

Markov Chain Order Estimation and χ^2 – *divergence* measure

A.R. Baigorri*
Mathematics Department
UnB

C.R. Gonçalves[†]
Mathematics Department
UnB

P.A.A. Resende[‡]
Mathematics Department
UnB

March 01, 2012

Abstract

We use the χ^2 – *divergence* as a measure of diversity between probability densities and review the basic properties of the estimator $\Delta_2(\cdot, \cdot)$. We define a few objects which capture relevant information from the sample of a Markov Chain to be used in the definition of a couple of estimators i.e. the *Local Dependency Level* and *Global Dependency Level* for a Markov chain sample. After exploring their properties we propose a new estimator for the Markov chain order. Finally we show a few tables containing numerical simulation results, comparing the performance of the new estimator with the well known and already established AIC, BIC and EDC estimators.

*baig@unb.br

[†]The author was partially supported by PROCAD/CAPES, CASADINHO/CAPES and PRONEX/FAPDF — catia@mat.unb.br

[‡]pa@mat.unb.br

1 Introduction

A Markov Chain is a discrete stochastic process $\mathbb{X} = \{X_n\}_{n \geq 0}$ with state space E , cardinality $|E| < \infty$ for which there is a $k \geq 1$ such that for $n \geq k$, $(x_1, \dots, x_n) \in E^n$

$$P(\mathbf{X}_1 = x_1, \dots, \mathbf{X}_n = x_n) = P(\mathbf{X}_1 = x_1, \dots, \mathbf{X}_k = x_k) \prod_{i=k+1}^n Q(x_i | x_{i-k}, \dots, x_{i-1})$$

for suitable transition probabilities $Q(\cdot|\cdot)$. The class of processes that holds the above condition for a given $k \geq 1$ will be denoted by \mathcal{M}_k , and \mathcal{M}_0 will denote the class of i.i.d. processes. The order of a process in $\cup_{i=0}^{\infty} \mathcal{M}_i$ is the smallest integer κ such that $\mathbb{X} = \{X_n\}_{n \geq 0} \in \mathcal{M}_\kappa$.

Along the last few decades there has been a great number of research on the estimation of the order of a Markov Chains, starting with M.S. Bartlett [6], P.G. Hoel [16], I.J. Good [15], T.W. Anderson & L.A. Goodman [4], P. Billingsley [7], [8] among others, and more recently, H. Tong [24], G. Schwarz [22], R.W. Katz [17], I. Csiszar and P. Shields [11], L.C. Zhao et al [25] had contributed with new Markov chain order estimators.

Since 1973, H. Akaike [1] entropic information criterion, known as AIC, has had a fundamental impact in statistical model evaluation problems. The AIC has been applied by Tong, for example, to the problem of estimating the order of autoregressive processes, autoregressive integrated moving average processes, and Markov chains. The Akaike-Tong (AIC) estimator was derived as an asymptotic approximate estimate of the Kullback-Leibler information discrepancy and provides a useful tool for evaluating models estimated by the maximum likelihood method. Later on, Katz derived the asymptotic distribution of the estimator and showed its inconsistency, proving that there is a positive probability of overestimating the true order no matter how large the sample size. Nevertheless, AIC is the most used and succesfull Markov chain order estimator used at the present time, mainly because it is more efficient than BIC for small sample.

The main consistent estimator alternative, the BIC estimator, does not perform too well for relatively small samples, as it was pointed out by Katz [17] and Csiszar & Shields [11]. It is natural to admit that the expansion of the Markov Chain complexity (size of the state space and order) has significant influence on the sample size required for the identification of the unknown order, even though, most of the time it is difficult to obtain sufficiently large samples.

In this notes we'll use a different entropic object called χ^2 - *divergence*, and

48 study its behaviour when applied to samples from random variables with
49 multinomial empirical distributions

$$50 \quad \mathcal{X} = \{X_i\}_{1 \leq i \leq r}$$

51 derived from a Markov Chain sample. Finally, we shall propose a new
52 strongly consistent Markov Chain order estimator more efficacious than the
53 already established AIC and BIC, which it shall be exhibited through the
54 outcomes of several numerical simulations.

55 In Section 2 we succinctly review the concept of *f-divergence* and its
56 properties. In Section 3, the χ^2 -divergence estimator is defined reviewing
57 some results concerning its convergence, as well as we briefly elaborate about
58 the Law of Iterated Logarithm (LIL) for our particular situation. In Section 4
59 the Markov chain sample is brought to attention, some notation introduced
60 and the estimators *Local Dependency Level* and *Global Dependency Level*,
61 which are the groundsill of the consistent Markov chain order estimator,
62 subsequently defined. Finally, in Section 4 we describe the procedures used
63 and the results obtained in an exploratory numerical simulations.

64 2 Entropy and f-divergences

65 2.1 Definitions and Notations

66 An *f-divergence* is a function that measures the discrepancy between two
67 probability distributions P and Q . The divergence is intuitively an average
68 of the function f of the odds ratio given by P and Q .

69 These divergences were introduced and studied independently by Csiszar,
70 Csiszar&Shields and Ali&Silvey among others ([10], [12], [3]) and sometimes
71 are referred as *Ali-Silvey distances*.

72 **Definition 2.1.** Let P and Q be discrete probability densities with sup-
73 port $S(P) = S(Q) = E = \{1, \dots, m\}$. For $f(t)$ convex function defined for
74 $t > 0$, $f(1) = 0$, the *f-divergence* for the distributions P and Q is

$$74 \quad D_f(P||Q) = \sum_{a \in A} Q(a) f\left(\frac{P(a)}{Q(a)}\right).$$

75 Here we take $0f(\frac{0}{0}) = 0$, $f(0) = \lim_{t \rightarrow 0} f(t)$, $0f(\frac{a}{0}) = \lim_{t \rightarrow 0} tf(\frac{a}{t}) =$
76 $a \lim_{u \rightarrow \infty} \frac{f(u)}{u}$. ♦

77 For example:

$$78 \quad f(t) = t \log(t) \Rightarrow D_f(P\|Q) = D(P\|Q) = \sum_{a \in A} P(a) \log \left(\frac{P(a)}{Q(a)} \right),$$

$$79 \quad f(t) = (1 - t^2) \Rightarrow D_f(P\|Q) = \sum_{a \in A} \frac{(P(a) - Q(a))^2}{Q(a)},$$

80 which are called *relative entropy* and χ^2 - *divergence*, respectively. From
81 now on the χ^2 - *divergence* shall be denote by $D_2(P\|Q)$.

