

# Intervalles de confiance valides en présence de sélection de modèle

François Bachoc\*, Hannes Leeb et Benedikt M. Pötscher

University of Vienna

(Full) linear model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{U}$$

- $\mathbf{X}$  of size  $n \times p$
- $p < n$
- $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$
- $\boldsymbol{\beta}$  of size  $p \times 1$
- $\mathbf{Y}$  observation vector

Least square estimator :

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

Standard variance estimator :

$$\hat{\sigma}^2 = \frac{1}{n-p} \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2$$

Working distribution  $P_{n,\boldsymbol{\beta},\sigma}$

## Linear submodels

Subsets  $M \subset \{1, \dots, p\}$  of the columns of  $\mathbf{X}$ . Give

$$\mathbf{Y} = \mathbf{X}[M]\mathbf{v} + \mathbf{U}$$

- $M$  of cardinality  $m$
- $\mathbf{X}[M]$  of size  $n \times m$  : only the columns of  $\mathbf{X}$  that are in  $M$
- $\mathbf{v}$  of size  $m \times 1$  : needs to be defined/estimated to give the **best representation** of the full linear model

Non-standard regression coefficient vector

$$\hat{\beta}_M^{(n)} = \underset{\mathbf{v}}{\operatorname{argmin}} \|\mathbf{X}\beta - \mathbf{X}[M]\mathbf{v}\|$$

$$\hat{\beta}_M^{(n)} = \beta[M] + (\mathbf{X}'[M]\mathbf{X}[M])^{-1} \mathbf{X}'[M]\mathbf{X}[M^c]\beta[M^c],$$

- $\beta[M]$  of size  $m \times 1$  : components of  $\beta$  in  $M$

Restricted least square estimator

$$\hat{\beta}_M = (\mathbf{X}'[M]\mathbf{X}[M])^{-1} \mathbf{X}'[M]\mathbf{Y}$$

## Model selection procedure

Data-driven selection of the model with  $\hat{M}(Y) = \hat{M}$

Ex. : BIC :

$$\hat{M}_{BIC}(Y) \in \underset{M}{\operatorname{argmin}} \|Y - \mathbf{X}[M]\hat{\beta}_M\|^2 + \log(n)|M|$$

- $|M|$  : cardinality of  $M$

Berk et al., 2013, *Annals of Statistics* consider the **non-standard target**

$$\beta_{\hat{M}}^{(n)}$$

as their target for confidence intervals

Comments :

- Model selector  $\hat{M}$  is "imposed"
- Objective : best coefficients in this imposed model
- Random target

Let  $\mathbf{x}_0$  be a fixed  $p \times 1$  vector and consider

$$y_0 = \mathbf{x}_0' \boldsymbol{\beta} + u_0$$

- $u_0 \sim \mathcal{N}(0, \sigma^2)$

We consider the **design-dependent non-standard target**

$$\mathbf{x}_0' [\hat{M}] \boldsymbol{\beta}_{\hat{M}}^{(n)}$$

Optimality property : when  $\mathbf{x}_0$  is random and follows the empirical distribution given by the lines of  $\mathbf{X}$  :

$$\mathbb{E}_{n, \beta, \sigma} \left( \left[ y_0 - \mathbf{x}_0' [\hat{M}] \boldsymbol{\beta}_{\hat{M}}^{(n)} \right]^2 \right) \leq \mathbb{E}_{n, \beta, \sigma} \left( \left[ y_0 - \mathbf{x}_0' [\hat{M}] v(Y) \right]^2 \right),$$

for any function  $v(Y) \in \mathbb{R}^{|\hat{M}|}$ .

