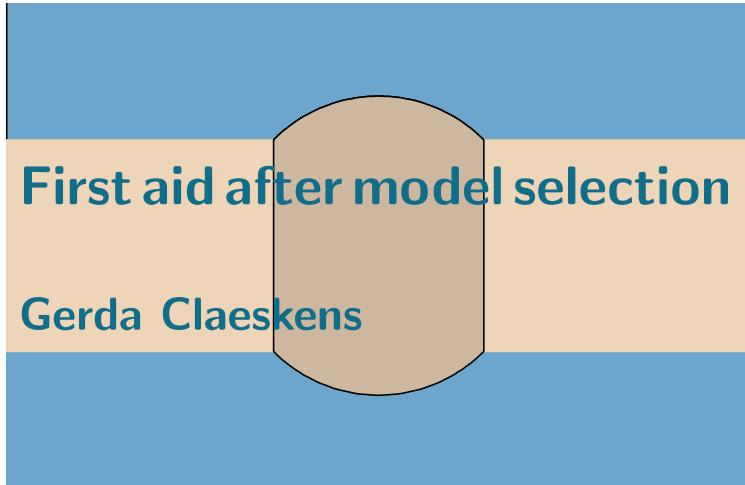


## Post-selection - team members

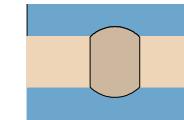


ORSTAT and Leuven Statistics Research Center, KU Leuven, Belgium  
26th Annual Meeting of the RSSB, 2018

Ali Charkhi (was at KU Leuven, now at ING)

Andrea C. Garcia-Angulo (PhD student, see also her poster!)

Jing Zhou (PhD student, see also her poster!)



## Model selection by AIC

- Model selection by Akaike's information criterion

$$\text{AIC}(M_k) = -2 \log \text{Likelihood}_{M_k}(\hat{\beta}_{M_k}) + 2p_{M_k}$$

Select model  $\hat{M}_{\text{aic}}$  for which  $\text{AIC}(\hat{M}_{\text{aic}})$  is the smallest amongst  $\text{AIC}(M_k), k = 1, \dots, K$

- Estimation in the selected model:

$$\mu(\hat{\beta}_{\hat{M}_{\text{aic}}}) = \sum_{k=1}^K \mu(\hat{\beta}_{M_k}) \cdot I(M_k = \hat{M}_{\text{aic}})$$

Selection is a special case of model averaging:

- Weights are only 0 or 1
- Weights are random, since data-driven.
- Randomness of weights influences further inference

## Confidence regions after AIC selection

In selected model  $\hat{M}_{\text{aic}}$ , estimate focus  $\mu$  by

$$\hat{\mu}(\hat{M}_{\text{aic}}) = \sum_{k=0}^K \hat{\mu}(M_k) \cdot I(M_k = \hat{M}_{\text{aic}}).$$

Confidence regions for  $\mu$  after selection?

Naive interval: pretend that the model has been given beforehand. Ignore randomness  $\hat{M}_{\text{aic}} = M_{\text{aic}}$ .

$$\hat{\mu}(M_{\text{aic}}) \pm 1.96 \cdot \text{standard error}(\hat{\mu}(M_{\text{aic}}))$$

Frequentist model averaging necessary to overcome low coverage of using naive method after AIC-selection.

(Hjort and Claeskens, 2003)

## Some recent literature on correct inference

Berk, Brown, Buja, Zhang, Zhao (2013): valid confidence intervals irrespective of the selection procedure. No overall true values for the parameters.

Bachoc, Leeb, Pötscher (2015): generalized above method to prediction intervals.

Efron (2014): calculate the variance of the bagging estimator  $\Rightarrow$  not based on the selected model for the original data.

Lee, Sun, Sun, Taylor (2016): obtained exact post-selection inference for lasso model selection for a given value of  $\lambda$ .

Schneider (2016): coverage of intervals based on thresholding estimators in high-dimensional linear regression models.

