

Variable Length Markov Chains

Peter Bühlmann¹ and Abraham J. Wyner²

ETH Zürich, Switzerland and University of Pennsylvania, USA

February 1999

Abstract

We study estimation in the class of stationary variable length Markov chains (VLMC) on a finite space. The processes in this class are still Markovian of higher order, but with memory of variable length yielding a much bigger and structurally richer class of models than ordinary higher order Markov chains. From a more algorithmic view, the VLMC model class has attracted interest in information theory and machine learning but statistical properties have not been explored very much. Provided that good estimation is available, an additional structural richness of the model class enhances predictive power by finding a better trade-off between model bias and variance and allows better structural description which can be of specific interest. The latter is exemplified with some DNA data.

A version of the tree-structured context algorithm, proposed by Rissanen (1983) in an information theoretical set-up, is shown to have new good asymptotic properties for estimation in the class of VLMC's, even when the underlying model increases in dimensionality: consistent estimation of minimal state spaces and mixing properties of fitted models are given.

We also propose a new bootstrap scheme based on fitted VLMC's. We show its validity for quite general stationary categorical time series and for a broad range of statistical procedures.

AMS 1991 subject classifications. Primary 62M05; secondary 60J10, 62G09, 62M10, 94A15

Key words and phrases. Bootstrap, categorical time series, central limit theorem, context algorithm, data compression, finite-memory sources, FSMX model, Kullback-Leibler distance, model selection, tree model.

Short title: Variable Length Markov Chain

¹Research supported in part by the Swiss National Science Foundation. Part of the work has been done while visiting the University of Heidelberg, Germany.

²Research supported by grant NSF DMS 9508933.

1 Introduction

One of the most general models for a stationary process $(X_t)_{t \in \mathbb{Z}}$ assuming no particular underlying mechanistic system is maybe a full Markov chain of high, but finite order. The only implicit assumption aside from stationarity then made is the finite memory of the process. We consider here exclusively the case where X_t takes values in a finite categorical space \mathcal{X} . And we always refer to a stationary full Markov chain of order k , whenever the transition mechanism carries no specific structure so that the state space is the entire \mathcal{X}^k . Probabilistically a nice model, such full Markov chains aren't very appropriate from the estimation point of view. Let us illustrate two main problems; to be more specific, we momentarily take for illustrative purposes cardinality $|\mathcal{X}| = 4$, e.g., $\mathcal{X} = \{A, C, G, T\}$ being the letters of a DNA string. But all the discussed problems below apply to any finite space \mathcal{X} .

Problem 1: The class of all finite order \mathcal{X} -valued full Markov chains is not structurally rich, implying that there aren't many members in the class. This structural poorness particularly implies that any kind of parsimonious representation of the state space is not possible. The table below additionally demonstrates such structural poorness in terms of the dimension of full Markov chain models (the number of free parameters) as a function of their orders k , i.e., $\text{Dim.} = (|\mathcal{X}| - 1)|\mathcal{X}|^k$ with cardinality $|\mathcal{X}| = 4$.

k	0	1	2	3	4	5	10
Dim.	3	12	48	192	768	3072	$\approx 3.1 \cdot 10^6$

There are no models 'in between', e.g., it is impossible to fit a model with say 72 parameters. Such a very 'discontinuous' increase in dimensionality of the model does not allow a good trade-off between bias (being low with many parameters) and variance (being low with few parameters) of a predictor.

Problem 2: As seen from the table above, the curse of dimensionality is heavily present when fitting higher order models: the dimensionality increases exponentially with the order k leading then to highly variable estimates.

A practical example where the table and problems 1 and 2 above apply is with full Markov chain modeling of DNA sequences, cf. Prum et al. (1995) and Braun and Müller (1998). The class of models and the estimator which we study in this paper will lead to an alternative and for many purposes better statistical description of DNA sequences. We give in section 3.3 a real-data example from this field of applications. Other examples of applications where our modeling is potentially attractive are precipitation analysis (Guttorp, 1995), flood analysis (Brillinger, 1995), or analysis of discrete directional data and repeated patterns of behavioral events (Raftery and Tavaré, 1994).

Problems 1 and 2 above can be cured with a very simple idea: the memory of a stationary Markov chain is allowed to be of variable length, being a function of the values from the past. More precisely, the time-homogeneous transition probabilities $\mathbb{P}[X_t = x_t | X_{t-1} = x_{t-1}, X_{t-2} = x_{t-2}, \dots]$ are functions depending only on a variable number ℓ of lagged values $\mathbb{P}[X_t = x_t | X_{t-1} = x_{t-1}, \dots, X_{t-\ell} = x_{t-\ell}]$, where $\ell = \ell(x_{t-1}, x_{t-2}, \dots)$ is itself a function of the past. If $\ell(x_{t-1}, x_{t-2}, \dots) \equiv k$ for all x_{t-1}, x_{t-2}, \dots , we obtain the full Markov chain model of order k . For variable $\ell(\cdot)$ with $\sup\{\ell(x_{t-1}, x_{t-2}, \dots); x_{t-1}, x_{t-2}, \dots\} = k$, we have an (embedding) Markov chain of order k , but with an additional well interpretable *structure of a variable length memory*: the structure implies that some transition probabil-

ities of the embedding Markov chain are lumped together. We call such a process variable length Markov chain (VLMC). It is closely related to models in information theory like ‘tree models’, ‘FSMX models’ or ‘finite-memory sources’, cf. Rissanen (1986), Weinberger et al. (1992, 1995) and also Feder et al. (1992). If one is able to choose in a data-driven way an appropriate member in the class of VLMC’s, there is nothing to lose but only to gain in comparison with the class of ordinary full Markov chains of higher order. We give in this paper new results for VLMC’s, in particular also addressing the problem how to select in a data-driven way an asymptotically correct member in the extremely large class of all VLMC’s.

The merits of our results are in different areas. On a theoretical level we offer a better understanding of problems 1 and 2 above. With this goal in mind the study of VLMC’s and their estimation is ‘per se’ an interesting task. The practical merits of our results, which offer new insights into VLMC’s, apply generally to problems involving categorical or binary time series (see the examples mentioned above). We also offer advances in statistical methodology; specifically through a new bootstrap scheme, based on VLMC’s, for categorical time series which is a very attractive alternative to the more general block-wise bootstrap (Künsch, 1989). Moreover, fitted VLMC’s can be used as an excellent exploratory tool for the dynamics of a categorical time series. This is accomplished by representing structural dependencies graphically and compactly (we demonstrate this for some DNA data in section 3.3). Such explorative information could also be used to build a more specific parametric model in a second stage. Finally, we offer new insights into information theory, where our findings sharpen and extend existing results on VLMC’s and compression rates, see Rissanen (1983) and Weinberger et al. (1995).

The notion of a variable length memory in a Markov chain is particularly attractive when there is long memory in certain ‘directions’. In such cases, the *minimal* state space is drastically smaller than the embedding state space of a full Markov chain (having many equivalent states which are lumped together in the VLMC): the VLMC yields a *parsimonious* parameterization of the state space. The difficulty with this attractive notion is then the estimation of that minimal state space. This can be seen as a model selection problem. However, due to the extremely large number of VLMC sub-models of a higher order Markov chain, global model selection techniques like AIC, BIC or MDL cannot be used. But estimation of the minimal state space and the probability distribution of a VLMC can be done with a tree structured scheme, called the context algorithm (Rissanen, 1983), which acts hierarchically on *local* pairwise decisions. Weinberger et al. (1995) proved consistency and optimal compression rates for the context algorithm under the assumption that the true underlying process is a finite-dimensional VLMC.

We give an entirely new consistency result where the true underlying model is allowed to grow in dimensionality as sample size n increases. The growth rate is in probabilistic terms and can be as large as $O(n^{1/2}/\log(n)^s)$ for an arbitrary $s > 1/2$. This describes much better the performance of the context algorithm, because now with increasing dimensionality, consistency is much less obvious. As an important consequence, our result implies a non-trivial balance between over- and under-estimation of the true model. It can be loosely translated to the fact that a bias-variance tradeoff in a possibly very high dimensional problem is handled by the context algorithm in an appropriate way. Also, by allowing for asymptotically infinite dimensional models, the new results contribute to explore the approximation of a general, sufficiently ‘nice’ stationary process by an esti-

mated VLMC. Consistency of the context algorithm for estimating stationary processes or minimal state spaces does not require a pre-specified model structure. Thus, estimation with the context algorithm is robust against model-misspecification.

We then make use of the general consistency result described above to propose a novel resampling scheme, the VLMC bootstrap. We prove asymptotic validity of the VLMC bootstrap for a whole class of estimators and argue, by proving a mixing property for estimated VLMC's, why such a scheme works under very general conditions. The VLMC bootstrap is tailored for categorical time series, has a nice probabilistic interpretation and enjoys the advantage of being applicable as a simple plug-in rule, as Efron's (1979) original proposal for the independent case, which is more user-friendly than the blockwise bootstrap (Künsch, 1989). Based on the results in theory and from a small simulation study we conclude that the VLMC bootstrap is a new universal resampling tool for categorical time series which is often expected to be better than the blockwise bootstrap.

The paper is organized as follows. In section 2 we give the definition of VLMC's, section 3 describes the process of fitting such models and gives new asymptotic properties thereof. In section 4 we discuss the VLMC bootstrap, its asymptotic validity and present results from a simulation study, including a comparison with the blockwise bootstrap. All the proofs are given in section 5.

2 Variable length Markov chains

As a starting point, consider a stationary full Markov chain $(X_t)_{t \in \mathbb{Z}}$ of finite order k with values in a finite categorical space \mathcal{X} . In the sequel, we denote by $x_i^j = x_j, x_{j-1}, \dots, x_i$ ($i < j$, $i, j \in \mathbb{Z} \cup \{-\infty, \infty\}$) a string whose components are written in reverse order, $wu = (w_{|w|}, \dots, w_2, w_1, u_{|u|}, \dots, u_2, u_1)$ is the concatenation of the strings w and u . We usually denote by capital letters X random variables and by small letters x fixed deterministic values. Thus,

$$\mathbb{P}[X_1 = x_1 | X_{-\infty}^0 = x_{-\infty}^0] = \mathbb{P}[X_1 = x_1 | X_{-k+1}^0 = x_{-k+1}^0], \text{ for all } x_{-\infty}^0. \quad (2.1)$$

Note that by stationarity, the time indices $-\infty, \dots, 0, 1$ are irrelevant and can be replaced by other indices $-\infty, \dots, t-1, t$ for some $t \in \mathbb{Z}$. We now introduce the idea of a variable length memory which can also be seen as lumping together irrelevant states in the history x_{-k+1}^0 in formula (2.1). Only some values from the infinite history $x_{-\infty}^0$ of the variable X_1 are relevant: these can be thought of as a *context* for X_1 . To achieve a flexible model class, ranging from some type of sparse to full Markov chains, we let the length of a context depend on (the first few of) the actual values $x_{-\infty}^0$. We formalize this by defining below a VLMC. Related models have been introduced in information theory as tree models, FSMX models or finite-memory sources, cf. Rissanen (1986) and Weinberger et al. (1992, 1995).

Definition 2.1 *Let $(X_t)_{t \in \mathbb{Z}}$ be a stationary process with values $X_t \in \mathcal{X}$, $|\mathcal{X}| < \infty$. Denote by $c : \mathcal{X}^\infty \rightarrow \mathcal{X}^\infty$ a (variable projection) function which maps*

$$\begin{aligned} c : x_{-\infty}^0 &\mapsto x_{-\ell+1}^0, \text{ where } \ell \text{ is defined by} \\ \ell = \ell(x_{-\infty}^0) &= \min\{k; \mathbb{P}[X_1 = x_1 | X_{-\infty}^0 = x_{-\infty}^0] = \mathbb{P}[X_1 = x_1 | X_{-k+1}^0 = x_{-k+1}^0] \text{ for all } x_1 \in \mathcal{X}\}, \\ &\text{where } \ell \equiv 0 \text{ corresponds to independence.} \end{aligned}$$

Then, $c(\cdot)$ is called a context function and for any $t \in \mathbb{Z}$, $c(x_{-\infty}^{t-1})$ is called the context for the variable x_t .

The name *context* refers to the portion of the past that influences the next outcome. The definition of ℓ implicitly reflects the fact that the context-length of a variable X_t is $\ell = \ell(x_{-\infty}^{t-1}) = |c(x_{-\infty}^{t-1})|$, depending on the history $X_{-\infty}^{t-1} = x_{-\infty}^{t-1}$. By the projection structure of the context function $c(\cdot)$, the context-length $\ell(\cdot) = |c(\cdot)|$ determines $c(\cdot)$ and vice-versa.

Definition 2.2 Let $(X_t)_{t \in \mathbb{Z}}$ be a stationary process with values $X_t \in \mathcal{X}$, $|\mathcal{X}| < \infty$ and corresponding context function $c(\cdot)$ as given in Definition 2.1. Let $0 \leq k \leq \infty$ be the smallest integer such that

$$|c(x_{-\infty}^0)| = \ell(x_{-\infty}^0) \leq k \text{ for all } x_{-\infty}^0 \in \mathcal{X}^\infty.$$

Then $c(\cdot)$ is called a context function of order k , and if $k < \infty$, $(X_t)_{t \in \mathbb{Z}}$ is called a stationary variable length Markov chain (VLMC) of order k .

