# Probabilistic and Statistical Properties of Words: An Overview

GESINE REINERT[1] and SOPHIE SCHBATH[2] and MICHAEL S. WATERMAN[3]

### Abstract

In the following, an overview is given on statistical and probabilistic properties of words, as occurring in the analysis of biological sequences. Counts of occurrence, counts of clumps and renewal counts are distinguished, and exact distributions as well as normal approximations, Poisson process approximations and compound Poisson approximations are derived. Here, a sequence is modelled as a stationary ergodic Markov chain; a test for determining the appropriate order of the Markov chain is described. The convergence results take the error made by estimating the Markovian transition probabilities into account. The main tools involved are moment generating functions, martingales, Stein's method, and the Chen-Stein method. Similar results are given for occurrences of multiple patterns, and, as an example, the problem of unique recoverability of a sequence from SBH chip data is discussed. Special emphasis lies on disentangling the complicated dependence structure between word occurrences, due to self-overlap as well as due to overlap between words. The results can be used to derive approximate, and conservative, confidence intervals for tests.

**Key words:** Word counts, renewal counts, Markov model, exact distribution, normal approximation, Poisson process approximation, compound Poisson approximation, occurrences of multiple words, sequencing by hybridization; martingales, moment generating function, Stein's method, Chen-Stein method
.

## 1 Introduction

Statistical and probabilistic properties of words have been of considerable interest in many fields, such as reliability theory, and most recently in the analysis of biological sequences. Here we provide an overview of the state of this research.

Two main aspects of word occurrences in biological sequences are: where do they occur and how many times do they occur? An important problem, for instance, was to determine the statistical significance of a word frequency in a DNA sequence. The naive idea is the following: a word may be significantly rare in a DNA sequence because it disrupts replication or gene expression, (perhaps a counter-selection factor), whereas a significantly frequent word may have a fundamental activity with regard to genome stability. Well-known examples of words with exceptional frequencies in DNA sequences are certain biological palindromes corresponding to restriction sites avoided for instance in *E. coli* (Karlin *et al.* (1992)), and the Cross-over Hotspot Instigator sites in several bacteria (see Biaudet *et al.* (1998), Chedin *et al.* (1998), Sourice *et al.* (1998)). Several papers aim at identifying over- and under-represented words in a particular genome (for instance, Leung *et al.* (1996), Rocha *et al.* (1998)). Statistical methods to study the distribution of the word locations along a sequence and word frequencies have also been an active field of research.

Because DNA sequences are long, asymptotic distributions were proposed first. Exact distributions exist now, motivated by the analysis of genes and protein sequences. Unfortunately, exact results are not adapted in practice for long sequences because of heavy numerical calculation, but they allow the validation of the quality of the stochastic approximations when no approximation error can be provided. For example, BLAST is probably the best-known algorithm for DNA matching, and it relies on a Poisson approximation. Approximate *p*-values can be given; yet the applicability of the Poisson approximation needs to be justified.

Statistical properties of words only make sense with respect to some underlying probability model. DNA sequences are commonly modeled as stationary random sequences. Typical models are homogeneous

[1] King's College and, Statistical Laboratory, Cambridge CB2 1ST, UK.
[2] Unité de Biométrie, INRA, 78352 Jouy-en-Josas, France.
[3], Departments of Mathematics, Department of Biological Sciences, and Department of Computer Science, USC, Los Angeles CA 90089, USA.

$m$-order Markov chains (model M$m$) in which the probability of occurrence of a letter at a given position depends only on the $m$ previous letters in the sequence (and not on the position); the independent case is a particular case with $m = 0$. Hidden Markov models (HMMs) reveal however that the composition of a DNA sequence may vary over the sequence (Churchill (1989), Muri 1998), Durbin *et al.* (1998)) and can be studied with HMMs. However, no statistical properties of words have been yet derived in such heterogeneous models. DNA sequences code for amino acid sequences (proteins) by non-overlapping triplets called *codons*. The three positions of the codons have distinct statistical properties, so that for coding DNA we naturally think of three sequences where the successive letters come from the three codon positions, respectively. The three chains and their transition matrices are denoted as M$m$-3. In this paper, we will focus on the homogeneous models M$m$ and give existing results for M$m$-3.

Because these probabilistic models have to be fitted to the observed biological sequence, we will pay attention to the influence of the model parameter estimation on the statistical results. Some asymptotic results take care of this problem but the exact results require that the true model driving the observed sequence is known.

The choice of the Markov model order depends on the sequence length, because of the data requirements in estimation. One might be able to test hierarchical models using Chi-square tests to assign which of Markovian dependence is appropriate for the underlying sequence. From a practical point of view, it also depends on the composition of the biological sequence one wants to take into account. Indeed, if the sequence was generated from an $m$-order Markov chain, then the model M$m$ should predict well the $(m + 1)$-letter words.

In this paper, we are concerned firstly with the occurrences of a single pattern in a sequence. To begin, we discuss the underlying probabilistic models (Section 2). The main complication for word occurrences arises from overlaps of words. One might be interested either in overlapping occurrences or in particular non-overlapping ones (Section 3). After presenting results for the statistical distribution of word locations along the sequence (Section 4), we focus on the distribution of the number of overlapping occurrences (Section 5) and the number of renewals (Section 6). In Section 7, we will study the occurrences of multiple patterns. Section 8 gives an example on how probabilistic and statistical considerations come into play for DNA sequence analysis. Namely, we analyze so-called SBH chips, a fast and effective method for determining a DNA sequence. These chips provide the $\ell$-tuple contents of a DNA sequence, where typically $\ell = 8, 10$ or $12$. A nontrivial combinatorial problem arises when determining the probability that a randomly chosen DNA sequence can be uniquely reconstructed from its $\ell$-tuple contents. Finally, Section 9, meant as an appendix, gives a compilation of more general techniques that are applied in this paper. Throughout we only consider finite words.

Necessarily, due to the abundance of literature, much of the existing work on probabilistic and statistical aspects of words had to be omitted. The present paper is intended not to serve as a complete literature survey (indeed even just a list of references would take up all the space of this volume), but rather to introduce the reader to the major aspects of this field, to provide some techniques and to warn of major pitfalls associated with the analysis of words.

For the same reason we completely omit the algorithmic aspect; an excellent starting point would be Waterman (1995) or Gusfield (1997); for a particular example see also Apostolico *et al.* (1998).

# 2 Probabilistic models for biological sequences

In this paper, a biological sequence is either a DNA sequence or a protein sequence, that is, a finite sequence of letters either in the 4-letter DNA alphabet $\{A, C, G, T\}$ or the 20-letter amino-acid alphabet. To model a biological sequence, we will consider models for random sequences of letters. Even if we observed a finite biological sequence $\underline{S} = s_1 s_2 \cdots s_n$, we consider for convenience an infinite random sequence $\underline{X} = (X_i)_{i \in \mathbb{Z}}$ on a finite alphabet $\mathcal{A}$, where $\mathbb{Z}$ is the set of integers. We present below two classes of Markov models widely used to analyze biological sequences and how to estimate their parameters according to the observed sequence. Then we give a classical Chi-square test to choose the appropriate order of the Markov model for a given sequence.

## 2.1 Models for random sequences of letters

The simplest model assumes that the letters $X_i$'s are independent and take on the value $a \in \mathcal{A}$ with probability $\mu(a) = 1/|\mathcal{A}|$, where $|\mathcal{A}|$ denotes the size of the alphabet. To refine this model, we can simply assume independent letters taking value in $\mathcal{A}$ with probability $(\mu(a))_{a \in \mathcal{A}}$ such that $\sum_{a \in \mathcal{A}} \mu(a) = 1$. This is called model M0. In practice, in particular for DNA sequences, this model is typically not very accurate.

Therefore, we consider a much more general homogeneous model, the model M$m$: an ergodic stationary $m$-order Markov chain on a finite alphabet $\mathcal{A}$ with transition matrix $\Pi = (\pi(a_1 \cdots a_m, a_{m+1}))_{a_1, \ldots, a_{m+1} \in \mathcal{A}}$ such that

$$\pi(a_1 \cdots a_m, a_{m+1}) = \mathbb{P}(X_i = a_{m+1} \mid X_{i-1} = a_m, \ldots, X_{i-m} = a_1).$$

In general, a stationary distribution $\mu$ of an ergodic stationary Markov chain with transition matrix $\Pi$ is defined as a solution of $\mu = \mu\Pi$. This implies that the above Markov chain has a unique stationary distribution $\mu$ on $\mathcal{A}^m$ defined by

$$
\begin{aligned}
\mu(a_1 \cdots a_m) &= \mathbb{P}(X_i \cdots X_{i+m-1} = a_1 \cdots a_m), \quad \forall i \in \mathbb{Z} \\
&= \sum_{b \in \mathcal{A}} \mu(ba_1 \cdots a_{m-1})\pi(ba_1 \cdots a_{m-1}, a_m).
\end{aligned}
$$

The model where the letters $\{X_i\}_{i \in \mathbb{Z}}$ are chosen independently with probabilities $p_1, p_2, \ldots, p_{|\mathcal{A}|}$, corresponds to the transition matrix $\Pi$ with identical rows $(p_1 \; p_2 \; \cdots \; p_{|\mathcal{A}|})$ and stationary distribution $\mu = (p_1, p_2, \ldots, p_{|\mathcal{A}|})$.

A coding DNA sequence is naturally read as successive non-overlapping 3-letter words called codons. These codons are then translated into amino acids via the genetic code to produce a protein sequence. Several different codons can code for the same amino acid, and often the first two letters of a codon suffice to determine the corresponding amino acid. Therefore, letters may have different importance depending on their position with respect to the codon partition. To distinguish the letter probabilities according to their position modulo 3 in the coding DNA sequence, we consider a stationary Markov chain with three distinct transition matrices $\Pi_1$, $\Pi_2$ and $\Pi_3$ such that, for $a_1, \ldots, a_{m+1} \in \mathcal{A}$ and $k \in \{1, 2, 3\}$

$$\pi_k(a_1 \cdots a_m, a_{m+1}) = \mathbb{P}(X_{3j+k} = a_{m+1} \mid X_{3j+k-1} = a_m, \ldots, X_{3j+k-m} = a_1).$$

This is model M$m$-3. The index $k \in \{1, 2, 3\}$ is called *phase* and represents the position of a letter inside a codon. By convention, the phase of a word is the phase of its last letter in the sequence; codons are then 3-letter words in phase 3.

The stationary distribution $\mu$ on $\mathcal{A}^m \times \{1, 2, 3\}$ is given by

$$
\begin{aligned}
\mu(a_1 \cdots a_m, k) &= \mathbb{P}(X_{3j+k-m+1} \cdots X_{3j+k} = a_1 \cdots a_m), \quad \forall j \in \mathbb{Z} \\
&= \sum_{b \in \mathcal{A}} \mu(ba_1 \cdots a_{m-1}, k-1)\pi_k(ba_1 \cdots a_{m-1}, a_m).
\end{aligned}
$$

## 2.2 Estimation of the model parameters

Modeling a biological sequence consists of choosing a probabilistic model (see previous paragraph) and then estimating the model parameters according to the unique realization that is the biological sequence. In the case of model M$m$, it means to estimate the transition probabilities $\pi(a_1 \cdots a_m, a_{m+1})$; their estimators are classically denoted by $\widehat{\pi}(a_1 \cdots a_m, a_{m+1})$.

We now derive the estimators that maximize the likelihood of the M1 model given the observed sequence; we will then give the maximum-likelihood estimators in models M$m$ and M$m$-3.

Assume $X_1 \cdots X_n$ is a stationary Markov chain on $\mathcal{A}$ with transition matrix $\Pi = (\pi(a, b))_{a, b \in \mathcal{A}}$ and stationary distribution $(\mu(a))_{a \in \mathcal{A}}$. The likelihood $L$ of the model is

$$L(\pi(a, b), a, b \in \mathcal{A}) = \mu(X_1) \prod_{a, b \in \mathcal{A}} (\pi(a, b))^{N(ab)},$$

where $N(ab)$ denotes the number of occurrences of the 2-letter word $ab$ in the random sequence $X_1 \cdots X_n$. To find the transition probabilities that maximize the likelihood, one maximizes the log likelihood

$$\log L(\pi(a, b), a, b \in \mathcal{A}) = \log \mu(X_1) + \sum_{a, b \in \mathcal{A}} N(ab) \log \pi(a, b).$$

One can separately maximize $\sum_{b \in \mathcal{A}} N(ab) \log \pi(a, b)$ for $a \in \mathcal{A}$, keeping in mind that $\sum_{b \in \mathcal{A}} \pi(a, b) = 1$. Let $a \in \mathcal{A}$ and choose $c \in \mathcal{A}$; we have

$$\sum_{b \in \mathcal{A}} N(ab) \log \pi(a, b) = \sum_{b \neq c} N(ab) \log \pi(a, b) + N(ac) \log \left(1 - \sum_{b \neq c} \pi(a, b)\right),$$

and for $b \neq c$

$$\frac{\partial}{\partial \pi(a,b)} \left( \sum_{b \in \mathcal{A}} N(ab) \log \pi(a,b) \right) = \frac{N(ab)}{\pi(a,b)} - \frac{N(ac)}{\pi(a,c)}.$$

All the partial derivatives equal to zero means that

$$\frac{N(ab)}{\pi(a,b)} = \frac{N(ac)}{\pi(a,c)} \quad \forall b \in \mathcal{A};$$

this implies in particular that

$$\frac{N(ab)}{\pi(a,b)} = \frac{\sum_{d \in \mathcal{A}} N(ad)}{\sum_{d \in \mathcal{A}} \pi(a,d)} = \sum_{d \in \mathcal{A}} N(ad) := N(a\bullet) \quad \forall b \in \mathcal{A}.$$

It follows that

$$\widehat{\pi}(a,b) = \frac{N(ab)}{N(a\bullet)} \quad \forall b \in \mathcal{A}.$$

Note that the second partial derivatives of the likelihood function are negative, assuring that we have indeed determined a maximum.

**Remark 1** *For convenience with the notation, the estimators used in the remainder of the paper will be $\widehat{\pi}(a,b) = N(ab)/N(a)$ since $N(a\bullet) = N(a)$ except for the last letter of the sequence for which the counts differ by 1.*

It is important to note that the estimators $\widehat{\pi}(a,b)$ are random variables. Assuming that the biological sequence is a realization of the random sequence, one can calculate a numerical value for the estimator of $\pi(a,b)$; that is

$$\widehat{\pi}^{\mathrm{obs}}(a,b) = \frac{N^{\mathrm{obs}}(ab)}{N^{\mathrm{obs}}(a\bullet)},$$

where $N^{\mathrm{obs}}(\cdot)$ denotes the observed count in the biological sequence. As we will see, some results are obtained assuming that the true parameters $\pi(a,b)$ are known and equal, in practice, to $N^{\mathrm{obs}}(ab)/N^{\mathrm{obs}}(a\bullet)$, and do not take care of the estimation. It is indeed a common practice to substitute the estimator for the corresponding parameter in distributional results, but sometimes it changes the distribution being studied, as it is illustrated in Waterman (1995) p. 313.

In the model M$m$, the maximum-likelihood estimator of $\pi(a_1 \cdots a_m, a_{m+1})$, $a_1, \ldots, a_{m+1} \in \mathcal{A}$, is

$$\widehat{\pi}(a_1 \cdots a_m, a_{m+1}) = \frac{N(a_1 \cdots a_m a_{m+1})}{N(a_1 \cdots a_m \bullet)},$$

and in the model M$m$-3, we have $\forall a_1, \ldots, a_{m+1} \in \mathcal{A}$, $\forall k \in \{1, 2, 3\}$,

$$\widehat{\pi}_k(a_1 \cdots a_m, a_{m+1}) = \frac{N(a_1 \cdots a_m a_{m+1}, k)}{\displaystyle\sum_{b \in \mathcal{A}} N(a_1 \cdots a_m b, k)}.$$

## 2.3 Test for the appropriate order of the Markov model

To test which Markov model would be appropriate for a given sequence of length $n$, the most straightforward test is a Chi-square test, which can be viewed as a generalized likelihood ratio test. Most well-known is the Chi-square test for independence, see, e.g., Rice (1995). In general, suppose we have a sample of size $n$ cross-classified in a table with $U$ rows and $V$ columns. For instance, we could have four rows labeled A, C, G, T, and four columns labeled A, C, G, T, and we count how often a letter from the row is followed by a letter from the column in the sequence.

First we test whether we may assume the sequence to consist of independent letters. To this purpose, recall that $N(ab)$ denotes the count in cell $(a,b)$, whereas $N(a\bullet)$ is the $a$th row count, and let $N(\bullet b)$ is the $b$th column count. Thus $N(ab)$ counts how often letter $a$ is followed by letter $b$ in the sequence. Let

$\pi(a, b)$ be the probability of cell $(a, b)$, let $\pi(a, \bullet)$ be the $a$th row marginal probability, and let $\pi(\bullet, b)$ be the $b$th column marginal probability. We test the null hypothesis of independence

$$H_0 : \pi(a, b) = \pi(a, \bullet)\pi(\bullet, b)$$

against the alternative that the $\pi(a, b)$ are free. Under $H_0$ the maximum-likelihood estimate of $\pi(a, b)$ is

$$\hat{\pi}(a, b) = \hat{\pi}(a, \bullet)\hat{\pi}(\bullet, b) = \frac{N(a\bullet)}{n-1}\frac{N(\bullet b)}{n-1},$$

and under the alternative, the maximum-likelihood estimate of $\pi(a, b)$ is simply, as there are $n - 1$ consecutive pairs in the sequence,

$$\hat{\pi}(a, b) = \frac{N(ab)}{n-1}.$$

The Pearson chi-square statistic is the sum of the square difference between observed and expected count, divided by the expected counts, namely

$$X^2 = \sum_{a=1}^{U}\sum_{b=1}^{V} \frac{(N(ab) - N(a\bullet)N(\bullet b)/n - 1)^2}{N(a\bullet)N(\bullet b)/n - 1}.$$

(We think of the letters in the alphabet being enumerated here.) Under the null hypothesis, $X^2$ follows asymptotically a chi-square distribution with $(U - 1)(V - 1)$ degrees of freedom. Thus we would reject the null hypothesis when $X^2$ is too large, compared to the corresponding chi-square distribution. As a rule of thumb, this test is applicable when the expected frequency in each row and column is at least 5. Applying this test to DNA counts, we thus would have to compare $X^2$ to a chi-square distribution with $(4-1)(4-1) = 9$ degrees of freedom. A typical cutoff level would be 5%, or, if one likes to be conservative, 1%. The corresponding critical values are 16.92 for 5 %, and 21.67 for 1 %. Thus, if $X^2 > 16.92$, we would reject the null hypothesis of independence at the 5 % level (meaning that, if we repeated this experiment many times, in about 5% of the cases we would reject the null hypothesis when it is true). If $X^2 > 21.67$, we could reject the null-hypothesis at the 1 % level (so in only about 1 % of all trials would we reject the null hypothesis when it is true). Otherwise we would not reject the null hypothesis.