## Asymptotics of AIC selection

AIC might overselect (Nested models: Woodroffe 1982, arc-sine laws) for  $n \rightarrow \infty$ ,  $\hat{p}_{\text{aic}} \in \{p_0, p_0 + 1, \dots, K\}$ ,

Minimize AIC  $\Leftrightarrow$  maximize AIC\*

$$\text{AIC}^*(M_k) = 2\{\ell_n(\hat{\theta}_{(k)}) - \ell_n(\vartheta)\} - 2k = 2\ell_{n,k}^* - 2k.$$

Joint convergence (likelihood ratio statistics)

$$2(\ell_{n,p_0}^*, \dots, \ell_{n,K}^*) \xrightarrow{d} \left( \underbrace{\sum_{i=1}^{a+p_0} Z_i^2}_{\chi_{a+p_0}^2}, \dots, \underbrace{\sum_{i=1}^{a+K} Z_i^2}_{\chi_{a+K}^2} \right)$$

with  $Z_1, \dots, Z_{a+K} \stackrel{i.i.d.}{\sim} N(0, 1)$

## Post-AIC-selection in nested models

Classical (easy) setting to start off with:

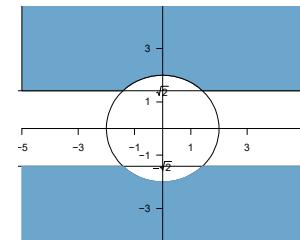
- ▶  $\mathcal{M}_{\text{nest}} = \{M_k; M_1 \subset M_2 \subset \dots \subset M_K\}$
- ▶ Maximum likelihood estimation in each model  $\hat{\theta}_{(k)} = (\hat{\theta}'_{(k)}, \mathbf{0}_{K-k}^t)^t$ .
- ▶ Minimal true model  $M_{p_0} \in \mathcal{M}_{\text{nest}}$ :  
Models with indices  $k < p_0$  are under-parametrized,  
Models with  $k > p_0$  are over-parametrized.
- ▶  $\text{AIC}(M_k) = -2\ell_n(\hat{\theta}_{(k)}) + 2(a+k)$

Select  $\hat{p}_{\text{aic}}$  if and only if  $\text{AIC}(M_{\hat{p}_{\text{aic}}}) \leq \text{AIC}(M_k) \quad \text{for all } k = 1, \dots, K$

## AIC selection region – Example

Three models:  $M_1 : \{\theta_1\}$ ;  $M_2 : \{\theta_1, \theta_2\}$ ;  $M_3 : \{\theta_1, \theta_2, \theta_3\}$

$a = 1$ ,  $K = 2$ , take  $M_1$  smallest true model,  $M_3$  selected,  $\hat{p}_{\text{aic}} = 2$



$$\begin{cases} \text{AIC}^*(M_3) > \text{AIC}^*(M_2) \\ \text{AIC}^*(M_3) > \text{AIC}^*(M_1) \end{cases}$$

Asymptotic representation:

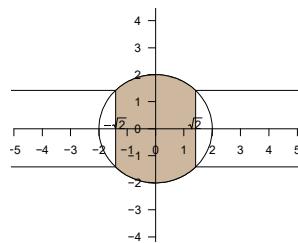
$$\begin{cases} Z_1^2 + Z_2^2 + Z_3^2 - 6 > Z_1^2 + Z_2^2 - 4 \\ Z_1^2 + Z_2^2 + Z_3^2 - 6 > Z_1^2 - 2 \end{cases}$$

$$\mathcal{A}_2(\mathcal{M}_{\text{nest}}) = \{z \in \mathbb{R}^3 : z_3^2 > 2, z_2^2 + z_3^2 > 4\}.$$

## AIC selection region – Example

Three models:  $M_1 : \{\theta_1\}$ ;  $M_2 : \{\theta_1, \theta_2\}$ ;  $M_3 : \{\theta_1, \theta_2, \theta_3\}$

$a = 1, K = 2$ , take  $M_1$  smallest true model,  $M_1$  selected,  $\hat{p}_{\text{aic}} = 0$



$$\begin{cases} \text{AIC}^*(M_1) > \text{AIC}^*(M_2) \\ \text{AIC}^*(M_1) > \text{AIC}^*(M_3) \end{cases}$$