We sometimes identify a VLMC $(X_t)_{t \in \mathbb{Z}}$ with its probability distribution P_c on $\mathcal{X}^\mathbb{Z}$. In the sequel, we often write for a probability measure P on $\mathcal{X}^\mathbb{Z}$, $P(x) = \mathbb{P}_P[X_1^m = x]$ ($x \in \mathcal{X}^m$) and $P(x|w) = P(xw)/P(w)$ ($x, w \in \cup_{m=1}^\infty \mathcal{X}^m$). The transition probabilities for P_c are denoted by $p(x_1|c(x_{-\infty}^0))$ and coincide with the conditional probabilities $P_c(x_1|c(x_{-\infty}^0))$. Clearly, a VLMC of order k is a Markov chain of order k , now having a *memory of variable length* ℓ . If the context function $c(\cdot)$ of order k is the full projection $x_{-\infty}^0 \mapsto x_{-k+1}^0$ for all $x_{-\infty}^0$, the VLMC is a full Markov chain of order k . Often the range space of the context function $c(\cdot)$ is not the full space \mathcal{X}^k , but also not the empty space. The class of context functions of length k is rich enough to obtain a broad class of Markov chains, including special sparse types. In particular, some context functions $c(\cdot)$ yield a substantial reduction in the number of states compared to a full Markov chain of the same order as the context function. Both of the latter phrases relate to solutions to problems 1 and 2 mentioned in section 1.

2.1 Tree representation of minimal state space

By requiring stationarity, a VLMC P_c is completely specified by its transition probabilities,

$$\mathbb{P}_{P_c}[X_1 = x_1 | X_{-\infty}^0 = x_{-\infty}^0] = p(x_1 | c(x_{-\infty}^0)), \quad x_{-\infty}^0 \in \mathcal{X}^\infty.$$

The states determining these transition probabilities are thus given by the values of the context function $c(\cdot)$. It is most convenient to represent these states, i.e., the minimal state space of the VLMC P_c , as a tree.

We consider trees with a root node on top, from which the branches are growing downwards, so that every internal node has at most $|\mathcal{X}|$ offsprings. Then, each value of a context function $c(\cdot) : \mathcal{X}^\infty \rightarrow \mathcal{X}^k$ can be represented as a branch (or terminal node) of such a tree. The context $w = c(x_{-\infty}^0)$ is represented by a branch, whose sub-branch on the top is determined by x_0 , the next sub-branch by x_{-1} and so on, and the terminal sub-branch by $x_{-\ell(x_{-\infty}^0)+1}$. As we will exemplify, context trees do not have to be complete, i.e., every internal node does not need to have exactly $|\mathcal{X}|$ offsprings.

Example 2.1 $|\mathcal{X}| = 2, k = 3$.

The function

$$c(x_{-\infty}^0) = \begin{cases} 0, & \text{if } x_0 = 0, x_{-\infty}^{-1} \text{ arbitrary} \\ 1, 0, 0, & \text{if } x_0 = 1, x_{-1} = 0, x_{-2} = 0, x_{-\infty}^{-3} \text{ arbitrary} \\ 1, 0, 1, & \text{if } x_0 = 1, x_{-1} = 0, x_{-2} = 1, x_{-\infty}^{-3} \text{ arbitrary} \\ 1, 1, & \text{if } x_0 = 1, x_{-1} = 1, x_{-\infty}^{-2} \text{ arbitrary} \end{cases}$$

can be represented by the tree τ_c on the left hand side in Figure 2.1. A ‘growing to the left’ sub-branch represents the symbol 0 and vice versa for the symbol 1.

Definition 2.3 Let $c(\cdot)$ be a context function of a stationary VLMC. The $(|\mathcal{X}|$ -ary) context tree τ and terminal node context tree τ^T are defined as

$$\begin{aligned} \tau &= \tau_c = \{w; w = c(x_{-\infty}^0), x_{-\infty}^0 \in \mathcal{X}^\infty\}, \\ \tau^T &= \tau_c^T = \{w; w \in \tau_c \text{ and } wu \notin \tau_c \text{ for all } u \in \mathcal{X}\}. \end{aligned}$$

The notion of a terminal context tree is convenient when formulating an estimation procedure for a context tree (minimal state space) (see section 3.1). Definition 2.3 says that only terminal nodes in the tree representation τ are considered as elements of the terminal node context tree τ^T , and states $w \in \tau_c$ do not need to be terminal nodes in τ_c . But we can reconstruct the context function $c(\cdot)$ from either τ_c or τ_c^T . The context tree τ_c is nothing else than the minimal state space of the VLMC P_c (we sometimes refer to the elements of τ_c as branches and sometimes as nodes in a tree). An internal node with $b < N = |\mathcal{X}|$ offsprings implicitly adds one complementary offspring, lumping the $N - b$ non-present offsprings together to a single new terminal node w_{new} which represents a single state in τ_c .

Example 2.2 $|\mathcal{X}| = 4, k = 2$.

The function

$$c(x_{-\infty}^0) = \begin{cases} 0, & \text{if } x_0 = 0, x_{-\infty}^{-1} \text{ arbitrary} \\ 1, & \text{if } x_0 = 1, x_{-\infty}^{-1} \text{ arbitrary} \\ 2, & \text{if } x_0 = 2, x_{-\infty}^{-1} \text{ arbitrary} \\ 3, & \text{if } x_0 = 3, x_{-1} \in \{0, 1, 2\}, x_{-\infty}^{-2} \text{ arbitrary} \\ 3, 3 & \text{if } x_0 = 3, x_{-1} = 3, x_{-\infty}^{-2} \text{ arbitrary} \end{cases}$$

can be represented by the tree τ_c on the right hand side in Figure 2.1. The rounded rectangle, which we usually don’t draw, symbolizes the absent nodes 0,1 and 2 in depth 2, which can be thought as a completion of the tree with nodes lumped together: in terms of transition probabilities it means that $p(x|3y)$ ($x \in \mathcal{X}$) is the same for all $y \in \{0, 1, 2\}$. The terminal node context tree is $\tau_c^T = \{0, 1, 2, 33\}$, whereas the context tree is $\tau_c = \{0, 1, 2, 3, 33\}$. The state 3 is represented by an internal node in the tree and hence is only an element of τ_c and not of τ_c^T . An alternative representation of the state 3 is given by the final complementary node, indicated by the rectangle, lumping the three non-present nodes together to a new terminal node.

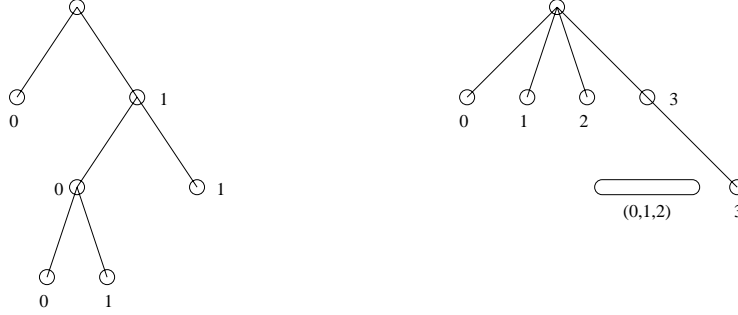


Figure 2.1: Tree representations of the context functions in Examples 2.1 and 2.2.

2.2 Semiparametric VLMC model and sequences of VLMC's

Rather than one finite order VLMC we consider a semiparametric model in the spirit of Ritov and Bickel (1990). The semiparametric VLMC model is

$$\mathcal{P} = \bigcup_{k=0}^{\infty} \mathcal{P}_k, \quad (2.2)$$

where \mathcal{P}_k is the set of stationary VLMC's of order k ,

$$\mathcal{P}_k = \{P_c; P_c \text{ a stationary VLMC of order } k\}.$$

Thus, every member of \mathcal{P} belongs to a nice (parametric) VLMC model whose order can be arbitrarily large. We will study in section 3 a consistent estimator for the semiparametric VLMC model \mathcal{P} without using additional structural information such as the underlying context function $c(\cdot)$. Since we do not specify any particular model structure in \mathcal{P} , the estimator is *robust* against model-misspecification in the set of all VLMC's (which is dense in the set of all stationary processes with respect to finite dimensional weak convergence).

The asymptotic analysis of such a robust estimator in \mathcal{P} is given for a framework where the underlying process changes with sample size, a so-called 'moving truth' model, see (2.3) below. The reasons are twofold. First, it is much more interesting to see whether an estimation technique is still consistent in such a situation: the 'non-moving' case, considered by Weinberger et al. (1995), is from an asymptotic point of view not so interesting, because the problem is finite (although high) dimensional. Secondly, the 'moving-truth' model has limiting elements on the boundary of the semiparametric model \mathcal{P} which are infinite-dimensional non-VLMC models. In this sense, the 'moving truth' model yields an interesting approximation for some general stationary \mathcal{X} -valued processes. The 'moving truth' is a sequence $(P_n)_{n \in \mathbb{N}}$ of VLMC's in \mathcal{P} from (2.2) from which the data are finite realizations in a triangular scheme,

$$\begin{aligned} X_{1,n}, \dots, X_{n,n} \text{ a finite realization of } P_n, \\ P_n \in \mathcal{P} \text{ from (2.2) with context function } c_n(\cdot), n \in \mathbb{N}. \end{aligned} \quad (2.3)$$

The transition probabilities corresponding to P_n are denoted by $p_n(\cdot|\cdot)$. In the sequel when writing data just as X_1, \dots, X_n , we usually think of a generating model as in (2.3).

3 Context algorithm and its consistency

Given data $X_{1,n}, \dots, X_{n,n}$ as in (2.3), the aim is to find the underlying context function $c_n(\cdot)$ (the minimal state space) and an estimate of P_n . A version of the context algorithm (Rissanen, 1983) will be used to solve the problem. Besides obvious uses of the numerical estimate of the probability distribution including a resampling scheme as given in section 4, the estimated context tree is an excellent exploratory tool for the dynamic structure of the underlying process, see section 3.3.

3.1 Context algorithm

We describe now the context algorithm for the aim mentioned above. The main strategy is as follows. First, a large context tree is grown, which represents an overfitted VLMC model. Since the value space \mathcal{X} is finite, there aren't any sophisticated problems with finding accurate splits of the predictor space and the construction of such a large tree turns out to be simple and computationally fast. Secondly, the algorithm employs a backward tree-pruning procedure by considering a local decision criterion. Thus, on this very basic level, the context algorithm has a similar architecture as many other tree fitting methods.

In the sequel we always make the convention that quantities involving time indices $\notin \{1, \dots, n\}$ equal zero (or are irrelevant). Let

$$N(w) = \sum_{t=1}^n 1_{[X_t^{t+|w|-1}=w]}, \quad w \in \cup_{m=1}^{\infty} \mathcal{X}^m, \quad (3.1)$$

denote the number of occurrences of the string w in the sequence X_1^n . Moreover, let

$$\hat{P}(w) = N(w)/n, \quad \hat{P}(u|w) = \frac{N(uw)}{N(w)}, \quad w, u \in \cup_{m=1}^{\infty} \mathcal{X}^m. \quad (3.2)$$

The algorithm below constructs the estimated context tree $\hat{\tau}$ as the biggest context tree (with respect to the order ' \preceq ' defined in Step 1 below) such that

$$\Delta_{wu} = \sum_{x \in \mathcal{X}} \hat{P}(x|wu) \log\left(\frac{\hat{P}(x|wu)}{\hat{P}(x|w)}\right) N(wu) \geq K \text{ for all } wu \in \hat{\tau}^T \text{ (} u \in \mathcal{X}\text{)} \quad (3.3)$$

with $K = K_n \sim C \log(n)$, $C > 2|\mathcal{X}| + 4$ a cut-off to be chosen by the user.

Step 1 Given \mathcal{X} -valued data X_1, \dots, X_n , fit a maximal context tree, i.e., search for the context function $c_{max}(\cdot)$ with terminal node context tree representation τ_{max}^T (see Definition 2.3), where τ_{max}^T is the biggest tree such that every element (terminal node) in τ_{max}^T has been observed at least twice in the data. This can be formalized as follows:

τ_{max}^T is such that $w \in \tau_{max}^T$ implies $N(w) \geq 2$, and such that for every τ^T , where $w \in \tau^T$ implies $N(w) \geq 2$, it holds that $\tau^T \preceq \tau_{max}^T$.

Here, $\tau_1 \preceq \tau_2$ means: $w \in \tau_1 \Rightarrow wu \in \tau_2$ for some $u \in \cup_{m=0}^{\infty} \mathcal{X}^m$ ($\mathcal{X}^0 = \emptyset$).

Set $\tau_{(0)}^T = \tau_{max}^T$.

Step 2 Examine every element (terminal node) of $\tau_{(0)}^T$ as follows (the order of examining is irrelevant, see Remark 3.3). Let $c(\cdot)$ be the corresponding context function of $\tau_{(0)}^T$ and let

$$wu = x_{-\ell+1}^0 = c(x_{-\infty}^0), \quad u = x_{-\ell+1}, \quad w = x_{-\ell+2}^0,$$

where wu is an element (terminal node) of $\tau_{(0)}^T$, which we compare with its pruned version $w = x_{-\ell+2}^0$ (if $\ell = 1$, the pruned version is the empty branch, i.e., the root node).

Prune $wu = x_{-\ell+1}^0$ to $w = x_{-\ell+2}^0$ if

$$\Delta_{wu} = \sum_{x \in \mathcal{X}} \hat{P}(x|wu) \log\left(\frac{\hat{P}(x|wu)}{\hat{P}(x|w)}\right) N(wu) < K,$$

with $K = K_n \sim C \log(n)$, $C > 2|\mathcal{X}| + 4$ and $\hat{P}(\cdot|\cdot)$ as defined in (3.2). Decision about pruning for every terminal node in $\tau_{(0)}^T$ yields a (possibly) smaller tree $\tau_{(1)} \preceq \tau_{(0)}^T$. Construct the terminal node context tree $\tau_{(1)}^T$.

Step 3 Repeat Step 2 with $\tau_{(i)}, \tau_{(i)}^T$ instead of $\tau_{(i-1)}, \tau_{(i-1)}^T$ ($i = 1, 2, \dots$) until no more pruning is possible. Denote this maximal pruned context tree (not necessarily of terminal node type) by $\hat{\tau} = \tau_{\hat{c}}$ and its corresponding context function by $\hat{c}(\cdot)$.

Step 4 If interested in probability distributions, estimate the transition probabilities $p(x_1|c(x_{-\infty}^0))$ by $\hat{P}(x_1|\hat{c}(x_{-\infty}^0))$, where $\hat{P}(\cdot|\cdot)$ is defined as in (3.2).