Let a nominal level  $1 - \alpha \in (0, 1)$  be fixed

We consider confidence intervals for  $\mathbf{x}'_0[\hat{M}]\beta_M^{(n)}$  of the form

$$CI = \mathbf{x}'_0[\hat{M}]\hat{\beta}_M \pm K \|\mathbf{s}_M\| \hat{\sigma},$$

with

$$\mathbf{s}'_M = \mathbf{x}'_0[M] (\mathbf{X}'[M]\mathbf{X}[M])^{-1} \mathbf{X}'[M]$$

## Interpretation

- "Constant"  $K$  does not depend on  $Y$  (but on  $\mathbf{X}$ ,  $\mathbf{x}_0$ ,  $\hat{M}$ )
- For **fixed**  $M$ ,

$$\mathbf{x}'_0[M]\hat{\beta}_M - \mathbf{x}'_0[M]\beta_M^{(n)} \sim \mathcal{N}(0, \|\mathbf{s}_M\|\sigma^2)$$

- Thus,  $K_{naive} = q_{S, n-p, 1-\alpha/2}$  (Student quantile) is valid when  $M$  is deterministic
- When  $\hat{M}$  is random,  $K$  needs to be larger (e.g. [Leeb et al. 2015, Statistical Science](#))

⇒ **Main issue** : choosing  $K$  ?

Observe that

$$\mathbf{x}'_0[\hat{M}]\hat{\beta}_{\hat{M}} - \mathbf{x}'_0[\hat{M}]\beta_M^{(n)} = \mathbf{s}'_{\hat{M}}(\mathbf{Y} - \mathbf{X}\beta)$$

Then, we have

$$\left| \frac{\mathbf{s}'_{\hat{M}}}{\|\mathbf{s}_{\hat{M}}\|\hat{\sigma}}(\mathbf{Y} - \mathbf{X}\beta) \right| \leq \max_{M \subseteq \{1, \dots, p\}} \left| \frac{\mathbf{s}'_M}{\|\mathbf{s}_M\|\hat{\sigma}}(\mathbf{Y} - \mathbf{X}\beta) \right|$$

Distribution of the upper-bound **independent** of  $\beta, \sigma \implies$  let  $K_1$  be its  $(1 - \alpha)$  quantile

The CI given by  $K_1$  satisfies

$$\inf_{\beta \in \mathbb{R}^p, \sigma > 0} P_{n, \beta, \sigma} \left( \mathbf{x}'_0[\hat{M}]\beta_M^{(n)} \in CI \right) \geq 1 - \alpha$$

$\implies$  **Uniformly valid** confidence interval

The constant  $K_1$  depends on all the components of  $\mathbf{x}_0$

It can happen that only  $\mathbf{x}_0[\hat{M}]$  is observed

- model selection for cost reason

We construct other constants (see the paper for details)

$$K_1 \leq K_2 \leq K_3 \leq K_4$$

(The CIs given by  $K_2, K_3, K_4$  are hence universally valid)

**Remark :** The case where only  $\mathbf{x}_0[\hat{M}]$  is observed motivates all the more the study of  $\mathbf{x}'_0[\hat{M}]\beta_{\hat{M}}^{(n)}$  as opposed to  $\mathbf{x}'_0\beta$



**Issue :** The target  $\mathbf{x}'_0[\hat{M}]\beta_{\hat{M}}^{(n)}$  depends on  $\mathbf{X}$

Issue is solved when lines of  $X$  and  $\mathbf{x}'_0$  are realizations from the same distribution  $\mathcal{L}$

Let, for  $\mathbf{x}' \sim \mathcal{L}$ ,  $\Sigma = \mathbb{E}(\mathbf{x}\mathbf{x}')$ . Then, define the **design-independent non-standard target** by

$$\mathbf{x}_0[\hat{M}]'\beta_{\hat{M}}^{(*)} = \mathbf{x}_0[\hat{M}]'\beta[\hat{M}] + \mathbf{x}_0[\hat{M}]' \left( \Sigma[\hat{M}, \hat{M}] \right)^{-1} \Sigma[\hat{M}, \hat{M}^c]\beta[\hat{M}^c],$$