Asymptotic representation:

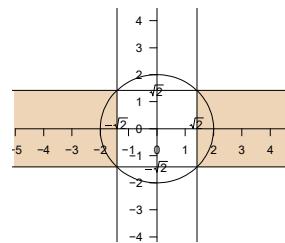
$$\begin{cases} Z_1^2 - 2 > Z_1^2 + Z_2^2 - 4 \\ Z_1^2 - 2 > Z_1^2 + Z_2^2 + Z_3^2 - 6 \end{cases}$$

$$\mathcal{A}_0(\mathcal{M}_{\text{nest}}) = \{z \in \mathbb{R}^3 : z_2^2 < 2, z_2^2 + z_3^2 < 4\}.$$

## AIC selection region – Example

Three models:  $M_1 : \{\theta_1\}$ ;  $M_2 : \{\theta_1, \theta_2\}$ ;  $M_3 : \{\theta_1, \theta_2, \theta_3\}$

$a = 1, K = 2$ , take  $M_1$  smallest true model,  $M_2$  selected,  $\hat{p}_{\text{aic}} = 1$



$$\begin{cases} \text{AIC}^*(M_2) > \text{AIC}^*(M_1) \\ \text{AIC}^*(M_2) > \text{AIC}^*(M_3) \end{cases}$$

Asymptotic representation:

$$\begin{cases} Z_1^2 + Z_2^2 - 4 > Z_1^2 + Z_2^2 + Z_3^2 - 6 \\ Z_1^2 + Z_2^2 - 4 > Z_1^2 - 2 \end{cases}$$

$$\mathcal{A}_1(\mathcal{M}_{\text{nest}}) = \{z \in \mathbb{R}^3 : z_3^2 < 2, z_2^2 > 2\}.$$

## Asymptotic distribution

$\mathbf{J}_p$ : Fisher information for model  $M_p$ .

$\tilde{\boldsymbol{\nu}}_{(p)} = (\nu_1, \dots, \nu_{a+p})^t$ : sub-vector of  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_{a+K})^t$  for model  $M_p$ .

$\boldsymbol{\vartheta}$  true value

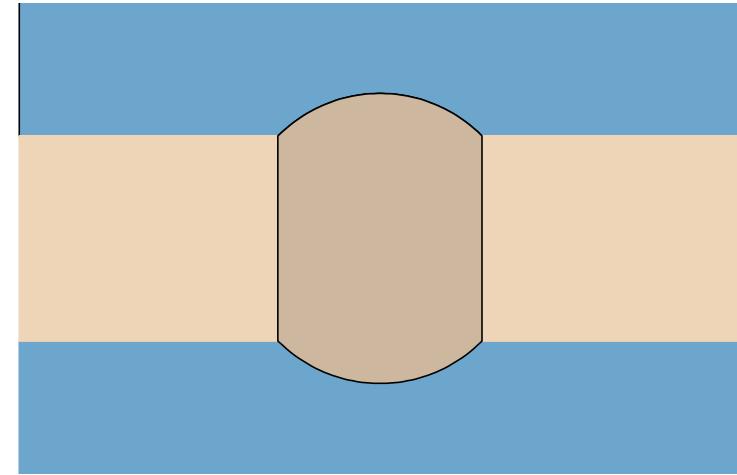
### Post-AIC distribution for nested models

For  $\mathcal{M}_{\text{nest}}$  with  $p_0$  denoting the true model order

$$\begin{aligned} F_p(t) &= \lim_{n \rightarrow \infty} P(n^{1/2}(\hat{\boldsymbol{\theta}}_{(p)} - \boldsymbol{\vartheta}) \leq t \mid \hat{p}_{\text{aic}} = p, \mathcal{M}_{\text{nest}}) \\ &= P\{\mathbf{J}_p^{-1/2}\tilde{\mathbf{Z}}_{(p)} \leq \tilde{\mathbf{t}}_{(p)} \mid \tilde{\mathbf{Z}}_{(p)} \in \mathcal{A}_p^{(s)}(\mathcal{M}_{\text{nest}})\} \cdot I(t \in \mathcal{T}_p), \end{aligned}$$

where  $\mathcal{A}_p^{(s)}(\mathcal{M}_{\text{nest}}) = \{\tilde{\mathbf{z}}_{(p)} \in \mathbb{R}^{a+p} : \bigcap_{j=p_0+1, \dots, p} [\sum_{i=j}^p (z_{a+i}^2 - 2) > 0]\}$  and  $\mathcal{T}_p = \mathbb{R}^{a+p} \times (\mathbb{R}^+)^{K-p}$ .