Remark 3.1. The pruning decision in Step 2 can be related to the Kullback-Leibler distance and to the likelihood ratio test. By definition,

$$\begin{aligned} \Delta_{wu} &= \sum_{x \in \mathcal{X}} \hat{P}(x|wu) \log\left(\frac{\hat{P}(x|wu)}{\hat{P}(x|w)}\right) N(wu) \\ &= D(\hat{P}(\cdot|wu) || \hat{P}(\cdot|w)) N(wu), \end{aligned} \tag{3.4}$$

where $N(wu)$ is defined in (3.1) and $D(P||Q) = \sum_{x \in \mathcal{X}} P(x) \log(P(x)/Q(x))$ is the Kullback-Leibler distance between two probability measures P and Q on \mathcal{X} .

Denote the estimated likelihood function (conditioned on the first state), based on context function $c(\cdot)$ by

$$\hat{P}_c(X_1^n) = \prod_{t=k+1}^n \hat{P}(X_t|c(X_{-\infty}^{t-1})), \tag{3.5}$$

where k is the order of $c(\cdot)$ and $\hat{P}(X_t|c(X_{-\infty}^{t-1}))$ is defined in (3.2).

Denote by $c(\cdot)$ the context function of a non-pruned context tree and by $c'(\cdot)$ the context function of the sub-tree, pruned at one terminal node $wu = x_{-\ell+1}^0$ to its parent node $w = x_{-\ell+2}^0$. By the multiplicative structure in (3.5), many terms cancel in the likelihood

ratio statistic and the only remaining term is at the node considered for pruning. One gets

$$\Delta_{wu} = \log\left(\frac{\hat{P}_c(X_1^n)}{\hat{P}_{c'}(X_1^n)}\right). \quad (3.6)$$

(If $c'(\cdot)$ is of lower order $k - 1$ than $c(\cdot)$, some minor edge effects due to conditioning on the first variables arise). Formula (3.6) says that our pruning criterion is nothing else than a likelihood ratio test, but now with a large acceptance region $[0, C \log(n)]$ for the pruned (sub-)tree. Our algorithm can be viewed as doing very many likelihood ratio tests.

Remark 3.2. The cut-off value $K_n \sim C \log(n)$ in Step 2 for the pruning decision is chosen by an asymptotic consideration. Clearly, by the interpretation as likelihood ratio tests, small cut-off values will result in larger context trees and overfitting occurs. Estimation of the cut-off value (the constant C in the present formulation) has been given in Bühlmann (1998): it aims for optimality with respect to some loss function, whose specification allows tailoring the procedure for particular aims, e.g., 0-1 prediction error loss. The cut-off value can be interpreted as a stepwise $1 - \alpha$ quantile for a multiple testing problem with $\alpha = \alpha_n \rightarrow 0$ ($n \rightarrow \infty$). The necessity for α_n converging to zero is explained in e.g. Rissanen (1989).

Remark 3.3. It does not matter which terminal node wu in Step 2 is examined first, second and so on: for every tree $\tau_{(i)}$ the order of testing the terminal nodes is irrelevant.

Remark 3.4. The pruning criterion Δ_{wu} does not need to be based on the Kullback-Leibler distance. The quantity $D(\hat{P}(\cdot|wu) || \hat{P}(\cdot|w))$ in (3.4) could be replaced by the squared L_1 -distance $\|\hat{P}(\cdot|wu) - \hat{P}(\cdot|w)\|_1^2$ (for a definition of $\|\cdot\|_1$ see assumption (A2) below). In this case, the cut-off in Step 2 of the context algorithm needs to satisfy $K_n \sim C \log(n)$, $C > 4|\mathcal{X}| + 8$. See proof of Theorem 5.1 (which is mainly in terms of $\|\cdot\|_1^2$) and Theorem 5.3.

Remark 3.5. The maximal tree τ_{max}^T in Step 1 is constructed on the basis of at least two occurrences of every terminal node in the sequence. The number two is a low enough value in practice which guarantees a sufficiently large initial tree and at least two observations to estimate transition probabilities associated to terminal nodes (states) in τ_{max}^T . It is easy to show under the assumptions in section 3.2 that $\mathbb{P}[N(w) \geq 2] \rightarrow 1$ ($n \rightarrow \infty$) for $w \in \tau_n = \tau_{c_n}$. Asymptotic properties of the algorithm remain unchanged when replacing the number two by any finite number.

Remark 3.6. The algorithm makes no a-priori length restriction for long contexts (i.e., deep nodes in the tree) such as $\frac{\log(n)}{\log(|\mathcal{X}|)}$ employed by Weinberger et al. (1995) which can be a severe restriction in practical applications.

Generally, the pruning in the context algorithm can be viewed as hierarchical backward selection. Dependence on some values further back in the history is weaker by considering deep nodes in the tree in a hierarchical way as less relevant. This hierarchic structure is a clear distinction to the CART algorithm (Breiman et al., 1984), where the tree architecture is binary and has no built in time structure.

3.2 Consistency under increasing dimensionality

We give two results, both dealing with consistency when the dimension of the underlying model is allowed to increase. The first one shows consistency for finding the minimal state spaces and the second one describes properties of the estimated probability distributions.

We consider a sequence of VLMC's $(P_n)_{n \in \mathbb{N}}$ with $P_n \in \mathcal{P}$, as described in (2.3) with context tree $\tau_n = \tau_{c_n}$, induced by the context function $c_n(\cdot)$. We make the following assumptions.

(A1) $(P_n)_{n \in \mathbb{N}}$ satisfies for some $r \in \mathbb{N}$,

$$\sup_{n \in \mathbb{N}} \sup_{A \subseteq \tau_n; w, w' \in \tau_n} |p_{Z_n}^{(r)}(A, w) - p_{Z_n}^{(r)}(A, w')| < 1 - 2\kappa, \text{ for some } \kappa > 0,$$

where $p_{Z_n}^{(r)}(A, w) = \mathbb{P}[Z_{r,n} \in A | Z_{0,n} = w]$ denotes the r -step transition kernel of the state process $Z_{t,n} = c_n(X_{-\infty,n}^{t,n})$ with $(X_{t,n})_{t \in \mathbb{Z}} \sim P_n$.

(A2) Let $b_n = \min_{w \in \tau_n} P_n(w)$ and $\epsilon_n = \min_{wu \in \tau_n, u \in \mathcal{X}} \|P_n(\cdot | wu) - P_n(\cdot | w)\|_1$ (with L_1 -distance $\|f\|_1 = \sum_{x \in \mathcal{X}} |f(x)|$ for some $f : \mathcal{X} \rightarrow \mathbb{R}$). Then,

$$\begin{aligned} b_n^{-1} &= O(\log(n)^{-(1/2+\beta)} n^{1/2}) \text{ for some } 0 < \beta < \infty \text{ } (n \rightarrow \infty), \\ \epsilon_n^{-1} &= O((\log(n)^{-(1+\delta)} n b_n)^{1/2}) \text{ for some } 0 < \delta < \infty \text{ } (n \rightarrow \infty). \end{aligned}$$

(A3) The minimal transition probabilities satisfy

$$\frac{1}{\min_{x \in \mathcal{X}, w \in \tau_n} p_n(x|w)} = O(n) \text{ } (n \rightarrow \infty).$$

Remark 3.7. The assumption about transition kernels in (A1) is related to the ergodicity coefficient for stationary Markov processes, cf. Iosifescu and Theodorescu (1969), Rajarshi (1990) or Doukhan (1994). It implies that the state processes $(Z_{t,n})_{t \in \mathbb{Z}}$ and hence also the VLMC's $(X_{t,n})_{t \in \mathbb{Z}}$ are geometrically ϕ -mixing with mixing coefficients bounded by

$$\sup_{n \in \mathbb{N}} \phi_n(k) \leq (1 - 2\kappa)^k \text{ for all } k \in \mathbb{N}.$$

Remark 3.8. Assumption (A2) about the minimum stationary probability bounds the size of the context tree as $|\tau_n| \leq b_n^{-1} = O(n^{1/2}/\log(n)^{1/2+\beta})$. Note that the number of transition probability parameters in the process P_n is $O(|\tau_n|)$. The above bound, in probabilistic terms, is a weak condition for the number of parameters and there is no explicit restriction on the order $c_n(\cdot)$ (the depth of the context tree τ_n).

Remark 3.9. For distinguishing a context wu from its parent node w in the context tree, assumption (A2) guarantees a minimal L_1 distance between the relevant conditional distributions.

In the special case with only one fixed VLMC $P = P_c$ with context tree τ_c , it is sufficient to only assume for the transition probabilities,

$$\min_{x \in \mathcal{X}, w \in \tau_c} p(x|w) > 0,$$

which implies assumptions (A1), (A2) with $b_n \equiv b > 0$, $\epsilon_n \equiv \epsilon > 0$ and (A3) with $O(1)$.

Theorem 3.1 *Consider data $X_{1,n}, \dots, X_{n,n}$ as in (2.3), where $c_n(\cdot)$ denotes the context function and τ_n the context tree of the process P_n , satisfying (A1)-(A3). Let $\hat{P}(\cdot|\cdot)$ be defined as in (3.2) and $\hat{c}(\cdot)$ the estimate in Step 3 of the context algorithm. Then,*

$$(i) \lim_{n \rightarrow \infty} \mathbf{P}[\hat{c}(\cdot) = c_n(\cdot)] = 1, \text{ or equivalently } \lim_{n \rightarrow \infty} \mathbf{P}[\hat{\tau} = \tau_n] = 1,$$

$$(ii) \sup_{x_{-\infty}^1 \in \mathcal{X}^\infty} |\hat{P}(x_1|\hat{c}(x_{-\infty}^0)) - p_n(x_1|c_n(x_{-\infty}^0))| = o_P(1) \text{ (} n \rightarrow \infty \text{)}.$$

A proof of Theorem 3.1 is given in section 5. There, more explicit bounds for the events of choosing too large or too small minimal state spaces are given. Theorem 3.1 explains why the context algorithm is a very powerful tool. Even if the dimensionality increases, the estimator $\hat{c}(\cdot)$ (or $\hat{\tau}$) neither chooses a too large nor a too small model asymptotically and is thus robust against model-misspecification with respect to sequences in the broad semiparametric class \mathcal{P} . The increase in dimensionality of the underlying model is restricted in probabilistic terms but allows a growth as fast as $O(n^{1/2}/\log(n)^s)$ for some $s > 1/2$. The problem is thus highly non-trivial: there is possible failure with simple estimation rules which consider models within a fixed increase in dimensionality, independent of the underlying process, but otherwise quite general of the order $O(n^r)$ for some $0 < r < 1/2$. This relates to a good bias-variance trade-off, even for some very high dimensional VLMC's and for general stationary processes which are on the boundary of \mathcal{P} . Theorem 3.1 describes the solution of a model selection problem which is impossible to deal with a global selection criterion, due to the extremely large number of possible models. For example, the number of all VLMC sub-model of a full \mathcal{X} -valued MC of order 4 with $|\mathcal{X}| = 4$ is $\approx 2 \cdot 10^{20}$. The selection criterion is based here on a hierarchical local criterion (the context algorithm) and interestingly, it works also in the case where the model dimension increases. In theory but never practically feasible, a minimum description length estimator might yield consistent state estimation as well, cf. Weinberger and Feder (1994) in a more general class of finite-state models.

The next result describes the construction and properties of the estimator \hat{P}_n for the underlying probability measure P_n . Define a metric for probability measures P, Q on \mathcal{X}^∞ ,

$$d(P, Q) = \sum_{m=1}^{\infty} 2^{-m} d_m(P \circ \pi_{1,\dots,m}^{-1}, Q \circ \pi_{1,\dots,m}^{-1}),$$

$$d_m(P \circ \pi_{1,\dots,m}^{-1}, Q \circ \pi_{1,\dots,m}^{-1}) = \sup_{x_1^m \in \mathcal{X}^m} |P(x_1^m) - Q(x_1^m)|, \quad (3.7)$$

where $\pi_{1,\dots,m} : x \mapsto x_1, \dots, x_m$, ($x \in \mathcal{X}^\infty$) is the coordinate function.

Theorem 3.2 *Consider data $X_{1,n}, \dots, X_{n,n}$ as in (2.3) with P_n satisfying (A1)-(A3). Then,*

(i) for $\hat{P}(\cdot|\cdot)$ as in (3.2) and $\hat{c}(\cdot)$ the estimate in Step 3 of the context algorithm,

$$\lim_{n \rightarrow \infty} \mathbb{P}[\text{the set } \{\hat{P}(\cdot|\hat{c}(x_{-\infty}^0)); x_{-\infty}^0 \in \mathcal{X}^\infty\} \text{ generates a unique stationary probability measure } \hat{P}_n \in \mathcal{P}] = 1,$$

(ii) for \hat{P}_n in (i) and $d(\cdot, \cdot)$ as in (3.7), $d(\hat{P}_n, P_n) = o_P(1)$ ($n \rightarrow \infty$).

(iii) the process \hat{P}_n in (i) satisfies

$$\mathbb{P}[\hat{P}_n \text{ is } \phi\text{-mixing with mixing coefficients satisfying } \phi_{\hat{P}_n}(k) \leq (1 - \kappa)^k \text{ for all } k \in \mathbb{N}_0] \rightarrow 1 \text{ (} n \rightarrow \infty \text{)}.$$

In particular, the bound for the mixing coefficients $\phi_{\hat{P}_n}(k)$ is non-random and the same for all $n \in \mathbb{N}$.

A proof of Theorem 3.2 is given in section 5. Statement (i) of Theorem 3.2 tells in a constructive way how to simulate the estimated underlying process, the fidi-convergence in (ii) is a minimal requirement for a reasonable estimator, whereas (iii) is important for simulation tasks like bootstrapping complicated statistics.

