

# Bootstrap test procedure for variance components in nonlinear mixed effects models in the presence of nuisance parameters and singular Fisher Information Matrix

T. GUEDON, C. BAEY, E. KUHN

## Abstract

We examine the problem of variance components testing in general mixed effects models using the likelihood ratio test. We account for the presence of nuisance parameters, i.e. the fact that some untested variances might also be equal to zero. Two main issues arise in this context leading to a non regular setting. First, under the null hypothesis the true parameter value lies on the boundary of the parameter space. Moreover, due to the presence of nuisance parameters the exact location of these boundary points is not known, which prevents from using classical asymptotic theory of maximum likelihood estimation. Then, in the specific context of nonlinear mixed-effects models, the Fisher information matrix is singular at the true parameter value. We address these two points by proposing a shrunk parametric bootstrap procedure, which is straightforward to apply even for nonlinear models. We show that the procedure is consistent, solving both the boundary and the singularity issues, and we provide a verifiable criterion for the applicability of our theoretical results. We show through a simulation study that, compared to the asymptotic approach, our procedure has a better small sample performance and is more robust to the presence of nuisance parameters. A real data application is also provided.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Proposed methodology</b>	<b>6</b>
2.1	Mixed effects models . . . . .	6
2.2	Variance components testing . . . . .	8
2.3	Testing procedure . . . . .	9
<b>3</b>	<b>Theoretical results</b>	<b>10</b>
3.1	Notations and theoretical setting . . . . .	10
3.2	Consistency of the bootstrap procedure in the identically distributed setting .	13
3.3	Extension to the non identically distributed setting . . . . .	16
3.4	Sufficient verifiable conditions for regularity assumptions . . . . .	17
<b>4</b>	<b>Experiments</b>	<b>18</b>
4.1	Simulation study . . . . .	18
4.2	Real data application . . . . .	22
<b>5</b>	<b>Discussion</b>	<b>23</b>
<b>6</b>	<b>Acknowledgment</b>	<b>24</b>
<b>7</b>	<b>Supplementary material</b>	<b>24</b>
7.1	Different parametrizations for mixed effects models . . . . .	25
7.2	Proof of proposition 1 . . . . .	25
7.3	Proof of proposition 2 . . . . .	26
7.4	Proof of proposition 4 . . . . .	27
7.5	Quadratic approximation of the log-likelihood . . . . .	30
7.6	Asymptotic distribution of the likelihood ratio test statistic . . . . .	32
7.7	Proof of proposition 3 . . . . .	34
7.8	Proof of theorem 1 . . . . .	36

7.9	Proof of proposition 5 . . . . .	43
7.10	Proof of proposition 6 . . . . .	45
7.11	Proof of theorem 2 . . . . .	45
7.12	Proof of proposition 7 . . . . .	48
7.13	Logistic growth model example . . . . .	59

# 1 Introduction

Mixed effects models are a powerful statistical tool to model longitudinal studies with repeated measurements or data with an underlying unknown latent structure as hierarchical data. There are many fields of applications, e.g. pharmacokinetic-pharmacodynamic [8], medicine [9], agriculture [39], ecology [7], psychology [29] or educational and social sciences [20]. These models allow to take into account two types of variabilities, between different individuals in a population and between several measurements made on the same individual, also called inter and intra variabilities. These are modeled by two types of effects: on the one hand, random effects that vary from one individual to another, and on the other hand, fixed effects, common to all individuals in the population (see [32], [15]).

From a modeling point of view, being able to distinguish among all effects those that can be modeled as fixed effects would allow to reduce the number of model parameters. This would also help to better identify the processes that are at the origin of the variability observed in the population. Two main approaches have been developed to tackle this task. On the one hand, some authors suggested methods based on variable selection, using a Bayesian procedure as in [13] or a penalized likelihood approach as in [27] or [22]. Specific selection criteria for mixed-effects models were also developed by [36, 23] and [16]. On the other hand, other authors focused on hypothesis testing for the nullity of some variance components of the random effects. Such a test is equivalent to comparing two nested models, and standard tools to address this question include the likelihood ratio, the score and the Wald tests statistics [37]. However two issues arise when testing the nullity of variance components, that prevent from using the usual asymptotic results of [38]. The first issue results from the true value of the variance parameter lying on the boundary of the parameter space, while the second is due to the singularity of the Fisher information matrix.

As far as the boundary issue is concerned, several authors studied the asymptotic of the likelihood ratio test statistic in this context. [14], [12], [33] derived the asymptotic distribution of the likelihood ratio test statistic in specific cases. [2], [34] gave a more general way of dealing with hypothesis testing when the true parameter is not constrained to be an

interior point of the ambient space. In the context of mixed effects models, [4] derived the asymptotic distribution of the likelihood ratio test statistic for testing that any subset of the variances of the random effects is null. However, most of these references assume that the untested parameters do not lie on the boundary of the parameter space. When this is the case, the asymptotic distribution of the likelihood ratio test statistic is intractable as it depends on the unknown location of these nuisance parameters on the boundary [33]. Therefore, procedures that do not involve the asymptotic distribution of the test statistic can be preferred. In particular, resampling methods such as bootstrap and permutations are powerful tools to address this issue. In addition, these methods are usually more robust in small samples context. [35] proposed a bootstrap-based score test in the context of generalized linear mixed models with one single random effect. [17] proposed a permutation-based test for any subset of the covariance matrix of the random effects in linear mixed models. The latter method is very easy to use in practice but is restricted to the context of linear models. The former requires the computation of the Fisher information matrix, which can be heavy in practice, especially in the context of nonlinear models. Moreover, the presence of nuisance parameters on the boundary of the parameter space is not considered in the aforementioned works, even though it can be a source of inconsistency for the proposed bootstrap procedures. Indeed, as discussed in [5], and highlighted in [3], the bootstrap is known to be inconsistent when the true parameter value is a boundary point. When estimating the expectation of a Gaussian distribution, restricted to be nonnegative, [3] proposed a parametric procedure that shrinks the parameter used to generate the bootstrap data near the boundary. Following this idea, [11] proposed a more general parametric bootstrap test procedure based on the Likelihood Ratio Test statistic with parameters lying on the boundary. Their method consists in shrinking the bootstrap parameter in order to accelerate its rate of convergence toward the boundary.

The second issue is the singularity of the Fisher information matrix that arises specifically in the context of mixed effect models, as discussed in [18]. Following the development of [25] to derive the new asymptotic distribution of the likelihood ratio test statistic, we show that this singularity issue is another source of inconsistency of the bootstrap procedure.

In this work we propose a shrunked parametric bootstrap test procedure for variance com-

ponents in nonlinear mixed effects models that addresses the two issues mentioned above. We show that given an appropriate choice of the bootstrap parameter, the procedure is consistent as the number of individuals grows to infinity. Our contribution is twofold: first, our procedure can be applied to linear, generalized linear and nonlinear models, and second, it takes into account the presence of nuisance parameters at unknown locations. We also provide a verifiable criterion to check the required regularity conditions. Finally, we illustrate our results on simulated and real data, exhibiting the good finite sample properties of the procedure and its applicability in practice.

Section 2 presents the mixed effects model that is considered in this article and describes the proposed shrunk parametric bootstrap procedure. Section 3 is dedicated to the theoretical study of the proposed procedure, and shows its consistency first in the independent and identically distributed setting and then in the more general independent but not identically distributed setting. Section 4 presents a simulation study that illustrates the performance of the test in different models and with varying sample sizes. A real data application is also presented. All the proofs and additional developments can be found in the supplementary material.

## 2 Proposed methodology

### 2.1 Mixed effects models

Consider  $N$  individuals each measured  $J_i < J$  times, where  $N$  and  $J_i$  are nonnegative integers. We denote by  $y_{ij}$  ( $i = 1, \dots, N; j = 1, \dots, J_i$ ) the  $j$ th observation of the  $i$ th individual and we define  $y_i = (y_{i1}, \dots, y_{iJ_i})$  and  $y_{1:N} = (y_1^T, \dots, y_N^T)$ . In the sequel,  $\mathcal{L}_p^+$  denotes the space of lower triangular matrices of size  $p \times p$  with positive diagonal coefficients,  $\mathbb{S}_+^p$  denotes the space of symmetric, positive semi-definite  $p \times p$  matrices,  $I_p$  is the identity matrix of size  $p \times p$ ,  $[A]_{ij}$  is the element on the  $i$ th line and  $j$ th column of matrix  $A$ , and  $\mathcal{N}(\mu, V)$  denotes the multivariate Gaussian distribution with expectation  $\mu \in \mathbb{R}^p$  and covariance matrix  $V$  of

size  $p \times p$ . We consider the following nonlinear mixed effects model

$$\begin{cases} y_{ij} = g(x_{ij}, \beta, \Lambda \xi_i) + \varepsilon_{ij} & \varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2) \\ \xi_i \sim \mathcal{N}(0, I_p) \end{cases}, \quad (1)$$

where  $(\xi_i)_{i=1, \dots, N}$  and  $(\varepsilon_{ij})_{i=1, \dots, N, j=1, \dots, J_i}$  are mutually independent random variables,  $g$  is a known nonlinear function,  $x_{ij}$  gathers all the covariates of the  $j$ th observation of the  $i$ th individual,  $\beta \in \mathbb{R}^b$  is the vector of fixed effects,  $\Lambda \in \mathcal{L}_p^+$  is a scaling parameter for random effect  $\xi_i$ , and  $\sigma^2$  is the positive noise variance.

**Remark 1** The covariance matrix of the scaled random effect  $b_i = \Lambda \xi_i$  is equal to  $\Gamma = \Lambda \Lambda^T$  which is positive semi-definite. Therefore, a natural choice for  $\Lambda$  is the lower triangular matrix in the Cholesky decomposition of the scaled random effects covariance matrix. Moreover, by constraining the diagonal coefficients of  $\Lambda$  to be nonnegative,  $\Lambda$  is uniquely defined. This decomposition is used for instance in [13].

**Remark 2** The definition of model (1) is slightly more general than the usual terminology of mixed effects models [32, p. 306] that defines  $y_{ij} = g(v_{ij}, \phi_i) + \varepsilon_{ij}$ , with  $\phi_i = A_{ij}\beta + B_{ij}b_i$  the  $i$ th individual parameter,  $\beta$  the vector of fixed effects associated with random effect  $b_i \sim \mathcal{N}(0, \Gamma)$ , and where  $v_{ij}$ ,  $A_{ij}$  and  $B_{ij}$  are known covariates. Model (1) covers this definition by taking  $x_{ij} = (v_{ij}, B_{ij}, A_{ij})$  and  $b_i = \Lambda \xi_i$ . A more detailed development of the differences between those two parameterizations is given in section (7.1) of the supplementary material.

Let us denote by  $\theta = (\beta, \Lambda, \sigma^2)$  the unknown vector of model parameters taking values in  $\Theta$ , by  $f_i(\cdot; \theta)$  the density of the  $i$ th individual response  $y_i$  given a parameter  $\theta \in \Theta$ , by  $f_i(y_i; \xi_i, \theta)$  the conditional density of  $y_i$  given the random effect  $\xi_i$  and a parameter  $\theta$ , and by  $\pi_p(\cdot)$  the density of the  $p$ -dimensional standard Gaussian density. With these notations we can define the log-likelihood of the model given the  $N$ -sample  $y_{1:N}$  by

$$l(\theta; y_{1:N}) = \log\{L_\theta(y_{1:N})\} = \log\left\{\prod_{i=1}^N f_i(y_i; \theta)\right\} = \sum_{i=1}^N \log\left\{\int f_i(y_i; \xi_i, \theta) \pi(\xi_i) d\xi_i\right\} \quad (2)$$

We consider the marginal likelihood defined as the complete likelihood integrated over the distribution of the random effects, since the random effects  $\xi_i$  are unobserved.

With our notations, the distribution of the random effects  $\xi_i$  does not depend on any parameter, contrary to the common case where the random effects are defined as the scaled version  $b_i$ . With the latter definition, the distribution of the latent variables depends on  $\Lambda$ , and is not defined on the entire parameter space since we only constrain  $\Gamma = \Lambda\Lambda^T$  to be positive semi-definite. When dealing with linear models, since the variance of the random effects adds up with the noise variance, the fact that some diagonal components in  $\Gamma$  are null is not an issue. However, the change of variable  $b_i = \Lambda\xi_i$  is a  $\mathcal{C}^1$  diffeomorphism if and only if the diagonal coefficients of  $\Lambda$  are nonnegative. Without this assumption the two parametrizations are no longer equivalent as illustrated in the supplementary material (see section 7.1). Our parametrization is similar to the so-called reparametrization trick proposed by [28] to train variational autoencoders with back-propagation.

## 2.2 Variance components testing

Let  $r \in \{1, \dots, p\}$  be the number of variances to be tested. Without loss of generality we assume that we test the nullity of the last  $r$  variances in  $\Gamma = \Lambda\Lambda^T$ . Therefore let us consider the following block matrix notation

$$\Lambda = \left( \begin{array}{c|c} \Lambda_1 & 0_{(p-r) \times r} \\ \hline \Lambda_{12} & \Lambda_2 \end{array} \right),$$

where  $\Lambda_1 \in \mathcal{L}_{p-r}^+$ ,  $\Lambda_2 \in \mathcal{L}_r^+$  and  $\Lambda_{12} \in \mathcal{M}_{r \times (p-r)}(\mathbb{R})$ . We write  $\theta_0$  the true parameter on which we consider the following test :

$$H_0 : \theta_0 \in \Theta_0 \quad \text{against} \quad H_1 : \theta_0 \in \Theta, \quad (3)$$



where

$$\begin{aligned}\Theta_0 &= \{\theta \in \mathbb{R}^q \mid \beta \in \mathbb{R}^b, \Lambda_1 \in \mathcal{L}_{p-r}^+, \Lambda_2 = 0, \Lambda_{12} = 0, \sigma^2 \in \mathbb{R}_*^+\}, \\ \Theta &= \{\theta \in \mathbb{R}^q \mid \beta \in \mathbb{R}^b, \Lambda \in \mathcal{L}_p^+, \sigma^2 \in \mathbb{R}_*^+\}.\end{aligned}$$

**Remark 3** We do not impose the diagonal of  $\Lambda_1$  to be strictly non-negative, which enables the case where some untested variances of the scaled random effects are in fact equal to zero. This will be discussed in more details in section (3.1) with the definition of nuisance parameters.

The likelihood ratio test statistic is defined as

$$\text{LRT}(y_{1:N}) = -2 \left( \sup_{\theta \in \Theta} l(\theta; y_{1:N}) - \sup_{\theta \in \Theta_0} l(\theta; y_{1:N}) \right).$$

In order to test (3) with a nominal level  $0 < \alpha < 1$ , we define the rejection region as  $R_\alpha = \{\text{LRT}(y_{1:N}) \geq q_\alpha\}$  with  $q_\alpha$  being the  $(1 - \alpha)$ th quantile of the distribution of  $\text{LRT}(y_{1:N})$ . Unfortunately, this distribution is often intractable.

In the following section, we detail the proposed shrunked parametric bootstrap procedure to test (3).

### 2.3 Testing procedure

We propose a parametric bootstrap procedure using a bootstrap parameter  $\theta_N^*$  and  $B \in \mathbb{N}^*$  bootstrap replications to test (3) with a type I error  $0 < \alpha < 1$ . As introduced in remark 3 and then detailed in section (3.1), some untested variances, at unknown locations, can be null. Therefore using  $\hat{\theta}_N = (\hat{\beta}_N, \hat{\Lambda}_N, \hat{\sigma}_N^2)$  the maximum likelihood estimator as a bootstrap parameter over  $\Theta$  would fail to asymptotically mimic the true distribution of the likelihood ratio test statistic. Indeed there are elements in  $\Lambda_0$  which are supposed to be zero, but that are non null in  $\hat{\Lambda}_N$ . Since we require that  $\theta_N^* \in \Theta_0$ , we use a shrinking parameter  $c_N$  to fix to zero the diagonal parameters of  $\Lambda$  that are smaller than  $c_N$ . The proposed algorithm is described

in algorithm 1, and the theoretical justification of this shrinking procedure is described in section 3.

---

**Algorithm 1** Shrunked parametric bootstrap for variance components testing

---

Input:  $c_N > 0$ ,  $B \in \mathbb{N}^*$ ,  $0 < \alpha < 1$   
Set  $\beta_N^* = \hat{\beta}_N$ ,  $\Lambda_N^* = \hat{\Lambda}_N$ , and  $\sigma_N^{*2} = \hat{\sigma}_N^2$   
Set  $\Lambda_{2,N}^* = \Lambda_{12,N}^* = 0$   
Set  $[\Lambda_{1,N}^*]_{mn} = [\hat{\Lambda}_{1,N}]_{mn} \mathbb{1}_{[\hat{\Lambda}_{1,N}]_{mn} > c_N}$   
For  $b = 1, \dots, B$   
    For  $i = 1, \dots, N$ , draw independently  $\varepsilon_i^{*,b} \sim \mathcal{N}(0, \sigma_N^{*2} I_{J_i})$  and  $\xi_i^{*,b} \sim \mathcal{N}(0, I_p)$   
    Build the  $i$ th value of the  $b$ th bootstrap sample  $y_i^{*,b} = g(x_i, \beta_N^*, \Lambda_N^* \xi_i^{*,b}) + \varepsilon_i^{*,b}$   
    Compute the likelihood ratio statistic  $\text{LRT}(y_{1:N}^{*,b})$   
Compute the bootstrap  $p$ -value as  $p_{boot} = \frac{1}{B} \sum_{b=1}^B \mathbb{1}_{\text{LRT}(y_{1:N}^{*,b}) > \text{LRT}(y_{1:N})}$   
Reject  $H_0$  if  $p_{boot} < \alpha$

---

The next section is dedicated to the asymptotic validity of this testing procedure.

## 3 Theoretical results

### 3.1 Notations and theoretical setting

In this section we are interested in the theoretical consistency of the bootstrap procedure presented in section 2.3. We consider the asymptotic as the number of individuals  $N$  grows to infinity, while the number of measurements per individual remains fixed and bounded by some value  $J$ . We denote by  $\theta_0 \in \Theta_0$  the true value of the parameter, such that the density of the response  $y_{1:N}$  is  $L_{\theta_0}(y_{1:N})$ . We denote by  $E\{T(y_{1:N})\}$  the expectation of any measurable function  $T$  of  $y_{1:N}$  if there is no confusion about the distribution of  $y_{1:N}$ . Otherwise we specify  $E_\theta\{T(y_{1:N})\}$  to emphasize that the expectation is with respect to the density  $L_\theta(y_{1:N})$ , for any  $\theta \in \Theta$ . As commonly used in the bootstrap literature, we denote by  $X^*$  the bootstrap version of a random variable  $X$ . We write  $E^*\{T(y_{1:N}^*)\} = E_{\theta_N^*}\{T(y_{1:N}^*) \mid y_{1:N}\}$ . Similarly, for any measurable subset  $A$  we write  $\text{pr}^*\{T(y_{1:N}^*) \in A\} = E_{\theta_N^*}\{\mathbb{1}_{T(y_{1:N}^*) \in A} \mid y_{1:N}\}$ .

