta}^{\mathbf{A}} \mathbf{J}_*^{-1} \mathbf{W}_1 + \frac{1}{2} \mathbf{u}^\top \mathbf{J}_{\vartheta}^{\mathbf{A}} \mathbf{u} + \lambda_0^a p(\boldsymbol{\delta}, \vartheta_0, \mathbf{A}^\top \mathbf{u}),$$

with  $(\mathbf{W}_1^\top, \mathbf{W}_2^\top)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I}^{\mathbf{A}})$  and, for  $\mathbf{u}_2 = (u_1, \dots, u_{d_2})^\top$ ,

$$p(\boldsymbol{\delta}, \vartheta_0, \mathbf{u}_2) = \sum_{i \in \mathcal{A}} \delta_i u_i \text{sign}(\vartheta_{0i}).$$

## Particular case with oracle property

Choosing a penalty term such that  $\lambda_0^a = 0$ , we obtain

$$\sqrt{n}\mathbf{A} \left( \hat{\boldsymbol{\vartheta}}_n^a - \boldsymbol{\vartheta}_0 \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\vartheta}}^{\mathbf{A}} \right),$$

where  $\boldsymbol{\Sigma}_{\boldsymbol{\vartheta}}^{\mathbf{A}}$  is the asymptotic distribution of the multistep QMLE that we would obtained **if we would know** which are the constant betas and **all the non active components of the model**.

After a first-step Lasso variable selection, the post-Lasso method which simply consists of a second-step estimation of the components selected from the first-step Lasso has the same asymptotic distribution.

This kind of property is interpreted as an oracle property by Fan and Li (2001) and Zou (2006), but the interpretation is strongly criticized by Leeb and Pötscher (2008) and Hansen (2016).

## Penalizing the constant does not work

The previously defined estimators do not allow to detect the irrelevant betas, *i.e.* the indices  $i$  such that  $\boldsymbol{\vartheta}_0^{(i)} = \mathbf{0}_{q+3}$ .

The naive solution would be to penalize all the  $\boldsymbol{\vartheta}^{(i)}$ 's coefficients, does not work because a constant beta  $\beta_{i,t+1} \equiv \beta$  can be written as

$$\beta_{i,t+1} = \varpi_i + \xi_i \frac{v_t x_{it}}{\mu_i^2 + g_{it}^2} + c_i \beta_{it} + \gamma_{1i} z_{1t} + \dots + \gamma_{qi} z_{qt},$$

with many possibilities for  $\boldsymbol{\vartheta}^{(i)}$ , in particular

$$\boldsymbol{\vartheta}^{(i)} = \boldsymbol{\vartheta}_1 := (\beta, \mathbf{0}'_{q+2})' \quad \text{or} \quad \boldsymbol{\vartheta}^{(i)} = \boldsymbol{\vartheta}_2 := (0, 0, 1, \mathbf{0}'_q)'$$

Note that the solution  $\boldsymbol{\vartheta}^{(i)} = \boldsymbol{\vartheta}_2$  would be favored by the penalized estimator if  $|\beta| > 1$ , because in this case  $\|\boldsymbol{\vartheta}_1\|_1 > \|\boldsymbol{\vartheta}_2\|_1$ .

Therefore, penalizing  $\varpi_i$  along with the other parameters  $\xi_i, c_i, \gamma_{1i}, \dots, \gamma_{qi}$  is likely to result in an **inconsistent estimator**.

## First detecting the zero $c_i$ 's and then the irrelevant betas

Having, in a first step, identified (some of) the  $\xi_{0i}$  and  $c_{0i}$  that are zero, we obtain the second-step model

$$\left\{ \begin{array}{l} \beta_{i,t+1} = \varpi_{0i} + \gamma_{01i}z_{1t} + \cdots + \gamma_{0qi}z_{qt}, \quad i = 1, \dots, p^1 \\ \beta_{i,t+1} = \varpi_{0i} + \xi_{0i} \frac{v_t x_{it}}{\mu_{0i}^2 + g_{it}^2} + \gamma_{01i}z_{1t} + \cdots + \gamma_{0qi}z_{qt}, \\ \quad i = p^1 + 1, \dots, p^1 + p^2 \\ \beta_{i,t+1} = \varpi_{0i} + \xi_{0i} \frac{v_t x_{it}}{\mu_{0i}^2 + g_{it}^2} + c_{0i}\beta_{it} + \gamma_{01i}z_{1t} + \cdots + \gamma_{0qi}z_{qt}, \\ \quad i = p^1 + p^2 + 1, \dots, p, \end{array} \right.$$

for  $0 \leq p^1 \leq p^1 + p^2 \leq p$ , with obvious convention. It is possible to shrink all the beta coefficients for  $i = 1, \dots, p^1 + p^2$  (but not the  $\varpi_{0i}$ 's for  $i = p^1 + p^2 + 1, \dots, p$ .)