3.3 DNA example

We now present an interesting and instructive application of VLMC estimation. In particular, we demonstrate the usefulness of an estimated context tree as an excellent graphical tool for detecting structure in the time series. Our data consists of three distinct sequences of DNA from the *Drosophila* genome. We point out that genetic data is a natural candidate for modeling by a VLMC since it has a finite alphabet \mathcal{X} , there is a ‘time’ index and its memory vanishes with increasing lag time. Furthermore, while it is known that the data is far from independent, higher order Markov models are not easily fitted because of the explosion in the number of parameters, cf. Braun and Müller (1998) and compare also with Problems 1 and 2 from section 1. This is compounded by an observed degree of non-stationarity which prohibits estimation over very long sequences. In this environment, it is of paramount importance that the model parameters are used with great economy in order to capture any significant dependent structure.

Our data began as a single 100-thousand base contiguous stretch of DNA from the *Drosophila* genome (Genbank number DS02740). Each base is one four possible DNA residues: Adenine, Cytosine, Guanine and Thymine, with abbreviations A, C, G and T, respectively. Using a variety of tools, biologists at Gerry Rubin’s lab at the University of California at Berkeley segmented the sequence into genes (which code for amino acid sequences) and non-genes (so-called junk DNA) which are ignored by the cell chemistry. Physically, the genes are spaced apart and separated by junk DNA which we term ‘inter-genes’. Moreover, the genes are further segmented into coding regions called exons and non-coding regions called introns. The cell’s engine for transcribing DNA first copies the gene (both intron and exon), it then splices out the intron sections. Each gene is in turn subdivided into alternating stretches of exon and intron. We form a single sequence of exons by concatenating all the exons (in the given order). Similarly, we form sequences of introns and inter-genes.

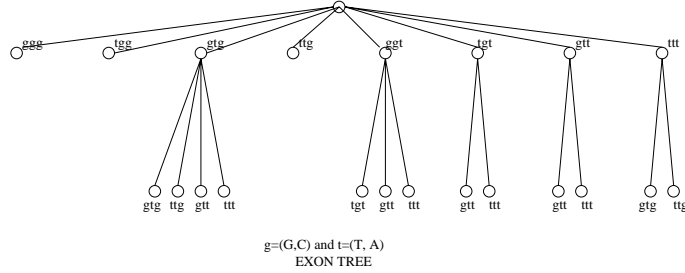


Figure 3.2: Triplet tree representation of the estimated minimal state space for exon sequence. The triplets are denoted in reverse order, e.g., the terminal node with concatenation (ggt)(ggt) describes the context $x_0 = g, x_{-1} = g, x_{-2} = t, x_{-3} = g, x_{-4} = t, x_{-5} = t$ for the variable x_1 .

Our goal is the application of the VLMC estimation algorithm to learn the dependence structure and to present the estimated minimal state space graphically as a tree, whose branches are the contexts. Application of the algorithm to each of the datasets suggests that complicated structures exist within the exons and the introns. On the other hand, the inter-genes showed no complex structure (a first order Markov model is a good fit). That exons exhibit such structure is not surprising due to constraints imposed by its coding function. The introns do not have a well understood function, but evidence of structure suggests that the intron is constrained in some way and is thus unable to freely mutate.

We also consider the sequences under a reduction of the 4-ary alphabet down to three possible binary alphabets, identifying: (1) G with C; (2) G with A; (3) G with T. Equivalences (1) and (2) have genetic meaning, the third has none (reducing the data to random bits). As expected, this final equivalence (3) produces sequences with no dependence structure. The most dramatic finding was produced by the exon sequence reduced to binary alphabet by identifying the base G with its bonding pair C (A is thus identified with T). The resulting context tree has branches of lengths 0, 3 and 6 only. Interestingly, we thus can represent it in terms of triplets, as shown in Figure 3.1. Because amino acids are known to be coded by triplets of DNA letters, the structure in Figure 3.1 has a beautiful biological interpretation. Our finding suggests that the triplet coding structure is strongly present despite the dramatic processing of the data (the actual code is not recoverable from this binary reduction). We point out, for emphasis, that the VLMC estimation algorithm learned the triplet structure on its own, a discovery made only for the reduced to binary alphabet coding sequences.

4 The VLMC bootstrap

Theorem 3.2 indicates, that the estimate \hat{P}_n of P_n can be used for resampling. Our proposal will be a bootstrap for stationary categorical time series. Since the semiparametric model (2.2) is dense (with respect to the metric in (3.7)) in the set of stationary \mathcal{X} -valued processes, this bootstrap is very general. It offers an attractive and often more accurate alternative to the model free blockwise bootstrap, which has been proposed by Künsch (1989). We proceed as follows.

Step 1 Given \mathcal{X} -valued data X_1, \dots, X_n , fit a VLMC as described in section 3.1, yielding a stationary probability measure \hat{P}_n on $\mathcal{X}^{\mathbb{Z}}$, see Theorem 3.2.

Step 2 Draw a finite realization

$$X_1^*, \dots, X_n^* \sim \hat{P}_n \circ \pi_{1, \dots, n}^{-1}.$$

The variables X_1^*, \dots, X_n^* are called the VLMC bootstrap sample, they are nothing else than one random sample from the fitted VLMC. In practice, one would choose some starting values, generate a longer random sample via the estimated transition probabilities $\hat{P}(x_1 | \hat{c}(x_{-\infty}^0))$ and then use the last n elements as our bootstrap sample. Such a device tries to avoid nonstationarity effects due to starting values in a simulated Markov chain. Of course, one could also draw bootstrap samples of size $m \neq n$, cf. Bickel et al. (1997).

Given an estimator $T_n = T_n(X_1, \dots, X_n)$, which is a measurable function of X_1, \dots, X_n , the bootstrapped estimator is defined by the plug-in rule $T_n^* = T_n(X_1^*, \dots, X_n^*)$. This plug-in rule, which is also the basis to Efron's (1979) bootstrap for the independent case, is very convenient in practice. Once the bootstrap sample is constructed, bootstrapping can be done with exactly the same computing tools or programs as for the original estimator T_n . This is not the case with the blockwise bootstrap (Künsch, 1989) if the estimator T_n is non-symmetric in the observations X_1, \dots, X_n , e.g., the estimators in (S1) and (S2) from section 4.2. Quantities induced by the resampling in Step 2 are denoted by an asterisk $*$.

4.1 Consistency of the VLMC bootstrap under increasing dimensionality

We present here an asymptotic result which justifies the use of the above defined VLMC bootstrap for estimators T_n being smooth functions of means. We will also discuss informally why the VLMC bootstrap should work in the more general framework of empirical processes, without giving the exact arguments.

We assume that we have observations $X_{1,n}, \dots, X_{n,n} \in \mathcal{X}$ from a sequence of VLMC's as given in (2.3). Consider the class of estimators, being smooth functions of means,

$$T_n = g\left\{(n-m+1)^{-1} \sum_{t=1}^{n-m+1} f(X_{t,n}^{t+m-1,n})\right\}, \quad 1 \leq m < \infty,$$

$$f = (f_1, \dots, f_v)' : \mathcal{X}^m \rightarrow \mathbb{R}^v, \quad g = (g_1, \dots, g_w)' : \mathbb{R}^v \rightarrow \mathbb{R}^w \text{ smooth.} \quad (4.1)$$

The function f is bounded, since $|\mathcal{X}| < \infty$. Examples include estimators of transition probabilities in finite state Markov chains of order $m-1$ or other functions of frequencies of tuples up to size m , such as the Z scores used in genetics, cf. Prum et al. (1995). We usually make the following assumption.

(B1) T_n is given by (4.1) with g having continuous partial derivatives in a neighborhood of $\theta_n = \mathbb{E}[f(X_{1,n}, \dots, X_{m,n})]$. Also, there exists an $n_0 \in \mathbb{N}$, such that for every $n \geq n_0$,

$$\left[\sum_{k=-n+1}^{n-1} Cov(f_i(X_{0,n}^{m-1,n}), f_j(X_{k,n}^{k+m-1,n})) \right]_{i,j=1}^v \text{ is positive definite.}$$

Remark 4.1. The assumption about positive definiteness of covariance matrices simplifies when assuming a limiting model P , where $\lim_{n \rightarrow \infty} d(P_n, P) = 0$ for the metric $d(\cdot, \cdot)$ defined in (3.7). Generally, P is not a VLMC anymore. It is then sufficient to assume

$$\begin{aligned} & \left| \sum_{k=-\infty}^{\infty} \text{Cov}(f_i(X_0^{m-1}), f_j(X_k^{k+m-1})) \right| < \infty, \quad i, j \in \{1, \dots, v\}, \\ & \left[\sum_{k=-\infty}^{\infty} \text{Cov}(f_i(X_0^{m-1}), f_j(X_k^{k+m-1})) \right]_{i,j=1}^v \text{ is positive definite,} \end{aligned}$$

where $(X_t)_{t \in \mathbb{Z}} \sim P$.

The following Theorem justifies the VLMC bootstrap for smooth functions of means.

Theorem 4.1 *Let $X_{1,n}, \dots, X_{n,n}$ be as in (2.3) with P_n satisfying (A1)-(A3). Assume also that (B1) holds. Let the VLMC bootstrap be defined as in section 4 and denote by $\theta_n^* = \mathbb{E}^*[f((X^*)_1^m)]$. Then,*

$$\sup_{x \in \mathbb{R}^w} |\mathbb{P}^*[n^{1/2}(T_n^* - g(\theta_n^*)) \leq x] - \mathbf{P}[n^{1/2}(T_n - g(\theta_n)) \leq x]| = o_P(1) \quad (n \rightarrow \infty).$$

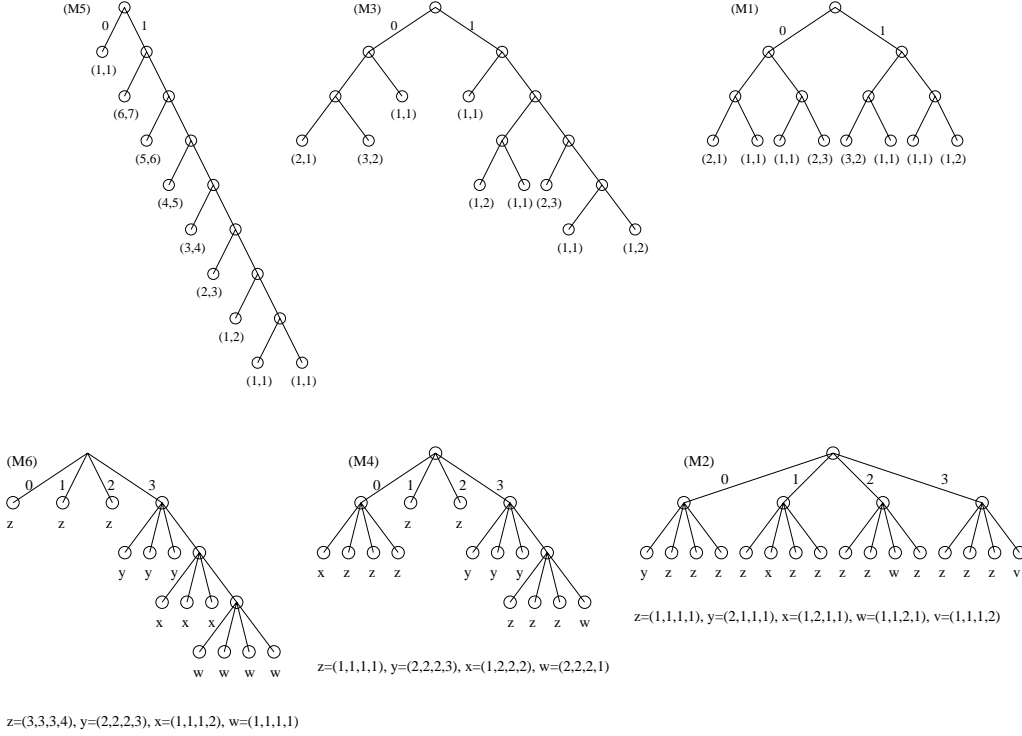
The proof of Theorem 4.1 is given in section 5. Our results can be generalized to consistency of the VLMC bootstrap for general empirical processes, due to the fact that the VLMC bootstrap for categorical time series satisfies a ϕ -mixing property with exponentially decaying mixing coefficients, see Theorem 3.2 (iii). These extensions are useful for studying the bootstrap consistency of estimators

$$T_n = T(\nu_n), \tag{4.2}$$

being a smooth functional of a general empirical measure ν_n . The class of estimators in (4.2) is considerably larger than the class in (4.1). It includes as examples the maximum likelihood estimators in generalized linear models of autoregressive type with quite general link functions, cf. Fahrmeir and Tutz (1994).

4.2 Simulations

We study here the VLMC bootstrap for variance estimation in various cases by simulation. We represent VLMC models by context trees and equip terminal nodes with tuples, describing the transition probabilities. A tuple $(i_0, \dots, i_{|\mathcal{X}|-1})$ corresponds to $p(j|w) = i_j / \sum_{j=0}^{|\mathcal{X}|-1} i_j$, $j \in \{0, \dots, |\mathcal{X}|-1\}$ (without loss of generality we let $\mathcal{X} = \{0, \dots, |\mathcal{X}|-1\}$). We consider the following models. (M1): full binary Markov chain of order 3; (M2): full 4-ary Markov chain of order 2; (M3): semi-sparse binary VLMC of order 5; (M4): semi-sparse 4-ary VLMC of order 3; (M5): sparse binary VLMC of order 8; (M6): sparse 4-ary VLMC of order 4. The precise specifications are given by the following trees and numbers:



We also consider

(M7): $X_t = 1_{[Y_t > 0]}$ where $(Y_t)_{t \in \mathbb{Z}}$ is a stationary nonlinear process

$$Y_t = (0.5 + 0.9 \exp(-2.354Y_{t-1}^2))Y_{t-1} - (0.8 - 1.8 \exp(-2.354Y_{t-1}^2))Y_{t-2} + Z_t$$

with $(Z_t)_{t \in \mathbb{Z}}$ an i.i.d. sequence, $Z_t \sim \mathcal{N}(0, 1)$ and Z_t independent from Y_s for all $s < t$. The process $(Y_t)_{t \in \mathbb{Z}}$ is also known as ‘Exponential AR(2)’.