82 Observe that the triangular inequality is not satisfied in general, so that
83 $D_2(P\|Q)$ defines no distance in the strict sense.

84 A basic theorem about *f-divergences* is the following approximation by the
85 $D_2(P\|Q)$.

86 **Theorem 2.1.** (*Csiszar & Shields [12]*) If f is twice differentiable at $t=1$ and
87 $f''(1) > 0$ then for any Q with support $S(Q) = A$ and P **close** to Q

$$88 \quad D_f(P\|Q) \sim \frac{f''(1)}{2} D_2(P\|Q).$$

89 Formally, $D_f(P\|Q)/D_2(P\|Q) \rightarrow f''(1)/2$ as $P \xrightarrow{D} Q$ ♦

90 The χ^2 -square divergence $D_2(P\|Q)$ test is well known statistical test proce-
91 dure close related to the chi-square distribution. See [19] for thorough and
92 detailed references.

93 3 Derived Markov Chains

94 Let $\mathbf{X}_1^n = (X_1, \dots, X_n)$ be a sample from a multiple stationary Markov chain
95 $\mathbb{X} = \{X_n\}_{n \geq 1}$ of unknown order κ . Assume that \mathbb{X} take values on a finite
96 state space $E = \{1, 2, \dots, m\}$ with transition probabilities given by

$$97 \quad p(x_{\kappa+1}|x_1^\kappa) = P(X_{n+1} = x_{n+1} | X_{n-\kappa+1}^n = x_1^\kappa) > 0 \quad (1)$$

98 where $x_1^\kappa = x_1^j x_{j+1}^\kappa = (x_1, \dots, x_\kappa) \in E^\kappa$.

99 Following Doob [13], from the process \mathbb{X} we can derive a first order MC,
 100 $\mathbb{Y}^{(\kappa)} = \{Y_n^{(\kappa)}\}_{n \geq 0}$ by setting $Y_n^{(\kappa)} = (X_n, \dots, X_{n+\kappa-1})$ so that for $v = (i_1, \dots, i_\kappa)$
 101 and $w = (i'_1, \dots, i'_\kappa)$

$$102 \quad P(Y_{n+1}^{(\kappa)} = w | Y_n^{(\kappa)} = v) = \tilde{p}_{vw} = \begin{cases} p(i'_\kappa | i_1 \dots i_\kappa), & i'_j = i_{j+1}, \quad j = 1, \dots, (\kappa - 1) \\ 0, & \text{otherwise.} \end{cases}$$

103 Clearly $\mathbb{Y}^{(\kappa)}$ is a first order and homogeneous MC that from now on shall be
 104 called the derived process, which by (1) is irreducible and positive recurrent
 105 MC having unique stationary distribution, say Π_κ . It is well known, see
 106 [[13]-Chap. 5.3], that the derived Markov Chains $\mathbb{Y}^{(l)}$, $l \geq \kappa$ is irreducible
 107 and aperiodic, consequently ergodic.

108 There exists an equilibrium (stationary) distribution $\Pi_\kappa(\cdot)$ satisfying for any
 109 initial distribution ν on E^κ

$$110 \quad \lim_{n \rightarrow \infty} |P_\nu(Y_n^{(\kappa)} = x_1^\kappa) - \Pi_\kappa(x_1^\kappa)| = 0,$$

111 and

$$112 \quad \Pi_\kappa(x_1^\kappa) = \sum_{z_1^\kappa} \Pi_\kappa(z_1^\kappa) p(x_\kappa | z_1^\kappa) = \sum_x \Pi_\kappa(x x_1^{\kappa-1}) p(x_\kappa | x x_1^{\kappa-1}).$$

113 Likewise, for $\mathbb{Y}^{(l)}$, $l > \kappa$

$$114 \quad \Pi_l(x_1^l) = \Pi_\kappa(x_1^\kappa) p(x_{\kappa+1} | x_1^\kappa) \dots p(x_l | x_{l-\kappa}^{l-1}) = \sum_x \Pi_l(x x_1^{l-1}) p(x_l | x x_{l-\kappa}^{l-1}). \quad (2)$$

115 which shows that Π_l defined above, is a stationary distribution for $\mathbb{Y}^{(l)}$. For
 116 the sake of notation's simplicity we'll use, from now on

$$117 \quad \boxed{\Pi(a_1^l) = \Pi_l(a_1^l), \quad l \geq \kappa.} \quad (3)$$

118 Now, let us return to $\mathbf{X}_1^n = (X_1, X_2, \dots, X_n)$ and define

$$119 \quad N(x_1^l | \mathbf{X}_1^n) = \sum_{j=1}^{n-l+1} 1(X_j = x_1, \dots, X_{j+l-1} = x_l) \quad (4)$$

120 i.e. the number of occurrences of x_1^l in X_1^n . If $l = 0$ we take $N(\cdot | \mathbf{X}_1^n) = n$. The
 121 sums are taken over positive terms $N(x_1^{l+1} | \mathbf{X}_1^n) > 0$, or else, we convention
 122 $0/0$ or $0.\infty$ as 0 .