## General optimization problem

Consider the optimization problem

$$\boldsymbol{\vartheta}(\lambda) = \arg \min_{\boldsymbol{\vartheta} \in \Theta} Q_\lambda(\boldsymbol{\vartheta}), \quad Q_\lambda(\boldsymbol{\vartheta}) = Q(\boldsymbol{\vartheta}) + \lambda p(\boldsymbol{\vartheta}),$$

with  $p(\boldsymbol{\vartheta}) = \sum_{i \in S} \delta_i |\vartheta_i|$ , where  $\lambda \geq 0$ ,  $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_d)^\top$ ,  $\Theta$  is a convex compact subset of  $\mathbb{R}^d$ ,  $\delta_1, \dots, \delta_d$  are given relative shrinkage coefficients with  $\delta_i \geq 0$  and  $S = \{i : \delta_i > 0\} \neq \emptyset$ .

Assume that  $Q(\cdot)$  is two times continuously differentiable but **not necessarily convex**. The Newton-Raphson method may not work since  $Q_\lambda(\cdot)$  is **not differentiable everywhere**.

## Nonlinear shooting (NLShoot) algorithm

► Illustration ► LQA

For the Lasso-LSE of linear regressions Fu (1998) proposes the "shooting algorithm", which is a coordinate-wise descent algorithm.

We first define set of negative  $T_i^-(\boldsymbol{\vartheta})$  and positive  $T_i^+(\boldsymbol{\vartheta})$  points, real functions  $Q_\lambda^{(i)}(\cdot; \boldsymbol{\vartheta})$ , and

$$\tilde{\vartheta}_i = \arg \min_{\vartheta \in T_i^-(\boldsymbol{\vartheta}) \cup T_i^+(\boldsymbol{\vartheta}) \cup \{0\}} Q_\lambda^{(i)}(\vartheta; \boldsymbol{\vartheta}).$$

We then propose the following generalized shooting algorithm: start with an initial value  $\boldsymbol{\vartheta} = \boldsymbol{\vartheta}^0$  and

replace the  $i$ -th coordinate of  $\boldsymbol{\vartheta}$  by  $\tilde{\vartheta}_i$  for  $i = 1, 2, \dots, d, 1, 2, \dots$

We show that the cluster point(s) of this NLShoot algorithm are stationary points of  $Q_\lambda(\cdot)$ .

## Monte-Carlo design

The explanatory  $p_{tv} + p_{cst}$  variables follow a DCC-GARCH model. For each replication, the penalized model is evaluated on a grid of 15 equidistant values of  $\lambda$  between 0 and  $\bar{\lambda}$ , using the nonlinear shooting (NLSshoot) algorithm initialized with the LQA algorithm.

The next table reports the number of time-varying betas, of constant betas, the percentage of correctly selected models, of correctly identified time-varying betas, of constant conditional betas, of correctly identified non active parameters and active parameters.

In Step 1, only the  $\xi$  and  $c$  parameters are penalized. In Step 2, the intercepts  $\varpi$ 's of the conditional betas identified as constant or for which the  $c$  parameter is 0 in Step 1 are also penalized.

# 100 replications for sample of 4,000 observations

More than 1 billion of possible models when  $p = 10$

$p_{tv}$	$p_{cst}$	% correct $p_{tv}$	% correct $p_{cst}$	% correct 0's	% correct !0's
Step 1					
3	3	93.3	66.0	82.3	94.3
6	2	94.0	45.5	71.7	95.8
2	6	88.4	78.3	88.0	90.6
5	5	92.4	66.4	81.7	93.9
7	3	94.8	52.3	75.0	96.1
3	7	88.0	76.3	86.6	90.0
Step 2					
3	3	92.6	100.0	97.4	95.9
6	2	93.0	99.5	96.2	94.0
2	6	77.5	99.3	97.5	91.4
5	5	88.8	99.0	95.6	94.0
7	3	92.0	99.3	95.6	95.3
3	7	87.0	97.9	95.6	94.9

## The data set

We regress 30 daily Dow Jones stock returns from January 5th, 2010 to June 29th, 2023 (3,394 values) on nine sector ETFs (potentially more than 1 billion of models for each stock).