If the null hypothesis of independence cannot be rejected at an appropriate level (say, 5 %), then one would fit an independent model. However, if the null hypothesis does get rejected, one would test for a higher-order dependence. The next step would thus be to test for a first-order Markov chain. Here we would proceed as above, but now regarding how often a transition from $a$ to $b$ is followed by a transition from $b$ to $c$, where $a, b, c \in \mathcal{A}$. Thus, for a DNA sequence, where $\mathcal{A} = \{A, C, G, T\}$, we would have $U = 4^2$ rows, where we record the 2-letter words $ab$, and $V = 4$ columns, for the transitions to $c$. Let $N(ab, c)$ be the count in cell $(ab, c)$, the number of transitions from $ab$ to $c$ (the number of times that the 2-letter word $ab$ is followed by $c$ in the sequence). (Previously this is also denoted by $N(abc)$; however, here we would like to keep the separation between $ab$ and $c$ to refer to the cell sizes.) Let $\pi(ab, c)$ be the probability of cell $(ab, c)$ and let $\pi(a, b)$ and all other quantities be as above. We test the null hypothesis of a first-order Markov chain

$$H_0 : \pi(ab, c) = \pi(a, b)\pi(b, c)$$

against the alternative that the $\pi(ab, c)$ are free. Under $H_0$ the maximum-likelihood estimate of $\pi(ab, c)$ is

$$\hat{\pi}(ab, c) = \frac{N(ab)}{N(a\bullet)}\frac{N(bc)}{N(b\bullet)},$$

and under the alternative, the maximum-likelihood estimate of $\pi(ab, c)$ is

$$\hat{\pi}(ab, c) = \frac{N(ab, c)}{n-2}.$$

Enumerating the letters $\{A, C, G, T\}$ as 1, 2, 3, 4, the Pearson chi-square statistic is

$$X^2 = \sum_{a=1}^{4}\sum_{b=1}^{4}\sum_{c=1}^{4} \frac{\left(N(ab, c) - (n - 1)N(a, b)N(b, c)/(N(a, \bullet)N(b, \bullet))\right)^2}{(n - 1)N(a, b)N(b, c)/(N(a, \bullet)N(b, \bullet))}.$$

This quantity now has to be compared to a chi-square statistic with $(16 - 1)(4 - 1) = 45$ degrees of freedom. If this hypothesis is rejected, one would test for a higher-order Markov chain in the analogous way.

# 3 Overlapping and non-overlapping occurrences

Statistical inference is often based on independence assumptions. Even if the sequence letters are independent and identically distributed, the different random indicators of word occurrences are not independent due to overlaps. For example, if $w = \texttt{ATAT}$ occurs at position $i$ in the sequence, then another occurrence of $w$ is much likely to occur at position $i + 2$ than if $w$ did not occur at position $i$, and an occurrence of $w$ at position $i + 1$ is not possible. Many of the arguments needed for a probabilistic and statistical analysis of word occurrences deal with disentangling this overlapping structure.

Let $w = w_1 \cdots w_\ell$ be a word of length $\ell$ on a finite alphabet $\mathcal{A}$. Two occurrences of $w$ may overlap in a sequence if and only if $w$ is periodic, meaning that there exists $p \in \{1, \dots, \ell - 1\}$ such that $w_i = w_{i+p}$, $i = 1, \dots, \ell - p$. A word may have several periods. The set $\mathcal{P}(w)$ of the periods of $w$ is defined by

$$\mathcal{P}(w) := \{p \in \{1, \dots, \ell - 1\} : w_i = w_{i+p}, \forall i = 1, \dots, \ell - p\}.$$

A word $w$ is not periodic if and only if $\mathcal{P}(w)$ is empty. For instance, $\texttt{AACAA}$ is periodic and admits two periods, 3 and 4. There are 4 occurrences of $\texttt{AACAA}$ in the sequence $\texttt{TG}\underline{\texttt{AACAAACAACAA}}\texttt{TAG}\ \underline{\texttt{AACAAAA}}$, starting respectively at positions 3, 7, 10 and 18. The first 3 occurrences overlap and form a clump. A clump of $w$ in a sequence is a maximal set of overlapping occurrences of $w$ in the sequence. By definition two clumps of $w$ in a sequence cannot overlap. A clump composed of exactly $k$ overlapping occurrences of $w$ is called a $k$-clump of $w$. There are 2 clumps of $\texttt{AACAA}$ in the previous sequence, the first one is a 3-clump starting at position 3 and the second one is a 1-clump starting at position 18. Let $\mathcal{C}_k(w)$ be the set of the concatenated words composed of exactly $k$ overlapping occurrences of $w$. For example, $\mathcal{C}_1(\texttt{AACAA}) = \{\texttt{AACAA}\}$ and $\mathcal{C}_2(\texttt{AACAA}) = \{\texttt{AACAACAA}, \texttt{AACAAACAA}\}$.

Renewals are another type of non-overlapping occurrences of interest that require scanning the sequence from one end to the other: the first occurrence of $w$ in the sequence is a renewal and a given occurrence of $w$ is a renewal if and only if it does not overlap a previous renewal. Renewals of $w$ do not overlap in a sequence. In the above example, there are 3 renewals of $\texttt{AACAA}$ starting at position 3, 10 and 18.

Depending on the problem, one could be interested in studying the overlapping occurrences of $w$ in a sequence, or in restricting attention to non-overlapping occurrences: the beginnings of clumps, the beginnings of $k$-clumps or the renewals. We now introduce notation related to occurrences of a word $w = w_1 \cdots w_\ell$, of a clump of $w$, of a $k$-clump of $w$, of a renewal of $w$ in a sequence, and to the corresponding counts.

## 3.1 Occurrence and number of overlapping occurrences

An occurrence of $w$ starts at position $i$ in the sequence $\underline{X} = (X_i)_{i \in \mathbb{Z}}$ if and only if $X_i \cdots X_{i+\ell-1} = w_1 \cdots w_\ell$. Let $Y_i(w)$ be the associated random indicator

$$Y_i(w) := \mathbb{I}\{w \text{ starts at position } i \text{ in } \underline{X}\}. \tag{1}$$

For convenience in some sections, $Y_i(w)$ will be the random indicator that an occurrence of $w$ ends at position $i$ in $\underline{X}$; it will be made precise in that case.

In the stationary $m$-order Markovian model, the expectation of $Y_i(w)$, that represents the probability that an occurrence of $w$ occurs at a given position in the sequence, is denoted by $\mu_m(w)$ and is given by

$$\mu_m(w) = \mu(w_1 \cdots w_m)\pi(w_1 \cdots w_m, w_{m+1}) \cdots \pi(w_{\ell-m} \cdots w_{\ell-1}, w_\ell). \tag{2}$$

When there is no ambiguity, the index $m$ referring to the order of the model will be omitted.

The number of overlapping occurrences of $w$ in the sequence $(X_i)_{i=1,\dots,n}$, simply called count of $w$, is defined by $N(w) = N_n(w) = \sum_{i=1}^{n-\ell+1} Y_i(w)$ (or $N(w) = \sum_{i=\ell}^{n} Y_i(w)$ if $Y_i(w)$ is associated with an occurrence of $w$ ending at position $i$).

## 3.2 Clump and declumped counts

A clump of $w$ starts at position $i$ in the infinite sequence $\underline{X}$ if and only if there is an occurrence of $w$ starting at position $i$ that does not overlap a previous occurrence of $w$. It follows that

$$
\begin{aligned}
\widetilde{Y}_i(w) \quad &:= \quad \mathbb{I}\{\text{a clump of } w \text{ starts at position } i \text{ in } \underline{X}\} \\
&= \quad Y_i(w)(1 - Y_{i-1}(w)) \cdots (1 - Y_{i-\ell+1}(w)).
\end{aligned}
\tag{3}
$$

6

Often $\widetilde{Y}_i(w)$ is zero, depending on the overlapping structure of $w$. If we define $\mathcal{P}'(w)$ as the set of the *principal* periods of $w$, namely the periods that are not strictly multiples of the minimal period $p_0(w)$ of $w$, then it turns that

$$\widetilde{Y}_i(w) = Y_i(w) - \sum_{p \in \mathcal{P}'(w)} Y_{i-p}(w^{(p)}w), \tag{4}$$

where $w^{(p)}$ denotes the prefix of $w$ of length $p$, $w^{(p)} = w_1 \cdots w_p$, and $w^{(p)}w$ is the concatenated word $w_1 \cdots w_p w_1 \cdots w_\ell$. If $p \in \mathcal{P}(w)$ then $w^{(p)}$ is called a *root* of $w$; if $p \in \mathcal{P}'(w)$, $w^{(p)}$ is called a *principal root* of $w$. Equation (4) is obtained from the two following steps: (i) note that an occurrence of $w$ starting at position $i$ overlaps a previous occurrence of $w$ if and only if it is directly preceded by an occurrence of a principal root of $w$, meaning that a principal root $w^{(p)}$, $p \in \mathcal{P}'(w)$, occurs at position $i - p$, (ii) note that the events $E_p = \{Y_{i-p}(w^{(p)}) = 1\}$, $p \in \mathcal{P}'(w)$, are disjoint. To prove (ii), we assume that two different principal roots $w^{(p)}$ and $w^{(q)}$ occur simultaneously at position $i-p$ and $i-q$. If so, the minimal root $w^{(p_0)}$ of $w$ could be decomposed into $w^{(p_0)} = xy = yx$ where $x$ and $y$ are two nonempty words. Proposition 1.3.2 from Lothaire (1983) says that two words commute if and only if they are powers of the same word. Thus, we would obtain the contradiction that the minimal root is not minimal (see Schbath (1995a) for more details).

It follows from Equation (4) that the probability $\widetilde{\mu}(w)$ that a clump of $w$ starts at a given position in $\underline{X}$ is given by

$$\widetilde{\mu}(w) = \mu(w) - \sum_{p \in \mathcal{P}'(w)} \mu(w^{(p)}w). \tag{5}$$

The number $\widetilde{N}(w)$ of clumps of $w$ in the finite sequence $X_1 \cdots X_n$ (or the declumped count) may be different from the sum $\widetilde{N}_{\inf}(w) = \sum_{i=1}^{n-\ell+1} \widetilde{Y}_i(w)$ because of a possible clump of $w$ that would start in $\underline{X}$ before position 1 and would stop after position $\ell - 1$. The difference $\widetilde{N}(w) - \widetilde{N}_{\inf}(w)$ is either equal to 0 or equal to 1. In fact, it can be shown that $\mathbb{P}(\widetilde{N}(w) \neq \widetilde{N}_{\inf}(w)) \leq (\ell-1)(\mu(w) - \widetilde{\mu}(w))$ (see Reinert and Schbath (1998)).

## 3.3 $k$-clump and number of $k$-clumps

A $k$-clump of $w$ starts at position $i$ in $\underline{X}$ if and only if there is an occurrence of a concatenated word $c \in \mathcal{C}_k(w)$ starting at position $i$ that does not overlap any other occurrence of $w$ in the sequence $\underline{X}$. As we proceeded for a clump occurrence, an occurrence of $c \in \mathcal{C}_k(w)$ is a $k$-clump of $w$ in $\underline{X}$ if and only if it is not directly preceded by any principal root $w^{(p)}$ of $w$ and it is not directly followed by any suffix $w_{(q)} = w_{\ell-q+1} \cdots w_\ell$ with $q \in \mathcal{P}'(w)$. Some straightforward calculation yields the following expression (Schbath (1995a)):

$$\widetilde{Y}_{i,k}(w) \quad := \quad \mathbb{1}\{\text{a } k\text{-clump of } w \text{ starts at position } i \text{ in } \underline{X}\} \tag{6}$$

$$= \sum_{c \in \mathcal{C}_k(w)} \left( Y_i(c) - \sum_{p \in \mathcal{P}'(w)} Y_{i-p}(w^{(p)}c) - \sum_{q \in \mathcal{P}'(w)} Y_i(cw_{(q)}) + \sum_{p,q \in \mathcal{P}'(w)} Y_{i-p}(w^{(p)}cw_{(q)}) \right).$$

It follows that the probability for a $k$-clump to start at a given position is given by

$$\widetilde{\mu}_k(w) = \sum_{c \in \mathcal{C}_k(w)} \mu(c) - 2 \sum_{c' \in \mathcal{C}_{k+1}(w)} \mu(c') + \sum_{c'' \in \mathcal{C}_{k+2}(w)} \mu(c'').$$

This formula can be improved. Note that $\mathcal{C}_{k+1}(w) = \{w^{(p)}c, c \in \mathcal{C}_k(w), p \in \mathcal{P}'(w)\}$ and $\mu(w^{(p)}c) = \mu(c)\frac{\mu(w^{(p)}c)}{\mu(c)} = \mu(c)\frac{\mu(w^{(p)}w)}{\mu(w)}$. By denoting

$$A(w) = \sum_{p \in \mathcal{P}'(w)} \frac{\mu(w^{(p)}w)}{\mu(w)},$$

we have that $\sum_{c' \in \mathcal{C}_{k+1}(w)} \mu(c') = A(w) \sum_{c \in \mathcal{C}_k(w)} \mu(c)$, and it follows that

$$
\begin{aligned}
\widetilde{\mu}_k(w) &= (1 - A(w))^2 \sum_{c \in \mathcal{C}_k(w)} \mu(c) \\
&= (1 - A(w))^2 A(w) \sum_{c \in \mathcal{C}_{k-1}(w)} \mu(c) \\
&\vdots \\
&= (1 - A(w))^2 A(w)^{k-1} \mu(w). \quad (7)
\end{aligned}
$$

As for the declumped count, the number of $k$-clumps of $w$ in the finite sequence may be different from the sum $\widetilde{N}_{\inf}^{(k)}(w) = \sum_{i=1}^{n-\ell+1} \widetilde{Y}_{i,k}(w)$ because of possible end effects. This difference can be controlled in probability. Moreover, possible end effects may lead to a difference between the count $N(w)$ and $\sum_{k>0} k \widetilde{N}_{\inf}^{(k)}(w)$ but this can also be controlled (Reinert and Schbath (1998)).

### 3.4   Renewal and renewal count

A renewal of $w$ starts at position $i$ in $X_1 \cdots X_n$ if and only if there is an occurrence of $w$ starting at position $i$ that either is the first one or does not overlap a previous renewal of $w$. Let $\mathbb{I}_i(w)$ be the associated random indicator:

$$
\begin{aligned}
\mathbb{I}_i(w) &= \mathbb{1}\{\text{a renewal of } w \text{ starts at position } i \text{ in } X_1 \cdots X_n\} \\
&= Y_i(w) \prod_{j=i-\ell+1}^{i-1} (1 - \mathbb{I}_j(w)) \quad (8)
\end{aligned}
$$

with the convention that $\mathbb{I}_j(w) = 0$ if $j < 1$. Thus, for $i \leq \ell$, a renewal occurrence of $w$ at position $i$ is exactly a clump occurrence of $w$ at $i$ in the finite sequence. Renewal count makes then extensive use of the linear ordering in the sequence: it is defined by $R(w) = R_n(w) = \sum_{i=1}^{n-\ell+1} \mathbb{I}_i(w)$.

## 4   Word locations along a sequence

Here we are concerned with the length of the gaps between word occurrences. First we describe how to get the exact distribution, and then we give asymptotic results.

### 4.1   Exact distribution of the length between word occurrences

Let $w = w_1 \cdots w_\ell$ be a word of length $\ell$ on a finite alphabet $\mathcal{A}$. We assume that $X_1 \cdots X_n$ is a stationary first-order Markov chain on $\mathcal{A}$ with transition matrix $\Pi = (\pi(a,b))_{a,b \in \mathcal{A}}$ and stationary distribution $(\mu(a))_{a \in \mathcal{A}}$. Here we are interested in the statistical distribution of the distance $D$ between two successive occurrences of $w$ and, more precisely, in the probabilities $f(d) = \mathbb{P}(D = d) = \mathbb{P}(w \text{ occurs at } i+d \text{ and there is no occurrence of } w \text{ between } i+1 \text{ and } i+d-1 \mid w \text{ occurs at } i), d \geq 1$. In this section, we say that a word $w$ occurs at position $i$ if an occurrence of $w$ ends at position $i$; it happens with probability $\mu(w)$ given in (2).