123 Now we define the empirical random variables $\mathbf{X}_{i\alpha}$, for $i \in E$ and $\alpha \in E^\eta$.

124 **Definition 3.1.** For $\alpha = (a_1, \dots, a_\eta) = a_1^\eta \in E^\eta$ and $i \in E$, let $X_{i\alpha}$ be
 125 the random variable taking values in E , extracted from the MC sample \mathbf{X}_1^n ,
 126 defined as

$$127 \quad P(X_{i\alpha} = l) = \frac{N(i a_1^\eta l | \mathbf{X}_1^n)}{N(i a_1^\eta | \mathbf{X}_1^n)}, \quad l \in E. \quad (5)$$

128 with

$$129 \quad \mathbf{X}_{i\alpha} = \left(X_{i\alpha}^{(1)}, \dots, X_{i\alpha}^{(n_{i\alpha})} \right) \text{ its sample of size } n_{i\alpha}.$$

130 Observe that for $i, j \in E$

$$131 \quad \mathbf{O}_n^\alpha(i, j) = N(i a_1^\eta | \mathbf{X}_1^n)$$

132 where \mathbf{O}_n^α is the empirical random variables that describe the $X_{i\alpha}$, $1 \leq i \leq m$
 133 observed frequencies. Likewise, we define the expected frequencies

$$134 \quad \mathbf{E}_n^\alpha(i, j) = \frac{\sum_l O_n^\alpha(i, l) \sum_l O_n^\alpha(l, j)}{\sum_{kl} O_n^\alpha(k, l)}$$

135 and the respective probability functions

$$136 \quad \mathbf{P}_{\mathbf{O}_n^\alpha}(i, j) = \frac{\mathbf{O}_n^\alpha(i, j)}{N(a_1^\eta | \mathbf{X}_1^n)}, \quad i, j \in E$$

$$137 \quad \mathbf{P}_{\mathbf{E}_n^\alpha}(i, j) = \frac{\mathbf{E}_n^\alpha(i, j)}{N(a_1^\eta | \mathbf{X}_1^n)}, \quad i, j \in E.$$

138 Finally the χ^2 -square divergence

$$\begin{aligned} 139 \quad \hat{\Delta}_2(\mathbf{P}_{\mathbf{O}_n^\alpha} \| \mathbf{P}_{\mathbf{E}_n^\alpha}) &= n \sum_{i=1}^r \sum_{j=1}^m \frac{(\mathbf{P}_{\mathbf{O}_n^\alpha}(i, j) - \mathbf{P}_{\mathbf{E}_n^\alpha}(i, j))^2}{\mathbf{P}_{\mathbf{E}_n^\alpha}(i, j)} \\ 140 \quad &= n \Delta_2(\mathbf{P}_{\mathbf{O}_n^\alpha} \| \mathbf{P}_{\mathbf{E}_n^\alpha}). \end{aligned} \quad (6)$$

141 Now we derive a version of the Law of Iterated Logarithm, significant for the
 142 establishment of subsequent results about the convergence of $\hat{\Delta}_2(\mathbf{P}_{\mathbf{O}_n^\alpha} \parallel \mathbf{P}_{\mathbf{E}_n^\alpha})$.

143 **Lemma 3.1.** [18](Theorems 17.0.1 & 17.2.2) Let $\mathbb{X} = \{X_n\}_{n>0}$ be a er-
 144 godic Markov chain with finite state space E and stationary distribution Π ,
 145 $g : E \longrightarrow \mathbb{R}$, $S_n(g) = \sum_{j=1}^n g(X_j)$ and

$$146 \quad \sigma_g^2 = E_\pi(g^2(X_1)) + 2 \sum_{j=2}^n E_\pi(g(X_1)g(X_j))$$

147 then:

148 (a) If $\sigma_g^2 = 0$, then a.s. $\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}}[S_n(g) - E_\pi(S_n(g))] = 0$.

149 (b) If $\sigma_g^2 > 0$, then a.s.

$$150 \quad \limsup_{n \rightarrow \infty} \frac{S_n(g) - E_\pi(S_n(g))}{\sqrt{2 \sigma_g^2 n \log(\log(n))}} = 1$$

151 and

$$152 \quad \liminf_{n \rightarrow \infty} \frac{S_n(g) - E_\pi(S_n(g))}{\sqrt{2 \sigma_g^2 n \log(\log(n))}} = -1,$$

153 (E_Π : expectation with initial distribution Π ; a.s. : almost surely). \blacklozenge

154 **Lemma 3.2.** [14](Lemma 2) If $\mathbb{Y}^{(\kappa)}$ is ergodic then for $\eta \geq \kappa - 1$, $\alpha = a_1^\eta$
 155 and $i \alpha j = (i, a_1, \dots, a_\eta, j) = i a_1^\eta j \in E^{\eta+2}$ we have a.s.

$$156 \quad \limsup_{n \rightarrow \infty} \frac{(N(i a_1^\eta j \mid \mathbf{X}_1^n) - N(i a_i^\eta \mid \mathbf{X}_1^n) p(j \mid i a_1^\eta))^2}{n \log(\log(n))} = 2 \Pi(i a_1^\eta j)(1 - p(j \mid i a_1^\eta)). \quad \blacklozenge$$

157 **Theorem 3.3.** Let us refer to (6) for the definition of $\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \parallel P_{\mathbf{E}_n^\alpha})$, as well
 158 as the beginning of the present section for complementary definitions and
 159 references related to the following result:

160 If $\kappa \leq \eta$, there exist $\mathcal{L} < \infty$ so that for every $\alpha = i_1^\eta \in E^\eta$

$$161 \quad P \left(\limsup_{n \rightarrow \infty} \left[\frac{\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \| P_{\mathbf{E}_n^\alpha})}{2 \log(\log(n))} \right] \leq \mathcal{L} \right) = 1. \quad (7)$$

162 If $\eta = \kappa - 1$, there exist a_1^η & $i, j, k \neq i$ such that, $p(j | i a_1^\eta) \neq p(j | k a_1^\eta)$,
 163 consequently

$$164 \quad P \left(\limsup_{n \rightarrow \infty} \left[\frac{\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \| P_{\mathbf{E}_n^\alpha})}{2 \log(\log(n))} \right] = \infty \right) = 1. \quad \blacklozenge \quad (8)$$

165 **Proof:** The following proof shall be divided in the next two cases.