The quantized binary process $(X_t)_{t \in \mathbb{Z}}$ in (M7) is non-Markovian, although the \mathbb{R} -valued $(Y_t)_{t \in \mathbb{Z}}$ is Markov of order 2. It is interesting to see whether the VLMC bootstrap provides a good finite-sample approximation to this model which is not a VLMC. This is also an interesting test case to make a fair comparison between the VLMC and blockwise bootstrap (Künsch, 1989).

As sample sizes, we choose $n = 1000$ and $n = 2000$. We consider the following statistics.

- (S1) $T_n = \hat{P}_n(1|0) = N_n(1, 0)/N_n(0)$ for binary models (M1), (M3) and (M5).
- (S2) $T_n = N_n(1, 3, 3)$, the frequency of the word $(x_3, x_2, x_1) = (1, 3, 3)$ as defined in (3.1), for 4-ary models (M2), (M4) and (M6).
- (S3) $T_n = \bar{X}_n = n^{-1} \sum_{t=1}^n X_t$, the relative frequency of symbol 1, for the binary model (M7).

The variance estimates are

$$\hat{\sigma}_n^2 = n \text{Var}^*(T_n^*) \text{ for } \sigma_n^2 = n \text{Var}(T_n) \text{ in (S1) and (S3),}$$

$$\hat{\sigma}_n^2 = \text{Var}^*(T_n^*) \text{ for } \sigma_n^2 = \text{Var}(T_n) \text{ in (S2),}$$

	σ_n^2	$\mathbb{E}[\hat{\sigma}_n^2] - \sigma_n^2$	$Var(\hat{\sigma}_n^2)$	$relMSE(\hat{\sigma}_n^2)$
(M1,S1) 95%	0.81	-0.04 (0.010)	0.02	0.033 (0.0035)
(M1,S1) 98%	0.81	-0.10 (0.011)	0.02	0.053 (0.0040)
(M1,S1) 99.9%	0.81	-0.26 (0.009)	0.02	0.127 (0.0038)
(M3,S1) 95%	0.67	-0.02 (0.009)	0.01	0.031 (0.0029)
(M3,S1) 98%	0.67	-0.05 (0.009)	0.01	0.036 (0.0030)
(M3,S1) 99.9%	0.67	-0.17 (0.005)	0.01	0.078 (0.0027)
(M5,S1) 95%	0.528	0.007 (0.0054)	0.006	0.021 (0.0027)
(M5,S1) 98%	0.528	-0.005 (0.0031)	0.002	0.007 (0.0006)
(M5,S1) 99.9%	0.528	0.003 (0.0027)	0.002	0.005 (0.0005)
(M2,S2) 95%	14.5	-0.5 (0.21)	9.1	0.045 (0.0051)
(M2,S2) 98%	14.5	0.1 (0.17)	5.8	0.028 (0.0027)
(M2,S2) 99.9%	14.5	0.0 (0.14)	3.8	0.018 (0.0019)
(M4,S2) 95%	14.1	-0.3 (0.18)	6.4	0.032 (0.0044)
(M4,S2) 98%	14.1	-0.4 (0.17)	5.5	0.029 (0.0042)
(M4,S2) 99.9%	14.1	-0.5 (0.12)	2.8	0.015 (0.0014)
(M6,S2) 95%	11.2	0.0 (0.15)	4.8	0.038 (0.0043)
(M6,S2) 98%	11.2	-0.1 (0.13)	3.1	0.025 (0.0029)
(M6,S2) 99.9%	11.2	-0.3 (0.10)	2.0	0.017 (0.0026)
(M7,S3) 95%	0.80	-0.03 (0.009)	0.02	0.029 (0.0029)
(M7,S3) 98%	0.80	-0.11 (0.009)	0.02	0.043 (0.0030)
(M7,S3) 99.9%	0.80	-0.24 (0.005)	0.01	0.100 (0.0032)

Table 4.1: VLMC bootstrap variance estimates, sample size $n = 1000$.

	σ_n^2	$\mathbb{E}[\hat{\sigma}_n^2] - \sigma_n^2$	$Var(\hat{\sigma}_n^2)$	$relMSE(\hat{\sigma}_n^2)$
(M1,S1) 95%	0.82	-0.01 (0.009)	0.01	0.022 (0.0021)
(M1,S1) 98%	0.82	-0.02 (0.007)	0.01	0.016 (0.0018)
(M1,S1) 99.9%	0.82	-0.14 (0.011)	0.02	0.065 (0.0048)
(M3,S1) 95%	0.67	0.00 (0.006)	0.01	0.014 (0.0013)
(M3,S1) 98%	0.67	-0.03 (0.006)	0.01	0.017 (0.0016)
(M3,S1) 99.9%	0.67	-0.09 (0.007)	0.01	0.042 (0.0033)
(M5,S1) 95%	0.518	0.007 (0.0038)	0.003	0.011 (0.0012)
(M5,S1) 98%	0.518	0.002 (0.0031)	0.002	0.007 (0.0008)
(M5,S1) 99.9%	0.518	0.009 (0.0025)	0.001	0.005 (0.0004)
(M2,S2) 95%	12.9	1.2 (0.16)	4.9	0.038 (0.0038)
(M2,S2) 98%	12.9	1.4 (0.15)	4.4	0.039 (0.0042)
(M2,S2) 99.9%	12.9	1.9 (0.13)	3.2	0.040 (0.0025)
(M4,S2) 95%	14.7	-0.7 (0.15)	4.3	0.022 (0.0022)
(M4,S2) 98%	14.7	-0.6 (0.11)	2.6	0.014 (0.0013)
(M4,S2) 99.9%	14.7	-1.0 (0.09)	1.7	0.012 (0.0011)
(M6,S2) 95%	11.5	-0.1 (0.14)	4.0	0.030 (0.0036)
(M6,S2) 98%	11.5	-0.3 (0.10)	1.9	0.015 (0.0015)
(M6,S2) 99.9%	11.5	-0.5 (0.08)	1.4	0.012 (0.0016)
(M7,S3) 95%	0.81	-0.03 (0.008)	0.01	0.018 (0.0017)
(M7,S3) 98%	0.81	-0.06 (0.007)	0.01	0.022 (0.0018)
(M7,S3) 99.9%	0.81	-0.23 (0.006)	0.01	0.089 (0.0036)

Table 4.2: VLMC bootstrap variance estimates, sample size $n = 2000$.

	σ_n^2	$\mathbb{E}[\hat{\sigma}_n^2] - \sigma_n^2$	$Var(\hat{\sigma}_n^2)$	$relMSE(\hat{\sigma}_n^2)$
(M5,S1) $\ell = 10$	0.528	0.106 (0.0057)	0.007	0.065 (0.0053)
(M5,S1) $\ell = 20$	0.528	0.058 (0.0071)	0.010	0.049 (0.0066)
(M5,S1) $\ell = 30$	0.528	0.030 (0.0084)	0.014	0.054 (0.0065)
(M7,S3) $\ell = 10$	0.80	-0.17 (0.005)	0.01	0.051 (0.0027)
(M7,S3) $\ell = 20$	0.80	-0.09 (0.008)	0.01	0.035 (0.0027)
(M7,S3) $\ell = 30$	0.80	-0.07 (0.011)	0.02	0.045 (0.0038)

Table 4.3: Blockwise bootstrap variance estimates, sample size $n = 1000$.

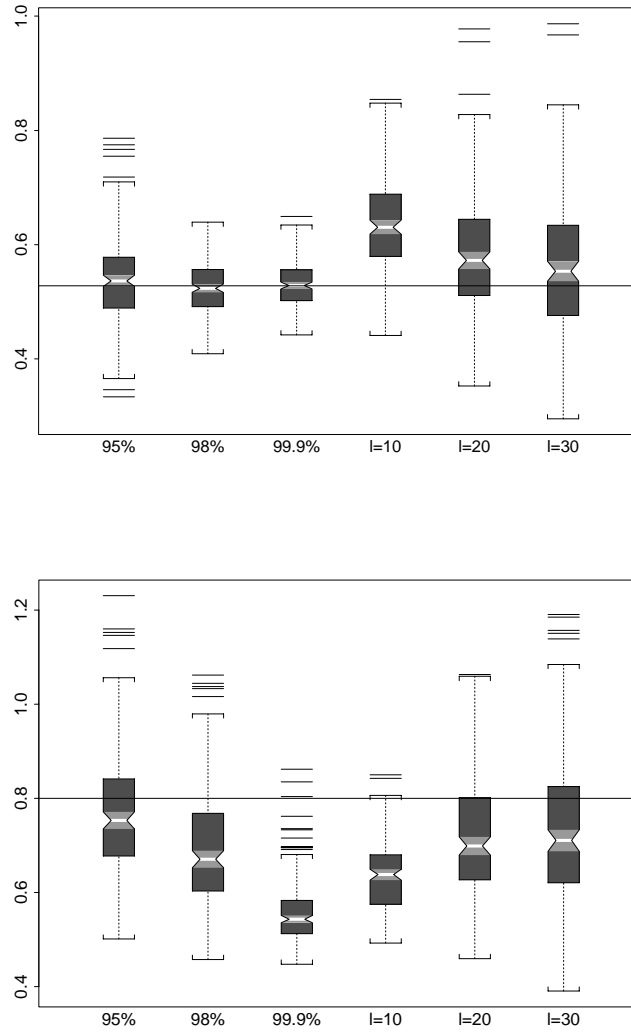


Figure 4.3: Boxplots of bootstrap variance estimates with $n = 1000$, for case (M5,S1) on top and for case (M7,S3) at the bottom. VLMC bootstrap estimates are denoted with their $\chi_{1,\alpha}^2/2$ -quantiles as cut-off values, blockwise bootstrap estimates are denoted with their blocklengths l , the line denotes the true variance.

based on the VLMC bootstrap with 500 resamples (note the different standardizations).

The results are given in Table 4.1 and 4.2. Moments of the bootstrap variances $\hat{\sigma}_n^2$ are estimated with 200 simulations over the different models, the true value of σ_n^2 is estimated with 1000 simulations. The relative mean square error is given by $relMSE(\hat{\sigma}_n^2) = \mathbb{E}|\hat{\sigma}_n^2 - \sigma_n^2|^2/\sigma_n^4$. Estimated standard errors for the bias and $relMSE$ are given in parentheses. We tried different cut-off values K which act as tuning parameter, see also Remark 3.2. We report them as $\chi_{|\mathcal{X}|-1;\alpha}^2/2$ quantiles with different levels α , corresponding to the asymptotic distribution of *one* log-likelihood ratio statistic in (3.6): this is sometimes more intuitive than the numerical value of C in the cut-off $K_n \sim C \log(n)$ from theory. On the other hand we point out the danger that direct interpretation of the cut-off as a $\chi_{|\mathcal{X}|-1;\alpha}^2/2$ quantile with α fixed and not depending on the sample size contradicts the very essence of the context algorithm in (3.3). For the binary models (M1), (M3), (M5) we used the cut-offs $\chi_{1,\alpha}^2/2$, for the 4-ary models (M2), (M4), (M6) the cut-offs $\chi_{3,\alpha}^2/2$, both denoted in short by $\alpha 100\%$.

The results are promising in that the relative mean square error is most often smaller than 5%. Though there are some exceptions, we found that often the performance is better for sparse models. This indicates that the algorithm adapts to sparseness; it is exactly in these cases, where other methods are more likely to fail.

For comparison, we also tried the blockwise bootstrap (Künsch, 1989) in the cases (M5,S1) and (M7,S3) for sample size $n = 1000$ with different blocklengths l , see Table 4.3. A graphical representation is given in Figure 4.1. The comparison is at least fair in (M7,S3) where the model is not a VLMC (and by the structure of $(Y_t)_{t \in \mathbb{Z}}$, the quantized series $(X_t)_{t \in \mathbb{Z}}$ doesn't allow sparse approximation). In this case, both bootstraps have similar performance, the best tuned VLMC bootstrap is about 15% better (in terms of relative means square error) than the best tuned blockwise bootstrap. In case (M5,S1) the blockwise bootstrap exhibits a serious bias and a large variability. The VLMC bootstrap is far better for this sparse VLMC (M5). We conclude that the VLMC bootstrap is at least as good as the blockwise bootstrap and enjoys the important practical advantage of being defined as a plug-in rule, see above.

The role of the cut-off as tuning parameter of the VLMC bootstrap is found as follows: the absolute value of the bias of the bootstrap variance estimator increases and the variance decreases with increasing cut-off parameter. This is expected since a larger cut-off parameter leads to a lower dimensional fitted VLMC model, by design of the context algorithm. Note that in some simulation examples, this behavior is not significantly visible. For the blocklength in the blockwise bootstrap in turn, the general asymptotic behavior is observed, i.e., that the bias decreases, whereas the variance increases with growing l .

5 Proofs

We first recall some notation. We usually denote by $w, u, v \in \cup_{m=0}^{\infty} \mathcal{X}^m$ sequences (written in reverse ‘time’) $w = (w_{|w|}, \dots, w_2, w_1)$, the concatenation is $wu = (w_{|w|}, \dots, w_2, w_1, u_{|u|}, \dots, u_1) \in \cup_{m=0}^{\infty} \mathcal{X}^m$ ($w, u \in \cup_{m=0}^{\infty} \mathcal{X}^m$). Transition probabilities in a context tree τ are indexed by $w \in \tau$ and abbreviated by $p_w(\cdot) = p(\cdot|w)$, estimated transition probabilities are denoted by $\hat{P}_w(x) = N(xw)/N(w)$, $N(\cdot)$ as defined in (3.1). We also abbreviate by $P_w(x) = P(xw)/P(w)$ for general $w \in \cup_{m=1}^{\infty} \mathcal{X}^m$, $x \in \mathcal{X}$ (w not necessarily a context in

τ) and P a stationary probability measure on $\mathcal{X}^{\mathbb{Z}}$ with $P(x) = \mathbb{P}_P[X_1^m = x]$ ($x \in \mathcal{X}^m$). We recall that for any $w = w'u$ ($u \in \mathcal{X}$) we have defined $\Delta_w = D(\hat{P}_{w'u} || \hat{P}_w)N(w)$. When looking at a sequence $(P_n)_{n \in \mathbb{N}}$ of VLMC's, we sometimes drop the index n .