The probability $f(d)$ can be obtained via a recursive formula (Robin and Daudin (1999)), as first proposed for independent and uniformly distributed letters by Blom and Thorburn (1982). It is clear that, if $1 \leq d \leq \ell - 1$ and $d \notin \mathcal{P}(w)$, then $f(d) = 0$. If $d \in \mathcal{P}(w)$ or if $d \geq \ell$, then we decompose the event

$$
E = \{w \text{ occurs at } i+d\}
$$

into the disjoint events

$$
E_1 = \{w \text{ occurs at } i+d \text{ and there is no occurrence of } w \text{ between } i+1 \text{ and } i+d-1\}
$$

and

$$
E_2 = \{w \text{ occurs at } i+d \text{ and there are some occurrences of } w \text{ between } i+1 \text{ and } i+d-1\}.
$$

Thus $\{E_1 \mid w \text{ at } i\}$ has probability $f(d)$. Moreover $E_2$ is itself decomposed as $E_2 = \cup_{h=1}^{d-1} E_2(h)$, where

$$
\begin{aligned}
E_2(h) \quad = \quad &\{\text{there is no occurrence of } w \text{ between } i+1 \text{ and } i+h-1, \\
&w \text{ occurs at } i+h \text{ and } i+d\}
\end{aligned}
$$

are again disjoint events.

If $1 \le d \le \ell-1$ and $d \in \mathcal{P}(w)$, then $\mathbb{P}(E \mid w \text{ at } i) = \mu(w)/\mu(w^{(\ell-d)})$. Moreover, if there are occurrences at positions $i+h$ and $i+d$, for some $h < d$, then the occurrences necessarily overlap, and this is only possible for $d - h \in \mathcal{P}(w)$; in this case, $\mathbb{P}(E_2(h) \mid w \text{ at } i) = f(h)\mu(w)/\mu(w^{(\ell-d+h)})$. Thus, we have

$$
\frac{\mu(w)}{\mu(w^{(\ell-d)})} = f(d) + \sum_{\substack{1 \le h \le d-1 \\ d-h \in \mathcal{P}(w)}} f(h)\frac{\mu(w)}{\mu(w^{(\ell-d+h)})}.
$$

If $d \ge \ell$, then $\mathbb{P}(E \mid w \text{ at } i) = \Pi^{d-\ell+1}(w_\ell, w_1)\mu(w)/\mu(w_1)$. If there is an occurrence at positions $i+h$ and $i+d$, for some $h < d$, then we distinguish two cases depending on the possible overlap between the occurrences at $i+h$ and $i+d$: if $d-\ell+1 \le h \le d-1$, they overlap and we use the previous calculation; if $1 \le h \le d-\ell$, they do not overlap and $\overline{\mathbb{P}}(E_2(h) \mid w \text{ at } i) = f(h)\Pi^{d-\ell-h+1}(w_\ell, w_1)\mu(w)/\mu(w_1)$. Thus, from

$$
\mathbb{P}(E \mid w \text{ at } i) = \mathbb{P}(E_1 \mid w \text{ at } i) + \sum_{h=1}^{d-1} \mathbb{P}(E_2(h) \mid w \text{ at } i)
$$

we get

$$
\begin{aligned}
\Pi^{d-\ell+1}(w_\ell, w_1)\frac{\mu(w)}{\mu(w_1)} \quad = \quad & f(d) + \sum_{1 \le h \le d-\ell} f(h)\Pi^{d-\ell-h+1}(w_\ell, w_1)\frac{\mu(w)}{\mu(w_1)} \\
& + \sum_{\substack{d-\ell+1 \le h \le d-1 \\ d-h \in \mathcal{P}(w)}} f(h)\frac{\mu(w)}{\mu(w^{(\ell-d+h)})}.
\end{aligned}
$$

This is the proof of the next theorem, see Robin and Daudin (1999).

**Theorem 1** *The distribution $f(d) = \mathbb{P}(D = d)$ of the distance $D$ between two successive occurrences of a word $w$ in a Markov chain is given by the following recursive formula:*

*If $1 \le d \le \ell-1$ and $d \notin \mathcal{P}(w)$, then $f(d) = 0$.*
*If $1 \le d \le \ell-1$ and $d \in \mathcal{P}(w)$,*

$$
f(d) = \frac{\mu(w)}{\mu(w^{(\ell-d)})} - \sum_{\substack{1 \le h \le d-1 \\ d-h \in \mathcal{P}(w)}} f(h)\frac{\mu(w)}{\mu(w^{(\ell-d+h)})}.
$$

*If $d \ge \ell$,*

$$
\begin{aligned}
f(d) = \Pi^{d-\ell+1}(w_\ell, w_1)\frac{\mu(w)}{\mu(w_1)} - \sum_{1 \le h \le d-\ell} f(h)\Pi^{d-\ell-h+1}(w_\ell, w_1)\frac{\mu(w)}{\mu(w_1)} \\
- \sum_{\substack{d-\ell+1 \le h \le d-1 \\ d-h \in \mathcal{P}(w)}} f(h)\frac{\mu(w)}{\mu(w^{(\ell-d+h)})}.
\end{aligned}
$$

Since $D$ is the distance between two successive occurrences of $w$, note that, even if $d \in \mathcal{P}(w)$, $f(d)$ can be null. For instance, by taking $w = \texttt{AAA}$, we have $\mathcal{P}(\texttt{AAA}) = \{1, 2\}$, and $f(1) = \mu(\texttt{AAA})/\mu(\texttt{AA}) = \pi(\texttt{A}, \texttt{A})$, $f(2) = \pi^2(\texttt{A}, \texttt{A}) - f(1)\pi(\texttt{A}, \texttt{A}) = 0$.

Note that the recurrence formula on $f(d)$ is not a "finite" recurrence since calculating $f(d)$ requires the calculation of $f(d-1)$, ..., $f(1)$, requiring substantial numerical calculations for large $d$. One can approach this computation problem by using the generating function defined by $\Phi_D(t) := \mathbb{E}(t^D) = \sum_{d \ge 1} f(d)t^d$. Indeed, the key argument is that the $\Phi_D(t)$ expression is a rational function (Theorem 2 and Remark 2; see Robin and Daudin (1999)) of the form $P(t)/Q(t)$ so the coefficient $f(d)$ of $t^d$ can be expressed with a recurrence formula whose order is the degree of the polynomial $Q(t)$ (see Section 9.3).

9

**Theorem 2** *The generating function of $D$ is*

$$\Phi_D(t) = 1 - \mu^{-1}(w) \left( \sum_{\substack{u=0 \\ u \in \mathcal{P}(W) \cup \{0\}}}^{\ell-1} \frac{t^u}{\mu(w^{(\ell-u)})} + \frac{1}{\mu(w_1)} \sum_{u \geq 1} \Pi^u(w_\ell, w_1) t^{\ell+u-1} \right)^{-1}.$$

**Remark 2** *If the transition matrix $\Pi$ is diagonalizable, there exists $\alpha_i$, $\beta_i \in \mathbb{C}$, $i = 2 \cdots |\mathcal{A}|$, such that*

$$\frac{1}{\mu(w_1)} \sum_{u \geq 1} \Pi^u(w_\ell, w_1) t^{\ell+u-1} = \frac{t^\ell}{1-t} \left( 1 + \frac{1-t}{\mu(w_1)} \sum_{i=2}^{|\mathcal{A}|} \frac{\alpha_i}{1 - t\beta_i} \right)$$

*implying that the above expression is a rational function with a pole equal to 1.*

**Remark 3** *Since $\Phi_D(t) = \sum_{d \geq 1} f(d) t^d$, we have the general following properties:*

$$\begin{aligned} \mathbb{E}(D) &= \Phi_D'(1) = \mu^{-1}(w) \\ Var(D) &= \Phi_D''(1) + \Phi_D'(1)(1 - \Phi_D'(1)). \end{aligned}$$

*Successive derivations of $\Phi_D(t)$ are obtained using the decomposition stated in the previous remark.*

**Proof.**
The proof of Theorem 2 is not complicated since one just has to develop the sum $\sum_{d \geq 0} f(d) t^d$ with $f(d)$ given by Theorem 1, but it is very technical. We thus only give the main lines of the calculation. By replacing $f(d)$ given by Theorem 1 in $\sum_{d \geq 0} f(d) t^d$, we obtain a sum of five term

$$\Phi_D(t) = K_1 - K_2 + K_3 - K_4 - K_5$$

with

$$\begin{aligned} K_1 &= \sum_{\substack{d=1 \\ d \in \mathcal{P}(W)}}^{\ell-1} \frac{\mu(w)}{\mu(w^{(\ell-d)})} t^d \\[2mm] K_2 &= \sum_{\substack{d=1 \\ d \in \mathcal{P}(W)}}^{\ell-1} \sum_{\substack{h=1 \\ d-h \in \mathcal{P}(W)}}^{d-1} f(h) \frac{\mu(w)}{\mu(w^{(\ell-d+h)})} t^d \\[2mm] &= \sum_{h=1}^{\ell-2} f(h) \sum_{\substack{d=h+1 \\ d-h \in \mathcal{P}(W)}} \frac{\mu(w)}{\mu(w^{(\ell-d+h)})} t^d \\[2mm] &= \sum_{h=1}^{\ell-2} f(h) \sum_{\substack{u=1 \\ u \in \mathcal{P}(W)}} \frac{\mu(w)}{\mu(w^{(\ell-u)})} t^{h+u} \\[2mm] K_3 &= \sum_{d \geq \ell} \Pi^{d-\ell+1}(w_\ell, w_1) \frac{\mu(w)}{\mu(w_1)} t^d \\[2mm] &= \frac{\mu(w)}{\mu(w_1)} t^{\ell-1} \sum_{u \geq 1} \Pi^u(w_\ell, w_1) t^u \end{aligned}$$

10

$$K_4 = \sum_{d \geq \ell} \sum_{h=1}^{d-\ell} f(h) \Pi^{d-\ell-h+1}(w_\ell, w_1) \frac{\mu(w)}{\mu(w_1)} t^d$$

$$= \frac{\mu(w)}{\mu(w_1)} t^{\ell-1} \sum_{z \geq 1} \sum_{h=1}^{z} f(h) \Pi^{z-h+1}(w_\ell, w_1) t^{z+1}$$

$$= \frac{\mu(w)}{\mu(w_1)} t^{\ell-1} \sum_{h \geq 1} f(h) t^h \sum_{z \geq h} \Pi^{z-h+1}(w_\ell, w_1) t^{z-h+1}$$

$$= \frac{\mu(w)}{\mu(w_1)} t^{\ell-1} \Phi_D(t) \sum_{u \geq 1} \Pi^u(w_\ell, w_1) t^u$$

and

$$K_5 = \sum_{d \geq \ell} \sum_{\substack{h=d-\ell+1 \\ d-h \in \mathcal{P}(W)}}^{d-1} f(h) \frac{\mu(w)}{\mu(w^{(\ell-d+h)})} t^d$$

$$= \sum_{z \geq 1} \sum_{\substack{h=z \\ z+\ell-h-1 \in \mathcal{P}(W)}}^{z+\ell-2} f(h) \frac{\mu(w)}{\mu(w^{(h-z+1)})} t^{z+\ell-1}$$

$$= \sum_{h=1}^{\ell-1} f(h) \sum_{\substack{z=1 \\ z+\ell-h-1 \in \mathcal{P}(W)}}^{h} \frac{\mu(w)}{\mu(w^{(h-z+1)})} t^{z+\ell-1}$$

$$+ \sum_{h \geq \ell} f(h) t^h \sum_{\substack{z=h-\ell+2 \\ z+\ell-h-1 \in \mathcal{P}(W)}}^{h} \frac{\mu(w)}{\mu(w^{(h-z+1)})} t^{z-h+\ell-1}$$

$$= \sum_{h=1}^{\ell-1} f(h) \sum_{\substack{u=\ell-h \\ u \in \mathcal{P}(W)}}^{\ell-1} \frac{\mu(w)}{\mu(w^{(\ell-u)})} t^{h+u} + \sum_{h \geq \ell} f(h) t^h \sum_{\substack{u=1 \\ u \in \mathcal{P}(W)}}^{\ell-1} \frac{\mu(w)}{\mu(w^{(\ell-u)})} t^u.$$

Grouping $K_1 - K_2 - K_5$ and $K_3 - K_4$ leads to

$$\Phi_D(t) = (1 - \Phi_D(t)) \left( \sum_{\substack{u=1 \\ u \in \mathcal{P}(W)}}^{\ell-1} \frac{\mu(w)}{\mu(w^{(\ell-u)})} t^u + \frac{\mu(w)}{\mu(w_1)} t^{\ell-1} \sum_{u \geq 1} \Pi^u(w_\ell, w_1) t^u \right),$$

hence

$$\Phi_D(t) = 1 - \left( 1 + \sum_{\substack{u=1 \\ u \in \mathcal{P}(W)}}^{\ell-1} \frac{\mu(w)}{\mu(w^{(\ell-u)})} t^u + \frac{\mu(w)}{\mu(w_1)} t^{\ell-1} \sum_{u \geq 1} \Pi^u(w_\ell, w_1) t^u \right)^{-1}.$$

Using $\mu(w)/\mu(w^{(\ell)}) = 1$ establishes the theorem.
$\square$


The distance $D$ between two successive occurrences of $w$ can be seen as the distance between the $j$-th and $(j+1)$-th occurrence of $w$ in the sequence, since we use an homogeneous model. It may be useful to study the distance $D^{(r)}$ between the $j$-th and $(j+r)$-th occurrence of $w$, called $r$-scan by Karlin and colleagues (e.g. Dembo and Karlin (1992)). The distance $D^{(r)}$ is the sum of $r$ independent and identically distributed random variables with the same distribution as $D$. So we have

$$\Phi_{D^{(r)}}(t) = \left( \Phi_D(t) \right)^r.$$

We get the exact distribution of $D^{(r)}$ from the Taylor expansion of $\Phi_{D^{(r)}}(t)$: the probability $\mathbb{P}(D^{(r)} = d)$ is the coefficient of $t^d$ in the series.

## 4.2 Asymptotic distribution of $r$-scans

In the preceding paragraph, we presented how to get the exact distribution of an $r$-scan $D^{(r)}$, the distance between a word occurrence and the $(r-1)$-th next one, in a stationary Markov chain. When analyzing a biological sequence, assume we observe $(h+1)$ occurrences of a given motif, so that we observe $h$ distances $D_1, \ldots, D_h$ between occurrences of the motif. Thus we observe $(h-r+1)$ so-called $r$-scans $D_i^{(r)} = \sum_{j=i}^{i+r-1} D_j$. To detect poor and rich regions with this motif, one is interested in studying the significance of the smallest and the largest $r$-scans, or more generally in the $k$th smallest $r$-scan, denoted by $m_k$, and the $k$th largest $r$-scan, denoted by $M_k$. In this section, we present a Poisson approximation for the statistical distribution of the extreme value $m_k$ obtained by Dembo and Karlin (1992) using the Chen-Stein method. A similar result also exists for $M_k$ by following an identical setup, so it will not be explained in detail here.

We start defining the Bernoulli variables that will be used in the Chen-Stein method (see Section 9.1):

$$W_i^-(d) := \mathbb{1}\{D_i^{(r)} \le d\}, \ d \ge 0.$$

Denote by

$$W^-(d) = \sum_{i=1}^{h-r+1} W_i^-(d)$$

the number of $r$-scans less or equal to $d$. Note the duality principle

$$\{W^-(d) < k\} = \{m_k > d\}, \ d \ge 0.$$

We now use Theorem 13 to get a Poisson approximation for the distribution of $W^-(d)$. To apply this theorem, we first need to choose a neighborhood of dependence for each indicator variable; ideally the indicator variables with index not from the neighborhood of dependence are independent of that indicator variable. Secondly there are three quantities to bound, called $b_1$, $b_2$, and $b_3$, given in (28), (29), and (30). Piecing this together gives a bound on the total variation distance between the distributions. Here we proceed as follows.

For $i \in \{1, \ldots, h-r+1\}$, we choose the neighborhood $B_i = \{j \mid |i-j| < r\}$. Let $Z_{\lambda^-}$ be the Poisson variable with expectation $\lambda^-$, where

$$\begin{aligned} \lambda^- &= \mathbb{E}W^-(d) \\ &= (h-r+1)\mathbb{E}W_i^-(d) \\ &= (h-r+1)\mathbb{P}(D^{(r)} \le d). \end{aligned}$$

Theorem 13 gives that

$$\begin{aligned} d_{\mathrm{TV}}\left(\mathcal{L}(W^-(d)), \mathcal{L}(Z_{\lambda^-})\right) &\le \frac{1-\mathrm{e}^{-\lambda^-}}{\lambda^-}\left(\sum_{i=1}^{h-r+1} \sum_{j \in B_i} \mathbb{E}W_i^-(d)\mathbb{E}W_j^-(d)\right. \\ &\left. + \sum_{i=1}^{h-r+1} \sum_{j \in B_i \setminus \{i\}} \mathbb{E}(W_i^-(d)W_j^-(d))\right). \end{aligned}$$

Here, $d_{\mathrm{TV}}$ denotes the total variation distance; see Section 9.1.