166 **Case I:** $0 \leq \kappa \leq \eta$.

167 From ([25], Lemma 3.1) and by Definition ?? we can calculate

$$168 \quad \mathbf{O}_n^\alpha(i, j) - \mathbf{E}_n^\alpha(i, j) = N(i a_1^\eta j | \mathbf{X}_1^n) - \frac{N(i a_1^\eta | \mathbf{X}_1^n) N(a_1^\eta j | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)}$$

169 or, in the limit

$$170 \quad \lim_{n \rightarrow \infty} \left(\mathbf{O}_n^\alpha(i, j) - \mathbf{E}_n^\alpha(i, j) \right)^2 = \lim_{n \rightarrow \infty} \left(N(i a_1^\eta j | \mathbf{X}_1^n) - N(i a_1^\eta | \mathbf{X}_1^n) p(j | i a_1^\eta) \right)^2$$

$$\begin{aligned} 171 \quad & \limsup_{n \rightarrow \infty} \frac{\left(\mathbf{O}_n^\alpha(i, j) - \mathbf{E}_n^\alpha(i, j) \right)^2}{n \log(\log(n)) \mathbf{P}_{\mathbf{E}_n^\alpha}(i, j)} = \\ 172 \quad & = \limsup_{n \rightarrow \infty} \left[\frac{\left(N(i a_1^\eta j | \mathbf{X}_1^n) - N(i a_1^\eta | \mathbf{X}_1^n) p(j | i a_1^\eta) \right)^2}{n \log(\log(n))} \frac{1}{\mathbf{P}_{\mathbf{E}_n^\alpha}(i, j)} \right]. \end{aligned}$$

174 Similarly

$$\begin{aligned}
175 \quad & \lim_{n \rightarrow \infty} \mathbf{P}_{\mathbf{E}_n^\alpha}(i, j) = \lim_{n \rightarrow \infty} \frac{\mathbf{E}_n^\alpha(i, j)}{N(a_1^\eta | \mathbf{X}_1^n)} = \lim_{n \rightarrow \infty} \left(\frac{N(i a_1^\eta | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)} \frac{N(a_1^\eta j | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)} \right) = \\
176 \quad & = \lim_{n \rightarrow \infty} \left(\frac{N(i a_1^\eta | \mathbf{X}_1^n)}{n} \frac{n}{N(a_1^\eta | \mathbf{X}_1^n)} \frac{N(a_1^\eta j | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)} \right) = \Pi(i a_1^\eta) \frac{1}{\Pi(a_1^\eta)} p(j | a_1^\eta) = \\
177 \quad & = \theta(i, j) > 0.
\end{aligned}$$

178 By (1) and Lemma 3.2 we have that $\min_{i,j} \theta(i, j) > 0$ with

$$179 \quad \mathcal{L} = \min_{i,j} \theta(i, j) \sum_{i=1}^m \sum_{j=1}^m \Pi(i a_1^\eta j) (1 - p(j | i a_1^\eta)) \leq 1$$

$$180 \quad P \left(\limsup_{n \rightarrow \infty} \frac{\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \| P_{\mathbf{E}_n^\alpha})}{2 \log(\log(n))} \leq \mathcal{L} \right) = 1.$$

181 **Case II:** $\eta = \kappa - 1$.

182 In accordance with the following

$$\begin{aligned}
183 \quad & \lim_{n \rightarrow \infty} \frac{N(a_1^\eta | \mathbf{X}_1^n)}{n} = \lim_{n \rightarrow \infty} \sum_{a \in E} \frac{N(a a_1^\eta | \mathbf{X}_1^n)}{n} = \sum_{a \in E} \Pi(a a_1^\eta) \text{ a.s.} \\
184 \quad & \lim_{n \rightarrow \infty} \frac{N(i a_1^\eta | \mathbf{X}_1^n)}{n} = \Pi(i a_1^\eta) \text{ a.s.} \\
185 \quad &
\end{aligned}$$

186 we can obtain, as in previous case

$$\begin{aligned}
187 \quad & \lim_{n \rightarrow \infty} \mathbf{P}_{\mathbf{E}_n^\alpha}(i, j) = \lim_{n \rightarrow \infty} \frac{\mathbf{E}_n^\alpha(i, j)}{N(a_1^\eta | \mathbf{X}_1^n)} = \lim_{n \rightarrow \infty} \left(\frac{N(i a_1^\eta | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)} \frac{N(a_1^\eta j | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)} \right) = \\
188 \quad & = \frac{\Pi(i a_1^\eta)}{\sum_{a \in E} \Pi(a a_1^\eta)} \frac{\Pi(a_1^\eta j)}{\sum_{a \in E} \Pi(a a_1^\eta)} \neq 0, \\
189 \quad &
\end{aligned}$$

190 and

$$\begin{aligned}
191 \quad & \lim_{n \rightarrow \infty} \mathbf{P}_{\mathbf{O}_n^\alpha}(i, j) = \lim_{n \rightarrow \infty} \frac{\mathbf{O}_n^\alpha(i, j)}{N(a_1^\eta | \mathbf{X}_1^n)} = \lim_{n \rightarrow \infty} \left(\frac{N(i a_1^\eta j | \mathbf{X}_1^n)}{N(a_1^\eta | \mathbf{X}_1^n)} \right) = \\
192 \quad & = \frac{\Pi(i a_1^\eta) p(j | i a_1^\eta)}{\sum_{a \in E} \Pi(a a_1^\eta)} \neq 0.
\end{aligned}$$

Clearly, if $\eta = \kappa - 1$, there exist $\alpha = a_1^\eta$ & $i, j \in E$ so that

$$\lim_{n \rightarrow \infty} (\mathbf{P}_{\mathbf{O}_n^\alpha}(i, j) - \mathbf{P}_{\mathbf{E}_n^\alpha}(i, j)) \neq 0$$

since, otherwise, it should imply that

$$p(j | i a_1^\eta) = \frac{\Pi(a_1^\eta j)}{\sum_{a \in E} \Pi(a a_1^\eta)}$$

i.e. $p(j | i a_1^\eta)$ does not depend on $i \in E$, contradicting the assumption that the order $\kappa > \eta$.

$$P \left(\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \| P_{\mathbf{E}_n^\alpha}) = n O(1) \right) = 1$$

and (8) is proved. \checkmark

3.1 Local and Global Dependency Level

Herein we define the Local Dependency Level and the Global Dependency Level.