For example, the BP series, which belongs to the Oil & Gas sector, has a penalized ACB model that includes only 3 ETFs:

- XLB (Material) and XLF (Financial services) with cst betas;
- XLE (Energy) with  $(\varpi, \xi, c) = (0.005, 0.009, 0.995)$ ;
- and an error term with vol  $(\omega, \alpha, \beta) = (0.079, 0.089, 0.846)$ .

In general, the ETF to which the stock is linked always appears as relevant regressor, and the error term of the regression is always conditionally heteroscedastic.

## Conclusion

- The Lasso can be used to solve the non-identifiability of our beta time series model, at the price of some bias.
- The adaptive version is able to suppress the bias.
- The Lasso Multi-Step QMLE is tractable even though the optimization is nonlinear, nonconvex, and involves a large number of parameters.
- The penalized estimator allows to select the constant betas in a first step, and the relevant regressors in a second term.

Thank you!

## The autoregressive conditional beta (ACB) model

We focus on simple GARCH models for simplicity, but there are no major conceptual or practical issues in considering more general time series models for the regressors  $x_{it}$  and for the error term  $v_t$  of the regression.

The term  $\frac{v_t x_{it}}{\mu_{0i}^2 + g_{it}^2}$  can be interpreted as an update variable for the dynamics of the betas.

This term has been obtained by using the **score-driven (SD) approach** ▶ SD approach with the aim to obtain a simple beta model, analogous to an ARMA or a GARCH for the first conditional moments.

## Interpretation of the update variable $\frac{v_t x_{it}}{\mu_{0i}^2 + g_{it}^2}$

The presence of the product  $v_t x_{it}$  allows dynamic monitoring of the orthogonality condition between the error term of the regression  $v_t$  and the regressors  $x_{it}$ .

Indeed, if  $v_t x_{it} \simeq 0$  and there is no exogenous variable ( $q = 0$ ), then we have the constant solution  $\beta_{it} \equiv \beta_i$ . Now, if  $\beta_{it}$  is negative and very small, then  $v_t x_{it} = y_t x_{it} - \sum_{j \neq i} \beta_{jt} x_{it} x_{jt} - \beta_{it} x_{it}^2$  tends to be positive. Therefore if  $\beta_{it}$  is small and  $\xi_i > 0$ , then  $\beta_{i,t+1}$  tends to be larger than  $\beta_{it}$ .

Also note that since the volatility  $g_{it}^2$  is in the denominator of the update variable, the time variation of the beta is smaller in periods of high volatility in the regressors. ◀



## Score Driven (GAS and Beta-t-GARCH) models

Assume that  $y_t$  follows a conditional density  $p(y_t|f_t, \mathcal{F}_{t-1}, \theta)$ , where  $f_t$  is a time-varying parameter of interest.

Harvey and Chakravarty (2008) and Creal, Koopman and Lucas (2012) proposed the score-driven (SD) model for  $f_t$ :

$$f_{t+1} = \varpi + \xi S(f_t) \frac{\partial \log p(y_t|f_t, \mathcal{F}_{t-1}, \theta)}{\partial f_t} + c f_t,$$

where  $\varpi$ ,  $\xi$  and  $c$  are unknown parameters and  $S(f_t)$  is the inverse of the conditional information matrix.

→ The **scaled score** is the updating mechanism in this approach.

## Examples of SD models

The SD model for  $m_t$  in the location model

$$y_t = m_t + \eta_t, \quad \eta_t \text{ iid } \mathcal{N}(0, \sigma^2),$$

is the **ARMA** model

$$m_{t+1} = \varpi + \xi(y_t - m_t) + cm_t.$$

The SD model for  $g_t^2$  in the scale model

$$y_t = g_t \eta_t, \quad \eta_t \text{ iid } \mathcal{N}(0, 1),$$

is the **GARCH** model

$$g_{t+1}^2 = \varpi + \xi(y_t^2 - g_t^2) + cg_t^2.$$

→ The SD approach provides benchmark models for Gaussian conditional distributions with location-scale parameters.