Indeed, the neighborhood $B_i$ has been chosen so that $W_i^-(d)$ is independent of $W_j^-(d)$, $\forall j \notin B_i$, leading to $b_3 = 0$. For $j > i$, we have

$$\begin{aligned} \mathbb{E}(W_i^-(d)W_j^-(d)) &= \mathbb{P}(D_i^{(r)} \le d, D_j^{(r)} \le d) \\ &= \mathbb{P}(D_j^{(r)} \le d \mid D_i^{(r)} \le d)\mathbb{P}(D_i^{(r)} \le d) \\ &= \mathbb{P}(D_{j-i+1}^{(r)} \le d \mid D_1^{(r)} \le d)\mathbb{P}(D^{(r)} \le d). \end{aligned}$$

Therefore,

$$\sum_{i=1}^{h-r+1} \sum_{j \in B_i \setminus \{i\}} \quad \mathbb{E}(W_i^-(d)W_j^-(d))$$

$$\leq 2(h-r+1)\mathbb{P}(D^{(r)} \leq d) \sum_{s=2}^{r} \mathbb{P}(D_s^{(r)} \leq d \mid D_1^{(r)} \leq d)$$

$$\leq 2\lambda^- \sum_{s=2}^{r} \mathbb{P}(D_s^{(r)} \leq d \mid D_1^{(r)} \leq d).$$

It can be shown that

$$\mathbb{P}(D_s^{(r)} \leq d \mid D_1^{(r)} \leq d) \leq \mathbb{P}\left( \sum_{i=r+1}^{s+r-1} D_i \leq d \right) = \mathbb{P}(D^{(s-1)} \leq d).$$

We finally get

$$d_{\mathrm{TV}}\left( \mathcal{L}(W^-(d)), \mathcal{L}(Z_{\lambda^-}) \right) \quad \leq \quad \left( (2r-1)\mathbb{P}(D^{(r)} \leq d) + 2\sum_{s=1}^{r-1} \mathbb{P}(D^{(s)} \leq d) \right)(1 - \mathrm{e}^{-\lambda^-}).$$

From the duality principle,

$$|\mathbb{P}(m_k > d) - \mathbb{P}(Z_{\lambda^-} < k)| \quad \leq \quad \left( (2r-1)\mathbb{P}(D^{(r)} \leq d) + 2\sum_{s=1}^{r-1} \mathbb{P}(D^{(s)} \leq d) \right)(1 - \mathrm{e}^{-\lambda^-}).$$

This approximation is very useful for the comparison between the expected distribution of the $r$-scans and the one observed in the biological sequence. It has been applied in Karlin and Macken (1991) to the *E. coli* genome by approximating the $r$-scan distribution given in Section 4.1 by a sum of $r - 1$ independent exponential random variables.

# 5    Word count distribution

Let again $w = w_1 \cdots w_\ell$ be a word of length $\ell$ on a finite alphabet $\mathcal{A}$ and $\underline{X} = (X_i)_{i \in \mathbb{Z}}$ be a random sequence on $\mathcal{A}$. This section is devoted to the statistical distribution of the count $N(w)$ of $w$ in the sequence $X_1 \cdots X_n$. First we state how to compute the exact distribution in the model M1, using recursion techniques. For long sequences, however, asymptotic results are obtainable, and, in general, easier to handle. Here the appropriate asymptotic regime depends crucially on the length $\ell$ of the target word relative to the sequence length $n$. For very short words, the law of large numbers can be applied to approximate the word count by the expected word count. This being a very crude estimate, one can easily improve on it by employing the Central Limit Theorem, stating that the word count distribution is asymptotically normal. This approximation will be satisfactory when the words are not too long. For rare words, as a rule of thumb words of length $\ell \asymp \log n$, a compound Poisson approximation will give better results. For the latter, the error made in the approximation can be bounded in terms of the sequence length, the word length, and word probabilities, so that it is possible to assess when a compound Poisson approximation will be a good choice. Moreover, the error bound can be incorporated to give conservative confidence intervals, as it will be explained below.

## 5.1    Exact distribution

If $\underline{X}$ is a stationary first-order Markov chain, the exact distribution of the count $N(w)$ can be easily obtained using the distribution of the successive positions $(T_j)_{j \geq 1}$ of the $j$-th occurrence of $w$ in $X_1 \cdots X_n$. Indeed, we have the following duality principle:

$$\{N(w) \geq j\} = \{T_j \leq n\}.$$

The exact distribution of $T_j$ can be obtained as in Section 4.1, by deriving the Taylor expansion of the generating function $\Phi_{T_j}(t)$ of $T_j$. If $j = 1$, the generating function $\Phi_{T_1}(t)$ can be obtained as $\Phi_D(t)$ (see Theorem 2). We just state the result of Robin and Daudin (1999):

$$\Phi_{T_1}(t) = \frac{t^\ell}{1-t}\left(\sum_{\substack{u=0 \\ u \in \mathcal{P}(W) \cup \{0\}}}^{\ell-1} \frac{t^u}{\mu(w^{(\ell-u)})} + \frac{1}{\mu(w_1)}\sum_{u \geq 1} \Pi^u(w_\ell, w_1)t^{\ell+u-1}\right)^{-1}.$$

Now, $T_j - T_1$ is a sum of $j - 1$ independent and identically distributed random variables with the same distribution as $D$. So we have $\Phi_{T_j}(t) = \Phi_{T_1}(t)\left(\Phi_D(t)\right)^{j-1}$. The coefficient of $t^a$ in the Taylor expansion of $\Phi_{T_j}(t)$ is then equal to $\mathbb{P}(T_j = a) = g_j(a)$. Using the duality principle, we get

$$\mathbb{P}(N(w) = j) = \sum_{a=\ell}^{n}(g_j(a) - g_{j+1}(a)).$$

This generalizes the exact result that Gentleman and Mullin (1989) obtained for the case that the sequence is composed of i.i.d. letters, where each letter occurs with equal probability. In this case, Gentleman (1994) also gives an algorithm for calculating the word frequency distribution. Moreover, in the Markov case the exact distribution of the count can also be obtained by other techniques: Kleffe and Langbecker (1990) used an automaton built on the pattern structure matrix, whereas Régnier (1998) and Régnier (1999) used a language decomposition approach with combinatorial methods.

## 5.2 The weak law of large numbers

As a crude first approximation, the weak law of large numbers states that the observed counts will converge towards the expected counts. Indeed we may use Chebyshev's inequality to bound the expected deviation of the observed counts from the expected number of occurrences. This approximation is valid only for relatively short words, and in this case a normal approximation gives more information. Such an approximation will be derived in the following subsection.

## 5.3 Asymptotic distribution: the Gaussian regime

We assume that $\underline{X} = (X_i)_{i \in \mathbb{Z}}$ is a stationary $m$-order Markov chain on $\mathcal{A}$, $0 \leq m \leq \ell - 2$, with transition probabilities $\pi(a_1 \cdots a_m, a_{m+1})$ and stationary distribution $\mu(a_1 \cdots a_m)$, $a_1, \ldots, a_{m+1} \in \mathcal{A}$. In this subsection, $N(w) = \sum_{i=\ell}^{n} Y_i(w)$ and

$$Y_i = Y_i(w) = \mathbb{1}\{w \text{ ends at position } i \text{ in } \underline{X}\}.$$

If the model is known, the asymptotic normality of $(N(w) - \mathbb{E}N(w))/\sqrt{n}$ directly follows from a Central Limit Theorem for Markov chains and the variance of $N(w)$ is given by Kleffe and Borodovsky (1992). When $m = 1$, the expectation and variance of $N(w)$ are

$$\begin{aligned}
\mathbb{E}N(w) &= (n - \ell + 1)\mu_1(w) \\
\text{Var}N(w) &= \mathbb{E}N(w) + 2\sum_{p \in \mathcal{P}(w)} \mathbb{E}N(w^{(p)}w) - \mathbb{E}^2 N(w) \\
&\quad + \frac{2}{\mu(w_1)}\mu_1^2(w)\sum_{d=1}^{n-2\ell+1}(n - 2\ell + 2 - d)\Pi^d(w_\ell, w_1)
\end{aligned}$$

where $\mu_1(w)$ is given in Eq. (2).

In the problem of finding exceptional words in biological sequences, the model is unknown and its parameters are estimated from the observed sequence. The expected mean of $N(w)$ is not available and is approximated with an estimator $\widehat{N}_m(w)$. In this paragraph, we get both the asymptotic normality of $(N(w) - \widehat{N}_m(w))/\sqrt{n}$ and the asymptotic variance. This is not a trivial problem since the estimation changes fundamentally the variance expression.

The expected mean of $N(w)$ is given by $\mathbb{E}N(w) = (n - \ell + 1)\mu(w)$ where $\mu(w) = \mu_m(w)$ is the probability that an occurrence of $w$ ends at a given position in the sequence (see Eq. (2)). Estimating each parameter by its maximum likelihood estimator gives an estimator $\widehat{N}_m(w)$ of $\mathbb{E}N(w)$:

$$\widehat{N}_m(w) = \frac{N(w_1 \cdots w_{m+1}) \cdots N(w_{\ell-m} \cdots w_\ell)}{N(w_2 \cdots w_{m+1}) \cdots N(w_{\ell-m} \cdots w_{\ell-1})}. \tag{9}$$

Let us first consider the maximal model ($m = \ell - 2$) that is mainly used to find exceptional words (Brendel *et al.* (1986), Leung *et al.* (1996), Rocha *et al.* (1998)). We introduce the following notation: $w^- = w_1 \cdots w_{\ell-1}$ (prefix of $w$ with length $\ell - 1$), $^-w = w_2 \cdots w_\ell$ (suffix of $w$ with length $\ell - 1$) and $^-w^- = w_2 \cdots w_{\ell-1}$. Under the maximal model, the estimator of $N(w)$ is

$$\widehat{N}_{\ell-2}(w) = \frac{N(w_1 \cdots w_{\ell-1})N(w_2 \cdots w_\ell)}{N(w_2 \cdots w_{\ell-1})} = \frac{N(w^-)N(^-w)}{N(^-w^-)};$$

moreover, the asymptotic normality of $(N(w) - \widehat{N}_{\ell-2}(w))/\sqrt{n}$ and the asymptotic variance can be obtained in an elegant way using martingale techniques. (For an introduction to martingales, see, e.g., Chung (1974)). Indeed, $\widehat{N}_{\ell-2}(w)$ is a natural estimator of $N(w^-)\pi(^-w^-, w_\ell)$, and $N(w) - N(w^-)\pi(^-w^-, w_\ell)$ is approximately a martingale as it is shown below.

We introduce the martingale $M_n = \sum_{i=\ell}^n (Y_i - \mathbb{E}(Y_i \mid \mathcal{F}_{i-1}))$ with $\mathcal{F}_i = \sigma(X_1, \ldots, X_i)$; it is easy to verify that $\mathbb{E}(M_n \mid \mathcal{F}_{n-1}) = M_{n-1}$. Moreover, we have

$$\begin{aligned}
\mathbb{E}(Y_i \mid \mathcal{F}_{i-1}) &= \mathbb{P}(w^- \text{ ends at } i-1 \text{ and } w_\ell \text{ occurs at } i \mid \mathcal{F}_{i-1}) \\
&= \mathbb{I}\{w^- \text{ ends at } i-1\}\pi(^-w^-, w_\ell),
\end{aligned}$$

and

$$\sum_{i=\ell}^n \mathbb{E}(Y_i \mid \mathcal{F}_{i-1}) = \big(N(w^-) - \mathbb{I}\{w^- \text{ ends at } n\}\big)\pi(^-w^-, w_\ell).$$

Therefore,

$$\frac{1}{\sqrt{n}}M_n = \frac{1}{\sqrt{n}}\big(N(w) - N(w^-)\pi(^-w^-, w_\ell)\big) - \frac{1}{\sqrt{n}}\mathbb{I}\{w^- \text{ ends at } n\}\pi(^-w^-, w_\ell). \tag{10}$$

Note that $n^{-1/2}\mathbb{I}\{w^- \text{ ends at } n\}\pi(^-w^-, w_\ell)$ tends to zero as $n \to \infty$. The next proposition establishes the asymptotic normality of $M_n/\sqrt{n}$.

**Proposition 1** *Let $V = \mu(w^-)\pi(^-w^-, w_\ell)(1 - \pi(^-w^-, w_\ell))$. We have*

$$\frac{1}{\sqrt{n}}M_n \xrightarrow{\mathcal{D}} \mathcal{N}(0, V) \text{ as } n \to \infty.$$

**Proof.**

This is an application of Theorem 17 for the one-dimensional random variable $\xi_{n,i} = n^{-1/2}(Y_i - \mathbb{E}(Y_i \mid \mathcal{F}_{i-1}))$. Three conditions have to be satisfied. Condition (i) holds from $\mathbb{E}(\xi_{n,i} \mid \mathcal{F}_{i-1}) = 0$. We then have to check that $\sum_{i=\ell}^n \text{Var}(\xi_{n,i} \mid \mathcal{F}_{i-1})$ converges to $V$ as $n \to \infty$. Since $Y_i$ is a 0-1 random variable, we have

$$\begin{aligned}
\text{Var}(Y_i \mid \mathcal{F}_{i-1}) &= \mathbb{E}(Y_i \mid \mathcal{F}_{i-1}) - \big(\mathbb{E}(Y_i \mid \mathcal{F}_{i-1})\big)^2 \\
&= \mathbb{I}\{w^- \text{ ends at } i-1\}\pi(^-w^-, w_\ell)\big(1 - \pi(^-w^-, w_\ell)\big).
\end{aligned}$$

We thus obtain

$$\begin{aligned}
\sum_{i=\ell}^n \text{Var}(\xi_{n,i} \mid \mathcal{F}_{i-1}) &= \frac{1}{n}\sum_{i=\ell}^n \text{Var}(Y_i \mid \mathcal{F}_{i-1}) \\
&= \frac{1}{n}N(w^-)\pi(^-w^-, w_\ell)(1 - \pi(^-w^-, w_\ell)) \\
&\quad - \frac{1}{n}\mathbb{I}\{w^- \text{ ends at } i-1\}\pi(^-w^-, w_\ell)(1 - \pi(^-w^-, w_\ell)) \\
&\longrightarrow V \text{ as } n \to \infty;
\end{aligned}$$

15

the convergence follows from the Law of Large Numbers: $N(w^-)/n \to \mu(w^-)$.

Finally, $|\xi_{n,i}| \leq \frac{2}{\sqrt{n}}$, so that $\forall \varepsilon > 0$, $\forall n > 4/\varepsilon^2$, $\mathbb{P}(|\xi_{n,i}| > \varepsilon) = 0$, establishing condition (iii). Using Theorem 17 proves the proposition. $\square$

Proposition 1 and Equation (10) also yield that

$$\frac{1}{\sqrt{n}}\big(N(w) - N(w^-)\pi(^-w^-, w_\ell)\big) \xrightarrow{\mathcal{D}} \mathcal{N}(0, V) \text{ as } n \to \infty.$$

We initially wanted to prove such convergence for

$$T_n = \frac{1}{\sqrt{n}}\big(N(w) - N(w^-)\widehat{\pi}(^-w^-, w_\ell)\big),$$

where

$$\widehat{\pi}(^-w^-, w_\ell) = \frac{N(^-w)}{N(^-w^-)}.$$

To this purpose, we decompose $T_n$ as follows:

$$
\begin{aligned}
T_n &= \frac{1}{\sqrt{n}}\big(N(w) - N(w^-)\pi(^-w^-, w_\ell)\big) - \frac{1}{\sqrt{n}}N(w^-)\big(\widehat{\pi}(^-w^-, w_\ell) - \pi(^-w^-, w_\ell)\big) \\
&= \frac{1}{\sqrt{n}}\big(N(w) - N(w^-)\pi(^-w^-, w_\ell)\big) - \frac{1}{\sqrt{n}}\frac{N(w^-)}{N(^-w^-)}\big(N(^-w) - N(^-w^-)\pi(^-w^-, w_\ell)\big) \\
&= \frac{1}{\sqrt{n}}M_n - \frac{1}{\sqrt{n}}\frac{N(w^-)}{N(^-w^-)}M'_n + \mathrm{o}(1),
\end{aligned}
\tag{11}
$$

where $M'_n$ is the martingale $M'_n = \sum_{i=\ell}^n \big(Y_i(^-w) - \mathbb{E}(Y_i(^-w) \mid \mathcal{F}_{i-1})\big)$. Now, using Theorem 17 gives

$$\frac{1}{\sqrt{n}}\begin{pmatrix} M_n \\ M'_n \end{pmatrix} \longrightarrow \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}; \begin{pmatrix} V & V_{12} \\ V_{21} & V_{22} \end{pmatrix}\right) \tag{12}$$

with

$$V_{21} = V_{12} = \lim_{n\to\infty} \frac{1}{n}\sum_{i=\ell}^n \mathbb{E}\left(\big(Y_i - \mathbb{E}(Y_i \mid \mathcal{F}_{i-1})\big)\big(Y_i(^-w) - \mathbb{E}(Y_i(^-w) \mid \mathcal{F}_{i-1})\big)\right)$$

and

$$V_{22} = \lim_{n\to\infty} \frac{1}{n}\sum_{i=\ell}^n \mathrm{Var}(Y_i(^-w) \mid \mathcal{F}_{i-1}).$$

With the same technique as for the derivation of $V$, as $Y_i Y_i(^-w) = Y_i$, we get $V_{21} = V_{12} = V$ and $V_{22} = \mu(^-w^-)\pi(^-w^-, w_\ell)(1 - \pi(^-w^-, w_\ell))$. Note that the Law of Large Numbers guarantees that, almost surely,

$$\frac{N(w^-)}{N(^-w^-)} \to \frac{\mu(w^-)}{\mu(^-w^-)} \text{ as } n \to \infty. \tag{13}$$

From (11)–(13), we are now able to deduce that $T_n$ converges in distribution to $\mathcal{N}(0, \sigma_{\ell-2}^2(w))$ with

$$
\begin{aligned}
\sigma_{\ell-2}^2(w) &= V_{11} - 2\frac{\mu(w^-)}{\mu(^-w^-)}V_{12} + \left(\frac{\mu(w^-)}{\mu(^-w^-)}\right)^2 V_{22} \\
&= \mu(w^-)\left(1 - \frac{\mu(w^-)}{\mu(^-w^-)}\right)\pi(^-w^-, w_\ell)(1 - \pi(^-w^-, w_\ell)) \\
&= \frac{\mu(w)}{\mu(^-w^-)}\big(\mu(^-w^-) - \mu(w^-)\big)(1 - \pi(^-w^-, w_\ell)) \\
&= \frac{\mu(w)}{\mu(^-w^-)}\big(\mu(^-w^-) - \mu(w^-) - \mu(^-w) + \mu(w)\big) \\
&= \frac{\mu(w)}{\mu(^-w^-)^2}\big(\mu(^-w^-) - \mu(^-w)\big)\big(\mu(^-w^-) - \mu(w^-)\big).
\end{aligned}
$$

We have just proved the following theorem.