Definition 3.2. Let $\mathbf{X}_n = \{X_i\}_{i=1}^n$ be a sample of a Markov chain \mathbb{X} of order $\kappa \geq 0$ and $\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \| P_{\mathbf{E}_n^\alpha})$ with $\alpha = a_1^\eta$, $\eta \geq 0$ as previously defined.

Let us assume that V is a χ^2 random variable with $(m-1)^2$ degrees of freedom where \mathcal{P} is the continuous strictly decreasing function $\mathcal{P} : \mathbb{R}^+ \rightarrow [0, 1]$

$$\mathcal{P}(x) = P(V \geq x), \quad x \in \mathbb{R}^+.$$

We define the Local Dependency Level $\widehat{LDL}_n(a_1^\eta)$, for $\alpha = a_1^\eta$ as

$$\widehat{LDL}_n(a_1^\eta) = \frac{\hat{\Delta}_2(P_{\mathbf{O}_n^\alpha} \| P_{\mathbf{E}_n^\alpha})}{2 \log(\log(n))},$$

and the Global Dependency Level $\widehat{GDL}_n(\eta)$ as

$$\widehat{GDL}_n(\eta) = \mathcal{P} \left(\sum_{a_1^\eta \in E^\eta} \left(\frac{N(a_1^\eta | \mathbf{X}_1^n)}{n} \right) \widehat{LDL}_n(a_1^\eta) \right). \quad \blacklozenge$$

213 Observe that, if the hypothesis \mathbf{H}_0^α is true, then $\forall a_1^\eta, \eta \geq \kappa$,

$$214 \quad P\left(\liminf_{n \rightarrow \infty} \left(\widehat{GDL}_n(\eta)\right) \geq \mathcal{P}(\mathcal{L})\right) = 1 \quad (9)$$

215 and for $\eta = \kappa - 1$

$$216 \quad P\left(\lim_{n \rightarrow \infty} \left(\widehat{GDL}_n(\eta)\right) = \mathcal{P}(\infty) = 0\right) = 1. \quad (10)$$

217 By (9) and (10) it is clear that, for n sufficiently large,

$$218 \quad P\left(\widehat{GDL}_n(\eta) \approx 0\right) = 1, \quad \eta = \kappa - 1,$$

219 and

$$220 \quad P\left(\widehat{GDL}_n(\eta) \approx \mathcal{P}(\mathcal{L})\right) = 1, \quad \eta \geq \kappa.$$

221 and consequently, for a multiple stationary Markov chain $\mathbb{X}_{n \geq 1}$ of order κ

$$222 \quad \kappa = 0 \Leftrightarrow \lim_{n \rightarrow \infty} \widehat{GDL}_n(\eta) = \mathcal{P}(\mathcal{L}), \quad \eta = 0, 1, \dots, B,$$

$$223 \quad \kappa = \max_{0 \leq \eta \leq B} \left\{ \eta : \lim_{n \rightarrow \infty} \widehat{GDL}_n(\eta) = 0 \right\} + 1.$$

224 Finally, let us define the Markov chain order estimator based on the infor-
225 mation contained in the vector GDL_n .

226 **Definition 3.3.** Given a fixed number $0 < B \in \mathbb{N}$, let us define the set
227 $\mathcal{S} = \{0, 1\}^{B+1}$ and the application $T : \mathcal{S} \rightarrow \mathbb{N}$

$$228 \quad T(s) = -1 \Leftrightarrow s_i = 1, \quad i = 0, 1, \dots, B$$

$$229 \quad T(s) = \max_{0 \leq i \leq B} \{i : s_i = 0, s_{i+1} = \mathcal{P}(\mathcal{L})\}, \quad s = (s_0, s_1, \dots, s_B). \quad \blacklozenge$$

230 **Definition 3.4.** Let $\mathbf{X}_n = \{X_i\}_{i=1}^n$ be a sample for the Markov chain \mathbb{X} of
 231 order κ , $0 \leq \kappa < B \in \mathbb{N}$ and $\{\widehat{GDL}_n(i)\}_{i=1}^B$ as above. We define the order's
 232 estimator $\widehat{\kappa}_{GDL}(\mathbf{X}_n)$ as

$$233 \quad \widehat{\kappa}_{GDL}(\mathbf{X}_n) = T(\sigma_n) + 1$$

234 with $\sigma_n \in \mathcal{S}$ so that $\forall s \in \mathcal{S}$

$$235 \quad \sum_{i=0}^B (\widehat{GDL}_n(i) - \sigma_n(i))^2 \leq \sum_{i=0}^B (\widehat{GDL}_n(i) - s(i))^2. \quad \blacklozenge$$

236 By (9),(10) and (3.1) it is clear that, for n large enough, $\{GDL_n(i)\}_{i=1}^B$
 237 satisfies the hypothesis of therefore, the order estimator converges almost
 238 surely to its value, i.e.,

$$239 \quad \boxed{P\left(\lim_{n \rightarrow \infty} \widehat{\kappa}_{GDL}(\mathbf{X}_n) = \kappa\right) = 1, \quad \kappa = 0, 1, 2, \dots, B.} \quad (11)$$