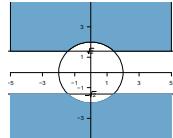
Simplified region for nested models due to independence of  $\tilde{\mathbf{Z}}_{(p)}$  and  $(Z_{p+1}, \dots, Z_K)$ .



# Asymptotic density

Truncated multivariate normal density

Denote  $\phi_p(\cdot | \mathcal{A}; \Sigma)$  the density of  $\Sigma^{-1/2} \tilde{\mathbf{Z}}_{(p)}$  where  $\tilde{\mathbf{Z}}_{(p)} \sim N_{a+p}(\mathbf{0}, I_{a+p})$  is truncated such that  $\tilde{\mathbf{Z}}_{(p)} \in \mathcal{A}$ .



**Corollary:** Post-AIC density for nested models

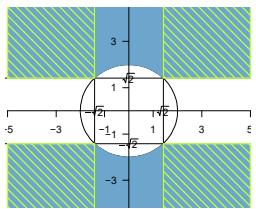
The limiting density function of  $n^{1/2}(\hat{\theta}_{(\hat{p}_{\text{aic}})} - \vartheta)$  conditional on the AIC-selection with  $\hat{p}_{\text{aic}} = p$  from the set of nested models  $\mathcal{M}_{\text{nest}}$  is

$$f_p(t) = \phi_p(\tilde{\mathbf{t}}_{(p)} | \mathcal{A}_p^{(s)}(\mathcal{M}_{\text{nest}}); \mathbf{J}_p^{-1}) \cdot I(t \in \mathcal{T}_p).$$

When the true model is selected,  $\hat{p}_{\text{aic}} = p_0$ ,

$$f_{p_0}(t) = \phi_{p_0}(\tilde{\mathbf{t}}_{(p_0)}) \cdot I(t \in \mathcal{T}_p)$$

## Inference depends on the set of models



Nested models:

$$M_1 : \{\theta_1\}; M_2 : \{\theta_1, \theta_2\}; M_3 : \{\theta_1, \theta_2, \theta_3\}$$

AIC selects  $M_3$ :

$$\text{AIC}^*(M_3) > \text{AIC}^*(M_1) \text{ and } \text{AIC}^*(M_3) > \text{AIC}^*(M_2)$$

$$\mathcal{A}_M(\mathcal{M}_{\text{nest}}) = \{z \in \mathbb{R}^3 : z_3^2 > 2, z_2^2 + z_3^2 > 4\}.$$

$M_{\text{all}} = \text{All subsets of } M_3 : \{\theta_1, \theta_2, \theta_3\}$

AIC selects  $M_3$ :

$$\text{AIC}^*(M_3) > \text{AIC}^*(M') \text{ for all models } M' \neq M_3 \in \mathcal{M}$$

$$\mathcal{A}_M(\mathcal{M}_{\text{all}}) = \{z \in \mathbb{R}^3 : z_2^2 > 2, z_3^2 > 2, z_2^2 + z_3^2 > 4\}$$

Different area, different conditioning

- ⇒ different distribution
- ⇒ different quantiles
- ⇒ different confidence intervals

## Exact calculations for nested models

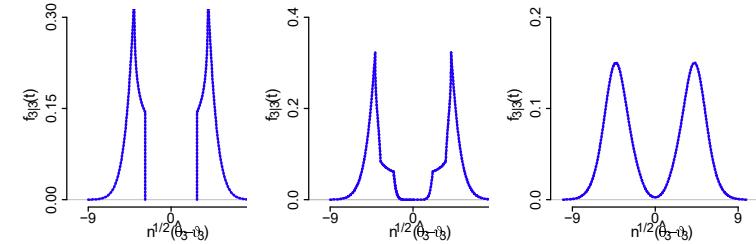
True value for parameters:  $\vartheta = (\vartheta_1, 0, 0)^t \rightarrow M_1$  true.