We want to show that the proposed bootstrap procedure is asymptotically valid which means that  $\text{LRT}(y_{1:N}^*)$  converges weakly in probability to the same limiting distribution as the one of  $\text{LRT}(y_{1:N})$ . More precisely, we want to show that, under some conditions, if there exists a random variable  $\text{LRT}_\infty$  such that  $\text{LRT}(y_{1:N})$  converges weakly to  $\text{LRT}_\infty$  then for every  $t \in \mathbb{R}$ , as  $N \rightarrow +\infty$ , it holds in probability that

$$\text{pr}^*\{\text{LRT}(y_{1:N}^*) \leq t\} \longrightarrow \text{pr}(\text{LRT}_\infty \leq t). \quad (4)$$

We also use the notations  $o_p(1)$  and  $O_p(1)$  for random sequences that respectively converge toward zero and are bounded in probability. More generally, this notation is used to compare two random sequences, using the definition of [37]. We also use their bootstrap versions  $o_{p^*}$  and  $O_{p^*}$  defined as follows: for a random quantity  $X_N^*$  computed on the bootstrap data,  $X_N^* = o_{p^*}(1)$  means that for any  $\varepsilon > 0$ ,  $\text{pr}^*(X_N^* > \varepsilon) \rightarrow 0$  in probability as  $N \rightarrow +\infty$ . Similarly  $X_N^* = O_{p^*}(1)$  means that for any  $\varepsilon > 0$  there exists a real  $M > 0$  and an integer  $N_0$  such that for all  $N > N_0$ , the event  $\{\text{pr}^*(\|X_N^*\| > M) < \varepsilon\}$  is arbitrary close to one in probability.

We now formalize what we call nuisance parameters. We suppose that, in addition to the last  $r$  tested variances,  $m$  untested variances are null. Without loss of generality we suppose that the last  $m + r$  variances of the individual parameters are null therefore  $\Lambda_0$  is of the form

$$\Lambda_0 = \left( \begin{array}{c|c|c} \Lambda_1^{nonuis} & 0_{(p-r-m) \times m} & 0_{(p-r-m) \times r} \\ \hline \Lambda_{12}^{nuis} & \Lambda_1^{nuis} & 0_{m \times r} \\ \hline \Lambda_{12,1} & \Lambda_{12,2} & \Lambda_2 \end{array} \right) = \left( \begin{array}{c|c|c} \Lambda_1^{nonuis} & 0_{(p-r-m) \times m} & 0_{(p-r-m) \times r} \\ \hline 0_{m \times (p-r-m)} & 0_{m \times m} & 0_{m \times r} \\ \hline 0_{r \times (p-r-m)} & 0_{r \times m} & 0_{r \times r} \end{array} \right).$$

It is important to notice that in real life applications the  $m$  rows inducing nuisance parameters are located at unknown positions in matrix  $\Lambda$ , and that the remaining  $p - m - r$  variances are strictly non-negative which is equivalent to the diagonal coefficients of  $\Lambda_1^{nonuis}$  being strictly non-negative.

We now split the parameter as  $\theta = (\psi, \delta, \lambda)$ , where  $\lambda$  stands for all the coefficients of  $\Lambda_2$  and  $\Lambda_{12,2}$ ,  $\delta$  represents the coefficients in  $\Lambda_1^{nuis}$  and  $\psi$  gathers all the remaining parameters.

The dimension of  $\lambda$  is  $d_\lambda = r(r+1)/2 + r(p-r-m)$ , the dimension of  $\delta$  is  $d_\delta = m(m+1)/2 + r \times m$  and the dimension of  $\psi$  is  $d_\psi = d_\theta - d_\lambda - d_\delta$ . Moreover,  $\theta_0 = (\psi_0, \delta_0, \lambda_0) = (\psi_0, 0_{d_\delta}, 0_{d_\lambda})$ . Before introducing our results we first state a set of general conditions on the model that will be required in this work.

**Assumption 1** *i)  $\Theta$  is compact, ii) the support of  $y \mapsto f_i(y; \theta)$  does not depend neither on  $\theta$ , nor on  $i$ , iii) the model is identifiable, iv) for all  $i \in \mathbb{N}$ ,  $y \in \mathbb{R}_i^J$ ,  $\xi \in \mathbb{R}^p$  the conditional likelihood  $\theta \mapsto f_i(y; \xi, \theta)$  is 4-times differentiable on the interior of  $\Theta$ , and directional derivatives exist on the boundary, v) each partial derivative of  $\theta \mapsto f_i(y; \xi, \theta)$  is bounded by a positive function which does not depend on  $\theta$  and is integrable with respect to the distribution of the random effects.*

**Remark 4** The compactness assumption is not verified for  $\Theta$ . However in practice it only requires that  $\sigma^2 \geq \rho$  for some non-negative number  $\rho$  and that each component of  $\theta$  is upper and lower bounded, which is reasonable in real data application. Assumptions ii) and iii) are usual in the context of maximum likelihood theory, iv) is needed to perform a Taylor expansion of the log likelihood and v) is needed to differentiate under the integral sign in (2).

The following proposition induces that if the Fisher Information Matrix exists, it will present blocks equal to zero, and will therefore be singular.

**Proposition 1** *Under assumption (1), for  $k = 0, 1$ , for all  $i \in \mathbb{N}$  and for all  $y \in \mathbb{R}^{J_i}$ ,  $\nabla_\delta^{2k+1} \log\{f_i(y; \theta_0)\} = 0_{d_\delta^{2k+1}}$  and  $\nabla_\lambda^{2k+1} \log\{f_i(y; \theta_0)\} = 0_{d_\lambda^{2k+1}}$ . In particular,  $\text{var}\{\nabla_\delta l(\theta; y_{1:N})\} = 0_{d_\delta \times d_\delta}$  and  $\text{var}\{\nabla_\lambda l(\theta; y_{1:N})\} = 0_{d_\lambda \times d_\lambda}$ .*

**Remark 5** If  $\theta \mapsto l(\theta; y_{1:N})$  admits higher order derivatives, the first part of Proposition 1 is true for every odds order derivatives. This comes from the null odds moments of the standard normal distribution of the random effects.

**Remark 6** As shown in the proof of proposition 1, in section 7.2 of the appendix, by considering the  $k$ th column  $[\Lambda]_{.k} = ([\Lambda]_{1k}, \dots, [\Lambda]_{pk})^T$  of  $\Lambda$ , for all  $j = 1, \dots, p$ ,  $\partial l(\theta; y_{1:N}) / \partial [\Lambda]_{jk} |_{[\Lambda]_{.k} = 0_p} = 0$ . That explains why the coefficients of  $\Lambda_{12,2}$  are part of the definition of  $\lambda$ .

### 3.2 Consistency of the bootstrap procedure in the identically distributed setting

We first deal with the simpler identically distributed case. In model (1) it corresponds to the case where  $(x_{ij})_{j=1,\dots,J_i}$  are common to every individual  $i$ . The next section is devoted to extending the results to the non identically distributed setting presented before.

Before studying the consistency of the test procedure, we first need to ensure the consistency of the restricted (respectively unrestricted) maximum likelihood estimator, i.e. computed over  $\Theta_0$  (respectively  $\Theta$ ). We first state the regularity conditions required for the asymptotic theory that follows.

#### Assumption 2

- i)*  $\sup_{\theta' \in \Theta} E_{\theta'} \{ \sup_{\theta \in \Theta} |\log f(y_i; \theta)|^2 \} < +\infty$
- ii)*  $\sup_{\theta' \in \Theta} E_{\theta'} \{ \sup_{\theta \in \Theta} \|\nabla_{\theta} \log f(y_i; \theta)\|^3 \} < +\infty$
- iii)*  $\sup_{\theta' \in \Theta} E_{\theta'} \{ \sup_{\theta \in \Theta} \|\nabla_{\theta}^2 \log f(y_i; \theta)\|^3 \} < +\infty$
- iv)*  $\sup_{\theta' \in \Theta} E_{\theta'} \{ \sup_{\theta \in \Theta} \|\nabla_{\theta}^3 \log f(y_i; \theta)\|^2 \} < +\infty$
- v)*  $\sup_{\theta' \in \Theta} E_{\theta'} \{ \sup_{\theta \in \Theta} \|\nabla_{\theta}^4 \log f(y_i; \theta)\|^2 \} < +\infty$

Assumption (2) *i)* is needed to ensure the consistency of the maximum likelihood estimators. Indeed it enables to derive a uniform law of large numbers. Assumptions *ii)* and *iii)* are similar to assumption (N8') in [26]. It is required to apply a central limit theorem to the score function, and the pseudo score function  $\tilde{S}_N(\theta)$  that appears in the quadratic expansion (see equation (21) in the appendix). Assumptions *iv)* and *v)* are needed to control the rest of the quadratic approximation. All the suprema are needed to control the consistency of the bootstrap distributions.

We now derive the consistency of the maximum likelihood estimators, following the result of [30].

**Proposition 2** *Under assumptions (1)–(2) i) :*

$$\arg \max_{\theta \in \Theta} l(\theta; y_{1:N}) = \theta_0 + o_p(1)$$

$$\arg \max_{\theta \in \Theta_0} l(\theta; y_{1:N}) = \theta_0 + o_p(1)$$

A natural choice for the bootstrap parameter  $\theta_N^*$  is the maximum likelihood estimator. However the bootstrap fails in presence of the nuisance parameters summarized in vector  $\delta$ . This is why care must be taken when choosing  $\delta_N^*$ . To explain and solve this issue we first need to derive the speed of convergence of the maximum likelihood estimator.

**Proposition 3** *Let  $\hat{\theta}_N = (\hat{\psi}_N, \hat{\delta}_N, \hat{\lambda}_N)$  and  $\tilde{\theta}_N = (\tilde{\psi}_N, \tilde{\delta}_N, 0_{d_\lambda})$  be respectively the unrestricted and restricted maximum likelihood estimators of  $\theta$ . Under assumptions (1) and (2)  $(\hat{\psi}_N, \tilde{\psi}_N) = O_p(N^{-1/2}), (\hat{\delta}_N, \tilde{\delta}_N, \hat{\lambda}_N) = O_p(N^{-1/4})$ .*

The usual way to derive the asymptotic distribution of the likelihood ratio statistic is to consider a quadratic approximation of the log-likelihood around the true value of the parameter, based on a second-order Taylor expansion. However in our case, due to the vanishing score property stated in proposition 1, this quadratic expansion is degenerate with respect to parameters  $\delta$  and  $\lambda$ . Using a reparametrization of parameter  $\theta$  as in [25], we obtain a new quadratic approximation based on a higher-order expansion of the log-likelihood. We then apply results from [2] and [34] to obtain an explicit formula for  $LRT_\infty$ . This new quadratic approximation involves a new matrix  $\tilde{I}(\theta_0)$  which plays the role of the Fisher information matrix, and which is defined explicitly in equation (21) of the supplementary material. This new matrix is no longer systematically degenerate, we can therefore state the usual assumption which must be verified case by case in real life applications.

**Assumption 3**  $\tilde{I}(\theta_0) \succ 0$

Before showing that the proposed test procedure is consistent, we first need to show that in the bootstrap world, if the bootstrap parameter is consistent, the bootstrap maximum likelihood estimators are, conditionally to the data, consistent.

**Proposition 4** *Under assumptions (1)–(2) i), if  $\theta_N^*$  the parameter used to generate the data is consistent, then :*

$$\arg \max_{\theta \in \Theta} l(\theta; y_{1:N}^*) = \theta_0 + o_p^*(1)$$

$$\arg \max_{\theta \in \Theta_0} l(\theta; y_{1:N}^*) = \theta_0 + o_p^*(1).$$

We can now state the main result of this article that guarantees the consistency of the bootstrap procedure.

**Theorem 1** *Under assumptions (1)–(3), if  $\theta_N^*$  is chosen such that  $\theta_N^* \in \Theta_0$ ,  $\theta_N^* = \theta_0 + o_p(1)$  and  $N^{1/4}\delta_N^* = o_p(1)$  then as  $N \rightarrow +\infty$ , it holds in probability that*

$$\text{pr}^* \{ \text{LRT}(y_{1:N}^*) \leq t \} \longrightarrow \text{pr}(\text{LRT}_\infty \leq t). \quad (5)$$

A way of choosing  $\theta_N^*$  that fulfills the hypothesis of theorem 1 is to follow the idea of [11] and shrinks the parameter toward 0. However here the rate of convergence of the shrinking parameter ( $c_N$ ) is not the same due to the singularity issue. The following lemma gives a procedure to choose  $\theta_N^*$  and justify the way it is chosen in algorithm 1.

**Proposition 5** *Let  $(c_N)_{N \in \mathbb{N}}$  be a sequence such that  $\lim_{N \rightarrow +\infty} c_N = 0$  and  $\lim_{N \rightarrow +\infty} N^{1/4} c_N = +\infty$ . Let  $\hat{\theta}_N = (\hat{\psi}_N, \hat{\delta}_N, \hat{\lambda}_N)$  be a maximum likelihood estimator (restricted or not) of  $\theta_0 = (\psi_0, 0_{d_\delta}, 0_{d_\lambda})$ . Under assumptions (1)–(3), by choosing  $\theta_N^* = (\psi_N^*, \delta_N^*, \lambda_N^*)$  such that:  $\forall k = 1, \dots, d_\psi$   $\psi_{N,k}^* = \hat{\psi}_{N,k} \mathbb{1}(\hat{\psi}_{N,k} > c_N)$ ,  $\forall k = 1, \dots, d_\delta$   $\delta_{N,k}^* = \hat{\delta}_{N,k} \mathbb{1}(\hat{\delta}_{N,k} > c_N)$  and  $\lambda_N^* = 0_{d_\lambda}$  then,  $\theta_N^*$  verifies the hypothesis of theorem 1.*

As we do not know which parameters are part of  $\delta$ , it is important to deal with every potential nuisance parameters. This is why we also consider a shrinkage bootstrap parameter for  $\psi$ . In the proof of this proposition we show that the shrinkage does not change the limit of the estimate, but only speeds up its convergence toward 0.

In addition to the boundary issue, the singularity is another source of inconsistency for the bootstrap procedure. As highlighted in the proof of theorem (1), this inconsistency comes

from the polluting random variables due to the asymptotic distribution of  $N^{\frac{1}{4}}\delta_N^*$ , that does not appear in the asymptotic distribution of  $LRT_\infty$ . The shrinkage enables to enforce that  $N^{\frac{1}{4}}\delta_N^* = o_p(1)$  and no longer  $O_p(1)$ .

### 3.3 Extension to the non identically distributed setting

As in the previous section we first derive the consistency of the maximum likelihood estimator. To do so, we need the regularity required in assumption (2) to hold uniformly over the different distributions of the individuals.

**Assumption 4** *We suppose that assumptions (2) i)–v) hold uniformly over the different individuals  $i \in \mathbb{N}$ .*

In addition to that, as discussed in [26] an additional assumption is required to ensure the unicity of the maximum of the asymptotic objective function.

**Assumption 5** *For every  $\theta \neq \theta_0$  :*

$$\lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \log \left\{ \frac{f_i(y_i; \theta)}{f_i(y_i; \theta_0)} \right\} \right] < 0$$

**Proposition 6** *Under assumptions (1), (4) and (5), propositions 2 and 4 still hold in the non identically distributed case.*

Following the same lines as in the last section the result of theorem 1 still holds.

**Theorem 2** *Under assumptions (1)–(5), if  $\theta_N^*$  is chosen such that  $\theta_N^* \in \Theta_0$ ,  $\theta_N^* = \theta_0 + o_p(1)$ , i.e.  $\lambda_N^* = 0$  and  $N^{\frac{1}{4}}\delta_N^* = o_p(1)$  then as  $N \rightarrow +\infty$ , it holds in probability that*

$$\text{pr}^* \{ \text{LRT}(y_{1:N}^*) \leq t \} - \text{pr}(LRT_N \leq t) = o_p(1). \quad (6)$$



### 3.4 Sufficient verifiable conditions for regularity assumptions

In this work we assume very strong regularity conditions on the model. Due to the integrated form of the likelihood in (2), it is often difficult to manipulate the quantities involved. [31] proposed some verifiable conditions under which the maximum likelihood estimators in non-linear mixed models is strongly consistent. However in his work he considered that the true parameter is an interior point of the parameter space, and that the Fisher information matrix is nonsingular. Furthermore in our context the conditions required are even more difficult to verify as we deal not only with maximum likelihood estimator consistency but also with likelihood ratio and bootstrap statistic consistency.

We propose an analytical sufficient criterion that only depends on the regularity of the known function  $g$  in model (1).

First we state a regularity condition on the derivatives of  $g$ .

**Assumption 6** For all  $k_1 = 0, \dots, 4$  and  $k_2 \in \mathbb{N}$

$$\sup_{i \in \mathbb{N}, j=1, \dots, J_i} \mathbb{E} \left\{ \sup_{\theta \in \Theta} \|\nabla_{\theta}^{k_1} g(x_{ij}, \beta, \Lambda \xi)\|^{k_2} \right\} < +\infty, \quad \xi \sim \mathcal{N}(0, I_p). \quad (7)$$

**Remark 7** This assumption seems very strong but in practice it only requires that the derivatives of  $g$  are not exponential in  $\|\xi\|^2$  which is verified for almost every commonly used models.

We now state the regularity condition on  $g$ , which is the proposed criterion to be verified case by case in real life applications.

**Assumption 7** for every  $\varepsilon > 0$ , there exists a compact set  $K \subset \mathbb{R}^p$  such that

$$\forall \xi \in \mathbb{R}^p \setminus K \quad \sup_{i \in \mathbb{N}, j=1, \dots, J_i} \sup_{\theta \in \Theta} \frac{\|g(x_{ij}, \beta, \Lambda \xi)\|}{\|\xi\|} \leq \varepsilon. \quad (8)$$

**Proposition 7** Suppose that assumption (1) holds, and that the function  $g$  verifies assumption (6)–(7), then assumption (4) is verified.

## 4 Experiments

### 4.1 Simulation study

We denote by  $\theta_0 = (\beta_0, \Lambda_0, \sigma_0^2)^T$  the true parameter used to generate the data. We use the notation  $\beta_k$  for the  $k$ th component of the fixed effects vector  $\beta$ , and we write  $\text{diag}(x_1, \dots, x_p)$  for a diagonal  $p \times p$  matrix, with a diagonal being equal to  $(x_1, \dots, x_p)^T$ . When it is not explicitly written we consider diagonal matrices  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)^T$ . The same way we write  $\xi_i = (\xi_{i1}, \dots, \xi_{ip})^T$  for the vector of random effects. We consider a linear and a nonlinear mixed effects models, with a varying number of random effects to account for the presence of nuisance parameters. Results were obtained using the `lme4` and `saemix` packages in R. Codes are available upon request from the first author.