## Applying the SD approach for the beta parameters

Let us define a SD model for  $\beta_{it}$  in the regression model

$$y_t = \beta_{1t}x_{1t} + \cdots + \beta_{pt}x_{pt} + v_t,$$

where  $v_t = g_t\eta_t$ ,  $g_t^2 = E_{t-1}v_t^2$ ,  $\eta_t$  iid  $\mathcal{N}(0, 1)$ . Let  $E_{t-1}x_{it} = \mu_{it}$  and  $E_{t-1}x_{it}^2 = g_{it}^2 + \mu_{it}^2$ .

We have  $l_t := \log p(y_t | \mathcal{F}_{t-1}, \theta) = \frac{-1}{2} \left\{ \frac{v_t^2}{g_t^2} + \log g_t^2 \right\}$ ,

$$\frac{\partial l_t}{\partial \beta_{it}} = \frac{v_t x_{it}}{g_t^2}, \quad S(\beta_{it}) = - \left( E_{t-1} \frac{\partial^2 l_t}{\partial^2 \beta_{it}} \right)^{-1} = \frac{g_t^2}{\mu_{it}^2 + g_{it}^2}.$$

Therefore the updating mechanism  $S(\beta_{it}) \frac{\partial l_t}{\partial \beta_{it}} = \frac{v_t x_{it}}{\mu_{it}^2 + g_{it}^2}$ .

## ACB with exogenous variables and special cases

For simplicity, assume  $\mu_{it} = \mu_i$  and a GARCH(1,1) volatility  $g_{it}$ .  
Let the ACB with  $q$  additional exogenous variables  $z_{1t}, \dots, z_{qt}$

$$\beta_{it+1} = \varpi_i + \xi_i \frac{v_t x_{it}}{\mu_i^2 + g_{it}^2} + c_i \beta_{it} + \gamma_{1i} z_{1t} + \dots + \gamma_{qi} z_{qt}.$$

- ① No GARCH effects in  $x_{it}$ :  $\alpha_i = \beta_i = 0 \rightarrow g_{it+1}^2 = \omega_i$ .

$$\beta_{it+1} = \varpi_i + \xi_i v_t x_{it} + c_i \beta_{it} + \gamma_{1i} z_{1t} + \dots + \gamma_{qi} z_{qt}.$$

- ②  $x_{it} = 1$  (intercept):  $\beta_{it}$  time varying unless  $\xi_i = 0$  and  $q = 0$

$$\beta_{it+1} = \varpi_i + \xi_i v_t + c_i \beta_{it} + \gamma_{1i} z_{1t} + \dots + \gamma_{qi} z_{qt}.$$

- ③ Assume that  $p = 1$ ,  $q = 0$  and that  $x_{1t} = 1$   
 $\rightarrow y_t = \varpi + (\xi - c)v_{t-1} + cy_{t-1} + v_t = \text{ARMA}(1,1)$ .

## Illustration of $\mathbf{A}(\lambda_0)$

Let

$$Q_\lambda(\boldsymbol{\vartheta}) = Q(\boldsymbol{\vartheta}) + \lambda \|\boldsymbol{\vartheta}\|_1$$

with  $\boldsymbol{\vartheta} = (\vartheta_1, \vartheta_2)^\top$ ,

$$Q(\boldsymbol{\vartheta}) = \mathcal{P}(\vartheta_1 - \vartheta_2), \quad \mathcal{P}(x) = \frac{x^4}{4} - \frac{2x^3}{3} - x^2 + 1.$$

Note that  $Q(\boldsymbol{\vartheta})$  is minimal for all  $\boldsymbol{\vartheta}$  such that  $\vartheta_2 = \vartheta_1 - 2$ .

- If  $\lambda > 0$  is very small but not zero,  $\boldsymbol{\vartheta}^* = \arg \min Q_\lambda(\boldsymbol{\vartheta})$  is unique and is close to  $(1, -1)$ ;
- if  $\lambda$  is very large  $\boldsymbol{\vartheta}^*$  is also unique and equal to  $(0, 0)$ ;
- if  $\lambda = 2.03$  the function reaches its minimum at 2 points
- $\lambda = 0$  the minimum is not well defined (*i.e.* not unique).