Models set:  $\mathcal{M}_{\text{nest}} = \{M_1, M_2, M_3\}$ .

Selected model:  $M_3 \Rightarrow \mathcal{A}_3 = \{(z_1, z_2, z_3) : z_3^2 > 2, z_2^2 + z_3^2 > 4\}$ .

$$\mathbf{J}^{-1/2}(\vartheta):$$

$$(a) \begin{pmatrix} 1.0 & 0 & 0 \\ 0 & 2.0 & 0 \\ 0 & 0 & 2.0 \end{pmatrix}, \quad (b) \begin{pmatrix} 1.0 & 0 & 0 \\ 0 & 2.0 & 0.5 \\ 0 & 0.5 & 2.0 \end{pmatrix}, \quad (c) \begin{pmatrix} 1.0 & 0.9 & 0.9 \\ 0.9 & 2.0 & 0.5 \\ 0.9 & 0.5 & 2.0 \end{pmatrix}.$$



Sample from a multivariate normal distribution subject to quadratic constraints. (Pakman and Paninski 2014)

## General likelihood models, any model set

For  $M_{\text{aic}} = M \in \mathcal{M}_O$ , overparametrized model, the confidence region

$$C(q_\alpha) = \left\{ \theta \in \mathbb{R}^{a+K} : n(\hat{\theta}_{(M)} - \tilde{\theta}_{(M)})^t \mathbf{J}_M(\hat{\theta}_{(M)} - \tilde{\theta}_{(M)}) \leq q_\alpha \right\},$$

where  $q_\alpha$  is determined by solving

$$\frac{P(\{\sum_{i \in M} Z_i^2 \leq q_\alpha\} \cap Z \in \mathcal{A}_M(\mathcal{M}_{\text{arb}} \cap \mathcal{M}_O))}{P(Z \in \mathcal{A}_M(\mathcal{M}_{\text{arb}} \cap \mathcal{M}_O))} = 1 - \alpha.$$

### Uniform result

Under some assumptions, incl. compact parameter space.

$$\lim_{n \rightarrow \infty} \inf_{\vartheta \in \Theta} P_\vartheta \{ \vartheta \in C(q_\alpha) \mid M_{\text{aic}} \in \mathcal{M}_O \} = 1 - \alpha.$$

When  $\mathcal{A}_M(\mathcal{M}_{\text{arb}})$  replaces  $\mathcal{A}_M(\mathcal{M}_{\text{arb}} \cap \mathcal{M}_O)$  to obtain  $\tilde{q}_\alpha$ ,  $\lim_{n \rightarrow \infty} \inf_{\vartheta \in \Theta} P_\vartheta \{ \vartheta \in C(\tilde{q}_\alpha) \mid M_{\text{aic}} \in \mathcal{M}_O \} \geq 1 - \alpha$ .

## About the uniformity result

- ▶ Limitation of AIC: selection of an overspecified model does not happen in a uniform way (Leeb, Pötscher, 2003).
- ▶ An overspecified model is needed. Hence result cannot be strengthened for AIC under this setup.
- ▶ For underparametrized model: need to work with pseudo-true values.
- ▶ 'Targets' differ from model to model.

G. Claeskens

KU LEUVEN

16/25

Assume (i) that there is a smallest model nested in all other models, (ii) similarity assumption of Vuong (1989), then construct selection region  $\mathcal{A}_{M_{\text{aic}}}(\mathcal{M})$  similar as before

### Uniform result under likelihood misspecification

$$\lim_{n \rightarrow \infty} \sup_{G_n \in \mathcal{G}_n} \sup_{t \in \mathbb{R}^{|M_{\text{aic}}|}} |P(n^{1/2}\{\hat{\theta}(M_{\text{aic}}) - \vartheta^*(M_{\text{aic}})\} \leq t \mid M_{\text{aic}}) - P(\Sigma^{1/2}Z \leq t \mid \mathcal{A}_{M_{\text{aic}}}(\mathcal{M}))| = 0$$