We first consider the linear case. We denote by  $m_1$  the linear model with two independent random effects, i.e. with  $g(x_{ij}, \beta, \Lambda \xi_i) = \beta_1 + \lambda_1 \xi_{i1} + (\beta_2 + \lambda_2 \xi_{i2}) x_{ij}$ . We set  $\beta_0 = (0, 7)^T$ ,  $\lambda_{01}^2 = 1.3, \lambda_{02} = 0$ . In this model, we consider the test  $H_0 : \lambda_2 = 0$  against  $H_1 : \lambda_2 \geq 0$ . We then denote by  $m_2$  the linear model with three independent random effects, i.e. with  $g(x_{ij}, \beta, \Lambda \xi_i) = \beta_1 + \lambda_1 \xi_{i1} + (\beta_2 + \lambda_2 \xi_{i2}) x_{ij} + (\beta_3 + \lambda_3 \xi_{i3}) x_{ij}^2$ . We set  $\beta_0 = (0, 7, 3)^T$ ,  $\lambda_{01}^2 = 1.3, \lambda_{02} = 0, \lambda_{03} = 0$ . In this model, we consider the test  $H_0 : \lambda_3 = 0$  against  $H_1 : \lambda_3 \geq 0$ , so that in this simulation  $\lambda_2$  is a nuisance parameter. For the choice of the shrinkage bootstrap parameter, we set  $c_N = 0.28 \approx 0.5 \times 20^{-0.2}$ . This choice is motivated by its theoretical definition, this parameter shrinks to zero the variances of the individual parameters  $\beta_2 + \lambda_2 \xi_{i2}$  with a relative standard deviation lower than 4%. In both settings, we set  $x_{ij} = j$ ,  $J = 5$  and  $\sigma_0 = 1.5$ . Finally, we denote by  $m_3$  the linear model with  $p = 8$  random effects and a varying number  $s$  of nuisance parameters, i.e. with  $g(x_{ij}, \beta, \Lambda \xi_i) = \sum_{k=1}^p x_{ijk} \lambda_k \xi_{ik}$ . Here we set  $N = 40$ ,  $J_i = 9$ ,  $\sigma^2 = 2$ , and every untested variance to 1. Finally we draw independently the covariates from a normal distribution with mean 2 and standard deviation 0.5. We want here to illustrate the effect of an increasing number of nuisance parameters on the performance of the test. We use three different values for the shrinkage parameter  $c_N \in \{0; 0.24; 0.9\}$ . We chose those values to consider three cases: first  $c_N = 0$  is equivalent to

the parametric bootstrap procedure without shrinkage, then  $c_N = 0.5 \times 40^{-0.2} \approx 0.24$  shrinks most of the nuisance parameters toward 0, and finally  $c_N = 0.9$  shrinks systematically the nuisance parameters (as if we were using the true model), but can also shrink some non-zero variances of the model. We consider the test  $H_0 : \lambda_1 = 0$  against  $H_1 : \lambda_1 \geq 0$ .

Next, we consider the nonlinear logistic model with three random effects denoted by  $m_4$ , where

$$g(x_{ij}, \beta, \Lambda \xi_i) = \frac{\beta_1 + \lambda_1 \xi_{i1}}{1 + \exp\left(-\frac{x_{ij} - (\beta_2 + \lambda_2 \xi_{i2})}{\beta_3 + \lambda_3 \xi_{i3}}\right)}. \quad (9)$$

We set  $\beta_0 = (200, 500, 150)^T$ ,  $\lambda_{01} = \lambda_{02} = 10$ ,  $\lambda_{03} = 0$  and  $\sigma_0^2 = 5^2$ . We set  $(x_{i1}, \dots, x_{iJ}) = (50, 287.5, 525, 762, 1000, 1100, 1200, 1300, 1400, 1500)$  for all  $i$ . In this model, we consider the test  $H_0 : \lambda_3 = 0$  against  $H_1 : \lambda_3 \geq 0$ .

First, we study the finite sample size properties of our procedure using models  $m_1$ ,  $m_2$  and  $m_4$ . We compute the empirical levels by generating  $K$  datasets under the null hypothesis as described in the previous paragraph, and by computing the proportion of these datasets for which we reject the null hypothesis, for a nominal level  $\alpha$  in  $\{0.01, 0.05, 0.10\}$  and a sample size  $N$  in  $\{10, 20, 30, 40, 100\}$  for  $m_1$ ,  $N$  in  $\{20, 30, 40\}$  for  $m_2$  and  $N = 40$  for  $m_4$ . Results are given in tables 1, 2 and 3. We compare the empirical level of the test associated with our bootstrap procedure with those obtained using the asymptotic distribution which is a 0.5 – 0.5 mixture between a Dirac distribution at zero and a chi-squared distribution with one degree of freedom [4]. We observe that the empirical levels obtained with our bootstrap procedure are closer to the nominal ones than those obtained with the asymptotic procedure, and that good results are already obtained for small values of  $N$  in the linear case. As expected, our procedure exhibits better small sample size properties than the asymptotic procedure, both in the linear and the nonlinear cases. It is noteworthy to mention that the existing non-asymptotic test procedures such as the one proposed by [17] can not be used in the latter case since they rely on explicit expressions for the parameter estimates, hence requiring the linearity assumption. We also observe that the presence of nuisance parameters also deteriorate the asymptotic results. It is not a surprise as it modifies the true asymptotic distribution. However we observe that the standard parametric bootstrap procedure is robust

Level $\alpha$	$N = 10$		$N = 20$		$N = 30$		$N = 40$		$N = 100$		max sd
	boot	asym	boot	asym	boot	asym	boot	asym	boot	asym	
1%	1.14	0.68	0.98	0.68	1.20	0.94	0.74	0.70	0.86	0.72	0.15
5%	5.20	3.64	5.22	3.82	5.74	4.30	4.86	3.94	5.26	4.50	0.33
10%	10.72	7.16	10.80	7.98	10.30	8.40	10.80	8.44	10.34	8.86	0.44

Table 1: Empirical levels (expressed as percentages) of the test that one variance is null in a linear model with two independent random effects, for  $K = 5000$  simulated datasets and  $B = 500$  bootstrap replicates. The last column gives the maximal standard deviation value obtained in each row

boot., parametric bootstrap procedure; asym., asymptotic procedure; sd, standard deviation .

Level $\alpha$	$N = 20$		$N = 30$		$N = 40$		max sd
	boot	asym	boot	asym	boot	asym	
1%	0.82	0.66	0.72	0.58	0.90	0.62	0.13
5%	4.46	3.54	3.96	3.28	4.14	3.08	0.29
10%	8.88	6.78	7.52	6.34	8.40	6.98	0.40

Table 2: Empirical levels (expressed as percentages) of the test that one variance is null in a linear model with three random effects and one nuisance parameter. The last column gives the maximal standard deviation value obtained in each row

boot., parametric bootstrap procedure; asym., asymptotic procedure; sd, standard deviation.

to the presence of a single nuisance parameter, and so the choice of  $c_N$  does not have a significant effect on this example. This must be due to the low number of nuisance parameters and the few number of parameters of the model.

We then study the empirical power of our procedure using models  $m_1$  and  $m_2$ . To this end, we consider a non diagonal matrix  $\Lambda$ , introducing a correlation between the components of the scaled random effects  $b_i$ . We denote by  $\rho_{kl}$  the correlation coefficient between the scaled random effects  $b_{ik}$  and  $b_{il}$ . We then consider increasing values of  $\lambda_2$  and  $\rho_{12}$  in  $m_1$ , and increasing values of  $\lambda_3$  and  $\rho_{13}$  in  $m_2$ . Results are given in figure 1.

As expected, we observe that, for fixed values of the correlation coefficient, the empirical power increases when the true value of the tested variance increases, and that, for fixed values of the variance, the power increases when the correlation coefficient increases. In  $m_1$ , since  $\beta_2 = 7$ , we obtain an empirical power of at least 70% for a relative standard deviation of 4.5% (i.e. when  $\lambda_2^2 = 0.1$ ). In  $m_2$ , since  $\beta_3 = 3$ , the empirical power is greater than 12.5% for a relative standard deviation of 4.7% (i.e. when  $\lambda_3^2 = 0.02$ ), and above 90% for a relative

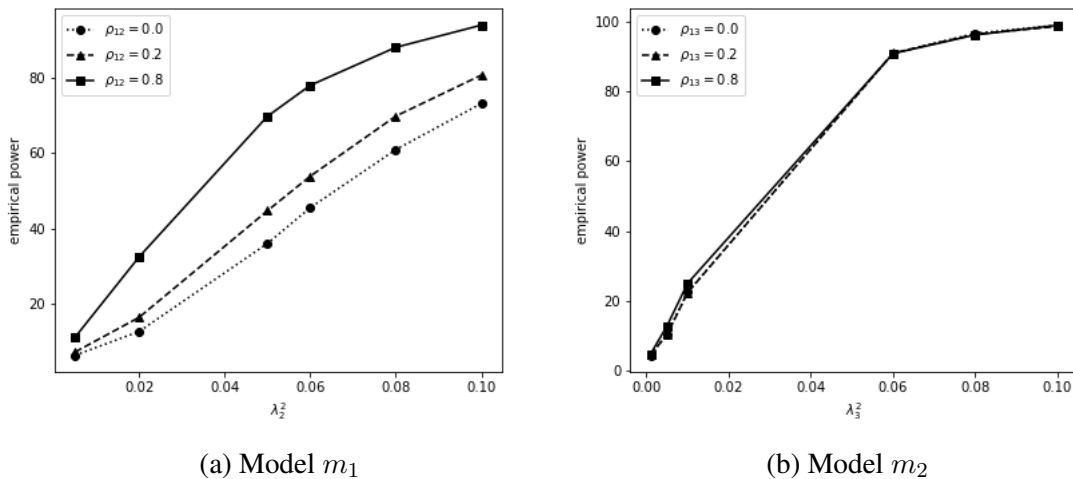


Figure 1: Empirical power of the test that one variance is null in a linear model with (a) two random effects and (b) three random effects, for varying value of the tested variance and of the correlation coefficient.

standard deviation of 10% (i.e. when  $\lambda_3^2 = 0.1$ ).

We then study the effect of shrinkage on the type I error using model  $m_3$ . Results are presented in figure 2 for a theoretical level of 5%. We see that the procedure is sensitive to extreme values of  $c_N$ . The performances of the shrunked bootstrap procedure are stable as the number of nuisance parameters increases, provided that the shrinkage parameter  $c_N$  is carefully chosen, whereas the performances of the regular bootstrap procedure with no shrinkage are downgraded in this context. Indeed choosing a value of  $c_N \approx 0.24$ , i.e. that shrinks most of the nuisance parameters, provides good results while choosing  $c_N = 0.9$ , i.e. of the same order of magnitude as the non-zero variances of the model (here,  $\lambda = 1$ ) deteriorates the results. On the other hand, neglecting the nuisance parameters, which corresponds to the case  $c_N = 0$ , also has an influence on the results and leads to an empirical level which is smaller than the theoretical one. These results suggest that in practice, users should be cautious when choosing the shrinkage parameter, and that smaller values of  $c_N$  could be preferred over too large values since the impact on the type I error seems higher in the latter case.

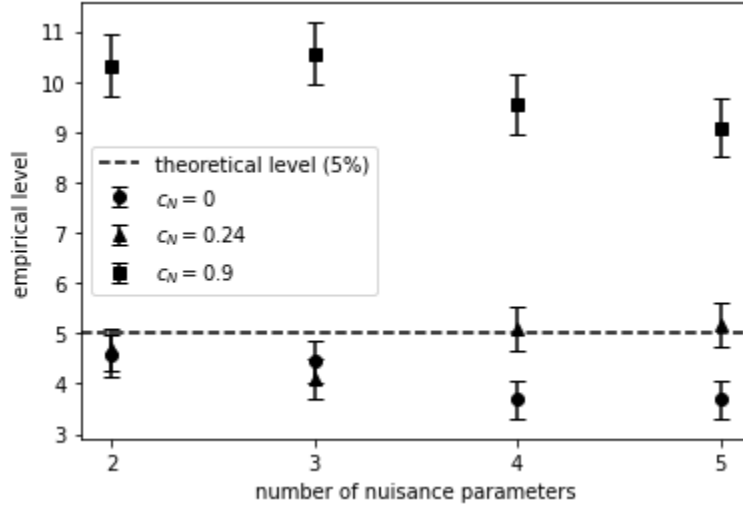


Figure 2: Comparison of the parametric bootstrap procedure and the shrunk parametric bootstrap procedure in  $m_3$ , using  $K = 2500$  datasets of size  $N = 30$  and  $B = 300$  bootstrap replicates for varying numbers of nuisance parameters, and different shrinkage parameters  $c_N$ .

Level $\alpha$	boot	asym	max sd
1%	0.80	0.80	0.28
5%	5.10	3.60	0.70
10%	10.30	7.00	0.96

Table 3: Comparison of the bootstrap procedure and the asymptotic procedure in the test in  $m_4$ , using  $K = 1000$  datasets of size  $N = 40$  and  $B = 300$  bootstrap replicates.

boot., parametric bootstrap procedure; asym., asymptotic procedure; sd, standard deviation.

## 4.2 Real data application

We then apply our procedure to a study of white-browed coucal growth rates, available as a Dryad package [21]. We use the logistic growth model defined in (9) to describe the evolution of the body mass of the nestlings as a function of their age. More precisely, if we denote by  $y_{ij}$  the body mass of nestling  $i$  at age  $t_j$  for  $i = 1, \dots, 292$ ,  $j = 1, \dots, N_i$ , we have that  $\beta_1 + \lambda_1 \xi_{i1}$  is the asymptotic nestling body mass,  $\beta_2 + \lambda_2 \xi_{i2}$  is the age (in days) at which nestling  $i$  reaches half its asymptotic body mass and  $\beta_3 + \lambda_3 \xi_{i3}$  is the growth rate of nestling  $i$ . When fitting the complete model the estimated scaling matrix is  $\hat{\Lambda} = \text{diag}(\sqrt{212.34}, \sqrt{0.89}, \sqrt{0.02})$ , which motivates the test that  $\lambda_2$  and  $\lambda_3$  are null.

In order to test for the presence of randomness in the inflexion point and in the growth

rate, we proceed sequentially. First, we test if the variances of both the inflexion point and the growth rate are null, i.e. we consider the test  $T_1 : H_0 : \lambda_2 = 0, \lambda_3 = 0$  against  $H_1 : \lambda_2 \geq 0, \lambda_3 \geq 0$ . If the null hypothesis in  $T_1$  is rejected, we perform two univariate tests, one for each variance tested in  $T_1$ . More precisely, we consider the tests  $T_2 : H_0 : \lambda_3 = 0$  against  $H_1 : \lambda_3 \geq 0$  and  $T_3 : H_0 : \lambda_2 = 0$  against  $H_1 : \lambda_2 \geq 0$ . If the null hypothesis is rejected in  $T_1$ , it means that at least one of the two tested variances is nonzero, thus we can then consider the univariate tests  $T_2$  and  $T_3$ , where  $\lambda_2$  and  $\lambda_3$  might be nuisance parameters. For each test, we consider two procedures, one where the estimate of the potential nuisance parameter ( $\lambda_2$  in  $T_2$  and  $\lambda_3$  in  $T_3$ ) is shrunk toward zero, and another one without shrinkage. This choice enables to consider the two possible cases :  $\lambda_2 > c_N$  and  $\lambda_2 < c_N$  for  $T_2$  and  $\lambda_3 > c_N$  and  $\lambda_3 < c_N$  for  $T_3$ . In practice the practitioner can choose the threshold according to his own level of significance of variability desired. For instance by saying that under X% of relative variability of a parameter, we consider it as a fixed effect. Table 4 compiles the results of the three tests.

Test	T1		T2		T3	
LRT	302.5		1.6		278.2	
procedure	no shrink	shrink	no shrink	shrink	no shrink	shrink
<i>p</i> -value	0	8.9	7.9	0	0	0

Table 4: Comparison of the *p*-value (in %) among the three tests  $T_1$ ,  $T_2$  and  $T_3$ , using  $B = 1000$  bootstrap replicates..

shrink., shrunked parametric bootstrap procedure; no shrink., regular parametric bootstrap procedure without shrinkage

We can see that the procedures lead to different *p*-values, and can thus, in practice, lead to different conclusions with respect to the null hypothesis depending on the type I error considered.

## 5 Discussion

This work can lead to several future developments, both from a theoretical and a practical point of view. For example, the choice of the shrinking parameter  $c_N$  is of interest for practi-

tioners that would want to apply our procedure. This issue can be tackled from two perspectives. From an expert-based point of view,  $c_N$  could be defined as a threshold, under which the variability of a random effect is not relevant in the task of interest. From a methodological point of view, it would be interesting to propose an automated procedure, following for example the idea of [6]. Our algorithm can be used indifferently for linear and nonlinear mixed effects models, however the computation time can be prohibitive in the latter case, especially for complex models. It would be interesting to find a criterion to optimize the choice of the bootstrap sample size. Another issue when dealing with complex models is the computation of the likelihood ratio test statistic. Indeed, it involves a ratio of likelihoods which are usually estimated separately using Monte Carlo approaches, leading to a biased estimate. This point is crucial since this quantity is calculated at each iteration of the bootstrap procedure.

## 6 Acknowledgment

This work was funded by the Stat4Plant project ANR-20-CE45-0012.

## 7 Supplementary material

We recall that we consider the following nonlinear mixed effects model, for any  $i = 1, \dots, N$  and any  $j = 1, \dots, J_i$ :

$$\begin{cases} y_{ij} = g(x_{ij}, \beta, \Lambda \xi_i) + \varepsilon_{ij} & \varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2) \\ \xi_i \sim \mathcal{N}(0, I_p) \end{cases} \quad (10)$$

We also recall that the log-likelihood of the model given the data  $y_{1:N}$  writes:

$$l(\theta; y_{1:N}) = \log\{L_\theta(y_{1:N})\} = \log\left\{\prod_{i=1}^N f_i(y_i; \theta)\right\} = \sum_{i=1}^N \log\left\{\int f_i(y_i; \xi_i, \theta) \pi(\xi_i) d\xi_i\right\}. \quad (11)$$

Finally we write

$$\hat{\theta}_N = \arg \sup_{\theta \in \Theta} l(\theta; y_{1:N}). \quad (12)$$



## 7.1 Different parametrizations for mixed effects models

The commonly used parametrization for nonlinear mixed effects models ([32] page 306) is for  $i = 1, \dots, N, j = 1, \dots, J_i$  :

$$\begin{cases} y_{ij} = g(v_{ij}, \phi_i) + \varepsilon_{ij} & \varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2) \\ \phi_i = A_{ij}\beta + B_{ij}b_i & b_i \sim \mathcal{N}(0, \Gamma) \end{cases} \quad (13)$$

where  $v_{ij}$ ,  $A_{ij}$  and  $B_{ij}$  are known covariates. With this definition, the log-likelihood is defined as follows :

$$l(\theta; y_{1:N}) = \log\{L_\theta(y_{1:N})\} = \log\left\{\prod_{i=1}^N f_i(y_i; \theta)\right\} = \sum_{i=1}^N \log\left\{\int f_i(y_i; b_i, \theta)\pi(b_i; 0, \Gamma)db_i\right\}, \quad (14)$$

where  $f_i(y_i; b_i, \theta)$  is the density of the conditional distribution of  $y_i$  given  $b_i$ , and  $\pi(b_i; 0, \Gamma)$  is the density of the  $p$ -dimensional centered Gaussian distribution with covariance  $\Gamma$ .

The change of variable  $b_i = \Lambda\xi_i$  is a  $\mathcal{C}^1$  diffeomorphism if and only if the diagonal coefficients of  $\Lambda$  are strictly nonnegative. Therefore in our setting where this condition is not verified these two parametrizations are not equivalent.

In particular taking  $\Lambda = 0$  in (2) is equivalent to considering a fixed-effects nonlinear model, while in (14), taking  $\Gamma \rightarrow 0$  makes  $l(\theta; y_{1:N}) \rightarrow -\infty$ .