Proof of Theorem 3.1.

We define first the events of under- and overestimation for sample size n ,

$$\begin{aligned} U_n &= \{\text{there exists } w \in \hat{\tau} \text{ with } wu \in \tau_n, wu \notin \hat{\tau} \text{ for some } u \in \cup_{m=1}^{\infty} \mathcal{X}^m\} \\ O_n &= \{\text{there exists } w \in \tau_n \text{ with } wu \in \hat{\tau}, wu \notin \tau_n \text{ for some } u \in \cup_{m=1}^{\infty} \mathcal{X}^m\}. \end{aligned}$$

Note that by formula (3.3) we can also characterize U_n and O_n in terms of the pruning criterion $\Delta_{wu} < K_n = C \log(n)$. The error event is

$$E_n = \{\hat{\tau} \neq \tau_n\} = U_n \cup O_n.$$

Theorem 5.1 *Assume that (A1) and (A2) with $\beta, \delta > 0$ hold. Then,*

$$\mathbb{P}[U_n] = O(\max\{n^{-\log(n)D_1\beta}, n^{-\log(n)D_2\delta}\}) \quad (n \rightarrow \infty)$$

for some constants $0 < D_1, D_2 < \infty$.

Proof: We partition the underestimation event U_n using the event

$$D_n = \{N(w) \geq \rho_n \text{ for every } w \in \tau_n\},$$

where ρ_n is a constant to be chosen later. Thus $\mathbb{P}[U_n] \leq \mathbb{P}[U_n \cap D_n] + \mathbb{P}[D_n^c]$. We will pursue a bound on $\mathbb{P}[U_n]$ by bounding both $\mathbb{P}[U_n \cap D_n]$ and $\mathbb{P}[D_n^c]$. First,

$$\begin{aligned} \mathbb{P}[U_n \cap D_n] &\leq \sum_{wu \in \tau_n, u \in \mathcal{X}} \mathbb{P}[\Delta_{wu} < C \log(n), N(wu) \geq \rho_n] \\ &= \sum_{wu \in \tau_n, u \in \mathcal{X}} \sum_{k=\rho_n}^n \sum_{j=k}^n \mathbb{P}[D(\hat{P}_{wu} || \hat{P}_w) < \frac{C \log(n)}{k}, N(wu) = k, N(w) = j]. \end{aligned} \quad (5.1)$$

It is well known, cf. Cover and Thomas (1991), that the divergence can be lower bounded by the the L_1 distance, $D(\hat{P}_{wu} || \hat{P}_w) \geq \frac{1}{2} \|\hat{P}_{wu} - \hat{P}_w\|_1^2$ and that $\|\hat{P}_{wu} - \hat{P}_w\|_1^2 = 2(\hat{P}_{wu}(A) - \hat{P}_w(A))^2$, where $A = \{x \in \mathcal{X}; \hat{P}_{wu}(x) > \hat{P}_w(x)\}$. Therefore,

$$\begin{aligned} &\mathbb{P}[D(\hat{P}_{wu} || \hat{P}_w) < \frac{C \log(n)}{k}, N(wu) = k, N(w) = j] \\ &\leq \mathbb{P}[(\hat{P}_{wu}(A) - \hat{P}_w(A))^2 < \frac{C \log(n)}{k}, N(wu) = k, N(w) = j]. \end{aligned} \quad (5.2)$$

Now because of assumption (A2), it must be that either $\hat{P}_{wu}(A)$ or $\hat{P}_w(A)$ is far from $P_{wu}(A)$ or $P_w(A)$, respectively. We formalize this by letting $\gamma_n^2(k) = \frac{C \log(n)}{k}$ and $\hat{P}_{wu}(x) = a$, $\hat{P}_w(x) = b$, $p_{wu}(x) = r$ and $p_w(x) = s$, where $x \in \mathcal{X}$. Our goal is to establish that if $|a - b|$ is small then either $|r - a|$ is large or $|s - b|$ is large. First assume, without loss of generality, that $r > s$. We have by (A2) that $r - s > \epsilon_n$. Now if $b < s$, then $|a - b| < \gamma_n(k)$ implies that $|a - r| > \epsilon_n - \gamma_n(k)$. Furthermore, if $b > r$, then it must be that $|s - b| > \epsilon_n$. Now if $s \leq b \leq r$ then either $s \leq b < s + \frac{r-s}{2}$, in which case

$|r - a| > \frac{\epsilon_n}{2} - \gamma_n(k)$ or $r - \frac{r-s}{2} \leq b \leq r$, in which case $|s - b| > \frac{\epsilon_n}{2}$. Taken together we have proved that if $|\hat{P}_{wu}(x) - \hat{P}_w(x)| < \gamma_n(k)$, then either $|\hat{P}_{wu}(x) - p_{wu}(x)| > \frac{\epsilon_n}{2} - \gamma_n(k)$ or $|\hat{P}_w(x) - p_w(x)| > \frac{\epsilon_n}{2} - \gamma_n(k)$. Thus, when applied to (5.2), we have proved that for

$$a_n(k) = \left(\frac{\epsilon_n}{2} - \gamma_n(k)\right)^2, \quad (5.3)$$

it must be that

$$\begin{aligned} & \mathbb{P}\left[D(\hat{P}_{wu}|\hat{P}_w) < \frac{C \log(n)}{k}, N(wu) = k, N(w) = j\right] \\ & \leq \mathbb{P}\left[\sum_{x \in A} |\hat{P}_{wu}(x) - p_{wu}(x)| \geq a_n(k)^{1/2}, N(wu) = k\right] \\ & + \mathbb{P}\left[\sum_{x \in A} |\hat{P}_w(x) - p_w(x)| \geq a_n(k)^{1/2}, N(w) = j\right] \\ & \leq |\mathcal{X}| \max_{x \in \mathcal{X}} \mathbb{P}[|\hat{P}_{wu}(x) - p_{wu}(x)| \geq a_n(k)^{1/2}, N(wu) = k] \\ & + |\mathcal{X}| \max_{x \in \mathcal{X}} \mathbb{P}[|\hat{P}_w(x) - p_w(x)| \geq a_n(k)^{1/2}, N(w) = j]. \end{aligned} \quad (5.4)$$

We will now choose $\rho_n = b_n n \geq n^{1/2} \log(n)^{1/2+\beta}$. Then, $\max_{k \geq \rho_n} \gamma_n^2(k) \leq \frac{C \log(n)}{nb_n}$. Thus, it follows by assumption on ϵ_n ,

$$\min_{k \geq \rho_n} a_n(k) = \min_{k \geq \rho_n} \left(\frac{\epsilon_n}{2} - \gamma_n(k)\right)^2 \geq \text{const.} \frac{\log(n)^{1+\delta}}{nb_n},$$

and hence

$$\min_{k \geq \rho_n} k a_n(k) \geq \text{const.} \log(n)^{1+\delta}.$$

We treat the two cases on the RHS of (5.4) simultaneously by denoting $v = wu$ or $v = w$, respectively. Let $p = P_v(x)$ and let $\hat{p} = \hat{P}_v(x)$. We would like to find an upper bound for the probability of the event $\{|p - \hat{p}|^2 > a_n(k), N(v) = k\}$. Since there are a random number of terms in the denominator of \hat{p} we cannot apply any large deviations bound directly. Instead we consider the extension of X_1^n to the infinite sequence $(X_t)_{t \in \mathbb{Z}}$. Define,

$$I_i = \{\text{the time of the } i^{\text{th}} \text{ occurrence of } v \text{ in } (X_t)_{t \in \mathbb{Z}}\}, \quad i \in \mathbb{N}.$$

Then let

$$W_i = X_{I_i+1}, \text{ the symbol that occurs after the } i^{\text{th}} \text{ occurrence of } v.$$

The sequence $(W_i)_{i \in \mathbb{N}}$ is a stationary ϕ -mixing sequence with mixing coefficients bounded by the same bound as the original sequence $(X_t)_{t \in \mathbb{Z}}$. The marginal probability distribution of W_1 on \mathcal{X} is equal to P_v . Let $Y_i = 1_{[W_i=x]}$. Now observe that

$$\left\{ \left| \sum_{i=1}^{N(v)} \frac{Y_i}{N(v)} - p \right|^2 > a_n(k), N(v) = k \right\} \subseteq \left\{ \left| \sum_{i=1}^k \frac{Y_i}{k} - p \right|^2 > a_n(k) \right\}.$$

Thus, we have established the upper bound,

$$\mathbb{P}[|\hat{p} - p|^2 > a_n(k), N(w) = k] \leq \mathbb{P}\left[\left|\sum_{i=1}^k \frac{Y_i}{k} - p\right|^2 > a_n(k)\right]. \quad (5.5)$$

At this point we are readily able to apply an exponential inequality.

Lemma 5.1 *Let $(Y_i)_{i \in \mathbb{N}}$ with $E[Y_1] = p$ be defined as above and $a_n(k)$ as in (5.3). Assume the conditions (A1) and (A2) with $\beta, \delta > 0$. Then, for $k \geq \rho_n = b_n n$,*

$$\sup_{0 < p < 1} \mathbb{P}\left[\left|\sum_{i=1}^k \frac{Y_i}{k} - p\right|^2 > a_n(k)\right] = O(n^{-\log(n)F\delta})$$

for some constant $0 < F < \infty$.

Proof: By assumption (A1), the process $(X_t)_{t \in \mathbb{Z}}$ has mixing coefficients $\phi(j) \leq (1 - 2\kappa)^j$, and the same bound applies also for the mixing coefficients of the process $(Y_i)_{i \in \mathbb{N}}$. Now apply Theorem 4 from Doukhan (1994, Ch. 1.4.2) with $q = \log(n)^{1+\delta}$ and note that $k \geq \rho_n = b_n n \geq \text{const.} n^{1/2} \log(n)^{1/2+\beta}$. Using that $ka_n(k) \geq \text{const.} \log(n)^{1+\delta}$ for all $k \geq \rho_n$ completes the proof. \square

Denote by $M_n = \exp(-F \log(n)^{1+\delta})$. A straightforward application of Lemma 5.1 to equation (5.5) proves that for $k, j \geq \rho_n$,

$$\begin{aligned} \max_{x \in \mathcal{X}} \mathbb{P}[(\hat{P}_{wu}(x) - p_{wu}(x))^2 \geq a_n(k), N(wu) = k] &= O(M_n), \\ \max_{x \in \mathcal{X}} \mathbb{P}[(\hat{P}_w(x) - p_w(x))^2 \geq a_n(k), N(w) = j] &= O(M_n). \end{aligned}$$

Thus, together with (5.1), (5.2) and (5.4),

$$\mathbb{P}[U_n \cap D_n] \leq |\mathcal{X}| \sum_{wu \in \tau_n} \sum_{k=\rho_n}^n \sum_{j=k}^n O(M_n). \quad (5.6)$$

By Remark 3.8, (5.6) and assumption (A2),

$$\begin{aligned} \mathbb{P}[U_n \cap D_n] &= O(|\tau_n| n^2 M_n) \\ &= O(b_n^{-1} n^2 \exp(-\text{const.} \log(n)^{1+\delta})) = O(n^{-\text{const.} \log(n)^\delta}). \end{aligned} \quad (5.7)$$

To complete the proof of Theorem 5.1, we need to bound $\mathbb{P}[D_n^c]$. Using the union bound we get,

$$\begin{aligned} \mathbb{P}[D_n^c] &\leq \sum_{w \in \tau_n} \mathbb{P}[N(w) < \rho_n] = \sum_{w \in \tau_n} \mathbb{P}[N(w) - nP_n(w) < \rho_n - nP_n(w)] \\ &\leq \sum_{w \in \tau_n} \mathbb{P}[N(w) - nP_n(w) < -b_n n/2] \leq \sum_{w \in \tau_n} \mathbb{P}[|N(w) - E[N(w)]| > b_n n/2]. \end{aligned}$$

We bound this quantity by the following exponential inequality.

Lemma 5.2 *Assume that (A1) and (A2) with $\beta, \delta > 0$ hold. Then,*

$$\max_{w \in \tau_n} \mathbf{P}[|N(w) - \mathbf{E}[N(w)]| \geq b_n n/2] = O(n^{-\log(n)G\beta})$$

for some constant $0 < G < \infty$.

Proof: Since $w \in \tau_n$ we can write

$$N(w) = \sum_{t=1}^n 1_{[Z_{t,n}=w]}, \quad Z_{t,n} = c(X_{-\infty,n}^{t,n}).$$

By assumption (A1), $(Z_{t,n})_{t \in \mathbb{Z}}$ is ϕ -mixing with mixing coefficients bounded by $\sup_{n \in \mathbb{N}} \phi_n(j) \leq (1 - 2\kappa)^j$, cf. Rajarshi (1990). Now apply Theorem 4 in Doukhan (1994, Ch. 1.4.2) with $q = \log(n)^{1+2\beta}$. Using that $nb_n \geq \text{const.} n^{1/2} \log(n)^{1/2+\beta}$ completes the proof. \square

By Lemma 5.2

$$\mathbf{P}[D_n^c] = O(b_n^{-1} n^{-\log(n)\text{const.}\beta}) = O(n^{-\log(n)\text{const.}\beta}).$$

where the last estimate follows from (A2). Together with (5.7) we complete the proof of Theorem 5.1. \square

We now consider the overestimation event $O_n = \{\text{there exists } w' \in \tau_n \text{ with } w = w'u \in \hat{\tau}, w = w'u \notin \tau_n \text{ for some } u \in \cup_{m=1}^{\infty} \mathcal{X}^m\}$. For a sequence w to be an element of $\hat{\tau}^T$, it is necessary that $N(w) > 1$ and $\Delta_w \geq C \log(n)$. Now Weinberger et. al. (1995) establish for any $w = w'u$ ($w' \in \tau_n, u \in \cup_{m=1}^{\infty} \mathcal{X}^m$) with $w = w'u \notin \tau_n$,

$$\mathbf{P}[\Delta_w \geq C \log(n+1)] \leq (n+1)^{2|\mathcal{X}|} (n+1)^{-C}.$$

In their algorithm, an overestimation event can only occur at any string w if $|w| \leq \frac{\log(n)}{\log(a)}$. Thus they establish that

$$\mathbf{P}[O_n] \leq \sum_{|w| \leq \frac{\log(n)}{\log(|\mathcal{X}|)}} (n+1)^{-C+2|\mathcal{X}|} \leq n^{-C+2|\mathcal{X}|+1}.$$

The last inequality follows since, for any m there are no more than $|\mathcal{X}|^m$ distinct sequences w with length $|w| = m$.