Define the set

Charkhi, Claeskens, 2018 Biometrika

$$C^*(q_\alpha) = \{\theta \in \mathbb{R}^{|M_{\text{aic}}|} : n\{\hat{\theta}(M_{\text{aic}}) - \theta(M_{\text{aic}})\}^t \Sigma_{M_{\text{aic}}}(\vartheta^*_{M_{\text{aic}}})^{-1}\{\hat{\theta}(M_{\text{aic}}) - \theta(M_{\text{aic}})\} \leq q_\alpha\},$$

$$\text{Then } \lim_{n \rightarrow \infty} \sup_{G_n \in \mathcal{G}_n} \sup_{\alpha \in [0,1]} |P_{G_n}\{\vartheta^*(M_{\text{aic}}) \in C^*(q_\alpha) \mid M_{\text{aic}}\} - (1-\alpha)| = 0.$$

G. Claeskens

KU LEUVEN

18/25

## All likelihood models misspecified

- ▶ Triangular array of observations  $\{Y_{ni}; i = 1, \dots, n, \text{ and } n \in \mathbb{N}_0\}$ ;
- ▶ True distribution of  $(Y_{n1}, \dots, Y_{nn})$  is  $G_n$
- ▶ Estimators maximizing  $\prod_{i=1}^n f_{ji}(y_i; \theta_j)$  converge to pseudo-true values  $\vartheta_n^*(M_j)$ ,  $j = 1, \dots, K$ ;
- ▶ When there is an estimator  $\hat{\Sigma}$  of  $\Sigma$  such that

$$\lim_{n \rightarrow \infty} \sup_{G_n \in \mathcal{G}_n} P_{G_n}(\|\hat{\Sigma}_n - \Sigma\| > \varepsilon) = 0,$$

then, with  $\mathcal{Z}_{m'} \sim N_{m'}(0, I_{m'})$

$$\lim_{n \rightarrow \infty} \sup_{G_n \in \mathcal{G}_n} \sup_{t \in \mathbb{R}^{m'}} |P(\hat{\Sigma}_n^{-1/2}n^{-1/2}(\hat{\theta}_{n,\mathcal{M}} - \vartheta_n^*, \mathcal{M}) \leq t) - P(\mathcal{Z}_{m'} \leq t)| = 0.$$

- ▶ White (1994) contains conditions. Sandwich estimator might overestimate  $\hookrightarrow$  conservative confidence intervals.

G. Claeskens

KU LEUVEN

17/25

## Simulation study – Linear Regression

Model to generate data

$$Y_i = \sum_{j=1}^{10} \theta_j x_{ji} + \varepsilon_i, \quad i = 1, \dots, n,$$

- ▶  $\varepsilon_i \sim N(0, 1)$
- ▶  $\boldsymbol{\theta}^t = (\theta_1, \dots, \theta_{10})^t = (2.25, -1.1, 2.43, -2.24, 2.5, \mathbf{0}_5^t)$ .
- ▶  $x_{1i} = 1$  and  $(x_{2i}, \dots, x_{10,i})^t \sim N(\mathbf{0}_9, \boldsymbol{\Omega})$  where  $\Omega_{ii} = 1$  and  $\Omega_{ij} = 0.25$ ,  $j \neq i$ .

Selection by AIC from  $\zeta_{\text{all}}^i$ : first  $i$  parameters are present in all models, all combinations of other parameters.

Focus model:  $M = (\theta_1, \dots, \theta_6, \theta_8)$ : Run simulation until  $M$  is selected 3000 times.

Report average confidence intervals over 3000 runs and coverage probabilities.