## 7.2 Proof of proposition 1

To prove the proposition we only have to prove that if the  $m$ th column of  $\Lambda$  is 0, then the odds orders of the partial derivatives of the log-likelihood with respect to the elements of this column are null regardless of the values of the  $y_{1:N}$ . We have, for all  $n = 1, \dots, p$ :

$$\frac{\partial f_i(y_i; \theta)}{\partial[\Lambda]_{nm}} \Big|_{\theta=\theta_0} = \int_{\xi_1} \dots \int_{\xi_p} \frac{\partial f_i(y_i; \xi, \theta_0)}{\partial[\Lambda]_{nm}} \pi(\xi_1)d\xi_1 \dots \pi(\xi_p)d\xi_p$$

Given the definition of model 1 :

$$f_i(y_i; \xi, \theta) = (2\pi\sigma^2)^{-\frac{J_i}{2}} \exp \left[ -\frac{\sum_{j=1}^{J_i} \{y_{ij} - g(x_{ij}, \beta, \Lambda\xi)\}^2}{2\sigma^2} \right]$$

therefore,

$$\begin{aligned} \frac{\partial f_i(y_i; \xi, \theta)}{\partial [\Lambda]_{nm}} &\propto \frac{\partial \Lambda \xi}{\partial [\Lambda]_{mn}} \nabla_{\Lambda \xi} \exp \left[ -\frac{\sum_{j=1}^{J_i} \{y_{ij} - g(x_{ij}, \beta, \Lambda\xi)\}^2}{2\sigma^2} \right] \\ &\propto \xi_m \nabla_{\Lambda \xi} \exp \left[ -\frac{\sum_{j=1}^{J_i} \{y_{ij} - g(x_{ij}, \beta, \Lambda\xi)\}^2}{2\sigma^2} \right] \end{aligned}$$

Evaluated at  $\theta = \theta_0$  the last term (the gradient) no longer depends on  $\xi_m$  as the  $m$ th column of  $\Lambda$  is null.

$$\begin{aligned} \frac{\partial f_i(y_i; \theta)}{\partial [\Lambda]_{nm}} \Big|_{\theta=\theta_0} &\propto \int_{\xi_m} \xi_m \pi(\xi_m) d\xi_m \\ &\times \int \nabla_{\Lambda \xi} \exp \left[ -\frac{\sum_{j=1}^{J_i} \{y_{ij} - g(x_{ij}, \beta, \Lambda\xi)\}^2}{\sigma^2} \right] \Big|_{\theta=\theta_0} \prod_{l \neq m} \pi(\xi_l) d\xi_l \end{aligned}$$

which is equal to 0 as the first term of the right hand side equation is expectation of a standard gaussian distribution. We can apply the same reasoning to every odds order derivatives.

### 7.3 Proof of proposition 2

Due to assumption (1) we have that :

$$\sup_{\theta \in \Theta} \mathbb{E} \{l(\theta; y_1)\} < \mathbb{E} \{l(\theta_0; y_1)\} \quad (15)$$

which comes from the identifiability of the model and the positivity of the Kullback-Leibler divergence. Assumption (2)i) enables to apply the uniform law of large number to the log-likelihood. Then the result follows from arguments as in [1] lemma A.1.

## 7.4 Proof of proposition 4

To prove the consistency of the bootstrap maximum likelihood estimator, we will use the same reasoning, as in the proof of proposition 2. The sketch of the proof is similar to the one of [11].

We first want to show that :

$$\sup_{\theta \in \Theta} \left| \frac{1}{N} l(\theta; y_{1:N}^*) - E \{l(\theta; y_1)\} \right| = o_{p^*}(1)$$

First of all we have that :

$$\sup_{\theta \in \Theta} \left| \frac{1}{N} l(\theta; y_{1:N}^*) - E \{l(\theta; y_1)\} \right| \leq \sup_{\theta \in \Theta} A_N^*(\theta) + \sup_{\theta \in \Theta} A_N(\theta)$$

where :

$$\begin{aligned} A_N^*(\theta) &= \left| \frac{1}{N} l(\theta; y_{1:N}^*) - E^* \{l(\theta; y_1^*)\} \right| \\ A_N(\theta) &= |E^* \{l(\theta; y_1^*)\} - E \{l(\theta; y_1)\}| \end{aligned}$$

We now want to apply the uniform law of large numbers to  $A_N(\theta)$ , and it's bootstrap version to  $A_N^*(\theta)$ . Therefore we shall show that both term converges toward 0 and that they are lipshitz.

We first consider  $A_N(\theta)$  for a given  $\theta \in \Theta$ .

$$\begin{aligned} E^* \{l(\theta; y_1^*)\} &= E \{l(\theta; y_1^*) | y_{1:N}\} \\ &= \int l(\theta; y_1) f(y_1; \theta_N^*) dy_1 \end{aligned}$$

Using assumption (1), we can state that there exist  $\theta^+$  between  $\theta_0$  and  $\theta_N^*$  such that,

$$\begin{aligned}
|f(y_1; \theta_N^*) - f(y_1; \theta_0)| &\leq \|\theta_0 - \theta_N^*\| \|\nabla_{\theta} f(y_1; \theta^+)\| \\
&\leq \|\theta_0 - \theta_N^*\| \|\nabla_{\theta} \log\{f(y_1; \theta^+)\}\| f(y_1; \theta^+)
\end{aligned}$$

therefore,

$$\begin{aligned}
A_N(\theta) &\leq \int |l(\theta; y_1)| |f(y_1; \theta_N^*) - f(y_1; \theta_0)| dy_1 \\
&\leq \int |l(\theta; y_1)| \|\theta_0 - \theta_N^*\| \|\nabla_{\theta} \log\{f(y_1; \theta^+)\}\| f(y_1; \theta^+) dy_1 \\
&\leq \|\theta_0 - \theta_N^*\| \int |l(\theta; y_1)| \|\nabla_{\theta} l(\theta^+; y_1)\| f(y_1; \theta^+) dy_1
\end{aligned}$$

Due to assumption (2)*i*–*ii*),  $|l(\theta; y_1)| \|\nabla_{\theta} l(\theta^+; y_1)\|$  is integrable with respect to the density  $f(y_1; \theta^+)$ . And finally using the elementary inequality ,

$$\frac{1}{2}(a^2 + b^2) \geq |ab|, \quad \forall a, b \in \mathbb{R} \tag{16}$$

we can state that

$$A_N(\theta) \leq \frac{1}{2} \|\theta_0 - \theta_N^*\| \sup_{\theta^+ \in \Theta} \int \sup_{\Theta \in \Theta} |l(\Theta; y_1)|^2 + \sup_{\theta_2 \in \Theta} \|\nabla_{\theta} l(\theta_2; y_1)\|^2 f(y_1; \theta^+) dy_1$$

Finally thanks to assumption (2), as  $N \rightarrow +\infty$ , it holds in probability that :

$$A_N(\theta) \rightarrow 0$$

We now consider :

$$A_N^*(\theta) = \left| \frac{1}{N} \sum_{i=1}^N l(\theta; y_i^*) - \mathbb{E}^* \{l(\theta; y_1^*)\} \right|$$

This quantity is a sum of conditionally independent and centered random variables. We

can't directly apply a law of large number as the parameter  $\theta_N^*$  and the index of the sum depends both on  $N$ .

For every real nonnegative number  $t$ , it holds almost surely that :

$$\text{pr}^*(A_N^* > t) \leq \text{pr}^*\left\{\frac{1}{N}\sum_i |l(y_i^*; \theta) - \text{E}^*\{l(\theta; y_1^*)\}| > t\right\} \leq \frac{\sup_{\theta' \in \Theta} \text{E}_{\theta'}\{\sup_{\theta \in \Theta} |l(\theta; y_1)|^2\}}{Nt^2}$$

by applying first triangular inequality and then Chebychev inequality, using assumption (2)i). And finally  $A_N^*(\theta) \rightarrow 0$  in probability, as  $N \rightarrow +\infty$ , which concludes the pointwise convergence. And we note that this result holds uniformly over  $\Theta$  so:

$$\sup_{\theta \in \Theta} \left| \frac{1}{N} l(\theta; y_{1:N}^*) - \text{E}\{l(\theta; y_1)\} \right| = o_p^*(1)$$

Let now use this result to show that the bootstrap maximum likelihood estimator is consistent.

Let  $\varepsilon > 0$ , using equation (15), there exists  $\delta > 0$  such that:

$$\inf_{\|\theta - \theta_0\| > \varepsilon} \text{E}\{l(\theta_0; y_1)\} - \text{E}\{l(\theta; y_1)\} \geq \delta$$

Let us introduce  $\theta_{mle}^B = \arg \max_{\theta \in \Theta} l(\theta; y_{1:N}^*)$ .

By writing  $V_\varepsilon = \{\theta \in \Theta : \|\theta - \theta_0\| > \varepsilon\}$ , we have that :

$$\begin{aligned} \text{pr}^*(\theta_{mle}^B \in V_\varepsilon) &\leq \text{pr}^*\left[\text{E}\{l(\theta_0; y_1)\} - \text{E}\{l(\theta_{mle}^B; y_1^*)\} \geq \delta\right] \\ &= \text{pr}^*\left[\text{E}\{l(\theta_0; y_1)\} - \frac{1}{N}l(\theta_{mle}^B; y_{1:N}^*) + \frac{1}{N}l(\theta_{mle}^B; y_{1:N}^*) - \text{E}\{l(\theta_{mle}^B; y_1^*)\} \geq \delta\right] \\ &\leq \text{pr}^*\left[\text{E}\{l(\theta_0; y_1)\} - \frac{1}{N}l(\theta_0; y_{1:N}^*) + \frac{1}{N}l(\theta_{mle}^B; y_{1:N}^*) - \text{E}\{l(\theta_{mle}^B; y_1^*)\} \geq \delta\right] \\ &\leq \text{pr}^*\left\{2 \sup_{\theta \in \Theta} \left| \frac{1}{N} l(\theta; y_{1:N}^*) - \text{E}\{l(\theta; y_1)\} \right| \geq \delta\right\} \\ &\leq o_p(1) \end{aligned}$$

Which concludes the proof that  $\theta_{mle}^B = \theta_0 + o_p^*(1)$ . The exact same proof still holds for the restricted bootstrap maximum likelihood estimator by replacing  $\Theta$  by  $\Theta_0$ .

## 7.5 Quadratic approximation of the log-likelihood

To derive the asymptotic distribution of  $LRT_N$ , we expand the log-likelihood around  $\theta_0$  (see [2] theorem 6). Under assumption (1), following the lines of [25], we can write:

$$\begin{aligned}
l(\theta; y_{1:N}) - l(\theta_0; y_{1:N}) &= (\psi - \psi_0)^T \nabla_{\psi} l(\theta_0; y_{1:N}) + \frac{1}{2} (\psi - \psi_0)^T \nabla_{\psi}^2 l(\theta_0; y_{1:N}) (\psi - \psi_0) \\
&+ (1/2) \sum_{i,j=1,\dots,d_{\delta}} \delta_i \delta_j \frac{\partial^2 l(\theta_0; y_{1:N})}{\partial \delta_i \partial \delta_j} + (3/3!) (\psi - \psi_0)^T \sum_{i,j=1,\dots,d_{\delta}} \delta_i \delta_j \frac{\partial^3 l(\theta_0; y_{1:N})}{\partial \delta_i \partial \delta_j \partial \psi} \\
&+ (1/2) \sum_{i,j=1,\dots,d_{\lambda}} \lambda_i \lambda_j \frac{\partial^2 l(\theta_0; y_{1:N})}{\partial \lambda_i \partial \lambda_j} + (3/3!) (\psi - \psi_0)^T \sum_{i,j=1,\dots,d_{\lambda}} \lambda_i \lambda_j \frac{\partial^3 l(\theta_0; y_{1:N})}{\partial \lambda_i \partial \lambda_j \partial \psi} \\
&+ (6/4!) \sum_{i,j=1,\dots,d_{\delta}} \sum_{k,l=1,\dots,d_{\lambda}} \delta_i \delta_j \lambda_k \lambda_l \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \delta_i \partial \delta_j \partial \lambda_k \partial \lambda_l} \\
&+ (1/4!) \sum_{i,j,k,l=1,\dots,d_{\delta}} \delta_i \delta_j \delta_k \delta_l \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \delta_i \partial \delta_j \partial \delta_k \partial \delta_l} \\
&+ (1/4!) \sum_{i,j,k,l=1,\dots,d_{\lambda}} \lambda_i \lambda_j \lambda_k \lambda_l \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \lambda_i \partial \lambda_j \partial \lambda_k \partial \lambda_l} + R_N(\theta)
\end{aligned}$$

With  $R_N(\theta)$  being the rest in the Taylor expansion. We define for all integers  $i, j$   $c_{ij} = \frac{1}{2}$  if  $i = j$  and 1 otherwise. We also define  $\mathcal{I}_{\lambda} = \{11, 22, \dots, d_{\lambda} - 1, d_{\lambda}\}$  and  $\mathcal{I}_{\delta} = \{11, 22, \dots, d_{\delta} - 1, d_{\delta}\}$  that respectively index  $v(\lambda) = (\lambda_i \lambda_j)_{ij \in \mathcal{I}_{\lambda}}$  and  $v(\delta) = (\delta_i \delta_j)_{ij \in \mathcal{I}_{\delta}}$ . For  $ij \in \mathcal{I}_{\lambda}$  (respectively  $\mathcal{I}_{\delta}$ ), we write  $\frac{\partial^2 l(\theta; y_{1:N})}{\partial v(\lambda)_{ij}} = \frac{\partial^2 l(\theta; y_{1:N})}{\partial \lambda_i \partial \lambda_j}$  (the same with  $\delta$ ). With these notations we have that:

$$\begin{aligned}
l(\theta; y_{1:N}) - l(\theta_0; y_{1:N}) &= (\psi - \psi_0)^T \nabla_\psi l(\theta_0; y_{1:N}) + \frac{1}{2} (\psi - \psi_0)^T \nabla_\psi^2 l(\theta_0; y_{1:N}) (\psi - \psi_0) \\
&+ \sum_{i \in \mathcal{I}_\lambda} v(\lambda)_i c_i \frac{\partial^2 l(\theta; y_{1:N})}{\partial v(\lambda)_i} + (1/3) (\psi - \psi_0)^T \sum_{i \in \mathcal{I}_\lambda} v(\lambda)_i c_i \frac{\partial^3 l(\theta_0; y_{1:N})}{\partial v(\lambda)_i \partial \psi} \\
&+ \sum_{i \in \mathcal{I}_\delta} v(\delta)_i c_i \frac{\partial^2 l(\theta; y_{1:N})}{\partial v(\delta)_i} + (1/3) (\psi - \psi_0)^T \sum_{i \in \mathcal{I}_\delta} v(\delta)_i c_i \frac{\partial^3 l(\theta_0; y_{1:N})}{\partial v(\delta)_i \partial \psi} \\
&+ 4 \times (6/4!) \sum_{i \in \mathcal{I}_\lambda, j \in \mathcal{I}_\delta} c_i c_j v(\lambda)_i v(\delta)_j \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial v(\lambda)_i \partial v(\delta)_j} \\
&+ (4/4!) \sum_{i \in \mathcal{I}_\delta, j \in \mathcal{I}_\delta} c_i c_j v(\delta)_i v(\delta)_j \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial v(\delta)_i \partial v(\delta)_j} \\
&+ (4/4!) \sum_{i \in \mathcal{I}_\lambda, j \in \mathcal{I}_\lambda} c_i c_j v(\lambda)_i v(\lambda)_j \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial v(\lambda)_i \partial v(\lambda)_j} + R_N(\theta)
\end{aligned}$$

We define the reparametrization of  $\theta$  :  $\phi(\theta) = (\psi, v(\delta), v(\lambda))$ .

By writing  $\tilde{\nabla}_{v(\lambda)} l(\theta; y_{1:N}) = \left( c_i \frac{\partial^2 l(\theta; y_{1:N})}{\partial v(\lambda)_i} \right)_{i \in \mathcal{I}_\lambda}$  (similarly for  $\tilde{\nabla}_{v(\delta)} l(\theta; y_{1:N})$ ) and  $\tilde{\nabla}_{v(\lambda)}^2 l(\theta; y_{1:N}) = (c_i c_j \frac{\partial^4 l(\theta; y_{1:N})}{\partial v(\lambda)_i \partial v(\lambda)_j})_{i,j=1,\dots,d_\lambda}$  (similarly for  $\tilde{\nabla}_{v(\delta)}^2 l(\theta; y_{1:N})$ ), we also define :

$$\tilde{S}_N(\theta_0) = \sqrt{N}^{-1} \left( \nabla_\psi l(\theta_0; y_{1:N})^T, \tilde{\nabla}_{v(\delta)} l(\theta_0; y_{1:N})^T, \tilde{\nabla}_{v(\lambda)} l(\theta_0; y_{1:N})^T \right)^T$$

$$\tilde{I}_N(\theta_0) = \begin{pmatrix} I_{N,\psi}(\theta_0) & I_{N,\psi,v(\delta)}(\theta_0) & I_{N,\psi,v(\lambda)}(\theta_0) \\ I_{N,\psi,v(\delta)}(\theta_0)^T & I_{N,v(\delta)}(\theta_0) & I_{N,v(\delta),v(\lambda)}(\theta_0) \\ I_{N,\psi,v(\lambda)}(\theta_0)^T & I_{N,v(\delta),v(\lambda)}(\theta_0)^T & I_{N,v(\lambda)}(\theta_0) \end{pmatrix}$$

Where  $I_{N,\psi}(\theta_0) = -\frac{1}{N} \nabla_\psi^2 l(\theta_0; y_{1:N})$ ,  $I_{N,\psi,v(\lambda)}(\theta_0) = \left( -\frac{c_i}{3N} \frac{\partial^3 l(\theta_0; y_{1:N})}{\partial \psi \partial v(\lambda)_i} \right)_{i \in \mathcal{I}_\lambda}$ ,  $I_{N,\psi,v(\delta)}(\theta_0) = \left( -\frac{c_i}{3N} \frac{\partial^3 l(\theta_0; y_{1:N})}{\partial \psi \partial v(\delta)_i} \right)_{i \in \mathcal{I}_\delta}$ ,  $I_{N,v(\lambda)}(\theta_0) = -\frac{1}{3N} \tilde{\nabla}_{v(\lambda)}^2 l(\theta_0; y_{1:N})$ ,  $I_{N,v(\delta)}(\theta_0) = -\frac{1}{3N} \tilde{\nabla}_{v(\delta)}^2 l(\theta_0; y_{1:N})$ ,  $I_{N,v(\lambda),v(\delta)}(\theta_0) = \left( -\frac{1}{N} c_i c_j \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial v(\lambda)_i \partial v(\delta)_j} \right)_{i \in \mathcal{I}_\lambda, j \in \mathcal{I}_\delta}$ .

Where the notation  $(a_i)_{i \in \mathcal{I}}$  stands for a  $\dim(a_i) \times \text{card}(\mathcal{I})$  matrix whose columns are  $a_i$  for  $i \in \mathcal{I}$ .

With these new notations we obtain the following quadratic approximation of the log-likelihood :

$$\begin{aligned} l(\theta; y_{1:N}) - l(\theta_0; y_{1:N}) &= \sqrt{N}(\phi(\theta) - \phi(\theta_0))^T \tilde{S}_N(\theta_0) \\ &\quad - \frac{1}{2} \sqrt{N}(\phi(\theta) - \phi(\theta_0))^T \tilde{I}_N(\theta_0) \sqrt{N}(\phi(\theta) - \phi(\theta_0)) + R_N(\theta) \end{aligned} \quad (17)$$

## 7.6 Asymptotic distribution of the likelihood ratio test statistic

We start from the expansion of the log-likelihood (17) derived in the last section, that we rewrite:

$$\begin{aligned} l(\theta; y_{1:N}) - l(\theta_0; y_{1:N}) &= \frac{1}{2} Z_N(\theta_0)^T \tilde{I}_N(\theta_0) Z_N(\theta_0) \\ &\quad - \frac{1}{2} (t_N(\theta) - Z_N(\theta_0))^T \tilde{I}_N(\theta_0) (t_N(\theta) - Z_N(\theta_0)) + R_N(\theta) \end{aligned}$$

where  $t_N(\theta) = \sqrt{N}(\phi(\theta) - \phi(\theta_0))$  and  $Z_N(\theta_0) = \tilde{I}_N(\theta_0)^{-1} \tilde{S}_N(\theta_0)$

**Remark 8** We consider the quantity  $\tilde{I}_N(\theta_0)^{-1}$  which implies that  $\tilde{I}_N(\theta_0)$  is non-singular which may not be always true. However under assumption (3), the probability that  $\tilde{I}_N(\theta_0)$  is non-singular tends to 1 as  $N \rightarrow +\infty$ .