**Theorem 3** *As $n \to \infty$, we have*

$$\frac{1}{\sqrt{n}} \left( N(w) - \widehat{N}_{\ell-2}(w) \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma_{\ell-2}^2(w))$$

*with*

$$\sigma_{\ell-2}^2(w) = \frac{\mu(w)}{\mu(^-w^-)^2} (\mu(^-w^-) - \mu(^-w))(\mu(^-w^-) - \mu(w^-)),$$

*and*

$$\frac{N(w) - \widehat{N}_{\ell-2}(w)}{\sqrt{n\widehat{\sigma}_{\ell-2}^2(w)}} \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1),$$

*where $n\widehat{\sigma}_{\ell-2}^2(w)$ is the plug-in estimator of $n\sigma_{\ell-2}^2(w)$:*

$$n\widehat{\sigma}_{\ell-2}^2(w) = \frac{\widehat{N}_{\ell-2}(w)}{N(^-w^-)^2} \left( N(^-w^-) - N(^-w) \right) \left( N(^-w^-) - N(w^-) \right).$$

In the nonmaximal models ($m < \ell - 2$), it is straightforward to extend the previous martingale approach to prove the asymptotic normality of $(N(w) - \widehat{N}_m(w))/\sqrt{n}$ and to derive the asymptotic variance. Indeed, for each value of $\ell - m$, the difference $N(w) - \widehat{N}_m(w)$ has to be decomposed as a linear combination of martingales, exactly as it was done for $T_n$. For instance, if $w = abcde$ and $m = 1$,

$$
\begin{aligned}
N(abcde) - \widehat{N}_1(abcde) &= N(abcde) - \frac{N(ab)N(bc)N(cd)N(de)}{N(b)N(c)N(d)} \\
&= N(abcde) - N(abcd)\frac{N(de)}{N(d)} + \frac{N(de)}{N(d)}\left(N(abcd) - N(abc)\frac{N(cd)}{N(c)}\right) \\
&\quad + \frac{N(de)N(cd)}{N(d)N(c)}\left(N(abc) - N(ab)\frac{N(bc)}{N(b)}\right).
\end{aligned}
$$

Another approach consists of using the $\delta$-method as proposed in Lundstrom (1990). The idea is to consider $N(w) - \widehat{N}_m(w)$ as $f(\underline{N})$ where $\underline{N}$ is the count vector $\underline{N} = (N(w), N(w_1 \cdots w_{m+1}), \dots, N(w_{\ell-m} \cdots w_\ell), N(w_2 \cdots w_{m+1}), \dots, N(w_{\ell-m} \cdots w_{\ell-1}))$ (see Equation (9)). There exists a covariance matrix $\Sigma$ such that

$$\frac{1}{\sqrt{n}}(\underline{N} - \mathbb{E}\underline{N}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma)$$

(see Waterman (1995) for an exposition). The next step is to use the $\delta$-method (Theorem 15) to transfer this convergence to $f(\underline{N})$:

$$\frac{1}{\sqrt{n}}(f(\underline{N}) - f(\mathbb{E}\underline{N})) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \nabla\Sigma\nabla^t),$$

where $\nabla = \left( \frac{\partial f(x_1, \dots, x_{2(\ell-m)})}{\partial x_j} \big|_{\mathbb{E}\underline{N}} \right)_{j=1, \dots, 2(\ell-m)}$ is the partial derivative vector of $f$. Since $f(\mathbb{E}\underline{N}) = 0$, we finally get

$$\frac{1}{\sqrt{n}}\left(N(w) - \widehat{N}_m(w)\right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \nabla\Sigma\nabla^t).$$

However, this method does not easily provide an explicit formula for the asymptotic variance since the function $f$ and its derivative depends on $\ell - m$.

The conditional approach of Prum *et al.* (1995) provides an alternative to the problem. Initially presented for the model M1, this method has been generalized to the general model M$m$, $0 \le m \le \ell - 2$ (Schbath (1995b)). The principle is to work conditionally on the sufficient statistic $\mathcal{S}_m$ of the model M$m$, namely the collection of counts $\{N(a_1 \cdots a_{m+1}), a_1, \dots, a_{m+1} \in \mathcal{A}\}$ and the first $m$ letters of the sequence. Using the technique developed by Cowan (1991), one can derive both the conditional expectation $\mathbb{E}(N(w) \mid \mathcal{S}_m)$ and the conditional variance of $N(w)$. The key arguments are first that the

conditional expectation is asymptotically equivalent to $\widehat{N}_m(w)$, leading to the asymptotic normality of $(N(w) - \mathbb{E}(N(w) \mid \mathcal{S}_m))/\sqrt{n}$, and second, that $n^{-1}\mathrm{Var}(N(w) \mid \mathcal{S}_m)$ has the limiting value $\sigma_m^2(w)$ with

$$
\begin{aligned}
\sigma_m^2(w) &= \mu(w) + 2 \sum_{p \in \mathcal{P}(w),\, p \leq \ell - m - 1} \mu(w^{(p)}w) + \mu(w)^2 \left( \sum_{a_1,\ldots,a_m} \frac{n(a_1 \cdots a_m \bullet)^2}{\mu(a_1 \cdots a_m)} \right. \\
&\quad \left. - \sum_{a_1,\ldots,a_{m+1}} \frac{n(a_1 \cdots a_{m+1})^2}{\mu(a_1 \cdots a_{m+1})} + \frac{1 - 2n(w_1 \cdots w_m \bullet)}{\mu(w_1 \cdots w_m)} \right) ,
\end{aligned}
\tag{14}
$$

where $n(\cdot)$ denotes the number of occurrences inside $w$, and $n(a_1 \cdots a_m \bullet)$ stands for $\sum_{b \in \mathcal{A}} n(a_1 \cdots a_m b)$. Since the conditional moment of order 4 of $N(w)/\sqrt{n}$ is bounded, it follows that

$$
\frac{1}{\sqrt{n}} \left( N(w) - \widehat{N}_m(w) \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma_m^2(w)).
$$

The overlapping structure of $w$ clearly appears in the limiting variance. It is an exercise to verify that the limiting variances given by Theorem 3 and Equation (14) with $m = \ell - 2$ are identical.

Both martingale and conditional approaches can be extended to the M$m$-3 model (see Introduction for definition and notation). When one wants to distinguish the occurrences of $w$ in a coding DNA sequence according to a particular phase $k \in \{1, 2, 3\}$ ($k$ represents the position of the word with respect to the codons), one is interested in the count $N(w, k)$ of $w$ in phase $k$ in $X_1 \cdots X_n$; recall that the word phase is the phase of its last letter. Here we state the result in the maximal model; see Schbath (1995b) for the general case.

**Theorem 4** *Assume $\underline{X} = (X_i)_{i \in \mathbb{Z}}$ is a stationary $(\ell - 2)$-order Markov chain on $\mathcal{A}$ with transition probabilities $\pi_k(a_1 \cdots a_{\ell-2}, b)$ and stationary distribution $\mu(a_1 \cdots a_{\ell-2}, k)$, $a_1, \ldots, a_{\ell-2}, b \in \mathcal{A}$, $k \in \{1, 2, 3\}$. As $n \to \infty$, we have*

$$
\frac{1}{\sqrt{n}} \left( N(w, k) - \frac{N(w^-, k-1)N(^-w, k)}{N(^-w^-, k-1)} \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma_{\ell-2}^2(w, k))
$$

*with*

$$
\sigma_{\ell-2}^2(w, k) = \frac{\mu(w, k)}{\mu(^-w^-, k-1)^2} \left( \mu(^-w^-, k-1) - \mu(^-w, k) \right) \left( \mu(^-w^-, k-1) - \mu(w^-, k-1) \right)
$$

*and*

$$
\begin{aligned}
\mu(w^-, k-1) &= \mu(w_1 \cdots w_{\ell-2}, k-2)\pi_{k-1}(w_1 \cdots w_{\ell-2}, w_{\ell-1}) \\
\mu(^-w, k) &= \mu(^-w^-, k-1)\pi_k(^-w^-, w_\ell) \\
\mu(w, k) &= \mu(w^-, k-1)\pi_k(^-w^-, w_\ell).
\end{aligned}
$$

Yet another approach is Stein's method for normal approximations, namely Theorem 14. Let us for the moment assume that we have the independent model M0. Put

$$
\begin{aligned}
\sigma^2 &= \mathrm{Var}N(w) \\
&= \sum_{i=\ell}^{n} \sum_{j=\ell}^{n} \left( \mathbb{E}(Y_i(w)Y_j(w)) - \mu(w)^2 \right) \\
&= 2 \sum_{i=\ell}^{n} \sum_{\ell \leq j \leq i} \left( \mathbb{E}(Y_i(w)Y_j(w)) - \mu(w)^2 \right) \\
&= 2 \sum_{i=\ell}^{n} \sum_{j=(i-\ell+1)\vee\ell}^{i} \left( \mathbb{E}(Y_i(w)Y_j(w)) - \mu(w)^2 \right) \\
&= 2 \sum_{p \in \mathcal{P}(w)\cup\{0\}} (n - \ell - p + 1)\left( \mu(w^{(p)}w) - \mu(w)^2 \right).
\end{aligned}
$$

18

Note that
$$\sigma \asymp \sqrt{n}.$$

Define
$$Z_i(w) = \frac{1}{\sigma}\left(Y_i(w) - \mu(w)\right), \quad i = \ell, \dots, n.$$

Put
$$W = \sum_{i=\ell}^{n} Z_i(w) = \frac{1}{\sigma}\{N(w) - (n - \ell + 1)\mu(w)\}.$$

Choose $\mathcal{S}_j = \{i : |i - j| \leq \ell - 1\}$, and $\mathcal{N}_j = \{i : |i - j| \leq 2\ell - 1\}$. With this neighborhood structure, we use Theorem 14 with $B = \frac{1}{\sigma}$, $D_1 = 2\ell - 1$, and $D_2 = 4\ell - 3$, and $\mathcal{H}$ the set of indicators of half-lines. This provides an explicit bound on the Kolmogorov-Smirnov distance to the normal distribution. Note that, due to the independence, the quantities $\chi_1, \chi_2$ and $\chi_3$ vanish.

**Theorem 5** *Assume the independent model M0. There are constants $a$ and $c$ such that, for any word $w$ of length $\ell$,*
$$|\mathbb{P}(W \leq x) - \Phi(x)| \leq c\left\{a\frac{1}{\sigma}(4\ell - 3) + \frac{n}{\sigma^{3/2}}(a + 1)(2\ell - 1)(4\ell - 3)\right\}.$$

Moreover the result could be generalized to the model M$m$, using a neighborhood of size proportional to $\log n$. As Theorem 14 would need some modifying (the quantities to bound should be in terms of conditional expectations with respect to sigma-fields rather than with respect to sums, to proceed as in Reinert and Schbath (1998)), the result as well as the multivariate generalization will not be presented here.

## 5.4   Asymptotic distribution: the Poisson regime

In the previous section, we showed that the count $N(w)$ of a word $w$ in a random sequence of length $n$ can be approximated by a Gaussian distribution for large $n$. This Gaussian approximation is in fact not good when the expected count $(n - \ell + 1)\mu(w)$ is very small, meaning that $w$ is a rare word. Poisson approximations are appropriate for counts of rare events. As an illustration, it is well-known that a sum of independent Bernoulli variables can be either approximated by a Gaussian distribution or a Poisson distribution, depending on the asymptotic behavior of the expected value.

When the sequence letters are independent, Poisson and compound Poisson approximations for $N(w)$ have been widely studied in the literature (Chryssaphinou and Papastavridis (1988a), Chryssaphinou and Papastavridis (1988b), Arratia *et al.* (1990), Godbole (1991), Hirano and Aki (1993), Godbole and Schaffner (1993), Fu (1993)). Markovian models under different conditions have then been considered (Rajarshi (1974), Godbole (1991), Godbole and Schaffner (1993), Hirano and Aki (1993), Geske *et al.* (1995), Schbath (1995a), Erhardsson (1997)), but few works concern general periodic words and provide explicit parameters of the limiting distribution. As we will see, a Poisson distribution is not satisfactory for periodic words because of possible overlaps; a compound Poisson distribution is proposed. Two classes of tools can be used: generating functions, which do not provide any approximation error, and the Chen-Stein method, which gives a bound for the total variation distance between the two distributions (see Section 9.1 for details). In this section, we chose to present the Chen-Stein approach under a first-order Markovian model with known parameters; generalizations to higher order and to estimated parameters are presented at the end of the section. No assumption is made on the overlapping structure of the word $w$. Our two basic references are Arratia *et al.* (1990) and Schbath (1995a); for the case that the sequence is composed of i.i.d. letters see also Apostolico *et al.* (1998).

We assume that $\underline{X} = (X_i)_{i \in \mathbb{Z}}$ is a stationary first-order Markov chain on $\mathcal{A}$, with transition probabilities $\pi(a, b)$ and stationary distribution $\mu(a)$, $a, b \in \mathcal{A}$. Let $w = w_1 \cdots w_\ell$ be a word of length $\ell$ on $\mathcal{A}$. Here, $Y_i = Y_i(w) = \mathbb{I}\{w \text{ starts at position } i \text{ in } \underline{X}\}$ and $\mu(w) = \mathbb{E}Y_i(w)$.

Applying Theorem 13 to the Bernoulli variables $Y_i$, we obtain a bound $b_1 + b_2 + b_3$ for the total variation distance between the distribution of $N(w)$ and the Poisson distribution with mean $(n - \ell + 1)\mu(w)$; the mean does not converge to 0 under the rare word assumption $n\mu(w) = O(1)$. Note that $n\mu(w) = O(1)$ also means $\ell = O(\log n)$. The main difficulty when applying Theorem 13 comes from the $b_2$ term and the

possible overlaps of periodic words. Indeed, let $w$ be a periodic word; its set of period $\mathcal{P}(w)$ is not empty. Take $B_i = \{i - 2\ell + 1, \dots, i + 2\ell - 1\}$ for the neighborhood of $i \in I = \{1, \dots, n - \ell + 1\}$; it guarantees that $b_1$ and $b_3$ tend to 0 as $n \to +\infty$. We get

$$b_2 := \sum_{i \in I} \sum_{j \in B_i \setminus \{i\}} \mathbb{E} Y_i Y_j = 2(n - \ell + 1) \sum_{p \in \mathcal{P}(w)} \mu(w^{(p)} w) + \mathrm{O}(n\ell\mu^2(w));$$

this quantity can be of order $\mathrm{O}(1)$ for small periods $p$. The Poisson approximation is however valid for the count of non-periodic words because the set of periods is empty. For periodic words, the crucial argument is to consider clumps that, by definition, cannot overlap. We first prove that the declumped count $\widetilde{N}(w)$ can be approximated by a Poisson distribution with mean $(n - \ell + 1)\widetilde{\mu}(w)$ (see Eq. (5)) by applying Theorem 13 to the Bernoulli variables $\widetilde{Y}_i(w)$ defined in (3). Then, we prove a compound Poisson approximation for $N(w)$ by applying Theorem 13 to the Bernoulli variables $\widetilde{Y}_{i,k}(w)$ defined in (6) and by using that $N(w)$ is asymptotically equivalent to $\sum_{i \in I} \sum_{k \geq 1} k\widetilde{Y}_{i,k}(w)$ in probability. For simplicity, the variables $\widetilde{Y}_i(w)$ and $\widetilde{Y}_{i,k}(w)$ are denoted by $\widetilde{Y}_i$ and $\widetilde{Y}_{i,k}$.