240 4 Numerical Simulations

241 In what follows we shall compare the non-asymptotic performance, mainly
 242 for small samples, of some of the most used Markov chains order estimators.
 243 Recalling the previous notations $\alpha = (a_1, \dots, a_{k+1}) = a_1^{k+1}$, $N(i a_1^{k+1} | \mathbf{X}_1^n)$
 244 as in (4) and denoting

$$245 \quad \hat{L}(\eta) = \Pi_{a_1^{\eta+1}} \left[\frac{N(i a_1^{\eta+1} | \mathbf{X}_1^n)}{N(i a_1^\eta | \mathbf{X}_1^n)} \right]^{N(i a_1^{\eta+1} | \mathbf{X}_1^n)}$$

246 the estimators of the Markov chain order κ , are defined, under the hypothesis:

$$247 \quad \boxed{\text{There exist a known } B \text{ so that } 0 \leq \kappa \leq B}$$

248 as

$$\begin{aligned} 249 \quad \widehat{\kappa}_{AIC} &= \operatorname{argmin}\{AIC(\eta); \eta = 0, 1, \dots, B\}, \\ 250 \quad \widehat{\kappa}_{BIC} &= \operatorname{argmin}\{BIC(\eta); \eta = 0, 1, \dots, B\}, \\ 251 \quad \widehat{\kappa}_{EDC} &= \operatorname{argmin}\{EDC(\eta); \eta = 0, 1, \dots, B\}, \end{aligned}$$

252 where

$$\begin{aligned}
253 \quad AIC(\eta) &= -2 \log \hat{L}(\eta) & + & & |E|^{\eta+1} 2 (|E| - 1), \\
254 \quad BIC(\eta) &= -2 \log \hat{L}(\eta) & + & & |E|^{\eta+1} 2 (|E| - 1) \left(\frac{\log(n)}{2} \right), \\
255 \quad EDC(\eta) &= -2 \log \hat{L}(\eta) & + & & |E|^{\eta+1} 2 (|E| - 1) \left(\frac{\log \log(n)}{2(|E| - 1)} \right), \\
256 \quad AIC(\eta) &\leq EDC(\eta) \leq BIC(\eta).
\end{aligned}$$

257 Clearly, for a given η , the order estimator $GDL(\eta)$, as well as $AIC(\eta)$ [24],
258 $BIC(\eta)$ [22] and $EDC(\eta)$ [25, 14] contain much of the information concerning
259 the sample's relative dependency, nevertheless numerical simulations as well
260 as theoretical considerations anticipates a great deal of variability for small
261 samples.

262 The following numerical simulation, based on an algorithm due to Raftery[21],
263 starts on with the generation of a Markov chain transition matrix, $\mathbf{Q} =$
264 $(q_{i_1 i_2 \dots i_\kappa; i_{\kappa+1}})$ with entries

$$265 \quad q_{i_1 i_2 \dots i_\kappa; i_{\kappa+1}} = \sum_{t=1}^{\kappa} \lambda_{i_t} R(i_{\kappa+1}, i_t), \quad 1 \leq i_t, i_{\kappa+1} \leq m. \quad (12)$$

266 where the matrix

$$267 \quad R(i, j), \quad 0 \leq i, j \leq m, \quad \sum_{i=1}^m R(i, j) = 1, \quad 1 \leq j \leq m$$

268 and the positive numbers

$$269 \quad \{\lambda_i\}_{i=1}^{\kappa}, \quad \sum_{i=1}^{\kappa} \lambda_i = 1$$

270 are arbitrarily chosen in advance.

271 Once the matrix $\mathbf{Q} = (q_{i_1 i_2 \dots i_\kappa; i_{\kappa+1}})$ is obtained, two hundreds replications of
272 the Markov chain sample of size n , space state E and transition matrix \mathbf{Q}
273 are generated to compare $GDL(\eta)$ performance against the standards, well
274 known and already established order estimators just mentioned above.

275 Katz(1981) [17] obtained the asymptotic distribution of $\hat{\kappa}_{AIC}$ and proved its
276 inconsistency showing the existence of a positive probability to overestimate
277 the order. See also Shibata(1976) [23].

278 On the contrary Schwarz (1978) [22] and Zhao(2001) [25] proved strong con-
 279 sistency for the estimators $\hat{\kappa}_{BIC}$ and $\hat{\kappa}_{EDC}$, respectively.

280 It is quite intuitive that the random information regarding the order of a
 281 Markov chain, is spread over an exponentially growing set of empirical dis-
 282 tributions Θ with $|\Theta| = m^{B+1}$, where \mathbf{B} is the maximum integer η , as in
 283 $\alpha = (i_1 i_2 \dots i_\eta)$. It seems reasonable to think that a small *viable* sample,
 284 i.e. samples able to retrieve enough information to estimate the chain order,
 285 should have size $n \approx O(m^{B+1})$. Keeping in mind that for the present nu-
 286 merical simulation, the maximum length to be used is $B = 5$, from now on
 287 the sample sizes for $|E| = 3$ and $|E| = 4$ should be $n \approx 1.500$ and $n \approx 5.000$,
 288 respectively.

289 Finally, after applying all estimators to each one of the replicated samples,
 290 the final results are registered in the form of tables.

291 **Case I: Markov Chain Examples with $\kappa = 0$, $|E| = 3$.**

292 Firstly, we choose the matrix $\{Q_1, Q_2, Q_3\}$ to produce samples with sizes
 293 $500 \leq n \leq 2.000$, originated from Markov chains of order $\kappa = 0$ with quite
 294 different probability distributions.

$$295 \quad Q_1 = \begin{bmatrix} 0.33 & 0.335 & 0.335 \\ 0.33 & 0.335 & 0.335 \\ 0.33 & 0.335 & 0.335 \end{bmatrix}, Q_2 = \begin{bmatrix} 0.05 & 0.475 & 0.475 \\ 0.05 & 0.475 & 0.475 \\ 0.05 & 0.475 & 0.475 \end{bmatrix}, Q_3 = \begin{bmatrix} 0.05 & 0.05 & 0.90 \\ 0.05 & 0.05 & 0.90 \\ 0.05 & 0.05 & 0.90 \end{bmatrix}.$$