The set a feasible values for  $t_N(\theta)$  is:

$$t_N(\Theta) = \{\sqrt{N}(\Theta_\psi - \psi_0)\} \times \{\sqrt{N}(v(\Theta_\lambda) - v(\lambda_0))\} \times \{(v(\Theta_\delta) - v(\delta_0))\} \quad (18)$$

$t_N(\Theta)$  does not depend on  $N$  ( it is a cartesian product whose terms are whether  $\mathbb{R}$  or  $[0, +\infty[$ ), and it is locally approximated by a cone (in fact it is a cone), we write it  $\mathcal{C}(\Theta)$ .



Therefore if we prove that  $R_N(\theta)$  is  $o_p(1)$  when evaluated at the maximum likelihood estimator, we can apply the result from [2]. For more details see [2] and paragraph 4.3.

In order to apply the theory of Andrews, one shall prove that  $\tilde{I}_N(\theta_0)$  converges in probability toward a nonnegative matrix  $\tilde{I}(\theta_0)$ , and that  $\tilde{S}_N(\theta_0)$  converges weakly to a random variable  $U(\theta_0)$ .

1) Let  $d_{\tilde{I}} \in \mathbb{N}$  such that  $\tilde{I}_N(\theta_0) \in \mathbb{R}^{d_{\tilde{I}} \times d_{\tilde{I}}}$ , let  $1 \leq m, n \leq d_{\tilde{I}}$ . We can write:

$$[\tilde{I}_N(\theta_0)]_{m,n} = \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)}(\theta_0)$$

where  $h_{m,n}^{(i)}(\theta_0)$  is of the form:

$$h_{m,n}^{(i)}(\theta_0) = c_{m,n} \frac{\partial^k \log f(y_i; \theta)}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} \Big|_{\theta=\theta_0}$$

with  $c_{m,n} \in \mathbb{R}$ ,  $k \in \{2, 4\}$ ,  $1 \leq i_1, \dots, i_k \leq d_\lambda + d_\psi + d_\delta$ .

With assumption (2) we can apply the law of large numbers to this empirical mean and therefore it holds in probability that:

$$\tilde{I}_N(\theta_0) \xrightarrow{N \rightarrow +\infty} [\mathbb{E}\{h_{m,n}^{(1)}(\theta_0)\}]_{1 \leq m, n \leq q_I} = \tilde{I}(\theta_0)$$

2) We now consider  $\tilde{S}_N(\theta_0)$ .

First the score is centered,

$$\mathbb{E}[\nabla_\psi \log f(x_i; \theta_0)] = 0$$

then, due to proposition (1),  $\forall m, n = 1, \dots, q_\lambda$ ,

$$\mathbb{E} \left[ \frac{\partial^2 \log f(y_i; \theta_0)}{\partial \lambda_m \partial \lambda_n} \right] = -\mathbb{E} \left[ \frac{\partial \log f(y_i; \theta_0)}{\partial \lambda_m} \frac{\partial \log f(y_i; \theta_0)}{\partial \lambda_n} \right] = 0$$

We apply the central limit theorem to  $\tilde{S}_N(\theta_0)$  which is a sum of independent and identically distributed centered random variables with finite variances, and therefore  $\tilde{S}_N(\theta_0)$  converges weakly to a random variable  $U(\theta_0)$

By doing the exact same quadratic approximation and development but considering the parameter space  $\Theta_0$ , combining (18)–(20), we use [2] to obtain :

$$LRT_\infty = \inf_{t \in \mathcal{C}(\Theta_0)} \|t - \tilde{I}(\theta_0)^{-1}U(\theta_0)\|_{\tilde{I}(\theta_0)} - \inf_{t \in \mathcal{C}(\Theta)} \|t - \tilde{I}(\theta_0)^{-1}U(\theta_0)\|_{\tilde{I}(\theta_0)} \quad (19)$$

Which leads to the expression of  $LRT_\infty$ .

## 7.7 Proof of proposition 3

To obtain an explicit form for  $R_N(\theta)$  we use the multivariate version of Taylor-Lagrange formula, which is for instance defined in [2] Theorem 6.

This way we have that  $R_N(\theta)$  is a sum of higher order derivatives with respect to  $\psi$  and the fourth crossed derivatives with respect to  $\lambda$ . Given assumption (2) all the derivatives of the log-likelihood are  $\mathcal{O}_p(N)$ , using Cauchy-Schwartz inequality and the fact that  $\|t_N(\theta)\|^2 = N(\|\psi - \psi_0\|^2 + \|v(\delta)\|^2 + \|v(\lambda)\|^2) = N(\|\psi - \psi_0\|^2 + \|\lambda\|^4 + \|\delta\|^4) \mathcal{O}(1)$  we have that:

$$\begin{aligned} |R_N(\theta)| &\leq \mathcal{O}_p(N)(\|\psi - \psi_0\|^3 + \|\psi - \psi_0\|^4 + \|\psi - \psi_0\|^2\|\lambda\|^2 + \|\psi - \psi_0\|^2\|\delta\|^2) \\ &\quad + \|\delta\|^4 \left| \sum_{m,n,o,p=1,\dots,d_\delta} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \delta_m \partial \delta_n \partial \delta_o \partial \delta_p} - \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \delta_m \partial \delta_n \partial \delta_o \partial \delta_p} \right| \\ &\quad + \|\lambda\|^4 \left| \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} - \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \right| \end{aligned}$$

and then:

$$\begin{aligned} |R_N(\theta)| &\leq \mathcal{O}_p(1)\|t_N(\theta)\|^2(o_p(1) + \frac{1}{N} \left| \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} - \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \right| \\ &\quad + \frac{1}{N} \left| \sum_{m,n,o,p=1,\dots,d_\delta} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \delta_m \partial \delta_n \partial \delta_o \partial \delta_p} - \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \delta_m \partial \delta_n \partial \delta_o \partial \delta_p} \right|) \end{aligned}$$

where  $\theta^+ = \theta_0 + t(\theta - \theta_0)$  for some  $0 < t < 1$ .

To show that the last two terms tend to zero we proceed as follows :

$$\begin{aligned} \frac{1}{N} \left| \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} - \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \right| = \\ \frac{1}{N} \left| \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} - \mathbb{E} \left[ \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_1)}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \right] \right| \\ + \mathbb{E} \left[ \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_{1:N})}{\partial \lambda_m \partial \lambda_1 \partial \lambda_o \partial \lambda_p} \right] - \mathbb{E} \left[ \frac{\partial^4 l(\theta_0; y_1)}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \right] \\ + \mathbb{E} \left[ \frac{\partial^4 l(\theta_0; y_1)}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \right] - \frac{\partial^4 l(\theta_0; y_{1:N})}{\partial \lambda_m \partial \lambda_n \partial \lambda_o \partial \lambda_p} \Big| \end{aligned}$$

We then apply triangular inequality to separate the 3 terms. The first and third terms are empirical means of centered random variables with bounded variances (assumption (2)v)). Therefore we can use each time Chebychev's inequality to obtain a weak law of large number and obtain the consistency toward 0. For the second term,

$$\left| \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta^+; y_1)}{\partial \lambda_m \partial \lambda_1 \partial \lambda_o \partial \lambda_p} - \sum_{m,n,o,p=1,\dots,d_\lambda} \frac{\partial^4 l(\theta; y_1)}{\partial \lambda_m \partial \lambda_1 \partial \lambda_o \partial \lambda_p} \right| \leq 2C \sup_{\theta \in \Theta} \|\nabla_{\theta}^4 l(\theta; y_1)\|$$

where  $C$  is the nonnegative constant that appears in the equivalence between the L1 and L2 norm. When we evaluate at  $\theta = \hat{\theta}_N$ ,  $\theta^+ = \theta_0 + t(\hat{\theta}_N - \theta_0) \rightarrow \theta_0$  in probability, as  $N \rightarrow +\infty$ . And by continuity of  $\nabla_{\theta}^4 l(\cdot; y_1)$ , and dominated convergence, the second term also converges toward 0. Finally we obtain :

$$|R_N(\hat{\theta}_N)| \leq o_p(1) \|t_N(\hat{\theta}_N)\|^2$$

$$\begin{aligned} \text{And then: } 0 &\leq l(\hat{\theta}_N; y_{1:N}) - l(\theta_0; y_{1:N}) \\ &\leq \|\tilde{S}_N(\theta_0)\| \|t_N(\hat{\theta}_N)\| - \frac{1}{2} \|t_N(\hat{\theta}_N)\|_{I_N(\theta_0)}^2 + o_p(\|t_N(\hat{\theta}_N)\|^2) \\ &\leq \|\tilde{S}_N(\theta_0)\| \|t_N(\hat{\theta}_N)\| - \frac{1}{2} (o_p(1) + a) \|t_N(\hat{\theta}_N)\|^2 \end{aligned}$$

where

$$a = \inf_{N > n_0} \inf_{x \neq 0} \frac{\|x\|_{\tilde{I}_N(\theta_0)}^2}{\|x\|^2}$$

where  $\|x\|_A$  stands for  $x^T A x$ , with  $A$  being a positive definite symmetric matrix.

By taking  $n_0$  large enough so that for every  $N > n_0$ ,  $\tilde{I}_N(\theta_0) \succ 0$  (assumption (3)) we have that  $0 < a < +\infty$ . The last inequality, shows that for  $N$  large enough, this polynomial of degree 2 in  $\|t_N(\hat{\theta}_N)\|$  is upper bounded (dominant coefficient negative) and lower bounded by 0. Which shows that

$$t_N(\hat{\theta}_N) = \mathcal{O}_p(1)$$

which concludes the proof, and :

$$R_N(\hat{\theta}_N) = o_p(1) \mathcal{O}_p(1) = o_p(1) \tag{20}$$

which is fundamental for the proof of theorem (1)

## 7.8 Proof of theorem 1

Now that we derived the expression of  $LRT_\infty$ , it remains to show that the bootstrap statistic also converges weakly in probability to this random variable. To do so we first derive a bootstrap quadratic approximation as in (21), where the expansion is done around  $\theta_N^*$ , as it is the true parameter of the bootstrap data. We obtain that :

$$\begin{aligned} l(\theta; y_{1:N}^*) - l(\theta_N^*; y_{1:N}^*) &= \frac{1}{2} Z_N^*(\theta_N^*)^T \tilde{I}_N^*(\theta_N^*) Z_N^*(\theta_N^*) \\ &\quad - \frac{1}{2} (t_N^*(\theta) - Z_N^*(\theta_N^*))^T \tilde{I}_N^*(\theta_N^*) (t_N^*(\theta) - Z_N^*(\theta_N^*)) + R_N^*(\theta) \end{aligned} \tag{21}$$

where the exponent  $*$  stands for "evaluated on the bootstrap data". The following proof will follow three main steps.

First we show that the bootstrap version of  $\tilde{I}_N(\theta_0)$  and of  $\tilde{S}_N(\theta_0)$  have the correct conditional limiting distribution which means that :

$$\begin{cases} \tilde{I}_N^*(\theta_N^*) - \tilde{I}_N(\theta_0) = o_{p^*}(1) \\ \forall t \in \mathbb{R}^{d_I} \quad \text{pr}^*\{\tilde{S}_N^*(\theta_N^*) < t\} - \text{pr}\{\tilde{S}_N(\theta_0) < t\} = o_p(1) \end{cases} \quad (22)$$

Second, we show that the new terms in the quadratic approximation that are supposed to be null converge toward 0 thanks to the shrinkage parameter.

Finally we show that the rest in the bootstrap quadratic approximation is a  $o_{p^*}(1)$ .

We first consider the case with no nuisance parameters i.e.  $\theta = (\psi, \lambda)$  to lighten the notations and work in two steps.

We first consider  $\tilde{I}_N(\theta_N^*)$ . We write :

$$[\tilde{I}_N(\theta_N^*)]_{m,n} = \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*)$$

where  $h_{m,n}^{(i)*}(\theta_N^*)$  is of the form:

$$h_{m,n}^{(i)*}(\theta_N^*) = c_{m,n} \frac{\partial^k \log f(y_i; \theta_N^*)}{\partial \theta_{i_1} \dots \partial \theta_{i_k}}$$

with  $c_{m,n} \in \mathbb{R}$ ,  $k \in \{2, 4\}$ ,  $1 \leq i_1, \dots, i_k \leq q_\lambda + q_\psi$ .

We want to show that:

$$\frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*) - \text{E}\{h_{m,n}^{(1)}(\theta_0)\} = o_{p^*}(1)$$

To show that we decompose the difference:

$$\frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*) - \text{E}[h_{m,n}^{(1)}(\theta_0)]$$

as:

$$\frac{1}{N} \sum_{i=1}^N \overbrace{h_{m,n}^{(i)*}(\theta_N^*) - \mathbb{E}^* [h_{m,n}^{(1)*}(\theta_N^*)]}^{(T1)} + \overbrace{\mathbb{E}^* [h_{m,n}^{(1)*}(\theta_N^*)] - \mathbb{E}^* [h_{m,n}^{(1)*}(\theta_0)]}^{(T2)} + \overbrace{\mathbb{E}^* [h_{m,n}^{(1)*}(\theta_0)] - \mathbb{E} [h_{m,n}^{(1)}(\theta_0)]}^{(T3)}$$

The term (T1) is a sum of centered random variables, we can apply Chebychev's inequality to have the convergence towards 0 (using assumption (2) to have that the variance is  $\mathcal{O}_p(1)$ ).

The term (T2) is controlled as follows :

$$\begin{aligned} |\mathbb{E}^* [(h_{m,n}^{(1)*}(\theta_N^*) - h_{m,n}^{(1)*}(\theta_0))]| &\leq \mathbb{E}^* [|(h_{m,n}^{(1)*}(\theta_N^*) - h_{m,n}^{(1)*}(\theta_0))|] \\ &\leq \sup_{\theta \in \Theta} \mathbb{E}_\theta [|(h_{m,n}^{(1)*}(\theta_N^*) - h_{m,n}^{(1)*}(\theta_0))|] \end{aligned}$$

as  $|(h_{m,n}^{(1)*}(\theta_N^*) - h_{m,n}^{(1)*}(\theta_0))| \leq 2 \sup_{\theta \in \Theta} |h_{m,n}^{(1)}(\theta)|$  almost surely, using assumption 2 and thanks to the consistency of  $\theta_N^*$  we can use the dominated convergence theorem to prove the convergence in probability toward 0.

Finally we deal with (T3) as follows :

$$\begin{aligned} |\mathbb{E}^* [h_{m,n}^{(1)*}(\theta_0)] - \mathbb{E} [h_{m,n}^{(1)}(\theta_0)]| &= \left| \int h_{m,n}^{(1)*}(\theta_0) \{f(y; \theta_N^*) - f(y; \theta_0)\} dy \right| \\ &\leq \int |h_{m,n}^{(1)*}(\theta_0)| |f(y; \theta_N^*) - f(y; \theta_0)| dy \end{aligned}$$

To show that this term tends toward 0 we use another time the equation (16), first we use a taylor expansion, there exist  $\theta^+$  between  $\theta_0$  and  $\theta_N^*$  such that :

$$\begin{aligned}
f(y; \theta_N^*) - f(y; \theta_0) &= (\theta_N^* - \theta_0)^T \nabla_{\theta} f(y; \theta^+) \\
&= (\theta_N^* - \theta_0)^T \nabla_{\theta} \log f(y; \theta^+) f(y; \theta^+)
\end{aligned}$$

therefore,

$$|f(y; \theta_N^*) - f(y; \theta_0)| \leq \|\theta_N^* - \theta_0\| \|\nabla_{\theta} \log f(y; \theta^+)\| f(y; \theta^+)$$

and,

$$\begin{aligned}
|\mathbb{E}^*[h_{m,n}^{(1)*}(\theta_0)] - \mathbb{E}[h_{m,n}^{(1)}(\theta_0)]| &\leq \|\theta_N^* - \theta_0\| \int |h_{m,n}^{(1)*}(\theta_0)| \|\nabla_{\theta} \log f(y; \theta^+)\| f(y; \theta^+) dy \\
&\leq \|\theta_N^* - \theta_0\| \int \frac{1}{2} \{ |h_{m,n}^{(1)*}(\theta_0)|^2 + \|\nabla_{\theta} \log f(y; \theta^+)\|^2 \} f(y; \theta^+) dy
\end{aligned}$$

Thanks to assumption (2), as  $\theta_N^* - \theta_0 = o_p(1)$  this last term is  $o_p(1)$  as  $N \rightarrow +\infty$ .

We then consider  $\tilde{S}_N^*(\theta_N^*)$ , which is  $\sqrt{N}$  times a sum of (conditionally) independent, centered, with finite variance random variables. Therefore, by proving that :

$$\mathbb{E}^*[\tilde{S}_N^*(\theta_N^*) \tilde{S}_N^*(\theta_N^*)^T] - \mathbb{E}[\tilde{S}_1(\theta_0)^T \tilde{S}_1(\theta_0)^T] = o_p(1) \quad (23)$$

and applying a multivariate version of the conditional central limit theorem of [10], it will conclude the second part of (13). Our third moment condition, and the consistency of the variance matrix of the score is much stronger than the Lindberg Feller conditions.

For each  $m, n$ ,  $[\tilde{S}_N^*(\theta_N^*) \tilde{S}_N^*(\theta_N^*)^T]_{m,n}$  can also be written as  $\frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*)$ , where  $h_{m,n}^{(i)*}(\theta_N^*) = [\tilde{S}_N^*(\theta_N^*)]_m \times [\tilde{S}_N^*(\theta_N^*)]_n$ . Due to assumption (2)ii)-iii), and the elementary equation (16),  $\mathbb{E}[h_{m,n}^{(i)*}(\theta_N^*)^{3/2}] < +\infty$ .

In order to prove (23), we once again split the sum :

$$\begin{aligned} & \mathbb{E}^*[h_{m,n}^{(1)*}(\theta_N^*)] - \mathbb{E}[h_{m,n}^{(1)}(\theta_0)] = \\ & \underbrace{\mathbb{E}^*[h_{m,n}^{(1)*}(\theta_N^*)] - \mathbb{E}^*[h_{m,n}^{(1)*}(\theta_0)]}_{(S1)} + \underbrace{\mathbb{E}^*[h_{m,n}^{(1)*}(\theta_0)] - \mathbb{E}[h_{m,n}^{(1)}(\theta_0)]}_{(S2)} \end{aligned}$$

We first deal with (S1) :

$$\begin{aligned} |\mathbb{E}^*[h_{m,n}^{(1)*}(\theta_N^*)] - \mathbb{E}^*[h_{m,n}^{(1)*}(\theta_0)]| & \leq \sup_{\theta \in \Theta} |\mathbb{E}_\theta[h_{m,n}^{(1)*}(\theta_N^*) - h_{m,n}^{(1)*}(\theta_0)]| \\ & \leq \sup_{\theta \in \Theta} \mathbb{E}_\theta[|h_{m,n}^{(1)*}(\theta_N^*) - h_{m,n}^{(1)*}(\theta_0)|] \\ & \leq 2 \sup_{\theta \in \Theta} \mathbb{E}_\theta[\sup_{\theta' \in \Theta} |h_{m,n}^{(1)*}(\theta')|] \\ & < +\infty \end{aligned}$$

Which enables, to apply dominated convergence to the first term, as  $\hat{\theta}_N$  is consistent.