**Poisson approximation for the declumped count** Our aim is to approximate the vector $\underline{\widetilde{Y}} = (\widetilde{Y}_i(w))_{i \in I}$ of Bernoulli variables by a vector $\underline{Z} = (Z_i)_{i \in I}$ with independent Poisson coordinates with mean $\mathbb{E}Z_i = \mathbb{E}\widetilde{Y}_i(w) = \widetilde{\mu}(w)$, where $\widetilde{\mu}(\cdot)$ is defined in (5). To apply Theorem 13, we choose the following neighborhood of $i \in I$:
$$B_i := \{j \in I : |j - i| \leq 3\ell - 3\}.$$

The neighborhood is such that, for $j$ not in $B_i$, there are no letters $X_h$ common to $\widetilde{Y}_i$ and $\widetilde{Y}_j$, and moreover, the $X_h$'s defining $\widetilde{Y}_i$ and those defining $\widetilde{Y}_j$ are separated by at least $\ell$ positions. It is important to consider a lag converging to infinity with $n$ since it leads to the exponential decay of the $b_3$ term given by Theorem 13 as we will see below. Deriving a bound for the total variation distance between $\underline{\widetilde{Y}}$ and $\underline{Z}$ consists of bounding the quantities $b_1, b_2$ and $b_3$ given in (28), (29) and (30). Bounding $b_1$ presents no difficulty:
$$b_1 := \sum_{i \in I} \sum_{j \in B_i} \mathbb{E}\widetilde{Y}_i \mathbb{E}\widetilde{Y}_j \leq (n - \ell + 1)(6\ell - 5)\widetilde{\mu}^2(w) = \mathrm{O}\left(\frac{\log n}{n}\right).$$

Since clumps of $w$ do not overlap in the sequence, $\widetilde{Y}_i \widetilde{Y}_j = 0$ for $|j - i| < \ell$. Therefore, we get

$$b_2 := \sum_{i \in I} \sum_{j \in B_i \setminus \{i\}} \mathbb{E}\widetilde{Y}_i \widetilde{Y}_j \leq 2 \sum_{i \in I} \sum_{j = i + \ell}^{i + 3\ell - 3} \mathbb{E}\widetilde{Y}_i \widetilde{Y}_j$$

using the symmetry of $B_i$. Now we have

$$\mathbb{E}\widetilde{Y}_i \widetilde{Y}_j \leq \mathbb{E}\widetilde{Y}_i Y_j = \widetilde{\mu}(w)\Pi^{j - i - \ell + 1}(w_\ell, w_1)\frac{\mu(w)}{\mu(w_1)}$$

and

$$b_2 \leq \frac{2}{\mu(w_1)}(n - \ell + 1)\widetilde{\mu}(w)\mu(w) \sum_{s=1}^{2\ell - 2} \Pi^s(w_\ell, w_1) = \mathrm{O}\left(\frac{\log n}{n}\right).$$

Bounding $b_3$ is a little more involved, but we give all the steps because the same technique is used for the compound Poisson approximation of the count and will not be described in detail then. By definition we have
$$b_3 := \sum_{i \in I} \mathbb{E}|\mathbb{E}(\widetilde{Y}_i - \mathbb{E}\widetilde{Y}_i \mid \sigma(\widetilde{Y}_j, j \notin B_i))|.$$

Since $\sigma(\widetilde{Y}_j, j \notin B_i) \subset \sigma(X_1, \dots, X_{i-2\ell+1}, X_{i+2\ell-1}, \dots, X_n)$, properties of conditional expectation and the Markov property give

$$
\begin{aligned}
b_3 \;\leq\; & \sum_{i \in I} \mathbb{E}|\mathbb{E}(\widetilde{Y}_i - \mathbb{E}\widetilde{Y}_i \mid X_{i-2\ell+1}, X_{i+2\ell-1})| \\
\leq\; & \sum_{i \in I} \sum_{x,y \in \mathcal{A}} |\mathbb{E}(\widetilde{Y}_i - \mathbb{E}\widetilde{Y}_i \mid X_{i-2\ell+1} = x, X_{i+2\ell-1} = y)| \\
& \times \mathbb{P}(X_{i-2\ell+1} = x, X_{i+2\ell-1} = y).
\end{aligned}
$$

To evaluate the right-hand term, we introduce the set of possible words of length $\ell - 1$ preceding a clump of $w$:

$$
\mathcal{G}(w) = \{ g = g_1 \cdots g_{\ell-1} : \text{ for all } p \in \mathcal{P}(w), g_{\ell-p} \cdots g_{\ell-1} \neq w^{(p)} \}. \tag{15}
$$

Thus a clump of $w$ starts at position $i$ in $(X_i)_{i \in \mathbb{Z}}$ if and only if one of the words $gw$, $g \in \mathcal{G}(w)$, starts at position $i - \ell + 1$. Therefore, we can write

$$
\widetilde{Y}_i(w) = \sum_{g \in \mathcal{G}(w)} Y_{i-\ell+1}(gw). \tag{16}
$$

This gives us

$$
\begin{aligned}
b_3 \;\leq\; & \sum_{i \in I} \sum_{x,y \in \mathcal{A}} \sum_{g \in \mathcal{G}(w)} |\mathbb{E}(Y_{i-\ell+1}(gw) - \mathbb{E}Y_{i-\ell+1}(gw) \mid X_{i-2\ell+1} = x, X_{i+2\ell-1} = y)| \\
& \times \mathbb{P}(X_{i-2\ell+1} = x, X_{i+2\ell-1} = y) \\
=\; & \sum_{i \in I} \sum_{x,y \in \mathcal{A}} \sum_{g \in \mathcal{G}(w)} |\mathbb{P}(X_{i-2\ell+1} = x, Y_{i-\ell+1}(gw) = 1, X_{i+2\ell-1} = y) \\
& -\mu(gw)\mathbb{P}(X_{i-2\ell+1} = x, X_{i+2\ell-1} = y)| \\
=\; & \sum_{i \in I} \sum_{x,y \in \mathcal{A}} \sum_{g \in \mathcal{G}(w)} \left| \mu(x)\Pi^\ell(x, g_1) \frac{\mu(gw)}{\mu(g_1)} \Pi^\ell(w_\ell, y) - \mu(gw)\mu(x)\Pi^{4\ell-2}(x, y) \right|.
\end{aligned}
$$

We now diagonalize the transition matrix. Let $(\alpha_t)_{t=1,\dots,|\mathcal{A}|}$ be the eigenvalues of $\Pi$ such that $|\alpha_1| \geq |\alpha_2| \geq \cdots \geq |\alpha_{|\mathcal{A}|}|$. The Perron-Frobenius Theorem (see, e.g., Karlin and Taylor (1975)) ensures that $\alpha_1 = 1$ and $|\alpha_2| < 1$; we abbreviate $\alpha_2$ by $\alpha$. $(1, 1, \dots, 1)^T$ is a right-eigenvector of $\Pi$ for the eigenvalue 1 whereas the vector of the stationary distribution $(\mu(a), a \in \mathcal{A})$ is a left-eigenvector of $\Pi$ for the eigenvalue 1. Let $D = \mathrm{Diag}(1, \alpha, \alpha_3, \cdots, \alpha_{|\mathcal{A}|})$. We decompose $\Pi = PDP^{-1}$ such that the first column of $P$ is $(1, 1, \dots, 1)^T$; then the first row of $P^{-1}$ is the vector of the stationary distribution $(\mu(a), a \in \mathcal{A})$. For all $t \in \{1, \dots, |\mathcal{A}|\}$, $I_t$ denotes the $|\mathcal{A}| \times |\mathcal{A}|$ matrix such that all its entries are equal to 0 except $I_t(t, t) = 1$, and we define $Q_t := PI_tP^{-1}$. We now use that $\Pi^\ell = PD^\ell P^{-1} = \sum_{t=1}^{|\mathcal{A}|} \alpha_t^\ell Q_t$ and $Q_1(a, b) = \mu(b)$, $\forall a, b \in \mathcal{A}$:

$$
\begin{aligned}
b_3 \;\leq\; & (n - \ell + 1)|\alpha|^\ell \sum_{g \in \mathcal{G}(w)} \mu(gw) \sum_{x,y \in \mathcal{A}} \mu(x) \left| \frac{1}{\mu(g_1)} \sum_{(t,t')} \frac{\alpha_t^\ell \alpha_{t'}^\ell}{\alpha^\ell} Q_t(x, g_1) Q_{t'}(w_\ell, y) \right. \\
& \left. - \sum_{t=1}^{|\mathcal{A}|} \frac{\alpha_t^{4\ell-2}}{\alpha^\ell} Q_t(x, y) \right| \\
=\; & (n - \ell + 1)|\alpha|^\ell \sum_{g \in \mathcal{G}(w)} \mu(gw) \sum_{x,y \in \mathcal{A}} \mu(x) \left| \frac{1}{\mu(g_1)} \sum_{(t,t') \neq (1,1)} \frac{\alpha_t^\ell \alpha_{t'}^\ell}{\alpha^\ell} Q_t(x, g_1) Q_{t'}(w_\ell, y) \right. \\
& \left. - \sum_{t=2}^{|\mathcal{A}|} \frac{\alpha_t^{4\ell-2}}{\alpha^\ell} Q_t(x, y) \right| \\
\leq\; & (n - \ell + 1)|\alpha|^\ell \sum_{g \in \mathcal{G}(w)} \mu(gw)\gamma(\ell, w_\ell),
\end{aligned}
$$

where

$$\gamma(\ell, a) = \max_{b \in \mathcal{A}} \sum_{x,y \in \mathcal{A}} \mu(x) \left| \frac{1}{\mu(b)} \sum_{(t,t') \neq (1,1)} \frac{\alpha_t^\ell \alpha_{t'}^\ell}{\alpha^\ell} Q_t(x,b) Q_{t'}(a,y) - \sum_{t=2}^{|\mathcal{A}|} \frac{\alpha_t^{4\ell-2}}{\alpha^\ell} Q_t(x,y) \right|.$$

Note that $\gamma(\ell, w_\ell) = \mathrm{O}(1)$. From (15) we have $\sum_{g \in \mathcal{G}(w)} \mu(gw) = \widetilde{\mu}(w)$ and

$$b_3 \leq (n - \ell + 1)\widetilde{\mu}(w)\gamma(\ell, w_\ell)|\alpha|^\ell = \mathrm{O}(|\alpha|^\ell).$$

We have proved the next theorem.

**Theorem 6** *Let $\underline{Z} = (Z_i)_{i \in I}$ be independent Poisson variables with expectation $\mathbb{E}Z_i = \mathbb{E}\widetilde{Y}_i(w) = \widetilde{\mu}(w)$. We have*

$$d_{TV}\left(\mathcal{L}(\underline{\widetilde{Y}}), \mathcal{L}(\underline{Z})\right) \leq (n - \ell + 1)(6\ell - 5)\widetilde{\mu}^2(w) + (n - \ell + 1)\widetilde{\mu}(w)\gamma(\ell, w_\ell)|\alpha|^\ell$$
$$+ \frac{2}{\mu(w_1)}(n - \ell + 1)\widetilde{\mu}(w)\mu(w) \sum_{s=1}^{2\ell-2} \Pi^s(w_\ell, w_1).$$

The declumped count $\widetilde{N}(w)$ can be approximated by $\widetilde{N}_{\inf}(w) := \sum_{i \in I} \widetilde{Y}_i(w)$ since

$$d_{\mathrm{TV}}\left(\mathcal{L}(\widetilde{N}(w)), \mathcal{L}(\widetilde{N}_{\inf}(w))\right) \leq \mathbb{P}(\widetilde{N}(w) \neq \widetilde{N}_{\inf}(w)) \leq (\ell - 1)(\mu(w) - \widetilde{\mu}(w))$$

(see Section 3). Using the triangle inequality leads to the following corollary:

**Corollary 1** *Let $Z$ be a Poisson variable with expectation $\mathbb{E}Z = (n - \ell + 1)\widetilde{\mu}(w)$. We have*

$$d_{TV}\left(\mathcal{L}(\widetilde{N}(w)), \mathcal{L}(Z)\right) \leq (n - \ell + 1)(6\ell - 5)\widetilde{\mu}^2(w) + (n - \ell + 1)\widetilde{\mu}(w)\gamma(\ell, w_\ell)|\alpha|^\ell$$
$$+ \frac{2}{\mu(w_1)}(n - \ell + 1)\widetilde{\mu}(w)\mu(w) \sum_{s=1}^{2\ell-2} \Pi^s(w_\ell, w_1) + (\ell - 1)(\mu(w) - \widetilde{\mu}(w)).$$

**Compound Poisson approximation for the count** To approximate the distribution of the count $N(w)$, we first use that $N(w)$ is asymptotically equivalent to $N_{\inf}(w) := \sum_{i=1}^{n-\ell+1} \sum_{k \geq 1} k\widetilde{Y}_{i,k}$ in probability (Reinert and Schbath (1998)):

$$d_{\mathrm{TV}}\left(\mathcal{L}(N(w)), \mathcal{L}(N_{\inf}(w))\right) \leq \mathbb{P}(N(w) \neq N_{\inf}(w)) \leq 2(\ell - 1)(\mu(w) - \widetilde{\mu}(w)).$$

Our goal is now to approximate the vector $(\widetilde{Y}_{i,k})_{(i,k) \in I}$, $I = \{1, \ldots, n - \ell + 1\} \times \{1, 2, \ldots\}$, of Bernoulli variables by a vector $(Z_{i,k})_{(i,k) \in I}$ with independent Poisson coordinates with expectation $\mathbb{E}Z_{i,k} = \mathbb{E}\widetilde{Y}_{i,k} = \widetilde{\mu}_k(w)$ where $\widetilde{\mu}_k(\cdot)$ is given in Equation (7). The neighborhood $B_{i,k}$ of $(i,k)$ is such that, for $(j,k')$ not in $B_{i,k}$, the letters $X_h$'s defining $\widetilde{Y}_{i,k}$ and those defining $\widetilde{Y}_{j,k}$ are separated by at least $\ell$ positions. Since $\widetilde{Y}_{i,k}$ can be described by at most $X_{i-\ell+1}, \ldots, X_{i+(k+1)(\ell-1)}$, we consider

$$B_{i,k} := \{(j, k') \in I : -(k' + 3)(\ell - 1) \leq j - i \leq (k + 3)(\ell - 1)\}.$$

We bound successively the quantities given in (28), (29) and (30). By definition

$$b_1 := \sum_{(i,k) \in I} \sum_{(j,k') \in B_{i,k}} \mathbb{E}\widetilde{Y}_{i,k} \, \mathbb{E}\widetilde{Y}_{j,k'}$$
$$\leq \sum_{i=1}^{n-\ell+1} \sum_{k \geq 1} \sum_{k' \geq 1} \sum_{j=i-(k'+3)(\ell-1)}^{i+(k+3)(\ell-1)} \widetilde{\mu}_k(w)\widetilde{\mu}_{k'}(w)$$
$$\leq (n - \ell + 1) \sum_{k \geq 1} \sum_{k' \geq 1} \left((k + k' + 6)(\ell - 1) + 1\right)\widetilde{\mu}_k(w)\widetilde{\mu}_{k'}(w).$$

22

From (5) and (7), we use that

$$\sum_{k\geq 1}\widetilde{\mu}_k(w) \;=\; \widetilde{\mu}(w)\,, \tag{17}$$

$$\sum_{k\geq 1}k\widetilde{\mu}_k(w) \;=\; \mu(w)\,, \tag{18}$$

and we obtain

$$b_1 \leq (n-\ell+1)\left(2(\ell-1)\widetilde{\mu}(w)\mu(w) + (6\ell-5)\widetilde{\mu}(w)^2\right).$$

The $b_2$ term involves products such as $\widetilde{Y}_{i,k}\widetilde{Y}_{j,k'}$ with $(j,k')\in B_{i,k}$. Since a $k$-clump of $w$ at position $i$ cannot overlap a $k'$-clump of $w$, many of these products are zero. To identify them, we need to describe in more detail the compound words $c\in\mathcal{C}_k(w)$ and $c'\in\mathcal{C}_{k'}(w)$ that may occur at positions $i$ and $j$. For this purpose, we introduce the set of words of length $\ell-1$ that can follow a clump of $w$:

$$\mathcal{D}(w) = \{d = d_1\cdots d_{\ell-1} : \forall p\in\mathcal{P}(w), d_1\cdots d_p \neq w_{\ell-p+1}\cdots w_\ell\}.$$

Therefore, we can write

$$\widetilde{Y}_{i,k}(w) = \sum_{g\in\mathcal{G}(w),c\in\mathcal{C}_k(w),d\in\mathcal{D}(w)} Y_{i-\ell+1}(gCd). \tag{19}$$

For convenience, we simply write $\sum_{gcd}$ for the sum over $g\in\mathcal{G}(w)$, $c\in\mathcal{C}_k(w)$, $d\in\mathcal{D}(w)$, and $\sum_{g'c'd'}$ for the sum over $g'\in\mathcal{G}(w)$, $c'\in\mathcal{C}_{k'}(w)$ and $d'\in\mathcal{D}(w)$. It gives us

$$
\begin{aligned}
b_2 \;:=&\; \sum_{(i,k)\in I}\sum_{(j,k')\in I\setminus\{(i,k)\}} \mathbb{E}\widetilde{Y}_{i,k}\widetilde{Y}_{j,k'}\\
=&\; \sum_{i=1}^{n-\ell+1}\sum_{k\geq 1}\sum_{k'\geq 1}\sum_{gcd}\sum_{g'c'd'}\sum_{j=i-(k'+3)(\ell-1)}^{i+(k+3)(\ell-1)} \mathbb{E}Y_{i-\ell+1}(gcd)Y_{j-\ell+1}(g'c'd').
\end{aligned}
$$

For $i-|c'| < j < i+|c|$, we have $Y_{i-\ell+1}(gcd)Y_{j-\ell+1}(g'c'd') = 0$ because clumps do not overlap. We then distinguish two cases:

(1) $g'c'd'$ at position $j-\ell+1$ overlaps $gcd$ at position $i-\ell+1$ (this is only possible over at most $2(\ell-1)$ letters); that is, for

$$j \in \{i-|c'|-2\ell+3,\ldots,i-|c'|\} \cup \{i+|c|,\ldots,i+|c|+2\ell-3\}\,;$$

let $b_{21}$ denote the associated term.

(2) $g'c'd'$ at position $j-\ell+1$ does not overlap $gcd$ at position $i-\ell+1$; that is, for

$$j \in \{i-(k'+3)(\ell-1),\ldots,i-|c'|-2\ell+2\} \cup \{i+|c|+2\ell-2,\ldots,i+(k+3)(\ell-1)\}\,;$$

let $b_{22}$ denote the associated term.