It is possible to prove a stronger result, eliminating the need for a length restriction on $|w|$. We just give an outline of such a proof.

Lemma 5.3 *Let swv be any possible string with $s \in \tau_n, w \in \cup_{m=0}^{\infty} \mathcal{X}^m, v \in \mathcal{X}$ and $swv \notin \tau_n$. Let $O_n(swv) = \{\Delta_{swv} \geq C \log(n), N(swv) > 1\}$. Denote by $p_{\min}(n) = \min_{x \in \mathcal{X}, w \in \tau_n} p_w(x)$ and by $\hat{\tau}_{\max}$ the maximal context tree in Step 1 of the context algorithm. Then, under the assumptions (A1)-(A3),*

$$\mathbf{P}[O_n(swv)] \leq \frac{1}{p_{\min}(n)} \mathbf{P}[sw \in \hat{\tau}_{\max}] n^{-C+2|\mathcal{X}|}.$$

A proof is given below.

Theorem 5.2 *Under the assumptions (A1)-(A3),*

$$\sum_{n=1}^{\infty} \mathbb{P}[O_n] \log(n) < \infty.$$

Proof: We apply Lemma 5.3 for swv ,

$$\mathbb{P}[O_n] \leq \sum_{swv} \mathbb{P}[O_n(swv)] = O(n^{-C+2|\mathcal{X}|+1}) \sum_{swv} \mathbb{P}[sw \in \hat{\tau}_{max}],$$

where the last estimate follows from (A3).

Let L be the number of sequences which occur at least twice in the data X_1^n . Then,

$$\sum_{swv} \mathbb{P}[sw \in \hat{\tau}_{max}] \leq |\mathcal{X}| \mathbb{E}[\sum_{sw} 1_{[sw \text{ occurs at least twice in } X_1^n]}] \leq |\mathcal{X}| \mathbb{E}[L] \leq |\mathcal{X}| n^2.$$

Therefore, since $C > 2|\mathcal{X}| + 4$ we complete the proof. \square

When defining the pruning criterion in Step 2 of the context algorithm in terms of the L_1 distance, we can sharpen Theorem 5.2. Let $\tilde{\Delta}_{wu} = \|\hat{P}_w(\cdot) - \hat{P}_{wu}(\cdot)\|_1^2$ and define $\tilde{O}_n = \{\text{there exists } w = w'u \text{ (} w' \in \tau_n, u \in \cup_{m=1}^{\infty} \mathcal{X}^m), \text{ such that } \tilde{\Delta}_w \geq C \log(n), N(w) > 1, \text{ and } w \notin \tau_n\}$.

Theorem 5.3 *Under the assumptions (A1)-(A3) but with cut-off in Step 2 of the context algorithm satisfying $K_n \sim C \log(n)$ for $C > 4|\mathcal{X}| + 8$,*

$$\sum_{n=1}^{\infty} \mathbb{P}[\tilde{O}_n] \log(n) < \infty$$

Proof: As used already in the proof of Theorem 5.1, $D(P||Q) \geq \frac{1}{2} \|P - Q\|_1^2$. Thus, $\tilde{\Delta}_w \leq 2\Delta_w$. \square

Proof of Lemma 5.3. Let $s \in \tau_n$ be a context and $su = swv$ with $w \in \cup_{m=0}^{\infty} \mathcal{X}^m$, $v \in \mathcal{X}$ and $swv \notin \tau_n$. Our aim is to bound the probability of overestimation at su . We begin by recalling several inequalities and definitions from Weinberger et. al. (1995). First, we fix a sequence x_1^n , being a realization from P_n . We can determine a probability law given by $Q_{su}(y_1^n | x_1^n)$ (on the set of sequences of length n), defined as follows:

$$\begin{aligned} \log(Q_{su}(y_1^n | x_1^n)) &= R_{sw}(y_1^n | S_s) + \sum_{x \in \mathcal{X}} \sum_{b \neq v} N_{y_1^n}(x | swb) \log(\hat{P}_{x_1^n}(x | sw)) \\ &\quad + \sum_{x \in \mathcal{X}} N_{y_1^n}(x | su) \log(\hat{P}_{x_1^n}(x | su)). \end{aligned}$$

where $R_{sw}(y_1^n | S_s)$, defined formally in Weinberger et al. (1995), is the sum of the log probability of all the symbols that occur in any context other than sw . An important observation is that for any sequence y_1^n with $N_{y_1^n}(sw) = 0$ the Q_{su} probability of y_1^n is the same as the P_n probability.

Now, for each x_1^n define $\sigma_{x_1^n}$ to be the set of all sequences y_1^n with $N_{y_1^n}(xsw) = N_{x_1^n}(xsw)$ and $N_{y_1^n}(xswv) = N_{x_1^n}(xswv)$ for all $x \in \mathcal{X}$. If $\Delta_{x_1^n}(swv) > C \log(n)$, it follows from (A9) in Weinberger et. al. (1995), that

$$P_n(\sigma_{x_1^n}) \leq Q_{su}(\sigma_{x_1^n}|x_1^n)n^{-C}. \quad (5.8)$$

At this point we need to introduce a new probability distribution given by Q' on the set of sequences of length n , closely related to Q_{su} . To that end, for every sequence y_1^t let x_0 be the symbol that occurs after the first occurrence of sw . Let b_0 be the symbol immediately preceding the first occurrence of sw . Thus x_0 occurs in the (extended) context swb_0 . If $b_0 \neq v$, we define

$$\log(Q'(y_1^n|x_1^n)) = \log(Q_{su}(y_1^t|x_1^n)) + \log(P_n(x_0|sw)) - \log(\hat{P}_{x_1^n}(x_0|sw)).$$

If $b_0 = v$ then we define

$$\log(Q'(y_1^n|x_1^n)) = \log(Q_{su}(y_1^t|x_1^n)) + \log(P_n(x_0|sw)) - \log(\hat{P}_{x_1^n}(x_0|swv)).$$

Thus, if $N_{y_1^n}(sw) < 2$ it must be that $P_n(y_1^n) = Q'(y_1^n|x_1^n)$. It also follows from the definition of Q' that

$$Q_{su}(y_1^n|x_1^n) \leq \frac{1}{p_{\min}(n)} Q'(y_1^n|x_1^n).$$

Therefore, together with (5.8) we have the bound,

$$P_n(\sigma_{x_1^n}) \leq Q'(\sigma_{x_1^n}) \frac{1}{p_{\min}(n)} n^{-C}.$$

The construction of $\sigma_{x_1^n}$ and the fact that $N_{x_1^n}(sw) > 1$ implies that

$$Q'(\sigma_{x_1^n}|x_1^n) \leq Q'(y_1^n; N_{y_1^n}(sw) > 1|x_1^n) = P_n(y_1^n; N_{x_1^n} > 1) = P_n(sw \in \hat{\tau}_{max}).$$

Furthermore, since there are at most $n^{2|\mathcal{X}|}$ distinct classes $\sigma_{x_1^n}$ it follows that

$$\mathbb{P}[O_n(swv)] = P_n(y_1^n; \Delta_{y_1^n}(swv) > C \log(n)) \frac{1}{p_{\min}(n)} \leq P_n(sw \in \hat{\tau}_{max}) n^{-C+2|\mathcal{X}|}.$$

□

Theorem 5.1 and 5.2 imply the assertion in Theorem 3.1 (i). The assertion in Theorem 3.1 (ii) follows from Theorem 3.1 (i) and along the lines of the proof of Theorem 3.1 (i): partition with the set D_n and use Lemma 5.1 and 5.2. □

Proof of Theorem 3.2.

Statements (i) and (iii) follow from the general formula (5.10) below, statement (ii) is an immediate consequence of Theorem 3.1 (ii).

We give here the analogon of assumption (A1) for the estimated process \hat{P}_n . The r -step transition kernel $p_{Z_n}^{(r)}(v, w) = \mathbb{P}[Z_{r,n} = v | Z_{0,n} = w]$ (for some $r \in \mathbb{N}$) for the state process

$Z_{t,n} = c(X_{-\infty,n}^{t,n})$ of a VLMC $(X_{t,n})_{t \in \mathbb{Z}}$ can be characterized by the transition probabilities $p_n(\cdot|\cdot)$ and the context function $c_n(\cdot)$, i.e.,

$$\begin{aligned} T(v|wX_{-\infty,n}^{-|w|,n}; r, p_n(\cdot|\cdot), c_n(\cdot)) &= p_{Z_n}^{(r)}(v, w) \\ &= \sum_{x_1^{r-1} \in \mathcal{X}^{r-1}, c_n(x_r \dots x_1 w X_{-\infty,n}^{-|w|,n})=v} \prod_{i=0}^{r-1} p_n(x_{r-i}|c(x_1^{r-i-1} w X_{-\infty,n}^{-|w|,n})). \end{aligned} \quad (5.9)$$

(Note that x_r is the first component of v). For every $n \in \mathbb{N}$, the process $(\hat{X}_{t,n})_{t \in \mathbb{Z}} \sim \hat{P}_n$ is a VLMC. We consider its r -step transition kernel for the states $\hat{P}^{(r)}(v, w) = \mathbb{P}_{\hat{P}_n}[\hat{Z}_{r,n} = v | \hat{Z}_{0,n} = w]$ ($r \geq 1$), where $\hat{Z}_{t,n} = \hat{c}_n(\hat{X}_{-\infty,n}^{t,n})$ is the state process of \hat{P}_n . This transition is characterized by

$$T(v|w\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}_n(\cdot|\cdot), \hat{c}_n(\cdot)) = \hat{P}^{(r)}(v, w) \quad (r \geq 1).$$

We now obtain an analogon of (A1) for \hat{P}_n . We have,

$$\begin{aligned} &|T(v|w\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}(\cdot|\cdot), \hat{c}(\cdot)) - T(v|w'\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}(\cdot|\cdot), \hat{c}(\cdot))| \\ &\leq |T(v|w\hat{X}_{-\infty,n}^{-|w|,n}; r, p(\cdot|\cdot), c(\cdot)) - T(v|w'\hat{X}_{-\infty,n}^{-|w|,n}; r, p(\cdot|\cdot), c(\cdot))| \\ &+ |T(v|w\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}(\cdot|\cdot), c(\cdot)) - T(v|w\hat{X}_{-\infty,n}^{-|w|,n}; r, p(\cdot|\cdot), c(\cdot))| \\ &+ |T(v|w'\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}(\cdot|\cdot), c(\cdot)) - T(v|w'\hat{X}_{-\infty,n}^{-|w|,n}; r, p(\cdot|\cdot), c(\cdot))| \\ &+ 2\mathbf{1}_{[\hat{c}(y) \neq c(y) \text{ for some } y \in \mathcal{X}^\infty]}. \end{aligned}$$

We now invoke (A1) for $T(\cdot|\cdot; r, p(\cdot|\cdot), c(\cdot))$ about the true underlying process. For the other terms we use the finiteness of r and \mathcal{X} , together with (5.9) and Theorem 3.1. We then obtain,

$$\begin{aligned} &\sup_{v, w, w' \in \hat{\tau}_n} |\hat{P}^{(r)}(v, w) - \hat{P}^{(r)}(v, w')| \\ &= \sup_{v, w, w' \in \hat{\tau}_n} |T(v|w\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}(\cdot|\cdot), \hat{c}(\cdot)) - T(v|w'\hat{X}_{-\infty,n}^{-|w|,n}; r, \hat{P}(\cdot|\cdot), \hat{c}(\cdot))| \\ &\leq 1 - 2\kappa + o_P(1), \end{aligned} \quad (5.10)$$

We then consider sets

$$A_n = \{\omega; \sup_{A \subseteq \tau_n; w, w' \in \tau_n} |\hat{P}^{(r)}(A, w) - \hat{P}^{(r)}(A, w')|(\omega) < 1 - \kappa \text{ and } \hat{c}(\cdot; \omega) = c_n(\cdot)\}.$$

On A_n , \hat{P}_n as constructed in Theorem 3.2(i) is uniquely determined, stationary and ϕ -mixing, with mixing coefficients bounded by

$$\phi_{\hat{P}_n}(k) \leq (1 - \kappa)^k \text{ for all } k \in \mathbb{N}_0 \text{ on the set } A_n,$$

cf. Rajarshi (1990, Lem. 2.1) or Doukhan (1994).

But by (5.10) and Theorem 3.1(i), $\mathbb{P}[A_n] \rightarrow 1$ as $n \rightarrow \infty$, which completes the proof of Theorem 3.2(i) and (iii). \square

Proof of Theorem 4.1.

We usually suppress the index n when writing X_t instead of $X_{t,n}$. Consider

$$U_n = (n - m + 1)^{-1} \sum_{t=1}^{n-m+1} f(X_t^{t+m-1}),$$

and denote by $\Sigma_n = \text{Cov}[U_n]$ the covariance matrix of U_n .