For (S2) we proceed as follows :

$$\begin{aligned} |\mathbb{E}^*[h_{m,n}^{(1)*}(\theta_0)] - \mathbb{E}[h_{m,n}^{(1)}(\theta_0)]| & = \left| \int h_{m,n}^{(1)*}(\theta_0) \{f(y; \theta_N^*) - f(y; \theta_0)\} dy \right| \\ & \leq \int |h_{m,n}^{(1)*}(\theta_0)| |f(y; \theta_N^*) - f(y; \theta_0)| dy \end{aligned}$$

To show that this term tends toward 0 we use the same reasoning as before, first we use a Taylor expansion: there exist  $\theta^+$  between  $\theta_0$  and  $\theta_N^*$  such that :

$$\begin{aligned} f(y; \theta_N^*) - f(y; \theta_0) & = (\theta_N^* - \theta_0)^T \nabla_\theta f(y; \theta^+) \\ & = (\theta_N^* - \theta_0)^T \nabla_\theta \log f(y; \theta^+) f(y; \theta^+) \end{aligned}$$

therefore,



$$|f(y; \theta_N^*) - f(y; \theta_0)| \leq \|\theta_N^* - \theta_0\| \|\nabla_{\theta} \log f(y; \theta^+)\| f(y; \theta^+)$$

and,

$$|\mathbb{E}^*[h_{m,n}^{(1)*}(\theta_0)] - \mathbb{E}[h_{m,n}^{(1)}(\theta_0)]| \leq \|\theta_N^* - \theta_0\| \int |h_{m,n}^{(1)*}(\theta_0)| \|\nabla_{\theta} \log f(y; \theta^+)\| f(y; \theta^+) dy$$

We can't directly use equation (16) here because  $h_{m,n}^{(1)}(\theta)$  doesn't admit second order moments ( $\tilde{S}_N(\theta)$  admits third order moments). Thanks to Holder's inequality using  $p = \frac{3}{2}$  and  $q = 3$  we have that :

$$\begin{aligned} & \|\theta_N^* - \theta_0\| \int |h_{m,n}^{(1)*}(\theta_0)| \|\nabla_{\theta} \log f(y; \theta^+)\| f(y; \theta^+) dy \\ & \leq \|\theta_N^* - \theta_0\| \left( \int \|\nabla_{\theta} \log f(y; \theta^+)\|^3 f(y; \theta^+) dy \right)^{\frac{1}{3}} \left( \int |h_{m,n}^{(1)*}(\theta_0)|^{\frac{3}{2}} f(y; \theta^+) dy \right)^{\frac{2}{3}} \\ & = o_p(1) \end{aligned}$$

the last inequality holds thanks to assumption (2)ii)–iii) that enables to state that the two integrals are finite. That concludes the proof that (23) holds.

We now deal with the nuisance parameters. The proof starts the same way as before. The issue is that the odd derivatives with respect to the parameter  $\delta$  are not 0 when evaluated at  $\theta_N^*$  because  $\delta_N^* \neq 0$ . We recall that  $\theta = (\psi, \delta, \lambda)$  and  $\theta_N^* = (\psi_N^*, \delta_N^*, 0)$

To lighten the calculations we write  $dx^* = x - x_N^*$ , for any quantity  $x$ .

After expanding the likelihood as (21), new terms appear in  $R_N^*(\theta)$ : the odds order derivatives with respect to  $\delta$  which would be 0 if  $\delta_N^* = 0$ .

We want to show that :  $d\delta^{*T} \nabla_{\delta} l(\theta_N^*; y_{1:N}^*) = o_{p^*}(1)$  and  $\sum_{ijk} d\delta_i^* d\delta_j^* d\delta_k^* \frac{\partial^3 l(\theta_N^*; y_{1:N}^*)}{\partial \delta_i \partial \delta_j \partial \delta_k} =$

$o_{p^*}(1)$ .

As said before the issue comes from the fact that the bootstrap parameter of  $\delta$  is not 0 . Indeed the "good" bootstrap parameter would be  $u_N^* = (\psi_N^*, 0, 0) = \theta_N^* - (0, \hat{\delta}_N, 0)$  ( the bootstrap parameter that we would use if we knew where were located the nuisance parameters)

We are now going to expand the terms around  $u_N^*$ , so that the odds order derivatives evaluated at  $u_N^*$  will be zero. And we will use the fact that  $\delta_N^*$  converges very fast to zero.

We write  $\theta^+ = t\theta_N^* + (1-t)u_N^*$ :

$$d\delta^{*T} \nabla_{\delta} l(\theta_N^*; y_{1:N}^*) = d\delta^{*T} \left( 0 + \nabla_{\delta}^2 l_N^*(u_N^*) \delta_N^* + \sum_{ij} \delta_{Ni}^* \delta_{Nj}^* \times 0 + \sum_{ijk} \delta_{Ni}^* \delta_{Nj}^* \delta_{Nk}^* \frac{\partial^4 l(\theta^+; y_{1:N}^*)}{\partial \delta_i \partial \delta_j \partial \delta_k \partial \delta} \right)$$

the first non zero term is a centered random variable with finite variance (assumption (2)iii), therefore by the central limit theorem it is  $\mathcal{O}_{p^*}(\sqrt{N})$ , and the last term is  $\mathcal{O}_{p^*}(1)$  by the law of large number. Therefore:

$$\begin{aligned} |d\delta^{*T} \nabla_{\delta} l(\theta_N^*; y_{1:N}^*)| &\leq \|d\delta^*\| \left( \mathcal{O}_p(\sqrt{N}) o_p(N^{-\frac{1}{4}}) + \mathcal{O}_p(N) o_p(N^{-\frac{3}{4}}) \right) \\ &\leq \|d\delta^*\| o_p(N^{\frac{1}{4}}) \end{aligned}$$

The same way we have:

$$\begin{aligned} \left| \sum_{ijk} d\delta_i^* d\delta_j^* d\delta_k^* \frac{\partial^3 l(\theta_N^*; y_{1:N}^*)}{\partial \delta_i \partial \delta_j \partial \delta_k} \right| &= 0 + |\delta_N^{*T} \sum_{ijk} d\delta_i^* d\delta_j^* d\delta_k^* \frac{\partial^4 l(\theta^+; y_{1:N}^*)}{\partial \delta \partial \delta_i \partial \delta_j \partial \delta_k}| \\ &\leq o_{p^*}(N^{-\frac{1}{4}}) \|d\delta^*\|^3 \times \mathcal{O}_{p^*}(N) \\ &\leq \|d\delta^*\|^3 o_{p^*}(N^{\frac{3}{4}}) \end{aligned}$$

By using the same reasoning as in the proof of proposition 1, evaluating the expansion of the bootstrap likelihood at the bootstrap maximum likelihood estimator (restricted or not)  $\hat{\theta}^*$  we

have that almost surely

$$\begin{aligned}
0 &\leq l(\hat{\theta}^*; y_{1:N}^*) - l(\theta_N^*; y_{1:N}^*) \\
&\leq \|\tilde{S}_N^*(\theta_N^*)\| \|t_N(\hat{\theta}^*)\| - \frac{1}{2} \|t_N(\hat{\theta}^*)\|_{I_N^*(\theta_N^*)}^2 + R_N^*(\hat{\theta}^*) \\
&\leq \|\tilde{S}_N^*(\theta_N^*)\| \|t_N(\hat{\theta}^*)\| - \frac{1}{2} (o_{p^*}(1) + a^*) \|t_N(\hat{\theta}^*)\|^2 + o_{p^*}(1) (\|t_N(\hat{\theta}^*)\|^{\frac{1}{2}} + \|t_N(\hat{\theta}^*)\|^{\frac{3}{2}})
\end{aligned}$$

Even if this quantity is no longer a polynomial, the dominant term remain the same, therefore this quantity is lower bounded by 0 and upper bounded in probability as a upper bounded function of  $\|t_N(\hat{\theta}^*)\|$ . Which implies that  $\|t_N(\hat{\theta}^*)\| = \mathcal{O}_{p^*}(1)$  (in this proof we don't show that it is not a  $o_{p^*}(1)$  but it is not important here as we already showed that the bootstrap score and the bootstrap FIM converge toward the correct limit). But the important is that we showed that  $R_N^*(\hat{\theta}^*) = o_{p^*}(1)$  which conclude the proof.

## 7.9 Proof of proposition 5

We recall that  $\hat{\theta}_N = (\hat{\psi}_N, \hat{\delta}_N, \hat{\lambda}_N) = \arg \max_{\theta \in \Theta} l(\theta; y_{1:N})$ , and  $(c_N)$  is a sequence defined as in proposition (5).

Consider  $\theta_N^* = (\psi_N^*, \delta_N^*, \lambda_N^*)$  such that  $\forall k = 1, \dots, d_\psi$   $\psi_{N,k}^* = \hat{\psi}_{N,k} \mathbf{1}(\hat{\psi}_{N,k} > c_N)$ ,  $\forall k = 1, \dots, d_\delta$   $\delta_{N,k}^* = \hat{\delta}_{N,k} \mathbf{1}(\hat{\delta}_{N,k} > c_N)$  and  $\lambda_N^* = 0_{d_\lambda}$ .

The proof of this proposition follows exactly the lines of [11] lemma 1. First the fact that  $N^{1/4} c_N \rightarrow +\infty$  as  $N \rightarrow +\infty$  implies also that  $\sqrt{N} c_N \rightarrow +\infty$ .

Let us establish a technical result that will then be applied to our proposition. Let  $(x_N)$  a real valued random sequence and  $x_0 \in \mathbb{R}$  such that  $r_N(x_N - x_0) = O_p(1)$ , for  $r_N$  being whether  $\sqrt{N}$  or  $N^{1/4}$ .

If  $x_0 = 0$ ,

$$\mathbf{pr}(x_N > c_N) = \mathbf{pr}(r_N x_N > r_N c_N) = \mathbf{pr}\{O_p(1) > r_N c_N\} = o(1)$$

The last equality holds because  $r_N c_N \rightarrow +\infty$ . Therefore  $\mathbf{1}(x_N > c_N) = o_p(1)$  and finally  $r_N x_N \mathbf{1}(x_N > c_N) = O_p(1) o_p(1) = o_p(1)$ .

If  $x_0 \neq 0$ ,

$$\mathbf{pr}(|x_N| > c_N) = \mathbf{pr}(|x_N - x_0 + x_0| > c_N) \geq \mathbf{pr}\{|x_N - x_0| - |x_0| > c_N\}$$

$$\mathbf{pr}\{|x_N - x_0| - |x_0| > c_N\} = \mathbf{pr}\{|x_0| - |x_N - x_0| > c_N\} + \mathbf{pr}\{|x_N - x_0| - |x_0| > c_N\}$$

First,  $|x_N - x_0| + c_N = o_p(1)$  and  $|x_0| > 0$  therefore  $\mathbf{pr}\{|x_0| - |x_N - x_0| > c_N\} \rightarrow 1$  as  $N \rightarrow +\infty$ . And  $r_N(|x_0| + c_N) \rightarrow +\infty$  so  $\mathbf{pr}\{|x_N - x_0| - |x_0| > c_N\} \rightarrow 0$  as  $N \rightarrow +\infty$ .

Therefore :

$$\mathbf{1}(x_N > c_N) - 1 = o_p(1) \tag{24}$$

Finally,

$$\begin{aligned} r_N \{x_N \mathbf{1}(x_N > c_N) - x_0\} &= r_N(x_N - x_0) \mathbf{1}(x_N > c_N) - x_0 \mathbf{1}(x_N \leq c_N) \\ &= r_N(x_N - x_0) \mathbf{1}(x_N > c_N) - x_0 \{1 - \mathbf{1}(x_N > c_N)\} \\ &= O_p(1) + o_p(1) \end{aligned}$$

using equation (24), it concludes with Slutsky's theorem that  $r_N \{x_N \mathbf{1}(x_N > c_N) - x_0\}$  and  $r_N(x_N - x_0)$  have the same limiting distribution.

Applying this result to  $x_N = \hat{\delta}_N$  and  $r_N = N^{1/4}$  we have that  $\delta_N^* = o_p(1)$ . the same holds for  $\hat{\psi}_N : \sqrt{N}(\hat{\psi}_N - \psi_0) = O_p(1)$  (if  $\psi_0 = 0$  then  $\sqrt{N}\hat{\psi}_N = o_p(1)$ ).

Finally it is obvious that  $\theta_N^* \in \Theta_0$  as  $\lambda = 0_{d_\lambda}$ . Which concludes the proof.

## 7.10 Proof of proposition 6

We verify that our hypothesis imply the conditions required in [26].

We show easily that the assumptions C(1), C(2), C(3'), C(4'), C(5) are verified. Assumptions C(1)–(2) are verified with assumption (1). Assumption C(3') is weaker than assumption (4). C(4') is equivalent to assumption (5). C(5) is verified as we the continuity of the likelihood with respect to  $\theta$ , for every  $y$  and the measurability with respect to  $y$  for every  $\theta$ . The result is then discussed for instance in [19] exercise 7.2.3.

## 7.11 Proof of theorem 2

This proof is very similar to the one of theorem (1).

First we have to derive the asymptotic distribution of the likelihood ratio test statistic. We start from the quadratic expansion (21). We first show that  $\tilde{I}_N(\theta_0)$  converges in probability toward a non random matrix  $\tilde{I}(\theta_0)$ .

As in the proof of theorem (1) we write

$$[\tilde{I}_N(\theta_0)]_{m,n} = \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)}(\theta_0)$$

where  $h_{m,n}^{(i)}(\theta_0)$  is of the form:

$$h_{m,n}^{(i)}(\theta_0) = c_{m,n} \frac{\partial^k \log f_i(y_i; \theta_0)}{\partial \theta_{i_1} \dots \partial \theta_{i_k}}$$

with  $c_{m,n} \in \mathbb{R}$ ,  $k \in \{2, 4\}$ ,  $1 \leq i_1, \dots, i_k \leq d_\psi + d_\lambda + d_\delta$ .

As a consequence of assumption (4), using Chebychev's inequality :

$$\frac{1}{N} \sum_{i=1}^N |h_{m,n}^{(i)}(\theta_0) - \mathbb{E}[h_{m,n}^{(i)}(\theta_0)]| = o_p(1)$$

which enables to define  $\tilde{I}(\theta_0) = \left[ \lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[h_{m,n}^{(i)}(\theta_0)] \right]_{m,n}$  which is a nonrandom matrix that is supposed to be positive definite (assumption (3)). Furthermore,

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [|h_{m,n}^{(i)}(\theta_0)|] &\leq \frac{1}{N} \sum_{i=1}^N \sup_{i \in \mathbb{N}} \mathbb{E} [|h_{m,n}^{(i)}(\theta_0)|] \\ &\leq \sup_{i \in \mathbb{N}} \mathbb{E} [|h_{m,n}^{(i)}(\theta_0)|] \\ &< +\infty \end{aligned}$$

which holds for every  $N \geq 0$ . This last inequality enables to invert the sum and the integral :

$$\begin{aligned} [\tilde{I}(\theta_0)]_{m,n} &= \lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[h_{m,n}^{(i)}(\theta_0)] \\ &= \mathbb{E} \left[ \lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)}(\theta_0) \right] \\ &= \lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)}(\theta_0) \end{aligned}$$

where the last equality holds as we consider a non random quantity.

We consider now  $\tilde{S}_N(\theta_0)$  which is a sum of centered random variables with finite variances, we want to apply theorem 6.5 of [24]. Assumption (2)ii)–iii) and assumption (4) enables to state that :

$$\lim_{N \rightarrow +\infty} \mathbb{E} [\tilde{S}_N(\theta_0) \tilde{S}_N(\theta_0)^T] < +\infty$$

which is a direct consequence of theorem A.5 of [26]. Furthermore, still thanks to assumption (2)ii)–iii) and assumption (4) equation (6.3) in [24] theorem 6.5 is verified for  $\delta = 1$ , and therefore  $\tilde{S}_N(\theta_0)$  is  $O_p(1)$  and converges in distribution toward a random variable that we call  $U(\theta_0)$ .

The next step of the proof is to prove (22).

We first deal with  $\tilde{I}_N^*(\theta_N^*)$ , we still write :

$$\left[\tilde{I}_N(\theta_N^*)\right]_{m,n} = \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*)$$

We proceed as in the proof of theorem (1), and we split :

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*) - \left[\tilde{I}(\theta_0)\right]_{m,n} &= \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*) - \lim_{N \rightarrow +\infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[h_{m,n}^{(i)}(\theta_0)] \\ &= \frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*) - \frac{1}{N} \sum_{i=1}^N \mathbb{E}[h_{m,n}^{(i)}(\theta_0)] + o(1) \end{aligned}$$

as

$$\begin{aligned} &\overbrace{\frac{1}{N} \sum_{i=1}^N h_{m,n}^{(i)*}(\theta_N^*) - \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_N^*)]}^{(U1)} + \overbrace{\frac{1}{N} \sum_{i=1}^N \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_N^*)] - \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_0)]}^{(U2)} \\ &\quad + \overbrace{\frac{1}{N} \sum_{i=1}^N \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_0)] - \mathbb{E} [h_{m,n}^{(i)}(\theta_0)]}^{(U3)} + o(1) \end{aligned}$$

The term (U1) is a sum of centered random variables with finite variance uniformly bounded over  $i \in \mathbb{N}$  and therefore is  $o_{p^*}(1)$ .

We deal with (U2) as before :

$$\begin{aligned} \left| \frac{1}{N} \sum_{i=1}^N \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_N^*)] - \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_0)] \right| &\leq \frac{1}{N} \sum_{i=1}^N \mathbb{E}^* [|h_{m,n}^{(i)*}(\theta_N^*) - h_{m,n}^{(i)*}(\theta_0)|] \\ &\leq \sup_{i \in \mathbb{N}} \sup_{\theta \in \Theta} \mathbb{E}_\theta [|h_{m,n}^{(i)*}(\theta_N^*) - h_{m,n}^{(i)*}(\theta_0)|] \end{aligned}$$

and for every  $i$ ,  $|h_{m,n}^{(i)*}(\theta_N^*) - h_{m,n}^{(i)*}(\theta_0)| \leq 2 \sup_{\theta' \in \Theta} |h_{m,n}^{(i)*}(\theta')|$ , thanks to assumption (4), we can apply dominated convergence.

For the term (U3), we apply the exact same reasoning as in the proof of theorem (1) to

show that

$$\left| \frac{1}{N} \sum_{i=1}^N \mathbb{E}^* [h_{m,n}^{(i)*}(\theta_0)] - \mathbb{E} [h_{m,n}^{(i)}(\theta_0)] \right|$$

is almost surely smaller than

$$\|\theta_N^* - \theta_0\| \sup_{i \in \mathbb{N}} \int \frac{1}{2} \left\{ |h_{m,n}^{(i)*}(\theta_0)|^2 + \|\nabla_{\theta} \log f_i(y; \theta^+)\|^2 \right\} f_i(y; \theta^+) dy$$

which is almost surely smaller than

$$\frac{\|\theta_N^* - \theta_0\|}{2} \left\{ \sup_{\theta^+ \in \Theta} \mathbb{E}_{\theta^+} \left[ \sup_{\theta \in \Theta} |h_{m,n}^{(i)}(\theta)|^2 \right] + \sup_{\theta^+ \in \Theta} \mathbb{E}_{\theta^+} \left[ \sup_{\theta \in \Theta} \|\nabla_{\theta} \log f_i(y; \theta)\|^2 \right] \right\}$$

which is  $o_p(1)$  du to the consistency of  $\theta_N^*$  and assumption (2). Which concludes the proof of theorem (2).