By symmetry, we have

$$b_{21} \leq 2\sum_{i=1}^{n-\ell+1}\sum_{k\geq 1}\sum_{k'\geq 1}\sum_{gcd}\sum_{g'c'd'}\sum_{j=i+|c|}^{i+|c|+2\ell-3} \mathbb{E}Y_{i-\ell+1}(gCd)Y_{j-\ell+1}(g'C'd')\,.$$

Summing over $k'$, $g'$, $c'$ and $d'$ gives

$$b_{21} \leq 2\sum_{i=1}^{n-\ell+1}\sum_{k\geq 1}\sum_{gcd}\sum_{j=i+|c|}^{i+|c|+2\ell-3} \mathbb{E}Y_{i-\ell+1}(gcd)\widetilde{Y}_j(w)\,;$$

23

now, summing over $d$ and using that $\widetilde{Y}_j(w) \le Y_j(w)$ leads to

$$b_{21} \le 2 \sum_{i=1}^{n-\ell+1} \sum_{k \ge 1} \sum_{gc} \sum_{j=i+|c|}^{i+|c|+2\ell-3} \mathbb{E}Y_{i-\ell+1}(gc)Y_j(w).$$

An occurrence of $gc$ at position $i-\ell+1$ does not overlap an occurrence of $w$ at position $j \ge i+|c|$; thus it follows that

$$\mathbb{E}Y_{i-\ell+1}(gc)Y_j(w) = \mu(gc)\Pi^{j-i-|c|+1}(w_\ell, w_1)\frac{\mu(w)}{\mu(w_1)},$$

and

$$b_{21} \le 2(n-\ell+1)\frac{\mu(w)}{\mu(w_1)} \sum_{s=1}^{2\ell-2} \Pi^s(w_\ell, w_1) \sum_{k \ge 1} \sum_{gc} \mu(gc).$$

Finally, note that

$$\sum_{k \ge 1} \sum_{gc} \mu(gc) = \sum_{k \ge 1} \sum_{k^* \ge k} \widetilde{\mu}_{k^*}(w) = \sum_{k^* \ge 1} k^* \widetilde{\mu}_{k^*}(w) = \mu(w),$$

which leads to

$$b_{21} \le 2(n-\ell+1)\frac{\mu^2(w)}{\mu(w_1)} \sum_{s=1}^{2\ell-2} \Pi^s(w_\ell, w_1) = O\left(\frac{\log n}{n}\right).$$

The $b_{22}$ term is easier to bound and we get

$$b_{22} \le 2(n-\ell+1)\frac{\widetilde{\mu}(w)}{\mu_{\min}}\left((\ell-2)\mu(w) + \widetilde{\mu}(w)\right) = O\left(\frac{\log n}{n}\right),$$

where $\mu_{\min}$ is the smallest value of $\{\mu(a), a \in \mathcal{A}\}$.

Finally, we have

$$b_2 \le 2(n-\ell+1)\frac{\mu^2(w)}{\mu(w_1)} \sum_{s=1}^{2\ell-2} \Pi^s(w_\ell, w_1) + 2(n-\ell+1)\frac{\widetilde{\mu}(w)}{\mu_{\min}}\left((\ell-2)\mu(w) + \widetilde{\mu}(w)\right).$$

Bounding $b_3$ consists of following the different steps previously described for the declumped count and using the decomposition (19) instead of (16). Since there is no interest in repeating this technical part, we just give the bound of $b_3$ and state the theorem:

$$b_3 \le (n-\ell+1)\widetilde{\mu}(w)\gamma_2(\ell)|\alpha|^\ell$$

with

$$\gamma_2(\ell) = \sum_{x,y \in \mathcal{A}} \mu(x) \max_{a,b \in \mathcal{A}} \left( \frac{1}{\mu(b)} \sum_{(t,t') \ne (1,1)} \left| \frac{\alpha_t^\ell \alpha_{t'}^\ell}{\alpha^\ell} Q_t(x,b)Q_{t'}(a,y) \right| + \sum_{t=2}^{|\mathcal{A}|} \left| \frac{\alpha_t^{5\ell-3}}{\alpha^\ell} Q_t(x,y) \right| \right).$$

**Theorem 7** *Let $(Z_{i,k})_{(i,k) \in I}$ be independent Poisson variables with expectation $\mathbb{E}Z_{i,k} = \mathbb{E}\widetilde{Y}_{i,k}(w) = \widetilde{\mu}_k(w)$. With the previous notation, we have*

$$d_{TV}\left( \mathcal{L}\left((\widetilde{Y}_{i,k}(w))_{(i,k) \in I}\right), \mathcal{L}\left((Z_{i,k})_{(i,k) \in I}\right) \right)$$

$$\le (n-\ell+1)\left( 2(\ell-1)\widetilde{\mu}(w)\mu(w) + (6\ell-5)\widetilde{\mu}(w)^2 \right) + (n-\ell+1)\widetilde{\mu}(w)\gamma_2(\ell)|\alpha|^\ell$$

$$+2(n-\ell+1)\frac{\mu^2(w)}{\mu(w_1)} \sum_{s=1}^{2\ell-2} \Pi^s(w_\ell, w_1) + 2(n-\ell+1)\frac{\widetilde{\mu}(w)}{\mu_{min}}\left((\ell-2)\mu(w) + \widetilde{\mu}(w)\right).$$

From the total variation distance properties, we have

$$d_{\text{TV}}\left(\mathcal{L}\left(\sum_{(i,k)\in I} k\widetilde{Y}_{i,k}\right), \mathcal{L}\left(\sum_{(i,k)\in I} kZ_{i,k}\right)\right) \leq d_{\text{TV}}\left(\mathcal{L}\left((\widetilde{Y}_{i,k}(w))_{(i,k)\in I}\right), \mathcal{L}\left((Z_{i,k})_{(i,k)\in I}\right)\right).$$

Since the $Z_{i,k}$'s are independent Poisson variables, $\sum_{(i,k)\in I} kZ_{i,k}$ is distributed like $\sum_{k\geq 1} kZ_k$ where the $Z_k$'s are independent Poisson variables with expectation $(n - \ell + 1)\widetilde{\mu}_k(w)$. Note that the latter is a compound Poisson distribution. Using the triangle inequality leads to the following corollary:

**Corollary 2** *Let $(Z_k)_{k\geq 1}$ be independent Poisson variables with expectation $\mathbb{E}Z_k = (n - \ell + 1)\widetilde{\mu}_k(w)$; CP denotes the compound Poisson distribution of $\sum_{k\geq 1} kZ_k$. With the previous notation, we have*

$$
\begin{aligned}
d_{TV}\left(\mathcal{L}(N(w)), CP\right) \quad &\leq (n - \ell + 1)\left(2(\ell - 1)\widetilde{\mu}(w)\mu(w) + (6\ell - 5)\widetilde{\mu}(w)^2\right) + (n - \ell + 1)\widetilde{\mu}(w)\gamma_2(\ell)|\alpha|^\ell \\
&\quad + 2(n - \ell + 1)\frac{\mu^2(w)}{\mu(w_1)}\sum_{s=1}^{2\ell - 2}\Pi^s(w_\ell, w_1) + 2(n - \ell + 1)\frac{\widetilde{\mu}(w)}{\mu_{min}}\left((\ell - 2)\mu(w) + \widetilde{\mu}(w)\right) \\
&\quad + 2(\ell - 1)(\mu(w) - \widetilde{\mu}(w)) \;= O\left(\frac{\log n}{n}\right).
\end{aligned}
$$

Also using the Chen-Stein method, Erhardsson (1997) obtained a different type of bound using a Markov chain coupling instead of the local approach. Related results for Poisson approximations using a Markov chain coupling have been announced by Roos and Stark (1999). The special case of runs of 1 in a random sequence of letters in the binary alphabet $\{0, 1\}$ has been extensively studied: Erdös and Rényi (1970) gave the asymptotic behavior of the longest run in a sequence of Bernoulli trials, and of the length of the longest segment that contains a proportion of 1 greater than a predescribed level $\alpha$. Their result was refined by Deheuvels *et al.* (1986). The compound Poisson approximation for counts of runs in the case where the sequence letters are independent has been considered by Eichelsbacher and Roos (1999), also employing the Chen-Stein method using recent results by Barbour and Utev (1998) (the limiting distribution is the same as the one given above, reduced to this special case). Very recently, Barbour and Xia (1999) obtained a more accurate limiting approximation for the case of runs of length 2; this approximation is based on a perturbation of a Poisson distribution.

Such a bound on the total variation distance between, for instance, the word count distribution and the associated compound Poisson distribution has the great advantage of providing confidence intervals (see Section 9.1). Indeed, using notation from Corollary 2, for all $t \in \mathbb{R}$, we have

$$\left|\mathbb{P}(N(w) \geq t) - \mathbb{P}\left(\sum_{k\geq 1} kZ_k \geq t\right)\right| \leq d_{\text{TV}}\left(\mathcal{L}(N(w)), \mathcal{L}\left(\sum_{k\geq 1} kZ_k\right)\right).$$

**Estimation of the parameters** When the transition probabilities are unknown and can only be estimated from the observed sequence, we need to evaluate the total variation distance between the word count distribution and the distribution of $\sum_{k\geq 1} kZ_k'$ where the $Z_k'$'s are independent Poisson variables with expectation $(n - \ell + 1)\widehat{\widetilde{\mu}}_k(w)$, where $\widehat{\widetilde{\mu}}_k(w)$ is the observed value of the plug-in maximum likelihood estimator of $\widetilde{\mu}_k(w)$. Similarly, we want to know the total variation distance between the declumped count, $\widetilde{N}(w)$, and the Poisson variable with expectation $(n - \ell + 1)\widehat{\widetilde{\mu}}(w)$. For this we use the triangle inequality and the fact that the total variation distance between two Poisson variables with expectation $\lambda$ and $\lambda'$ is less than $|\lambda - \lambda'|$:

$$
\begin{aligned}
d_{\text{TV}}\left(\mathcal{L}(\widetilde{N}(w)), \mathcal{P}o((n - \ell + 1)\widehat{\widetilde{\mu}}(w))\right) \quad &\leq \quad d_{\text{TV}}\left(\mathcal{L}(\widetilde{N}(w)), \mathcal{P}o((n - \ell + 1)\widetilde{\mu}(w))\right) \\
&\quad + (n - \ell + 1)|\widehat{\widetilde{\mu}}(w) - \widetilde{\mu}(w)|.
\end{aligned}
$$

Using the Law of Iterated Logarithm for Markov chains (Senoussi (1990)) and Equation (2), one can show that

$$\widehat{\mu}(w) = \mu(w)\left(1 + O\left(\frac{\ell\sqrt{\log\log n}}{\sqrt{n}}\right)\right) \qquad \text{almost surely (a.s.)}$$

(see Schbath (1995b)). Under the rare word condition $n\mu(w) = O(1)$, we get

$$n\widehat{\mu}(w) - n\mu(w) = O\left(\frac{\ell\sqrt{\log\log n}}{\sqrt{n}}\right) \quad \text{a.s.}$$

Now, using Equation (5), we obtain

$$n\widehat{\widetilde{\mu}}(w) - n\widetilde{\mu}(w) = O\left(\frac{\ell^2\sqrt{\log\log n}}{\sqrt{n}}\right) \quad \text{a.s.}$$

This quantity converges to zero as $n \to \infty$, because the rare word condition implies that $\ell = O(\log n)$. Thus,

$$d_{\mathrm{TV}}\left(\mathcal{L}(\widetilde{N}(w)), \mathcal{P}\mathrm{o}\big((n-\ell+1)\widehat{\widetilde{\mu}}(w)\big)\right) \leq d_{\mathrm{TV}}\left(\mathcal{L}(\widetilde{N}(w)), \mathcal{P}\mathrm{o}\big((n-\ell+1)\widetilde{\mu}(w)\big)\right) + O\left(\frac{\ell^2\sqrt{\log\log n}}{\sqrt{n}}\right).$$

The approximation follows from Corollary 2.

Note that we do not have an explicit bound for this additional error term. However, for long sequences the error term due to the maximum-likelihood estimation will be small compared to the bound on the (compound) Poisson approximation error.

Similarly, the total variation distance between the two compound Poisson distributions is bounded by

$$d_{\mathrm{TV}}\left(\mathcal{L}\left(\sum_{k\geq 1} kZ_k\right), \mathcal{L}\left(\sum_{k\geq 1} kZ_k'\right)\right) \leq \sum_{k\geq 1} |n\widehat{\widetilde{\mu}}_k(w) - n\widetilde{\mu}_k(w)|.$$

Using Equation (7), this quantity tends to zero as $n \to \infty$ when $n\mu(w) = O(1)$ (see Schbath (1995b)).

**Generalization to M$m$.** Let us now assume that the sequence $(X_i)_{i\in\mathbb{Z}}$ is a $m$-order Markov chain on the alphabet $\mathcal{A}$, with transition probabilities $\pi(a_1\cdots a_m, a_{m+1})$, $a_1,\cdots,a_{m+1} \in \mathcal{A}$. The basic idea is to rewrite the sequence over the alphabet $\mathcal{A}^m$ by defining

$$\mathbb{X}_i = X_i X_{i+1} \cdots X_{i+m-1},$$

so that the sequence $(\mathbb{X}_i)_{i\in\mathbb{Z}}$ is a first-order Markov chain on $\mathcal{A}^m$ with transition probabilities such that, for $\mathbb{A} = a_1\cdots a_m \in \mathcal{A}^m$ and $\mathbb{B} = b_1\cdots b_m \in \mathcal{A}^m$,

$$\Pi(\mathbb{A}, \mathbb{B}) = \begin{cases} \pi(a_1\cdots a_m, b_m) & \text{if } a_2\cdots a_m = b_1\cdots b_{m-1} \\ 0 & \text{else.} \end{cases}$$

Denote by $\mathbb{W} = \mathbb{W}_1 \cdots \mathbb{W}_{\ell-m+1}$ the word $w = w_1\ldots w_\ell$ written using the alphabet $\mathcal{A}^m$, so that $\mathbb{W}_j = w_j\ldots w_{j+m-1}$. The results presented below are valid for the number $N(\mathbb{W})$ of overlapping occurrences and the number $\widetilde{N}(\mathbb{W})$ of clumps of $\mathbb{W}$ in $\mathbb{X}_1\cdots\mathbb{X}_{n-m+1}$. Since an occurrence of $w$ at position $i$ in $X_1\cdots X_n$ corresponds to an occurrence of $\mathbb{W}$ at position $i-m+1$ in $\mathbb{X}_1\cdots\mathbb{X}_{n-m+1}$, we simply have $N(w) = N(\mathbb{W})$. In contrast, clumps of $\mathbb{W}$ in $\mathbb{X}_1\cdots\mathbb{X}_{n-m+1}$ are different from clumps of $w$ in $X_1\cdots X_n$ because $\mathbb{W}$ is less periodic than $w$, leading to $\widetilde{N}(\mathbb{W}) \neq \widetilde{N}(w)$. Let us take a simple example: $w = \mathtt{ATA}$ and $m = 2$. Put $\mathbb{A} = \mathtt{AT} \in \mathcal{A}^2$ and $\mathbb{B} = \mathtt{TA} \in \mathcal{A}^2$; we then have $\mathbb{W} = \mathbb{AB}$. The sequence $\mathtt{TATATATAT}$ contains a unique clump of $\mathtt{AT}$ whereas the associated sequence $\mathbb{BABABABA}$ contains 3 clumps of $\mathbb{AB}$. Indeed, $\mathbb{AB}$ has no period and $\mathtt{ATA}$ has one period. In fact, the periods of $\mathbb{W}$ are those periods of $w$ that are strictly less than $\ell-m+1$. Therefore, the Poisson approximation for the declumped count in a $m$-order Markov chain does not follow immediately from the case $m = 1$; a rigorous proof would require applying the Chen-Stein theorem with an adapted neighborhood and to bound the new quantities $b_1$, $b_2$ and $b_3$ in M$m$, but this has not been carried out yet.

Since $N(w) = N(\mathbb{W})$, Corollary 2 ensures that $N(w)$ can be approximated by a sum $\sum_{k\geq 1} kZ_k$, where $Z_k$ is a Poisson variable whose expectation is $(n-\ell+1)$ times the probability that a $k$-clump of $\mathbb{W}$ starts at a given position in $\mathbb{X}_1\cdots\mathbb{X}_{n-m+1}$. From Equation (7), we obtain

$$\mathbb{E}Z_k = (n-\ell+1)(1-A'(w))^2 A'(w)^{k-1}\mu(w)$$

26

with

$$A'(w) = \sum_{p \in \mathcal{P}'(w) \cup \{1, \ldots, \ell-m\}} \frac{\mu(w^{(p)}w)}{\mu(w)}.$$

An important consequence is that, in M$m$, the compound Poisson approximation for words that cannot overlap on more than $m - 1$ letters comes from a single Poisson approximation.

## 5.5 Large deviation approach

For long sequences, the probability that a given word occurs more than a certain number of times can be approximated using a Gaussian or a compound Poisson distribution (Sections 5.3 and 5.4). The aim of this section is to show that large deviation techniques can also be used to approximate the probability that a given word frequency deviates from its expected value. Let $w = w_1 \cdots w_\ell$ be a word of length $\ell$; recall that $\mu(w)$ denotes the probability that $w$ occurs at a given position in $X_1 \cdots X_n$. We now aim to provide good approximations for $\mathbb{P}(\frac{1}{n-\ell+1}N(w) \geq \mu(w) + b)$ and $\mathbb{P}(\frac{1}{n-\ell+1}N(w) \leq \mu(w) - b)$ with $0 < b < 1$.