R function PostAIC

G. Claeskens

KU LEUVEN

19/25

## CI for regression parameter

PoSI: valid inference irrespective of model selection method.  
 (Berk, Brown, Buja, Zhang, Zhao 2013 AoS)

$n$	method	$\theta_j$	$\zeta_{\text{all}}^3$	Cov.	$\zeta_{\text{all}}^6$	Cov.
100	PostAIC	$\theta_4$	[-2.54, -1.94]	0.99	[-2.46, -2.02]	0.94
		$\theta_6$	[-0.30, 0.31]	0.95	[-0.22, 0.22]	0.95
		$\theta_8$	[-0.30, 0.31]	0.95	[-0.29, 0.30]	0.95
	PoSI	$\theta_4$	[-2.58, -1.90]	1.00	[-2.54, -1.94]	0.99
		$\theta_6$	[-0.33, 0.34]	0.98	[-0.30, 0.30]	0.99
	Naive	$\theta_8$	[-0.34, 0.34]	0.98	[-0.29, 0.31]	0.95
		$\theta_4$	[-2.46, -2.02]	0.93	[-2.46, -2.02]	0.93
	Naive	$\theta_6$	[-0.22, 0.22]	0.66	[-0.22, 0.22]	0.94
		$\theta_8$	[-0.22, 0.22]	0.66	[-0.22, 0.23]	0.69

$\zeta_{\text{all}}^3$  : More conservative for  $\theta_4$ .

$\zeta_{\text{all}}^6$  : Almost exact for all parameters.

KU LEUVEN

## Extensions to other criteria

AIC-like criteria, overselection is exploited.

- Takeuchi's information criterion (1976)

$$\text{TIC}(M) = -2\ell_n(\hat{\theta}(M)) + 2\text{tr}\{\mathbf{Q}_M(\boldsymbol{\vartheta}^*)^{-1}\mathbf{J}_M(\boldsymbol{\vartheta}^*)\}.$$

Replace  $|M|$  with  $\text{tr}\{\mathbf{Q}_M(\boldsymbol{\vartheta}^*)^{-1}\mathbf{J}_M(\boldsymbol{\vartheta}^*)\}$

- Generalized information criterion (Konishi, Kitagawa 1996)

$$\text{GIC}(M) = -2\ell_n(\hat{\theta}(M)) + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\text{Infl}(Y_i) \frac{\partial}{\partial \boldsymbol{\theta}_M^t} \log f(Y_i; \hat{\theta}_M)\}.$$

- Mallows's  $C_p$  for linear regression (1973)

$$C_p(M) = \hat{\sigma}^{-2} \hat{\sigma}^2(M) + 2|M| - n$$

When  $n \rightarrow \infty$ ,  $C_p(M) - C_p(M^*) \sim \chi_q^2/q + 2q$  where  
 $q = |M^*| - |M|$

KU LEUVEN

## CI for a linear combination of the parameters $x^t\theta$

PoSIp (Bachoc, Leeb, Pötscher, 2015)

Smoothed bootstrap CI (Efron, 2014)

$\sigma$	$n$	method	$\zeta_{\text{all}}^3$		$\zeta_{\text{all}}^5$	
			length	cov.	length	cov.
1	30	PostAIC	3.11	0.97	2.61	0.95
		Bootstrap	3.67	0.92	3.32	0.92
		PoSIp	4.38	1.00	4.39	1.00
100	100	PostAIC	1.42	0.98	1.17	0.95
		Bootstrap	1.25	0.94	1.25	0.94
		PoSIp	1.83	1.00	1.83	1.00
3	30	PostAIC	11.76	0.98	7.82	0.94
		Bootstrap	11.46	0.92	9.95	0.92
		PoSIp	12.65	0.99	13.16	1.00
100	100	PostAIC	4.25	0.98	3.50	0.95
		Bootstrap	3.77	0.94	3.74	0.94
		PoSIp	5.47	1.00	5.48	1.00

KU LEUVEN

## Take home messages

- Inference needs to be adapted after model/variable selection
- Think carefully about the sets of models to include in a search.
- Also regularized estimation causes inference problems (lasso etc.)
- Regarding AIC:
  - Asymptotic AIC selection regions are formed by quadratic functions.
  - Stronger uniform results under misspecification (pseudo-true parameter values).
- Confidence curves give a fuller picture. Work in progress

Thank you!

KU LEUVEN