## 7.12 Proof of proposition 7

Recall that we want to show that :

$$\sup_{\theta' \in \Theta} \mathbb{E}_{\theta'} \left\{ \sup_{\theta \in \Theta} \|\nabla_{\theta}^k \log f_i(y_i; \theta)\|^{\gamma} \right\} < +\infty$$

with  $\gamma = 2$  for  $k = 0, 3, 4$  and  $\gamma = 3$  for  $k = 1, 2$ , for every  $i \in \mathbb{N}$ .

For a sake of clarity, we consider the following simplified notations :

$$g(x_{ij}, \beta, \Lambda \xi_i) = g_{\theta}^{ij}(\xi_i) \text{ and } g_{\theta}^i(\xi_i) = (g_{\theta}^{ij}(\xi_i))_{j=1, \dots, J_i}$$

$$f(y_i; \theta) = \mathbb{E}\{f(y_i; \xi, \theta)\} \propto \mathbb{E}[\exp\{-V(\theta, y_i, \xi_i)\}] \text{ with}$$

$$V(\theta, y_i, \xi_i) = \frac{\sum_j (y_{ij} - g_{\theta}^j(\xi_i))^2}{2\sigma^2} = \frac{\|y_i - g_{\theta}^i(\xi_i)\|^2}{2\sigma^2}$$

where the expectation is taken with respect to the random variable  $\xi$ .



As these quantities are individual, we get rid of the subscript  $i$  to lighten the notations.

As we consider the parameter space  $\Theta$  compact, the residual variance  $\sigma$  is restricted to lie in a segment  $[\sigma_{min}; \sigma_{max}]$  with  $0 < \sigma_{min} < \sigma_{max} < +\infty$ . That is why we consider the gaussian density up to a constant that won't change the reasoning. From now on we won't write  $\propto$  and we make the shortcut  $f(y; \xi, \theta) = e^{-V(\theta, y, \xi)}$ .

We also write  $E^Z\{\cdot\}$  when the expectation is taken with respect to the random variable  $Z$ .

We suppose now that  $\|g_\theta(\xi)\| \rightarrow +\infty$  as  $\|\xi\| \rightarrow +\infty$  which is the most complicated case.

We first deal with the simplest case  $k = 0$ , we want to show that

$$E_{\theta'} \left\{ \sup_{\theta \in \Theta} |\log f(y; \theta)|^2 \right\} < +\infty$$

Let  $\theta, \theta' \in \Theta$ , let  $M > 0$ ,

$$\begin{aligned} f(y; \theta) &= E^\xi \{ f(y; \xi, \theta) \} \\ &\geq E^\xi \{ f(y; \xi, \theta) \mathbf{1}(\xi \leq M) \} \\ &\geq E^\xi \{ e^{-V(\theta, y, \xi)} \mathbf{1}(\xi \leq M) \} \\ &= E^\xi \left\{ e^{-\frac{\|y - g_\theta(\xi)\|^2}{2\sigma}} \mathbf{1}(\xi \leq M) \right\} \\ &= E^\xi \left\{ e^{-\frac{\|y\|^2 + \|g_\theta(\xi)\|^2 - 2y^T g_\theta(\xi)}{2\sigma}} \mathbf{1}(\xi \leq M) \right\} \\ &\geq E^\xi \left\{ e^{-\frac{\|y\|^2 + \|g_\theta(\xi)\|^2 + 2\|y\| \|g_\theta(\xi)\|}{2\sigma}} \mathbf{1}(\xi \leq M) \right\} \end{aligned}$$

where the last inequality is a direct application of Cauchy Schwartz's inequality.

The quantity in the exponential is a polynomial in  $\|g_\theta(\xi)\|$  that goes to  $-\infty$  when  $\|\xi\| \rightarrow +\infty$ . Its minimal value is achieved at  $\alpha_M(\theta) = \sup_{\xi: \|\xi\| \leq M} \|g_\theta(\xi)\| = \|g_\theta(\tilde{\xi})\|$

Therefore we obtain that for every  $M > 0$ ,

$$f(y; \theta) \geq \text{pr}(\xi \leq M) e^{-\frac{(\|y\| + \alpha_M(\theta))^2}{2\sigma^2}} \quad (25)$$

yet, by writing  $\kappa_M = \text{pr}(\xi \leq M)$ ,

$$\begin{aligned}
1 &> f(y; \theta) > \kappa_M e^{-\frac{(\|y\| + \alpha_M(\theta))^2}{2\sigma^2}} \\
\Leftrightarrow 0 &> \log f(y; \theta) > \log(\kappa_M) - \frac{(\|y\| + \alpha_M(\theta))^2}{2\sigma^2} \\
\Leftrightarrow 0 < |\log f(y; \theta)|^2 &< \log(\kappa_M)^2 + \frac{(\|y\| + \alpha_M(\theta))^4}{4\sigma^2} - \frac{\log(\kappa_M)(\|y\| + \alpha_M(\theta))^2}{\sigma^2} \\
\Leftrightarrow 0 < |\log f(y; \theta)|^2 &< \log(\kappa_M)^2 + \frac{\{\|y\| + \alpha_M(\theta)\}^4}{4\sigma^4} \\
\Leftrightarrow 0 < |\log f(y; \theta)|^2 &< \log(\kappa_M)^2 + \frac{\{\|y\| + \alpha_M(\bar{\theta})\}^4}{4\sigma_{min}^4}
\end{aligned}$$

where  $\bar{\theta} = \text{argsup}_{\theta \in \Theta} \alpha(\theta) \in \Theta$  by continuity of  $\alpha(\cdot)$  (continuity of  $\theta \rightarrow g_\theta(\tilde{\xi})$ ) and compactness of  $\Theta$

We define the quantity  $P_M(y) = \log(\kappa_M)^2 + \frac{\{\|y\| + \alpha_M(\bar{\theta})\}^4}{4\sigma_{min}^4}$  that does not depend on  $\theta$ .

Therefore we have that :

$$\begin{aligned}
\mathbb{E}_{\theta'} \left\{ \sup_{\theta \in \Theta} |\log f(y; \theta)|^2 \right\} &\leq \mathbb{E}_{\theta'} \{ P_M(y) \} \\
&= \int_y P_M(y) f(y; \theta') dy \\
&= \int_y P_M(y) \int_\xi f(y; \xi, \theta') \pi_p(\xi) d\xi dy \\
&= \int_\xi \int_y P_M(y) f(y; \xi, \theta') \pi_p(\xi) dy d\xi \\
&= \int_\xi \int_u 2\sigma'^2 P_M(2\sigma'^2 u + g_{\theta'}(\xi)) \pi_J(u) \pi_p(\xi) du d\xi \\
&\leq \int_\xi \int_u 2\sigma_{max}^2 P_M(2\sigma_{max}^2 u + g_{\theta'}(\xi)) \pi_J(u) \pi_p(\xi) du d\xi
\end{aligned}$$

using first Fubini-Tonelli's theorem and then a change of variable :  $u = \frac{y - g_{\theta'}(\xi)}{2\sigma'^2}$ .

$P_M(2\sigma'^2 u + g_{\theta'}(\xi))$  is a polynomial of degree 4 in  $\|u\|$  and  $\|g_{\theta'}(\xi)\|$ . Thanks to the

assumption (6) of proposition (7) with  $k_1 = 0$  and  $k_2 = 4$ , we have that :

$$\sup_{\theta' \in \Theta} \mathbb{E}_{\theta'} \left\{ \sup_{\theta \in \Theta} |\log f(y; \theta)|^2 \right\} \leq \int_{\xi} \int_u 2\sigma_{max}^2 \sup_{\theta' \in \Theta} P_M(2\sigma_{max}^2 u + g_{\theta'}(\xi)) \Psi_J(u) \Psi_p(\xi) du d\xi$$

which concludes the first part of this proof.

We now consider the case  $k = 1, 2, 3, 4$ , and  $\gamma = 2, 3$  (the fact of considering  $\gamma = 2$  or  $3$  is the same, therefore we will consider  $\gamma = 3$  as it is stronger).

$$\|\nabla_{\theta}^k \log f(y; \theta)\|^3 \leq \sup_{\substack{I \in \{1, \dots, d_{\theta}\}^k \\ I = (i_1, \dots, i_k)}} d_{\theta}^k \left\| \frac{\partial^k \log f(y; \theta)}{\prod_{j=1}^k \partial \theta_{i_j}} \right\|^3$$

Let  $I_0 = (i_1, \dots, i_k)$  the subset of indexes where the sup in the right hand side is achieved.

We consider the quantity

$$\frac{\partial^k \log f(y; \theta)}{\prod_{j=1}^k \partial \theta_{i_j}}$$

as the derivatives of the comosition

$$\theta \mapsto f(y; \theta) \mapsto \log f(y; \theta)$$

and we use Faa di Bruno's formula to develop this expression :

$$\frac{\partial^k \log f(y; \theta)}{\prod_{j=1}^k \partial \theta_{i_j}} = \sum_{\Psi \in \mathcal{P}(\{i_1, \dots, i_k\})} \alpha_{\Psi} f(y; \theta)^{-|\Psi|} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}$$

where  $\mathcal{P}(K)$  stands for all partitions of a set  $K$ , and  $(\alpha_{\Psi})_{\Psi}$  are constants. A sufficient condition for this quantity to be  $\mathbb{L}^3$  is that each term of the sum is  $\mathbb{L}^3$ .

Let  $\Psi \in \mathcal{P}(\{i_1, \dots, i_k\})$ , we write  $m$  the cardinal of  $\Psi$ .

We recall that :

$$f(y; \theta) = \mathbb{E}^{\xi} \left\{ e^{-V(\theta, y, \xi)} \right\}$$

therefore, for every  $B \in \Psi$ ,

$$\begin{aligned} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b} &= \mathbb{E}^\xi \left\{ \frac{\partial^{|B|} f(y; \xi, \theta)}{\prod_{b \in B} \partial \theta_b} \right\} \\ &= \mathbb{E}^\xi \left[ P^B \{V(\theta, y, \xi)\} e^{-V(\theta, y, \xi)} \right] \end{aligned}$$

where  $P^B \{V(\theta, y, \xi)\}$  is a polynomial of degree  $m$  in  $y$  and the partial derivatives of  $g_\theta(\xi)$ .

$$\begin{aligned} |f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}| &= f(y; \theta)^{-m} \prod_{B \in \Psi} \left| \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b} \right| \\ &\stackrel{(Jensen)}{\leq} f(y; \theta)^{-m} \prod_{B \in \Psi} \mathbb{E}^\xi \left[ |P^B \{V(\theta, y, \xi)\}| e^{-V(\theta, y, \xi)} \right] \end{aligned}$$

the random variables  $\xi$  can be renamed  $\xi_B$  for each term of the product so that :

$$|f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}| \leq f(y; \theta)^{-m} \prod_{B \in \Psi} \mathbb{E}^{\xi_B} \left[ |P^B \{V(\theta, y, \xi_B)\}| e^{-V(\theta, y, \xi_B)} \right]$$

We introduce  $\Xi = (\xi_{B_1}^T, \dots, \xi_{B_m}^T)^T \sim \mathcal{N}(0, I_{m \times p})$ , so that we can write :

$$\begin{aligned} f(y; \theta)^{-m} \prod_{B \in \Psi} \mathbb{E}^{\xi_B} \left[ |P^B \{V(\theta, y, \xi_B)\}| e^{-V(\theta, y, \xi_B)} \right] \\ &= f(y; \theta)^{-m} \mathbb{E}^\Xi \left[ \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}| e^{-V(\theta, y, \xi_B)} \right] \\ &= f(y; \theta)^{-m} \mathbb{E}^\Xi \left( \left[ \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}| \right] e^{-\sum_{B \in \Psi} V(\theta, y, \xi_B)} \right) \end{aligned}$$

every terms in the product inside the expectation is nonnegative, therefore when rising this quantity to the power of 3 we can use Jensen by convexity on  $\mathbb{R}^+$  of  $x \mapsto x^3$ . And we find

that :

$$\begin{aligned} |f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}|^3 &\leq f(y; \theta)^{-3m} \mathbb{E}^\Xi \left( \left[ \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}| \right] e^{-\sum_{B \in \Psi} V(\theta, y, \xi_B)} \right)^3 \\ &\leq f(y; \theta)^{-3m} \mathbb{E}^\Xi \left( \left[ \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}|^3 \right] e^{-\sum_{B \in \Psi} 3V(\theta, y, \xi_B)} \right) \end{aligned}$$

using Jensen's inequality.

By using equation (25), we know that :

$$f(y; \theta) \geq \kappa_M e^{-\frac{(\|y\| + \alpha_M(\bar{\theta}))^2}{2\sigma^2}}$$

$\kappa_m$  is a constant that has no impact on the reasoning therefore we neglect it to lighten the notations. We write  $V(y, \sigma) = \frac{(\|y\| + \alpha_M(\bar{\theta}))^2}{2\sigma^2}$  so that we have :

$$|f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}|^3 \leq e^{3mV(y, \sigma)} \mathbb{E}^\Xi \left( \left[ \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}|^3 \right] e^{-\sum_{B \in \Psi} 3V(\theta, y, \xi_B)} \right)$$

We recall that the cardinal of  $\Psi$  is equal to  $m$  so :

$$|f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}|^3 \leq \mathbb{E}^\Xi \left( \left[ \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}|^3 \right] e^{-3\sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right) \quad (26)$$

Once more to lighten the notations we define  $M_\theta^\Psi(y, \Xi) = \prod_{B \in \Psi} |P^B \{V(\theta, y, \xi_B)\}|^3$

Let  $\theta' \in \Theta$ ,

$$\begin{aligned}
& \mathbb{E}_{\theta'}^y \left\{ \left| f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b} \right|_3 \right\} \\
& \leq \mathbb{E}_{\theta'}^y \left\{ \mathbb{E}^{\Xi} \left( M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right) \right\} \\
& = \mathbb{E}^{\Xi} \left\{ \mathbb{E}_{\theta'}^y \left( M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right) \right\} \\
& = \mathbb{E}^{\Xi} \left\{ \int_y M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} f(y; \theta') dy \right\} \\
& = \mathbb{E}^{\Xi} \left[ \int_y M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \mathbb{E}^{\xi} \{ e^{-V(\theta', y, \xi)} \} dy \right] \\
& = \mathbb{E}^{\Xi, \xi} \left\{ \int_y M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} e^{-V(\theta', y, \xi)} dy \right\} \tag{27}
\end{aligned}$$

using once more Fubini Tonelli's theorem.

We focus on the exponential :

$$\begin{aligned}
-3 \sum_{B \in \Psi} \{V(\theta, y, \xi_B) - V(y, \sigma)\} - V(\theta', y, \xi) &= -3 \sum_{B \in \Psi} \{V(\theta, y, \xi_B) - V(y, \sigma)\} - V(\theta', y, \xi) \\
&= -3 \sum_{B \in \Psi} \left\{ \frac{\|y - g_{\theta}(\xi_B)\|^2}{2\sigma^2} - \frac{\{\|y\| + \alpha(\bar{\theta})\}^2}{2\sigma^2} \right\} - \frac{\|y - g_{\theta'}(\xi)\|^2}{2\sigma'^2} \\
&= -3 \sum_{B \in \Psi} \frac{1}{2\sigma^2} \{ \|g_{\theta}(\xi_B)\|^2 - 2y^T g_{\theta}(\xi_B) + 2\|y\| \alpha(\bar{\theta}) - \alpha(\bar{\theta})^2 \} \\
&\quad - \frac{1}{2\sigma'^2} \{ \|y\|^2 + \|g_{\theta'}(\xi)\|^2 - 2y^T g_{\theta'}(\xi) \}
\end{aligned}$$

We use Cauchy-Schwartz's inequality to get rid of the scalar product and consider scalar quantities :

$$\begin{aligned}
-3 \sum_{B \in \Psi} \{V(\theta, y, \xi_B) - V(y, \sigma)\} - V(\theta', y, \xi) &= -3 \sum_{B \in \Psi} \{V(\theta, y, \xi_B) - V(y, \sigma)\} - V(\theta', y, \xi) \\
&\leq -3 \sum_{B \in \Psi} \frac{1}{2\sigma^2} \{\|g_{\theta}(\xi_B)\|^2 - 2\|y\| \|g_{\theta}(\xi_B)\| + 2\|y\| \alpha(\bar{\theta}) - \alpha(\bar{\theta})^2\} \\
&\quad - \frac{1}{2\sigma'^2} \{\|y\|^2 + \|g_{\theta'}(\xi)\|^2 - 2\|y\| \|g_{\theta'}(\xi)\|\}
\end{aligned}$$

Therefore by taking the integrated form of (27) we have that :

$$\begin{aligned}
\mathbb{E}_{\theta'}^y &\left\{ |f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b}|^3 \right\} \\
&\leq \int_{\Xi, \xi} \int_y M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} e^{-V(\theta', y, \xi)} dy e^{-\frac{\|\Xi\|^2}{2} - \frac{\|\xi\|^2}{2}} d\Xi d\xi \\
&\leq \int_{\Xi, \xi, y} M_{\theta}^{\Psi}(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma) - V(\theta', y, \xi) - \frac{\|\Xi\|^2}{2} - \frac{\|\xi\|^2}{2}} dy d\Xi d\xi \\
&\leq \int_{\Xi, \xi, y} M_{\theta}^{\Psi}(y, \Xi) e^{H(\theta, \theta', y, \Xi, \xi)} dy d\Xi d\xi
\end{aligned}$$

where

$$\begin{aligned}
H(\theta, \theta', y, \Xi, \xi) &= -3 \sum_{B \in \Psi} \frac{1}{2\sigma^2} \{\|g_{\theta}(\xi_B)\|^2 - 2\|y\| \|g_{\theta}(\xi_B)\| + 2\|y\| \alpha(\bar{\theta}) - \alpha(\bar{\theta})^2\} \\
&\quad - \frac{1}{2\sigma'^2} \{\|y\|^2 + \|g_{\theta'}(\xi)\|^2 - 2\|y\| \|g_{\theta'}(\xi)\| - \frac{\|\Xi\|^2}{2} - \frac{\|\xi\|^2}{2}\}
\end{aligned}$$

by splitting  $\|\Xi\|^2$  and rearranging the terms we get :

$$\begin{aligned}
H(\theta, \theta', y, \Xi, \xi) &= \|y\| \left\{ \frac{\|g_{\theta'}(\xi)\|}{\sigma'^2} - \frac{3m}{\sigma^2} \alpha(\bar{\theta}) \right\} + \frac{3m}{2\sigma^2} \alpha(\bar{\theta})^2 + \sum_{B \in \Psi} \frac{3}{\sigma^2} \|y\| \|g_{\theta}(\xi_B)\| \\
&\quad - \frac{1}{2\sigma'^2} \|y\|^2 - \frac{1}{2} \left\{ \|\xi\|^2 + \frac{1}{\sigma'^2} \|g_{\theta'}(\xi)\|^2 \right\} - \frac{1}{2} \sum_{B \in \Psi} \left\{ \frac{3}{\sigma^2} \|g_{\theta}(\xi_B)\|^2 + \|\xi_B\|^2 \right\}
\end{aligned}$$

As it is constant we can omit  $\frac{3m}{2\sigma^2} \alpha(\bar{\theta})^2$  in the development as it can be taken out from the integrals. Furthermore  $-\alpha(\bar{\theta})\|y\|/\sigma^2 < 0$  therefore it can be upper bounded by zero.