$\mathbf{A}(\lambda_0)$  always satisfied, except for  $\lambda_0 = 0$  and  $\lambda_0 = 2.03$

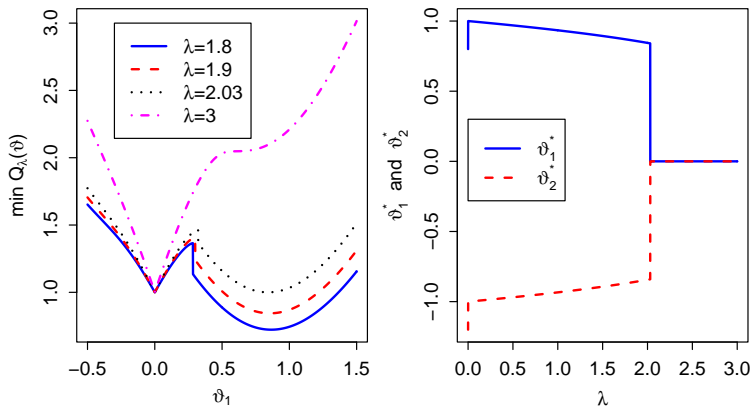


Figure: Left graph:  $\min Q_\lambda(\vartheta)$  as function of  $\vartheta_1$  for several values of  $\lambda$ .  
 Right graph:  $(\vartheta_1^*, \vartheta_2^*) = \arg \min Q_\lambda(\vartheta)$  as function of  $\lambda$ .

Example showing that we can have  $\lambda^i \neq \lambda^s$ 

Let

$$Q_\lambda(\boldsymbol{\vartheta}) = Q(\boldsymbol{\vartheta}) + \lambda|\vartheta_1|, \quad \boldsymbol{\vartheta} = (\vartheta_1, \vartheta_2)^\top,$$

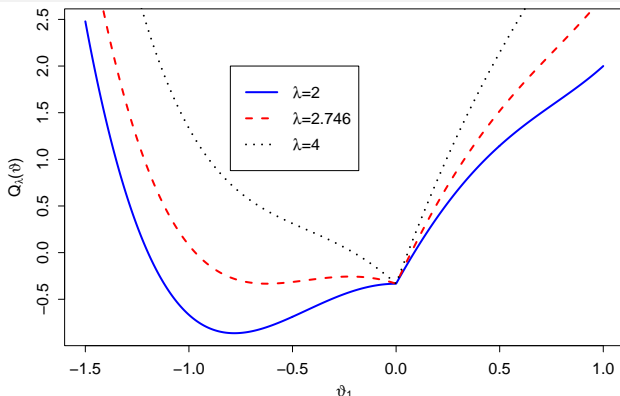
with  $Q(\boldsymbol{\vartheta}) = \sum_{i=1}^2 \vartheta_i^4 - \frac{2}{3}\vartheta_i^3 - 2\vartheta_i^2 + 2\vartheta_i + 1$ . Note that  $\boldsymbol{\vartheta}^* := \arg \min Q_\lambda(\boldsymbol{\vartheta}) = (\vartheta_1^*, -1)^\top$ , and that for  $\boldsymbol{\vartheta}^c := (0, -1)^\top$  we have

$$\lambda^i = \partial Q(\boldsymbol{\vartheta}^c) / \partial \vartheta_1 = 2.$$

However,  $\arg \min Q_2(\boldsymbol{\vartheta}) = (-0.781, -1)^\top \neq \boldsymbol{\vartheta}^c$ . We have  $\boldsymbol{\vartheta}^* = \boldsymbol{\vartheta}^c$  for  $\lambda \geq 2.746$ , which is in agreement with

$$\lambda^s = \sup_{\vartheta_1 \leq 0} 4\vartheta_1^3 - 2\vartheta_1^2 - 4\vartheta_1 + 2 = 3.032.$$

## Example where $\lambda^i < \lambda^s$



**Figure:** Function  $\vartheta_1 \mapsto Q_\lambda(\vartheta_1, -1) = Q(\vartheta_1, -1) + \lambda|\vartheta_1|$ . For  $\vartheta^c = (0, -1)$ , we have  $\partial Q(\vartheta^c)/\partial \vartheta_1 = 2$ , but  $\arg \min_{\vartheta} Q_\lambda(\vartheta) = \vartheta^c$  only for  $\lambda \geq 2.746$  and not for  $\lambda = \lambda^i = 2$ .



## Example of computation of $\lambda^i$ and $\lambda^s$

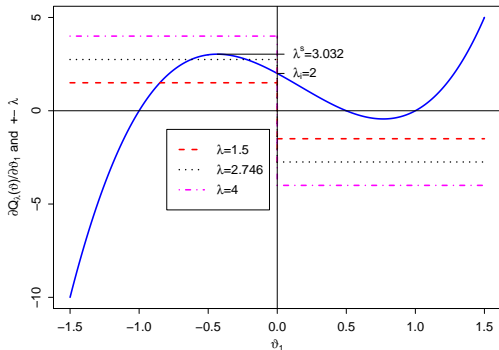


Figure: Graphical representation of  $\lambda^i$  and  $\lambda^s$ .