Then, we have for  $\mathbf{x}_0 \sim \mathcal{L}$ ,

$$\mathbb{E} \left( \left[ y_0 - \mathbf{x}'_0[\hat{M}]\beta_{\hat{M}}^{(*)} \right]^2 \right) \leq \mathbb{E} \left( \left[ y_0 - \mathbf{x}'_0[\hat{M}]\mathbf{v}(\mathbf{Y}) \right]^2 \right),$$

for any function  $\mathbf{v}(\mathbf{Y}) \in \mathbb{R}^{|\hat{M}|}$

Observe that

$$\begin{aligned} & \left( \mathbf{x}_0[\hat{M}]' \beta_{\hat{M}}^{(*)} - \mathbf{x}_0[\hat{M}]' \beta_{\hat{M}}^{(n)} \right) = \\ & \mathbf{x}'_0[\hat{M}] \left( \left( \mathbf{X}'[\hat{M}]\mathbf{X}[\hat{M}] \right)^{-1} \mathbf{X}'[\hat{M}]\mathbf{X}[\hat{M}^c] - \left( \Sigma[\hat{M}, \hat{M}] \right)^{-1} \Sigma[\hat{M}, \hat{M}^c] \right) \beta[\hat{M}^c] \end{aligned}$$

### Theorem

Assume that

$$\sqrt{n} [(\mathbf{X}'\mathbf{X}/n) - \Sigma] = O_p(1)$$

and that for any  $M$  with  $|M| < p$  and for any  $\delta > 0$ ,

$$\sup \left\{ P_{n,\beta,\sigma}(\hat{M} = M | X) : \beta \in \mathbb{R}^p, \sigma > 0, \|\beta[M^c]\| / \sigma \geq \delta \right\} = o_p(1)$$

Then, for  $CI$  obtained by  $K_1, K_2, K_3, K_4$ ,

$$\inf_{\beta \in \mathbb{R}^p, \sigma > 0} P_{n,\beta,\sigma} \left( \mathbf{x}'_0[\hat{M}] \beta_{\hat{M}}^{(*)} \in CI \mid X \right) \geq (1 - \alpha) + o_p(1)$$

For  $\alpha = 0.05$  and  $p = 10$  we evaluate

$$\inf_{\beta \in \mathbb{R}^p, \sigma > 0} P_{n, \beta, \sigma} \left( \mathbf{x}_0' [\hat{M}] \beta_{\hat{M}}^{(n, *)} \in CI \mid X \right),$$

for one realization of  $\mathbf{X}$

Results :

$n$	model selector	target							
		design-dependent				design-independent			
		$\mathbf{x}_0' [\hat{M}] \beta_{\hat{M}}^{(n)}$				$\mathbf{x}_0' [\hat{M}] \beta_{\hat{M}}^{(*)}$			
		$K_{naive}$	$K_1$	$K_3$	$K_4$	$K_{naive}$	$K_1$	$K_3$	$K_4$
20	AIC	0.84	0.99	1.00	1.00	0.79	0.97	0.99	0.99
20	BIC	0.84	0.99	1.00	1.00	0.74	0.96	0.98	0.98
20	LASSO	0.90	1.00	1.00	1.00	0.18	0.48	0.61	0.61
100	AIC	0.87	0.99	1.00	1.00	0.88	0.99	1.00	1.00
100	BIC	0.88	0.99	1.00	1.00	0.87	0.99	1.00	1.00
100	LASSO	0.88	0.99	1.00	1.00	0.87	0.99	1.00	1.00

## Conclusion :

- It is known that in the classical case (estimation of  $\beta$ ), it is difficult to construct valid post-model-selection confidence intervals
- Recently, alternative targets have been studied
- This removes some obstacles
- But naive procedures still fail

## Prospects :

- Asymptotics where  $d$  is large
- Generalized linear models

## The paper :

- ✎ **F. Bachoc, H. Leeb, B.M. Pötscher (2014+). Valid confidence intervals for post-model-selection predictors, <http://arxiv.org/abs/1412.4605>, submitted**

Thank you for your attention !