Then we define for  $\theta \in \Theta$ ,  $\xi \in \mathbb{R}^p$   $r_\theta(\xi)$  the ratio  $\frac{\|g_\theta(\xi)\|}{\|\xi\|}$  such that,

$$\begin{aligned}
H(\theta, \theta', y, \Xi, \xi) &\leq \frac{r_{\theta'}(\xi)}{\sigma'^2} \|\xi\| \|y\| + \sum_{B \in \Psi} \frac{3r_\theta(\xi_B)}{\sigma^2} \|\xi_B\| \|y\| \\
&\quad - \frac{1}{2\sigma'^2} \|y\|^2 - \frac{1}{2} \left(1 + \frac{r_{\theta'}(\xi)^2}{\sigma'^2}\right) \|\xi\|^2 - \frac{1}{2} \sum_{B \in \Psi} \left\{1 + \frac{3r_\theta(\xi_B)^2}{\sigma^2}\right\} \|\xi_B\|^2 \\
&\leq -\frac{1}{2} \left(\frac{y^T}{\sigma_{min}}, \tilde{\xi}^T, \tilde{\Xi}^T\right)^T \begin{pmatrix} 1 & -\frac{\sigma_{min} r_{\theta'}(\xi)}{\sigma_{max}^2} & \dots & -\frac{3\sigma_{min} r_\theta(\xi_B)}{\sigma_{max}^2} & \dots \\ -\frac{\sigma_{min} r_{\theta'}(\xi)}{\sigma_{max}^2} & 1 & 0 & \dots & 0 \\ \vdots & 0 & \ddots & 0 & \vdots \\ -\frac{3\sigma_{min} r_\theta(\xi_B)}{\sigma_{max}^2} & \vdots & 0 & \ddots & 0 \\ \vdots & 0 & \dots & 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{y}{\sigma_{min}} \\ \tilde{\xi} \\ \tilde{\Xi} \end{pmatrix}
\end{aligned}$$

Where  $\tilde{\xi} = \sqrt{\left(1 + \frac{r_{\theta'}(\xi)^2}{\sigma'^2}\right)} \times \xi$  and  $\tilde{\Xi} = \left(\sqrt{\left\{1 + \frac{3r_\theta(\xi_B)^2}{\sigma^2}\right\}} \times \xi_B\right)_{B \in \Psi}$  Let  $\Omega(\theta, \theta', y, \xi, \Xi)$  the symmetric matrix involved in the previous equation. In order to bound this last quantity we study the spectrum of this matrix. To do so we determine its characteristic. Let  $x \in \mathbb{R}$ . We call  $n = n$  By using the cofactor expansion formula for the first line we find :

$$\begin{aligned}
\det \{xI_n - \Omega(\theta, \theta', y, \xi, \Xi)\} &= (x-1)^n - \frac{\sigma_{min}^2 r_{\theta'}^2(\xi)}{\sigma_{max}^4} (x-1)^{n-2} - \sum_{B \in \Psi} \frac{9\sigma_{min}^2 r_\theta^2(\xi_B)}{\sigma_{max}^4} (x-1)^{n-2} \\
&= (x-1)^{n-2} \left[ (x-1)^2 - \left\{ \frac{\sigma_{min}^2 r_{\theta'}^2(\xi)}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 r_\theta^2(\xi_B)}{\sigma_{max}^4} \right\} \right]
\end{aligned}$$

therefore the spectrum of  $\Omega(\theta, \theta', y, \xi, \Xi)$  is

$$\left\{ 1; 1 - \sqrt{\frac{\sigma_{min}^2 r_{\theta'}^2(\xi)}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 r_\theta^2(\xi_B)}{\sigma_{max}^4}}; 1 + \sqrt{\frac{\sigma_{min}^2 r_{\theta'}^2(\xi)}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 r_\theta^2(\xi_B)}{\sigma_{max}^4}} \right\}$$



and we get :

$$\begin{aligned} H(\theta, \theta', y, \Xi, \xi) &\leq -\frac{1}{2} \left\{ 1 - \sqrt{\frac{\sigma_{min}^2 r_{\theta'}^2(\xi)}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 r_{\theta}^2(\xi_B)}{\sigma_{max}^4}} \right\} \left( \frac{\|y\|^2}{\sigma_{min}^2} + \|\tilde{\xi}\|^2 + \|\tilde{\Xi}\|^2 \right) \\ &\leq -\frac{1}{2} \left\{ 1 - \sqrt{\frac{\sigma_{min}^2 r_{\theta'}^2(\xi)}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 r_{\theta}^2(\xi_B)}{\sigma_{max}^4}} \right\} \left( \frac{\|y\|^2}{\sigma_{min}^2} + \|\xi\|^2 + \|\Xi\|^2 \right) \end{aligned}$$

where the last equality holds because  $\|\tilde{\xi}\| > \|\xi\|$  and  $\|\tilde{\Xi}\| > \|\Xi\|$ .

Let  $\varepsilon > 0$  small enough such that :

$$0 < 1 - \sqrt{\frac{\sigma_{min}^2 \varepsilon^2}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 \varepsilon^2}{\sigma_{max}^4}} < 1$$

Let  $\omega = \sqrt{\frac{\sigma_{min}^2 \varepsilon^2}{\sigma_{max}^4} + \sum_{B \in \Psi} \frac{9\sigma_{min}^2 \varepsilon^2}{\sigma_{max}^4}}$

Let  $K$  a compact set such that assumption (7) is verified, and  $K_M = \{y \in \mathbb{R}^J \leq M\}$  for some  $M > 0$

we finally get to :

$$\begin{aligned} &\mathbb{E}_{\theta'}^y \left\{ \left| f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b} \right|^3 \right\} \\ &\leq \int_{(\Xi, \xi, y) \in K^{m+1} \times K_M} M_{\theta}^{\Psi}(y, \Xi) e^{H(\theta, \theta', y, \Xi, \xi)} dy d\Xi d\xi \\ &+ \int_{(\Xi, \xi, y) \in \mathbb{R}^{(m+1)p+J} \setminus K^{m+1} \times K_M} M_{\theta}^{\Psi}(y, \Xi) e^{-\frac{1}{2}(1-\omega) \left( \frac{\|y\|^2}{\sigma_{min}^2} + \|\xi\|^2 + \|\Xi\|^2 \right)} dy d\Xi d\xi \end{aligned}$$

Now, to be precise, we restart from equation (26) and the following calculation leading to equation (27), to show that the sup can be considered inside the integral :

$$\begin{aligned}
& \mathbb{E}_{\theta'}^y \left\{ \sup_{\theta \in \Theta} \left| f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b} \right|^3 \right\} \\
& \leq \mathbb{E}_{\theta'}^y \left\{ \sup_{\theta \in \Theta} \mathbb{E}^\Xi \left( M_\theta^\Psi(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right) \right\} \\
& = \int_y \sup_{\theta \in \Theta} \left[ \mathbb{E}^\Xi \left\{ M_\theta^\Psi(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right\} \right] f_{\theta'}(y) dy \\
& (f_{\theta'}(y) > 0) = \int_y \sup_{\theta \in \Theta} \left[ \mathbb{E}^\Xi \left\{ M_\theta^\Psi(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right\} f_{\theta'}(y) \right] dy \\
& (Jensen) \leq \int_y \mathbb{E}^\Xi \left[ \sup_{\theta \in \Theta} \left\{ M_\theta^\Psi(y, \Xi) e^{-3 \sum_{B \in \Psi} V(\theta, y, \xi_B) - V(y, \sigma)} \right\} \right] f_{\theta'}(y) dy
\end{aligned}$$

where the second equality holds because  $f_{\theta'}(y)$  is nonnegative and does not depend on  $\theta$ .

And finally :

$$\begin{aligned}
& \mathbb{E}_{\theta'}^y \left\{ \sup_{\theta \in \Theta} \left| f(y; \theta)^{-m} \prod_{B \in \Psi} \frac{\partial^{|B|} f(y; \theta)}{\prod_{b \in B} \partial \theta_b} \right|^3 \right\} \\
& \leq \int_{(\Xi, \xi, y) \in K^{m+1} \times K_M} \sup_{\theta \in \Theta} M_\theta^\Psi(y, \Xi) e^{H(\theta, \theta', y, \Xi, \xi)} dy d\Xi d\xi \\
& + \int_{(\Xi, \xi, y) \in \mathbb{R}^{(m+1)p+J} \setminus K^{m+1} \times K_M} \sup_{\theta \in \Theta} M_\theta^\Psi(y, \Xi) e^{-\frac{1}{2}(1-\omega) \left( \frac{\|y\|^2}{\sigma_{min}^2} + \|\xi\|^2 + \|\Xi\|^2 \right)} dy d\Xi d\xi \quad (28) \\
& \leq \sup_{\substack{\theta \in \Theta \\ (\Xi, \xi, y) \in K^{m+1} \times K_M}} \text{Vol}(K^{m+1} \times K_M) M_\theta^\Psi(y, \Xi) e^{H(\theta, \theta', y, \Xi, \xi)} \\
& + \int_{(\Xi, \xi, y) \in \mathbb{R}^{(m+1)p+J} \setminus K^{m+1} \times K_M} \sup_{\theta \in \Theta} M_\theta^\Psi(y, \Xi) e^{-\frac{1}{2}(1-\omega) \left( \frac{\|y\|^2}{\sigma_{min}^2} + \|\xi\|^2 + \|\Xi\|^2 \right)} dy d\Xi d\xi
\end{aligned}$$

The first term of the last quantity is a suprema of a continuous function over a compact set and is therefore finite.

$M_\theta^\Psi(y, \Xi)$  is a polynomial with respect to  $\|y\|$  and the partial derivatives of  $\theta \mapsto g_\theta(\xi_B)$  for each  $B \in \Psi$ , therefore, thanks to assumption (6), considering  $\varepsilon$  small enough such that  $\varepsilon < \delta$  where  $\delta$  is defined in proposition (7), we finally get to the conclusion :

$$\mathbb{E}_{\theta'}^y \left\{ \sup_{\theta \in \Theta} |f(y; \theta)|^{-m} \prod_{B \in \Psi} \frac{|\partial^{|B|} f(y; \theta)|}{\prod_{b \in B} \partial \theta_b} \right\} < +\infty$$

In equation (28) it is important to notice that the last term no longer depends on  $\theta'$ , and that the first term in the sup is a continuous function of  $\theta'$ , therefore by compactness of  $\Theta$  the suprema can be taken with respect to  $\theta'$  and

$$\sup_{\theta' \in \Theta} \mathbb{E}_{\theta'}^y \left\{ \sup_{\theta \in \Theta} |f(y; \theta)|^{-m} \prod_{B \in \Psi} \frac{|\partial^{|B|} f(y; \theta)|}{\prod_{b \in B} \partial \theta_b} \right\} < +\infty$$

which concludes the proof.

### 7.13 Logistic growth model example

We consider here the logistic growth model that is commonly used in many fields. We show how to use the criteria given in proposition (7).

We recall the definition of the logistic growth function :

$$g(x_{ij}, \beta, \Lambda \xi_i) = \frac{\beta_1 + \lambda_1 \xi_{i1}}{1 + \exp\left(-\frac{x_{ij} - (\beta_2 + \lambda_2 \xi_{i2})}{\beta_3 + \lambda_3 \xi_{i3}}\right)}. \quad (29)$$

To verify that assumption (6) is verified, we need to calculate the derivatives of  $g$  with respect to  $\theta$ . To avoid heavy calculations we consider the parameter  $\lambda_3$  :

$$\begin{aligned} \frac{\partial g(x_{ij}, \beta, \Lambda \xi_i)}{\partial \lambda_3} &= -\xi_{i3} \left( \frac{j - (\beta_2 + \lambda_2 \xi_{i2})}{(\beta_3 + \lambda_3 \xi_{i3})^2} \right) \exp\left(-\frac{j - (\beta_2 + \lambda_2 \xi_{i2})}{\beta_3 + \lambda_3 \xi_{i3}}\right) \\ &\quad \times \frac{\beta_1 + \lambda_1 \xi_{i1}}{\left(1 + \exp\left(-\frac{j - (\beta_2 + \lambda_2 \xi_{i2})}{\beta_3 + \lambda_3 \xi_{i3}}\right)\right)^2} \end{aligned}$$

The only issue of this quantity occurs at  $\xi_{i3} = -\frac{\beta_3}{\lambda_3}$  that tends toward 0 as  $\xi_{i3} \rightarrow -\frac{\beta_3}{\lambda_3}$ . Finally this quantity is integrable with respect to the Gaussian distribution. By iterating the derivatives, we find expressions similar to this one and assumption (6) is verified.

Let  $\varepsilon > 0, M > 0$ , and we define the set :

$$K_M = \{\xi \in \mathbb{R}^3 : |\xi_1| \leq M, |\xi_2| \leq M|\xi_1|, |\xi_3| \leq M|\xi_1|\}$$

Therefore for every  $\xi \in \mathbb{R}^3 \setminus K_M$  :

$$\begin{aligned} \frac{\|g(x_{ij}, \beta, \Lambda\xi)\|^2}{\|\xi\|^2} &\leq \frac{\lambda_1^2 \xi_1^2}{\xi_1^2 + \xi_2^2 + \xi_3^2} \\ &\leq \frac{\lambda_1^2 \xi_1^2}{M^2 \xi_1^2} \\ &\leq \frac{\lambda_1^2}{M^2} \end{aligned}$$

By taking  $M > \frac{\lambda_1}{\varepsilon}$ , it verifies assumption (7).

## References

- [1] Donald WK Andrews. Tests for parameter instability and structural change with unknown change point. *Econometrica: Journal of the Econometric Society*, pages 821–856, 1993.
- [2] Donald WK Andrews. Estimation when a parameter is on a boundary. *Econometrica*, 67(6):1341–1383, 1999.
- [3] Donald WK Andrews. Inconsistency of the bootstrap when a parameter is on the boundary of the parameter space. *Econometrica*, pages 399–405, 2000.
- [4] Charlotte Baey, Paul-Henry Cournède, and Estelle Kuhn. Asymptotic distribution of likelihood ratio test statistics for variance components in nonlinear mixed effects models. *Computational Statistics & Data Analysis*, 135:107–122, 2019.
- [5] Rudolf Beran. Diagnosing bootstrap success. *Annals of the Institute of Statistical Mathematics*, 49(1):1–24, 1997.
- [6] Peter J Bickel and Anat Sakov. On the choice of m in the m out of n bootstrap and confidence bounds for extrema. *Statistica Sinica*, pages 967–985, 2008.

- [7] Benjamin M Bolker, Mollie E Brooks, Connie J Clark, Shane W Geange, John R Poulsen, M Henry H Stevens, and Jada-Simone S White. Generalized linear mixed models: a practical guide for ecology and evolution. *Trends in ecology & evolution*, 24(3):127–135, 2009.
- [8] Peter L Bonate. Nonlinear mixed effects models: theory. *Pharmacokinetic-pharmacodynamic modeling and simulation*, pages 233–301, 2011.
- [9] Helen Brown and Robin Prescott. *Applied mixed models in medicine*. John Wiley & Sons, 2015.
- [10] Alexander V Bulinski. Conditional central limit theorem. *Theory of Probability & Its Applications*, 61(4):613–631, 2017.
- [11] Giuseppe Cavaliere, Heino Bohn Nielsen, Rasmus Søndergaard Pedersen, and Anders Rahbek. Bootstrap inference on the boundary of the parameter space, with application to conditional volatility models. *Journal of Econometrics*, 2020.
- [12] D Chant. On Asymptotic Tests of Composite Hypotheses in Nonstandard Conditions. *Biometrika*, 61(2):291–298, 1974.
- [13] Zhen Chen and David B Dunson. Random effects selection in linear mixed models. *Biometrics*, 59(4):762–769, 2003.
- [14] Herman Chernoff. On the Distribution of the Likelihood Ratio. *The Annals of Mathematical Statistics*, 25(3):573–578, 1954.
- [15] Marie Davidian and David M Giltinan. *Nonlinear models for repeated measurement data*. Routledge, 2017.
- [16] Maud Delattre, Marc Lavielle, and Marie-Anne Poursat. A note on BIC in mixed-effects models. *Electronic Journal of Statistics*, 8(1):456 – 475, 2014.
- [17] R. Drikvandi, G. Verbeke, A. Khodadadi, and V. Partovi Nia. Testing multiple variance components in linear mixed-effects models. *Biostatistics*, 14(1):144–159, 2013.

- [18] Karl Oskar Ekvall and Matteo Bottai. Confidence regions near singular information and boundary points with applications to mixed models. *arXiv preprint arXiv:2103.10236*, 2021.
- [19] Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*. Cambridge university press, 2021.
- [20] Katherine R Gordon. How mixed-effects modeling can advance our understanding of learning and memory and improve clinical and educational practice. *Journal of Speech, Language, and Hearing Research*, 62(3):507–524, 2019.
- [21] Wolfgang Goymann, Ignas Safari, Christina Muck, and Ingrid Schwabl. Sex roles, parental care and offspring growth in two contrasting coucal species. *Royal Society Open Science*, 3(10):160463, 2016.
- [22] Andreas Groll and Gerhard Tutz. Variable selection for generalized linear mixed models by l1-penalized estimation. *Statistics and Computing*, 24:137–154, 2014.
- [23] Matthew J Gurka. Selecting the best linear mixed model under reml. *The American Statistician*, 60(1):19–26, 2006.
- [24] Bruce Hansen. *Econometrics*. Princeton University Press, 2022.
- [25] Kasahara Hiroyuki, Shimotsu Katsumi, et al. Testing the number of components in finite mixture models. Technical report, Institute of Economic Research, Hitotsubashi University, 2012.
- [26] Bruce Hoadley. Asymptotic properties of maximum likelihood estimators for the independent not identically distributed case. *The Annals of mathematical statistics*, pages 1977–1991, 1971.
- [27] Joseph G Ibrahim, Hongtu Zhu, Ramon I Garcia, and Ruixin Guo. Fixed and random effects selection in mixed effects models. *Biometrics*, 67(2):495–503, 2011.
- [28] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2020.

- [29] Lotte Meteyard and Robert A.I. Davies. Best practice guidance for linear mixed-effects models in psychological science. *Journal of Memory and Language*, 112:104092, 2020.
- [30] PAP Moran. The uniform consistency of maximum-likelihood estimators. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 70, pages 435–439. Cambridge University Press, 1971.
- [31] Lei Nie. Strong consistency of the maximum likelihood estimator in generalized linear and nonlinear mixed-effects models. *Metrika*, 63(2):123–143, 2006.
- [32] José Pinheiro and Douglas Bates. *Mixed-effects models in S and S-PLUS*. Springer science & business media, 2006.
- [33] Steven G Self and Kung-Yee Liang. Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association*, 82(398):605–610, 1987.
- [34] Mervyn J Silvapulle and Pranab Kumar Sen. *Constrained statistical inference: Inequality, order and shape restrictions*. John Wiley & Sons, 2005.
- [35] Sanjoy K Sinha. Bootstrap tests for variance components in generalized linear mixed models. *Canadian Journal of Statistics*, 37(2):219–234, 2009.
- [36] Florin Vaida and Suzette Blanchard. Conditional akaike information for mixed effects models. *Corrado Lagazio, Marco Marchi (Eds)*, page 101, 2005.
- [37] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [38] S. S. Wilks. The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses. *The Annals of Mathematical Statistics*, 9(1):60 – 62, 1938.
- [39] Xinbin Zhou, Gerard B.M. Heuvelink, Yusuke Kono, Tsutomu Matsui, and Takashi S.T. Tanaka. Using linear mixed-effects modeling to evaluate the impact of edaphic factors

on spatial variation in winter wheat grain yield in japanese consolidated paddy fields.  
*European Journal of Agronomy*, 133:126447, 2022.