$ E = 3 \quad \leftrightarrow \quad \kappa = 0 \quad \leftrightarrow \quad \lambda_i = 1/3, i = 1, 2, 3.$												
	Q_1				Q_1				Q_1			
	$n = 500$				$n = 1.000$				$n = 1.500$			
k	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl
296 0	75.5%	100%	100%	99%	80%	100%	100%	99.5%	71.5%	100%	100%	99%
1	24.5%			1%	18%			0.5%	22.5%			1%
2					2%				6%			
3												
4												

$ E = 3 \quad \leftrightarrow \quad \kappa = 0 \quad \leftrightarrow \quad \lambda_i = 1/3, i = 1, 2, 3.$												
	Q_2				Q_2				Q_2			
	$n = 1.000$				$n = 1.500$				$n = 500$			
k	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl
0	63.5%	100%	100%	99%	63%	100%	100%	99%	59%	100%	100%	99%
1	29%			1%	34.5%			1%	37%			1%
2	7.5%				2.5%				4%			
3												
4												

$ E = 3 \quad \leftrightarrow \quad \kappa = 0 \quad \leftrightarrow \quad \lambda_i = 1/3, i = 1, 2, 3.$												
	Q_3				Q_3				Q_3			
	$n = 1.000$				$n = 1.500$				$n = 2.000$			
k	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl
0	43%	100%	100%	98%	47%	100%	99.5%	96%	46%	100%	100%	97%
1	53%			2%	51.5%		0.5%	4%	50.5%			2%
2	4%				1.5%				3.5%			1%
3												
4												

Notice that for a fixed sample size $n = \{500, 1.000, 1.500, 2.000\}$, the order estimator $\hat{\kappa}_{AIC}$ steadily overestimate the real order $\kappa = 0$ with the excessiveness depending on the probability distribution of the Markov chain. Differently, the order estimators $\hat{\kappa}_{BIC}$, $\hat{\kappa}_{EDC}$ and $\hat{\kappa}_{GDL}$ show consistent performance, mainly obtaining the right order, free from the influence of the sample size and the generating matrix. Regarding $\hat{\kappa}_{BIC}$ and $\hat{\kappa}_{EDC}$ improved effect, most likely depends on their correcting factor, $\frac{\log(n)}{2}$ and $\left(\frac{\log \log(n)}{2(|E|-1)}\right)$ which tend to decrease the estimated order.

Case II: Markov Chain Examples with $\kappa = 3, |E| = 3$ and $\kappa = \{2, 3, 0\}, |E| = 4$

Secondly, we choose the matrix $\{Q_4, Q_5\}$ to produce samples with sizes $n \in \{500, 1.000, 1.500, 2.000\}$, originated from Markov chains for $|E| = 3$ of order $\kappa = 3$.

312

$$Q_4 = \begin{bmatrix} 0.05 & 0.05 & 0.90 \\ 0.05 & 0.90 & 0.05 \\ 0.90 & 0.05 & 0.05 \end{bmatrix}, \quad Q_5 = \begin{bmatrix} 0.475 & 0.475 & 0.05 \\ 0.475 & 0.05 & 0.475 \\ 0.05 & 0.475 & 0.475 \end{bmatrix}.$$

313

$ E = 3 \quad \leftrightarrow \quad \kappa = 3 \quad \leftrightarrow \quad \lambda_i = 1/3, i = 1, 2, 3.$												
	Q_4				Q_4				Q_4			
	$n = 1.000$				$n = 1.500$				$n = 2.000$			
k	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl
0												
1												
2		99.5%	88.5%	41%		76.5%	16.5%	5%		17%	0.5%	1%
3	100%	0.5%	11.5%	59%	100%	23.5%	83.5%	95%	100%	83%	99.5%	99%
4												

314

$ E = 3 \quad \leftrightarrow \quad \kappa = 3 \quad \leftrightarrow \quad \lambda_i = 1/3, i = 1, 2, 3.$												
	Q_5				Q_5				Q_5			
	$n = 1.000$				$n = 1.500$				$n = 2.500$			
k	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl
0		0.5%										
1		92.5%	69.5%	6.5%		54.5%	19.5%	1%				
2	16.5%	7%	30.5%	92%	2%	45.5%	80.5%	80.5%		100%	98.5%	8.5%
3	83.5%			1.5%	98%			18.5%	100%		1.5%	91.5%
4												

315 For $|E| = 3$, $\kappa = 3$ the estimator $\hat{\kappa}_{AIC}$ overestimate the order in a lesser
316 extent than the previous case, while $\hat{\kappa}_{BIC}$ and $\hat{\kappa}_{EDC}$ overweighted by the
317 respective constants $\frac{\log(n)}{2}$ and $\left(\frac{\log \log(n)}{2(|E|-1)}\right)$, underestimate the order more than
318 it was supposed to be. Concerning $\hat{\kappa}_{GDL}$, it rapidly converges to the right
319 order depending on the sample size n .

320 For $|E| = 4$ the greater complexity of a Markov chain of order $\kappa = 3$ impose
321 the use of larger sample size for estimators to accomplish some reliability.
322 Finally, we choose the matrix $\{Q_6, Q_7\}$ to produce samples with size $n =$
323 5.000, originated from Markov chains of order $\kappa \in \{2, 3, 0\}$ like in the previous
324 cases.

$$Q_6 = \begin{bmatrix} 0.05 & 0.05 & 0.05 & 0.85 \\ 0.05 & 0.05 & 0.85 & 0.05 \\ 0.05 & 0.85 & 0.05 & 0.05 \\ 0.85 & 0.05 & 0.05 & 0.05 \end{bmatrix}, \quad Q_7 = \begin{bmatrix} 0.05 & 0.05 & 0.05 & 0.85 \\ 0.05 & 0.05 & 0.05 & 0.85 \\ 0.05 & 0.05 & 0.05 & 0.85 \\ 0.05 & 0.05 & 0.05 & 0.85 \end{bmatrix}.$$