## NLShoot for a simple optimization

Let us use the NLshoot algorithm to compute

$$\arg \min_{\vartheta} P(\vartheta) + \lambda|\vartheta|, \quad P(x) = x^4 - \frac{2}{3}x^3 - 2x^2 + 2x.$$

NLShoot reduces to a single application of the algorithm. The figure below shows that, for  $\lambda = 2.5$ , the set  $\{\vartheta < 0 : P'(\vartheta) = \lambda\}$  contains 2 points, that 0 is also a point at which a generalized gradient is zero, and that the set  $\{\vartheta > 0 : P'(\vartheta) = -\lambda\}$  is empty. To find the minimum, it is then sufficient to compare the value of the penalized function at these 3 points.

## The objective function

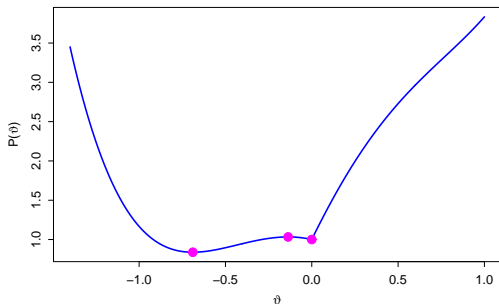


Figure: Function  $P(\vartheta) + 2.5|\vartheta|$  with a minimum at  $\vartheta = -0.689$ , and two other "critical points" at  $\vartheta = -0.137$  and  $\vartheta = 0$ .

## NLShoot in a simple (degenerate) case

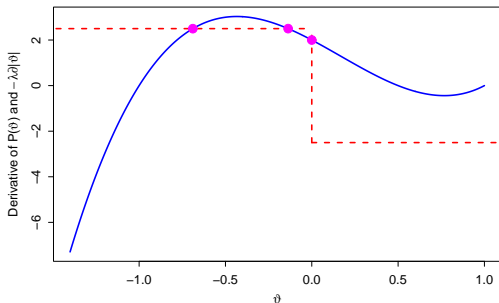


Figure: NLShoot finds the 3 critical points.



## FOC satisfied by $\widehat{\vartheta}_n$ and sparsity of the estimator

We have

$$0 \in \left\{ \frac{\partial \widetilde{O}_n(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\vartheta}}_n)}{\partial \boldsymbol{\vartheta}} \right\} + \lambda_n \partial p(\widehat{\boldsymbol{\vartheta}}_n).$$

Therefore, for all  $i \in S$ ,

- $\left| \frac{1}{n} \sum_{t=2}^n \frac{\partial \widetilde{\ell}_t(\widehat{\boldsymbol{\varphi}}_n)}{\partial \vartheta_i} \right| \leq \lambda_n$  if  $\widehat{\vartheta}_{ni} = 0$ ,
- $\frac{1}{n} \sum_{t=2}^n \frac{\partial \widetilde{\ell}_t(\widehat{\boldsymbol{\varphi}}_n)}{\partial \vartheta_i} = -\lambda_n \text{sign}(\widehat{\vartheta}_{ni})$  if  $\widehat{\vartheta}_{ni} \neq 0$ .



## The LQA algorithm

For penalty terms that are non convex and non differentiable, Fan and Li (2001) proposed the local quadratic approximation (LQA) algorithm. In our framework the LQA algorithm consists of repeatedly solving

$$\boldsymbol{\vartheta}^{(k+1)} = \arg \min_{\boldsymbol{\vartheta} \in \Theta} \left\{ Q(\boldsymbol{\vartheta}) + \lambda \sum_{i \in S} \delta_i \frac{\text{sign}(\vartheta_i^{(k)})}{2\vartheta_i^{(k)}} \vartheta_i^2 \right\}$$

until convergence, where  $\boldsymbol{\vartheta}^{(0)} \in \Theta$  is an initial value.

Fan and Li (2001) suggested to set  $|\vartheta_i^{(\cdot)}| \equiv 0$  and to remove this component from the optimization problem if  $|\vartheta_i^{(k)}| < \varepsilon$  for some small  $\varepsilon > 0$ , but the **choice of  $\varepsilon$**  matters. ◀