Lemma 5.4 *Assume (B1) with $(X_{t,n})_{t \in \mathbb{Z}} \sim P_n$ satisfying (A1). Then,*

(i) *there exists $n_0 \in \mathbb{N}$ such that $n\Sigma_n$ is positive definite for all $n \geq n_0$.*

(ii) *for $Z \sim \mathcal{N}_v(0, I)$,*

$$\sup_{x \in \mathbb{R}^v} |\mathbb{P}[\Sigma_n^{-1/2}(U_n - \theta_n) \leq x] - \mathbb{P}[Z \leq x]| = o(1) \quad (n \rightarrow \infty).$$

Proof: For every $n \in \mathbb{N}$, the process $(X_{t,n})_{t \in \mathbb{Z}}$ is ϕ_n -mixing whose mixing coefficients are bounded by

$$\sup_{n \in \mathbb{N}} \phi_n(k) \leq (1 - 2\kappa)^k \text{ for all } k \in \mathbb{N}, \quad (5.11)$$

cf. Rajarshi (1990, Lem. 2.1).

Bounding covariances in terms of mixing coefficients, cf. Doukhan (1994), and using the bound in (5.11) implies for $i, j \in \{1, \dots, v\}$,

$$(n - m + 1)(\Sigma_n)_{i,j} = \sum_{k=-n+1}^{n-1} \text{Cov}(f_i(X_0^{m-1}), f_j(X_k^{k+m-1})) + O(n^{-1}). \quad (5.12)$$

Hence, assertion (i) follows from the assumption in (B1).

Assertion (i), assumption (B1) and (5.12) allow us to write

$$\Sigma_n^{-1/2} = n^{1/2}, \quad \sup_{n \in \mathbb{N}} \max_{1 \leq i, j \leq v} |(\cdot, \cdot)_{i,j}| < \infty. \quad (5.13)$$

Now write

$$\Sigma_n^{-1/2}(U_n - \theta_n) = n^{-1/2} \sum_{t=1}^{n-m+1} \tilde{f}_n(X_t^{t+m-1})(1 + o(1)),$$

where $\tilde{f}_n(X_t^{t+m-1}) = n(f(X_t^{t+m-1}) - \theta_n)$.

By construction and (5.13),

$$\begin{aligned} \mathbb{E}[\tilde{f}_n(X_1^m)] &= 0, \\ \text{Cov}(n^{-1/2} \sum_{t=1}^{n-m+1} \tilde{f}_n(X_t^{t+m-1})) &\rightarrow I \quad (n \rightarrow \infty), \\ \sup_{n \in \mathbb{N}} \mathbb{E} \|\tilde{f}_n(X_1^m)\|^2 &< \infty. \end{aligned} \quad (5.14)$$

We can then apply Theorem 2.1 in Withers (1981) to $n^{-1/2} \sum_{t=1}^{n-m+1} \tilde{f}_n(X_t^{t+m-1})$. The conditions (version (A) or (B), note also the corrigendum in Vol. 63) are easily verified by invoking the mixing bound in (5.11) and (5.14). Thus,

$$n^{-1/2} \sum_{t=1}^{n-m+1} \tilde{f}_n(X_t^{t+m-1}) \Rightarrow \mathcal{N}_v(0, I),$$

and assertion (ii) follows by Polya's Theorem. \square

By the smoothness assumption about g we use a first order Taylor expansion,

$$n^{1/2}(T_n - g(\theta_n)) = n^{1/2}Dg(\tilde{\theta}_n)(U_n - \theta_n), \quad (5.15)$$

where $Dg(\theta) = [\frac{\partial g_i(u)}{\partial u_j}]_{i,j}$, ($1 \leq i \leq w$, $1 \leq j \leq v$) and $\|\tilde{\theta}_n - \theta_n\| \leq \|U_n - \theta_n\|$.

By (5.13) and Lemma 5.4 (ii), $U_n - \theta_n = o_P(1)$, so that

$$[Dg(\tilde{\theta}_n) - Dg(\theta_n)]_{i,j} = o_P(1), \quad 1 \leq i \leq w, \quad 1 \leq j \leq v.$$

This, together with (5.15), the boundedness of $n^{1/2}\Sigma_n^{1/2}$ (use (5.12)) and Lemma 5.4 (ii) implies

$$\sup_{x \in \mathbb{R}^w} |\mathbb{P}[n^{1/2}(T_n - g(\theta_n)) \leq x] - \mathbb{P}[n^{1/2}\Sigma_n^{1/2}Dg(\theta_n)Z \leq x]| = o(1) \quad (n \rightarrow \infty), \quad (5.16)$$

where $Z \sim \mathcal{N}_v(0, I)$.

We are going now to show the bootstrap analog of (5.16). By Theorem 3.2(i) and (iii), the bootstrap process $(X_t^*)_{t \in \mathbb{Z}}$ is with high probability stationary and geometrically ϕ -mixing with mixing coefficients denoted by $\phi_n^*(k) = \phi_{\hat{\rho}_n}(k)$ from Theorem 3.2(iii). Note that the distribution of $(X_t^*)_{t \in \mathbb{Z}}$ depends again on the sample size n .

Denote by $U_n^* = (n - m + 1)^{-1} \sum_{t=1}^{n-m+1} f((X_t^*)^{t+m-1})$ and let $\Sigma_n^* = \text{Cov}^*[U_n^*]$ be the covariance matrix of U_n^* with respect to the bootstrap distribution.

Lemma 5.5 *Assume the conditions of Theorem 4.1. Then,*

(i) $n(\Sigma_n^* - \Sigma_n)_{i,j} = o_P(1) \quad (n \rightarrow \infty), \quad i, j = 1, \dots, v.$

(ii) $\lim_{n \rightarrow \infty} \mathbb{P}[n\Sigma_n^* \text{ is positive definite}] = 1.$

(iii) for $Z \sim \mathcal{N}_v(0, I)$,

$$\sup_{x \in \mathbb{R}^v} |\mathbb{P}^*[(\Sigma_n^*)^{-1/2}(U_n^* - \theta_n^*) \leq x] - \mathbb{P}[Z \leq x]| = o_P(1) \quad (n \rightarrow \infty).$$

Proof: For any $i, j \in \{1, \dots, v\}$,

$$\begin{aligned} n(\Sigma_n^*)_{i,j} &= \sum_{k=-n+m}^{n-m} \text{Cov}^*(f_i((X^*)_0^{m-1}), f_j((X^*)_k^{k+m-1})) \left(1 - \frac{|k|}{n-m+1}\right) \\ &= \sum_{k=-M}^M \text{Cov}^*(f_i((X^*)_0^{m-1}), f_j((X^*)_k^{k+m-1})) \left(1 - \frac{|k|}{n-m+1}\right) + \Delta_{n,M}, \end{aligned} \quad (5.17)$$

where M is a finite constant.

By well known bounds of covariances in terms of mixing coefficients, cf. Doukhan (1994),

$$|\Delta_{n,M}| \leq 2const. \sum_{k=M+1}^{\infty} \phi_n^*(k).$$

Therefore by Theorem 3.2(iii),

$$\mathbb{P}[\lim_{M \rightarrow \infty} |\Delta_{n,M}| = 0] \rightarrow 1 \quad (n \rightarrow \infty). \quad (5.18)$$

By Theorem 3.2 (ii),

$$\max_{x_1^d \in \mathcal{X}^d} |\mathbb{P}^*[(X^*)_1^d = x_1^d] - \mathbb{P}[X_1^d = x_1^d]| = o_P(1) \quad (d \in \mathbb{N}). \quad (5.19)$$

This, the boundedness of f and the finiteness of M imply,

$$\begin{aligned} & \left| \sum_{k=-M}^M Cov^*(f_i((X^*)_0^{m-1}), f_j((X^*)_k^{k+m-1})) - \sum_{k=-M}^M Cov(f_i(X_0^{m-1}), f_j(X_k^{k+m-1})) \right| \\ &= o_P(1) \quad (n \rightarrow \infty). \end{aligned} \quad (5.20)$$

By the geometric ϕ -mixing property of $(X_t)_{t \in \mathbb{Z}}$, see (5.11), and the boundedness of f ,

$$\begin{aligned} & \left| \sum_{k=-M}^M Cov(f_i(X_0^{m-1}), f_j(X_k^{k+m-1})) - \sum_{k=-\infty}^{\infty} Cov(f_i(X_0^{m-1}), f_j(X_k^{k+m-1})) \right| \\ &= o(1) \quad (M \rightarrow \infty). \end{aligned} \quad (5.21)$$

Thus, by (5.17)-(5.21) we have shown assertion (i).

Assertion (ii) follows by (i) and Lemma 5.4 (i).

Assertion (iii) can be proved as Lemma 5.4 (ii); we now invoke the mixing bound in Theorem 3.2(iii) and use (i). \square

By (5.19) and the finiteness of $|\mathcal{X}|$ we have,

$$\theta^* - \theta_n = \sum_{x_1^m \in \mathcal{X}^m} f(x_1^m) (\mathbb{P}^*[(X^*)_1^m = x_1^m] - \mathbb{P}[X_1^m = x_1^m]) = o_P(1), \quad (5.22)$$

and hence by the continuous differentiability of g ,

$$[Dg(\tilde{\theta}_n^*) - Dg(\theta_n)]_{i,j} = o_P(1) \quad \text{for } \|\tilde{\theta}_n^* - \theta^*\| \leq \|U_n^* - \theta^*\|, \quad (1 \leq i \leq w, 1 \leq j \leq v). \quad (5.23)$$

A first order Taylor expansion, (5.23), Lemma 5.5 (iii) and the boundedness of $n\Sigma_n^* = O_P(1)$ imply

$$\sup_{x \in \mathbb{R}^w} |\mathbb{P}^*[n^{1/2}(T_n^* - g(\theta_n^*)) \leq x] - \mathbb{P}[n^{1/2}\Sigma_n^{1/2}Dg(\theta_n)Z \leq x]| = o_P(1) \quad (n \rightarrow \infty), \quad (5.24)$$

where $Z \sim \mathcal{N}_v(0, I)$.

By (5.16) and (5.24) we complete the proof of Theorem 4.1. \square

Acknowledgments. We thank Itai Zukerman for carrying out the computations. We also acknowledge interesting general comments from P. Bickel, H.-R. Künsch, R. Olshen, T. Speed and specific suggestions from F. Ferrari, the Editor and the referees.

References

- [1] Bickel, P.J., Götze, F. and van Zwet, W.R. (1997). Resampling fewer than n observations: gains, losses, and remedies for losses. *Statistica Sinica* **7** 1-32.
- [2] Breiman, L., Friedman, J.H., Olshen, R.A. and Stone, C.J. (1984). *Classification and Regression Trees*. Wadsworth.
- [3] Brillinger, D.R. (1995). Trend analysis: binary-valued and point cases. *Stochastic Hydrology and Hydraulics* **9** 207-213.
- [4] Braun, J.V. and Müller, H.-G. (1998). Statistical methods for DNA sequence. *Statist. Sci.* **13** 142-162.
- [5] Bühlmann, P. (1998). Model selection for variable length Markov chains and tuning the context algorithm. To appear in *Ann. Inst. Statist. Math.*
- [6] Cover, T.M. and Thomas, J.A. (1991). *Elements of Information Theory*. Wiley.
- [7] Doukhan, P. (1994). *Mixing. Properties and Examples*. *Lect. Notes in Stat.* **85**. Springer.
- [8] Efron, B. (1979). Bootstrap methods: another look at the jackknife. *Ann. Statist.* **7** 1-26.
- [9] Fahrmeir, L. and Tutz, G. (1994). *Multivariate Statistical Modelling Based in Generalized Linear Models*. Springer.
- [10] Feder, M., Merhav, N. and Gutman, M. (1992). Universal prediction of individual sequences. *IEEE Trans. Inform. Theory* **IT-38** 1258-1270.
- [11] Guttorp, P. (1995). *Stochastic Modeling of Scientific Data*. Chapman & Hall.
- [12] Iosifescu, M. and Theodorescu, R. (1969). *Random Processes and Learning*. Springer.
- [13] Künsch, H.R. (1989). The jackknife and the bootstrap for general stationary observations. *Ann. Statist.* **17** 1217-1241.
- [14] Prum, B., Rodolphe, F. and deTurckheim, E. (1995). Finding words with unexpected frequencies in deoxyribonucleic acid sequences. *J. Roy. Statist. Soc. B* **57** 205-220.
- [15] Raftery, A. and Tavaré, S. (1994). Estimation and modelling repeated patterns in high order Markov chains with the mixture transition distribution model. *Appl. Statist.* **43** 179-199.
- [16] Rajarshi, M.B. (1990). Bootstrap in Markov-sequences based on estimates of transition density. *Ann. Inst. Statist. Math.* **42** 253-268.
- [17] Rissanen, J. (1983). A universal data compression system. *IEEE Trans. Inform. Theory* **IT-29** 656-664.
- [18] Rissanen, J. (1986). Complexity of strings in the class of Markov sources. *IEEE Trans. Inform. Theory* **IT-32** 526-532.

- [19] Rissanen, J. (1989). Stochastic Complexity in Statistical Inquiry. World Scientific.
- [20] Ritov, Y. and Bickel, P.J. (1990). Achieving information bounds in non and semi-parametric models. Ann. Statist. **18** 925-938.
- [21] Weinberger, M.J. and Feder, M. (1994). Predictive stochastic complexity and model estimation for finite-state processes. J. Statist. Plann. Infer. **39** 353-372.
- [22] Weinberger, M.J., Lempel, A. and Ziv, J. (1992). A sequential algorithm for the universal coding of finite memory sources. IEEE Trans. Inform. Theory **IT-38** 1002-1014.
- [23] Weinberger, M.J., Rissanen, J.J. and Feder, M. (1995). A universal finite memory source. IEEE Trans. Inform. Theory **IT-41** 643-652.
- [24] Withers, C.S. (1981). Central limit theorems for dependent variables I. Z. Wahrsch. verw. Gebiete **57** 509-534 (Corr: **63** p555).

Seminar für Statistik
 ETH Zürich
 CH-8092 Zürich
 Switzerland
 E-mail: buhlmann@stat.math.ethz.ch

Department of Statistics
 University of Pennsylvania
 Philadelphia 19104
 USA
 E-mail: ajw@compstat.wharton.upenn.edu