$ E = 4 \quad \leftrightarrow \quad n = 5.000$												
$Q_6 \Leftrightarrow \lambda_i = 1/2, i = 1, 2.$					$Q_6 \Leftrightarrow \lambda_i = 1/3, i = 1, 2, 3.$				$Q_7 \Leftrightarrow \lambda_i = 1/3, i = 1, 2, 3.$			
$\kappa = 2$					$\kappa = 3$				$\kappa = 0$			
k	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl	Aic	Bic	Edc	Gdl
0									85%	100%	100%	100%
1									15%			
2	100%	100%	100%	100%		99%		4%				
3					100%	1%	100%	96%				
4												
5												
6												

For the order for $|E| = 4$, $\kappa = 0$, apparently $\hat{\kappa}_{AIC}$ keeps overestimating the order in some degree, while $\hat{\kappa}_{BIC}$ as in example $\kappa = 3$ severely underestimate the order, presumably due to the excessive weight of the correcting factors $\frac{\log(n)}{2}$. On the contrary $\hat{\kappa}_{EDC}$ and $\hat{\kappa}_{GDL}$ behaves quite well in same setting.

5 Conclusion

The pioneer research started with the contributions of Bartlett[6], Hoel[16], Good [15], Anderson & Goodman [4], Billingsley([7], [8]) among others, where they developed tests of hypothesis for the estimation of the order of a given Markov chain.

Later on these procedures were adapted and improved with the used of *Penalty Functions* (Tong[24], Katz[17]) together with other tools created in the realm of Models Selection (Akaike[1], Schwarz[22]). Since then, there have been a considerable number of subsequent contributions on this subject, several of them consisting in the enhancement of the already existing techniques (Csiszar[11], Zhao et al[25]).

In this notes we propose a new Markov chain order estimator based on a different idea which makes it behave in a quite different form. This estimator is

strongly consistent and more efficient than AIC (inconsistent), outperforming the well established and consistent BIC and EDC, mainly on relatively small samples.

References

- [1] H. Akaike, “A New Look at the Statistical Model Identification”, *IEEE Trans. Autom. Cont.*, vol. 19, 1998.
- [2] H. Akaike & G. Kitagawa,(Eds.), “The Practice of Time Series Analysis”, *Springer-Verlag, New York*, 1974.
- [3] M. S. Ali & D. Silvey, “A general class of coefficients of divergence of one distribution from another”, *Journal of the Royal Statistical Society, Ser. B, No. 28, pp. 131-140*, 1966.
- [4] T. W. Anderson & Leo A. Goodman, “Statistical Inference about Markov Chains”, *Annals of Mathematical Statistics*, vol. 28, 1957.
- [5] A. Banerjee, X. Guo & H. Wang, “The Optimallity of Conditional Expectation as a Bregman Predictor”, *IEEE Transactions on Information Theory*, vol. 51, No 7, 2005.
- [6] M. S. Bartlett, “The Frequency Goodness of Fit Test for Probability Chains”, *Math. Proc. Camb. Phil. Soc.*, vol. 9, 1951.
- [7] P. Billingsley, “Statistical Methods in Markov Chains”, *Annals of Mathematical Statistics*, Vol. 32, No 1, 1961.
- [8] P. Billingsley, “Statistical Inference for Markov Chains”, *University of Chicago Press, Chicago*, 1961.
- [9] W. Cochran, “The χ^2 Test of Goodness of Fit”, *Annals of Mathematical Statistics*, vol. 23, No 3, 1952.
- [10] I. Csiszar, “Information-type measures of difference of probability distributions and indirect observation”, *Studia Sci. Math. Hungar.*, Vol. 2, pp. 229-318, 1967.
- [11] I. Csiszar & P. Shields, “The Consistency of the BIC Markov Order Estimator”, *The Annals of statistics*, Vol 28, No 6, pp. 1601-1619, 2000.

- [12] I. Csiszar & P. Shields, “Information Theory and Statistics: A Tutorial”, *Foundations and Trends in Communications and Information Theory*, Vol. 1, No 4, pp. 417-528, 2004.
- [13] J.L. Doob, “Stochastic Processes(Wiley Publication in Statistics)”, *John Wiley & Sons. Inc.*, 1966.
- [14] , C.C.Y. Dorea, “Optimal Penalty Term for EDC Markov Chain Order Estimator”, *Annales de l’Institut de Statistique de l’Universite de Paris (l’ISUP)*, v. 52, p. 15-26, 2008.
- [15] I.J. Good, “The Likelihood Ratio Test for Markov Chains”, *Biometrika*, vol. 42, (3/4), pp. 531-533, 1955.
- [16] P. G. Hoel, “A test for Markov Chains”, *Biometrika*, vol. 41, 1954.
- [17] R. W. Katz, “On Some Criteria for Estimating the Order Of Markov Chains”, *Technometrics*, vol. 23(3), 1981.
- [18] S.P. Meyn & R.L. Tweedie, “Markov Chain and Stochastic Stability”, *Springer-Verlag, London*, 1993.
- [19] L. Pardo, “Statistical Inference Based on Divergence Measures”, *Taylor & Francis Group, LLC, New York*, 2006.
- [20] R Development Core Team, “R: A language and environment for statistical computing”, *R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>*, 2008.
- [21] A. E. Raftery, “A model for High-Order Markov Chains”, *Journal of the Royal Statistical Society, Serie B*, vol. 47, No. 3, 1985.
- [22] G. Schwarz, “Estimating the Dimension of a Model”, *The Annals of Statistics*, vol. 6, No 2, 1978.
- [23] R. Shibata, “Selection of the Order of an Autoregressive model by Akaike’s Information Criterion.”, *Biometrika*, vol. 6, 1976.
- [24] H. Tong, “Determination of the order of a Markov Chain by Akaikes Information Criterion”, *Journal of Applied Probability*, vol. 12, pp. 488-497, 1975.
- [25] H. C. Zhao, C.C.Y. Dorea & C.R. Gonçalves, “On Determination of the Order of a Markov Chain”, *Stat. Infer. for Stoc. Processes*, vol. 4, 2001.