We assume that $X_1 \cdots X_n$ is a stationary first-order Markov chain on a finite alphabet $\mathcal{A}$ with transition probabilities $\pi(a, b) > 0$, $a, b \in \mathcal{A}$. (Generalizing to M$m$ follows the same setup as in Section 5.4.) To use Theorem 16 for $\frac{1}{n-\ell+1}N(w)$, we need to consider the irreducible Markov chain $\mathbb{X}_1, \ldots, \mathbb{X}_{n-\ell+1}$ on $\mathcal{A}^\ell$, where $\mathbb{X}_i = X_i \cdots X_{i+\ell-1}$, with transition matrix $\mathbb{\Pi} = (\Pi(u, v))_{u,v \in \mathcal{A}^\ell}$ such that

$$\Pi(u_1 \cdots u_\ell, v_1 \cdots v_\ell) = \begin{cases} \pi(u_\ell, v_\ell) & \text{if } u_{j+1} = v_j, \ j = 1 \cdots \ell - 1, \\ 0 & \text{else.} \end{cases}$$

The count $N(w)$ can then be written as

$$\begin{aligned} N(w) &= \sum_{i=1}^{n-\ell+1} \mathbb{I}\{X_i \cdots X_{i+\ell-1} = w_1 \cdots w_\ell\} \\ &= \sum_{i=1}^{n-\ell+1} \mathbb{I}\{\mathbb{X}_i = w\} := \sum_{i=1}^{n-\ell+1} f(\mathbb{X}_i) \end{aligned}$$

with $f(u) = \mathbb{I}\{u = w\}$, $u \in \mathcal{A}^\ell$. We have $\mathbb{E}f(\mathbb{X}_1) = \mu(w)$. Let $I$ be the function $I(x) = \sup_{\theta \in \mathbb{R}}(\theta x - \log \lambda(\theta))$, $x \in \mathbb{R}$, with $\lambda(\theta)$ be the largest eigenvalue of the matrix $\mathbb{\Pi}_\theta = (\Pi_\theta(u, v))_{u,v \in \mathcal{A}^\ell}$ defined by

$$\Pi_\theta(u, v) = \begin{cases} e^\theta \Pi(u, v) & \text{if } v = w, \\ \Pi(u, v) & \text{else.} \end{cases}$$

Let $0 < b < 1$; applying Theorem 16 to the closed subset $[\mu(w) + b, +\infty]$ and the open subset $(\mu(w) + b, +\infty)$, we get

$$\lim_{n \to +\infty} \frac{1}{n - \ell + 1} \log \mathbb{P}\left(\frac{1}{n - \ell + 1}N(w) \geq \mu(w) + b\right) = -I(\mu(w) + b) \, ;$$

similarly we have

$$\lim_{n \to +\infty} \frac{1}{n - \ell + 1} \log \mathbb{P}\left(\frac{1}{n - \ell + 1}N(w) \leq \mu(w) - b\right) = -I(\mu(w) - b).$$

Denoting the observed count of $w$ in the biological sequence by $N^{\text{obs}}(w)$, we get for large $n$: if $N^{\text{obs}}(w) > (n - \ell + 1)\mu(w)$ and $b := \frac{N^{\text{obs}}(w)}{n-\ell+1} - \mu(w)$, then

$$\mathbb{P}(N(w) \geq N^{\text{obs}}(w)) \simeq \exp\left(-(n - \ell + 1)I\left(\frac{N^{\text{obs}}(w)}{n - \ell + 1}\right)\right),$$

if $N^{\text{obs}}(w) < (n - \ell + 1)\mu(w)$ and $b := \mu(w) - \frac{N^{\text{obs}}(w)}{n-\ell+1}$, then

$$\mathbb{P}(N(w) \leq N^{\text{obs}}(w)) \simeq \exp\left(-(n - \ell + 1)I\left(\frac{N^{\text{obs}}(w)}{n - \ell + 1}\right)\right).$$

Note that this approximation has been obtained assuming the transition probabilities $\pi(a,b)$, $a, b \in \mathcal{A}$ are known. Moreover, since $\lambda(\theta)$ is an eigenvalue of a $|\mathcal{A}|^\ell \times |\mathcal{A}|^\ell$ matrix, the word length $\ell$ is a limiting factor for the numerical calculation, even if $|\mathcal{A}| = 4$.

# 6   Renewal count distribution

As a particular case of non-overlapping occurrence counts, in this section we count renewals of a word $w = w_1 w_2 \ldots w_\ell$ in a random sequence $X_1 \cdots X_n$ as defined in Section 3. We then consider the renewal count $R_n(w) = \sum_{i=1}^{n-\ell+1} \mathbb{I}_i(w)$, where $\mathbb{I}_i(w)$ is the random indicator that a renewal of $w$ starts at position $i$ in $X_1 \cdots X_n$ (see (8)). Exact results can be found in Régnier (1999); in particular, a combinatorial approach and language decompositions are used to derive the moment generating function of the renewal count. Because those tools are very different from the ones used in this paper, we only present asymptotic results.

## 6.1   Gaussian approximation

When the letters $X_1, \ldots, X_n$ are independent and identically distributed, the asymptotic distribution of the renewal count was studied by Breen *et al.* (1985); Tanushev (1996) proved a Central Limit Theorem in the Markovian case. The main technique being generating functions, no bound on the rate of convergence is obtained. Here we present the result from Tanushev (1996).

Note that, once asymptotic mean and variance are established, the normal approximation follows from the Markov renewal Central Limit Theorem. First we derive the expected renewal count.

If the $\mathbb{I}_i(w)$'s had the same expectation, say $\mu_R(w)$, then $\mathbb{E}R_n(w) = (n - \ell + 1)\mu_R(w)$. This is the commonly used expectation (see Breen *et al.* (1985) or Tanushev (1996) for instance), but it ignores the end effect. For $i > \ell$, the $\mathbb{I}_i(w)$'s are effectively identically distributed by stationarity of the Markov process, but this is not the case for $1 \leq i \leq \ell$.

We start with the calculation of $\mu_R(w)$. Recall that $\mathcal{P}(w)$ is the set of periods of $w$ and that $w^{(p)} = w_1 w_2 \cdots w_p$ denotes the word composed of the first $p$ letters of $w$. We consider the overlap-matching polynomial $Q(z)$ associated with $w$ (see, e.g. Guibas and Odlyzko (1980), Li (1980), Biggins and Cannings (1987)) defined by

$$Q(z) = \sum_{p \in \mathcal{P}(w) \cup \{0\}} \frac{\mu(w)}{\mu(w^{(\ell-p)})} z^p.$$

When the Markov process is in stationarity, we have from renewal theory that

$$\mu_R(w) = \frac{\mu(w)}{Q(1)}. \tag{20}$$

To understand this formula, note that we can decompose the event {there is an occurrence of $w$ starting at position $i$}, $i > \ell$, as the disjoint union of {there is a renewal of $w$ starting at position $i$} and {there is a renewal of $w$ starting at position $j$ directly followed by the letters $w_{\ell-i+j+1} \cdots w_\ell$ and $j - i$ is a period of $w$}, for $j \in \{i - \ell + 1, \ldots, i - 1\}$. This can be written as follows

$$
\begin{aligned}
Y_i(w) &= \sum_{j=i-\ell+1}^{i} \mathbb{I}_j(w) Y_{j+\ell}(w_{\ell-i+j+1} \cdots w_\ell) \mathbb{I}\{i - j \in \mathcal{P}(w) \cup \{0\}\} \\
&= \sum_{p \in \mathcal{P}(w) \cup \{0\}} \mathbb{I}_{i-p}(w) Y_{i+\ell-p}(w_{\ell-p+1} \cdots w_\ell).
\end{aligned}
$$

Taking expectations on both sides thus gives

$$
\mu(w) = \sum_{p \in \mathcal{P}(w) \cup \{0\}} \mu_R(w) \mu(w_{\ell-p} \cdots w_\ell) \frac{1}{\mu(w_{\ell-p})}.
$$

28

Hence

$$\mu_R(w) \;=\; \frac{\mu(w)}{\sum_{p\in\mathcal{P}(w)\cup\{0\}} \pi(w_{\ell-p},w_{\ell-p+1})\cdots\pi(w_{\ell-1},w_\ell)},$$

which gives the result (20).

As previously noted, the first variables $\mathbb{I}_1(w)$, ..., $\mathbb{I}_\ell(w)$ are not identically distributed because of boundary effects. For the asymptotic results we are interested in in this section, this end effect may be ignored. Note however that Régnier (1999) provides the exact renewal count expectation, namely

$$\mathbb{E}R_n(w) = (n-\ell+1)\mu_R(w) + \mu_R(w)\frac{Q'(1)}{Q(1)}.$$

Calculating the asymptotic variance is a little more involved, relying much on the overlap-matching polynomial. To this purpose, similarly to Tanushev (1996) we define $A$ as the $|\mathcal{A}|\times|\mathcal{A}|$ matrix where each row is the vector $(\mu(a), a\in\mathcal{A})$ of the stationary distribution. With $\Pi$ denoting the Markovian transition matrix, put

$$Z = \sum_{k=1}^{\infty}(\Pi-A)^k. \tag{21}$$

Put

$$\sigma^2 = \mu_R^2(w)\left((1-2\ell) + 2\frac{Q'(1)}{Q(1)} + \frac{2Z(w_\ell,w_1)}{\mu(w_1)}\right).$$

Then Tanushev (1996) proved the following theorem, using generating functions. (It is a special case of his theorem, which in fact proves a multivariate approximation.)

**Theorem 8** *We have that, as $n\to\infty$,*

$$\frac{R_n(w)-n\mu_R(w)}{\sqrt{n}} \overset{\mathcal{D}}{\longrightarrow} \mathcal{N}(0,\sigma^2).$$

The theorem is much easier to prove in the i.i.d. case, which is carried out in Waterman (1995). In fact all that is needed is to establish the variance expression (that will not be presented here); everything else follows from the Renewal Central Limit Theorem. Note that we do not have a corresponding result when mean and standard deviation are estimated.

## 6.2 Poisson approximation

Similarly as with the declumped count, we can also derive a Poisson approximation for the renewal count under the rare word condition $n\mu(w) = O(1)$. Indeed this is very simple. Recall (1)

$$Y_i(w) := \mathbb{I}\{w \text{ starts at position } i \text{ in } \underline{X}\}.$$

We write, for $i > \ell$,

$$
\begin{aligned}
\mathbb{I}_i(w) \;=\;& Y_i(w)\prod_{j=i-\ell+1}^{i-1}(1-\mathbb{I}_j(w))\\
\;=\;& Y_i(w)\prod_{j=i-\ell+1}^{i-1}(1-Y_j(w)) + Y_i(w)\left(\prod_{j=i-\ell+1}^{i-1}(1-\mathbb{I}_j(w)) - \prod_{j=i-\ell+1}^{i-1}(1-Y_j(w))\right)\\
\;=\;& \widetilde{Y}_i(w) + Y_i(w)\left(\prod_{j=i-\ell+1}^{i-1}(1-\mathbb{I}_j(w)) - \prod_{j=i-\ell+1}^{i-1}(1-Y_j(w))\right) \tag{22}
\end{aligned}
$$

whereas $\mathbb{I}_i(w) = Y_i(w) \prod_{j=1}^{i-1}(1 - Y_j(w))$ if $1 \leq i \leq \ell$. Note that a renewal occurrence in the first $\ell$ positions is a clump occurrence observed in the finite sequence, and conversely. Thus we have

$$
\begin{aligned}
R_n(w) &= \sum_{i=1}^{n-\ell+1} \mathbb{I}_i(w) \\
&= \widetilde{N}(w) + \sum_{i=\ell+1}^{n-\ell+1} Y_i(w) \left( \prod_{j=i-\ell+1}^{i-1}(1 - \mathbb{I}_j(w)) - \prod_{j=i-\ell+1}^{i-1}(1 - Y_j(w)) \right).
\end{aligned}
$$

We have already derived a Poisson approximation for the number of clumps $\widetilde{N}(w)$ (see Section 5.4). Let us consider the difference

$$
R_n(w) - \widetilde{N}(w) = \sum_{i=\ell+1}^{n-\ell+1} Y_i(w) \left( \prod_{j=i-\ell+1}^{i-1}(1 - \mathbb{I}_j(w)) - \prod_{j=i-\ell+1}^{i-1}(1 - Y_j(w)) \right).
$$

For a summand to be nonzero, firstly we need that $Y_i(w) = 1$. Note that a renewal always implies an occurrence, so that

$$
\prod_{j=i-\ell+1}^{i-1}(1 - \mathbb{I}_j(w)) \geq \prod_{j=i-\ell+1}^{i-1}(1 - Y_j(w)).
$$

The product being always 0 or 1, the two products are different if and only if $\prod_{j=i-\ell+1}^{i-1}(1 - \mathbb{I}_j(w)) = 1$ and $\prod_{j=i-\ell+1}^{i-1}(1 - Y_j(w)) = 0$. This implies that there is no renewal between the positions $i - \ell + 1$ and $i - 1$, but that there must be an occurrence not only at position $i$ but also at some position $j$ between $i - \ell + 1$ and $i - 1$. This occurrence again cannot be a renewal, so that it must be part of a larger clump; repeating this argument we see that the occurrence at $i$ must be part of a clump that started before position $i - \ell + 1$. This implies that there had to be an occurrence of $w$ somewhere between $i - 2\ell + 2$ and $i - \ell$, and this occurrence is in the same clump as the occurrence at $i$. Thus

$$
\begin{aligned}
\mathbb{P}(R_n(w) \neq \widetilde{N}(w)) &\leq \sum_{i=\ell+1}^{n-\ell+1} \sum_{j=i-2\ell+2}^{i-\ell} \mathbb{E}Y_i(w)Y_j(w) \\
&\leq (n - 2\ell + 1)(\ell - 1)\mu(w)^2 \frac{1}{\mu(w_1)}.
\end{aligned} \tag{23}
$$

This quantity will be small under the asymptotic framework $n\mu(w) = O(1)$. Thus we may use the Poisson bound for the number of clumps derived above, and just add an error term of order $\log n/n$. Indeed this has been the idea behind the proof of Geske *et al.* (1995), although Geske *et al.* (1995) prove the result only for words having at most one principal period. Related results have been obtained by Chryssaphinou and Papastavridis (1988b).

# 7 Occurrences of multiple patterns

When characterizing protein families via short motifs, for instance, one is interested in the distribution of the joint occurrences of multiple patterns rather than single patterns. It is also the case to determine the statistical significance of the count of degenerated words like A(C or G)G(A or T). Asymptotic results, similar to the above approximations, are available for the distribution of joint occurrences and joint counts of multiple patterns and we will present them in this section. As we will see, the main new feature one has to consider are the possible overlaps between different words from the target family.

Thus we are interested in the family of $q$ words $\{w^1, \dots, w^q\}$, where $w^r = w_1^r w_2^r \cdots w_{\ell_r}^r$. For two words $w^1 = w_1^1 w_2^1 \cdots w_{\ell_1}^1$ and $w^2 = w_1^2 w_2^2 \cdots w_{\ell_2}^2$ on $\mathcal{A}$, we describe the possible overlaps between $w^1$ and $w^2$ by defining

$$
\mathcal{P}(w^1, w^2) := \{p \in \{1, \dots, \ell_1 - 1\} : w_i^2 = w_{i+p}^1, \forall i = 1, \dots, (\ell_1 - p) \wedge \ell_2\}.
$$

Thus $\mathcal{P}(w^1, w^2) \neq \emptyset$ means that an occurrence of $w^2$ can overlap an occurrence of $w^1$ from the right, and $\mathcal{P}(w^2, w^1) \neq \emptyset$ means that $w^2$ can overlap $w^1$ from the left. Note the lack of symmetry; for example, if $w^1 = $ AAAGAAGAA and $w^2 = $ AAGAATCA, we have $\mathcal{P}(w^1, w^2) = \{4, 7, 8\}$ and $\mathcal{P}(w^2, w^1) = \{7\}$. To avoid trivialities, we make the following assumption.

**(A1)** $\qquad\qquad\qquad\qquad \forall r \neq r'$, $w^r$ is not a substring of $w^{r'}$ .

Again we model the sequence $\{X_i\}_{i \in \mathbb{Z}}$ as a stationary ergodic Markov chain.

**Gaussian approximation for the joint distribution of multiple word counts** We assume the general model M$m$, $m \leq \min\{\ell_r - 2, r = 1, \ldots, q\}$. We state the asymptotic normality of the vector $n^{-1/2}(N(w^r) - \widehat{N}_m(w^r))_{r=1,\ldots,q}$:

$$\frac{1}{\sqrt{n}} \big(N(w^r) - \widehat{N}_m(w^r)\big)_{r=1,\ldots,q} \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_m).$$

To prove this result, we use a multivariate martingale central limit theorem. The estimated count $\widehat{N}_m(w^r)$ is given by (9). The novelty consists here of deriving the asymptotic covariance matrix $\Sigma_m = (\Sigma_m(w^r, w^{r'}))_{r,r'=1,\ldots,q}$.

Suppose all the words $w^r$ have the same length $\ell$ and $m = \ell - 2$ (the maximal model) then the martingale technique (see Section 5.3) leads to

$$\Sigma_{\ell-2}(w^r, w^{r'}) = \mu(w^r)\mu(w^{r'}) \left( \frac{\mathbb{I}\{w^r = w^{r'}\}}{\mu(w^r)} - \frac{\mathbb{I}\{(w^r)^- = (w^{r'})^-\}}{\mu((w^r)^-)} \right.$$
$$\left. - \frac{\mathbb{I}\{^-(w^r) =^- (w^{r'})\}}{\mu(^-(w^r))} + \frac{\mathbb{I}\{^-(w^r)^- =^- (w^{r'})^-\}}{\mu(^-(w^r)^-)} \right).$$

Note that when $r = r'$, this formula reduces to the asymptotic variance $\sigma_{\ell-2}^2(w^r)$ of Section 5.3.

More generally, for $r \neq r'$, the conditional approach (see Section 5.3) leads to

$$\Sigma_m(w^r, w^{r'}) = \sum_{\substack{p \in \mathcal{P}(w^r, w^{r'}) \\ p \leq \ell_r - m - 1}} \mu\left((w^r)^{(p)} w^{r'}\right) + \sum_{\substack{p \in \mathcal{P}(w^{r'}, w^r) \\ p \leq \ell_{r'} - m - 1}} \mu\left((w^{r'})^{(p)} w^r\right)$$
$$+ \mu(w^r)\mu(w^{r'}) \left( \sum_{a_1,\ldots,a_m} \frac{n(a_1 \cdots a_m \bullet)n'(a_1 \cdots a_m \bullet)}{\mu(a_1 \cdots a_m)} \right.$$
$$- \sum_{a_1,\ldots,a_{m+1}} \frac{n(a_1 \cdots a_{m+1})n'(a_1 \cdots a_{m+1})}{\mu(a_1 \cdots a_{m+1})} - \frac{n(w_1^{r'} \cdots w_m^{r'} \bullet)}{\mu(w_1^{r'} \cdots w_m^{r'})}$$
$$\left. + \frac{\mathbb{I}\{w_1^r \cdots w_m^r = w_1^{r'} \cdots w_m^{r'}\} - n'(w_1^r \cdots w_m^r \bullet)}{\mu(w_1^r \cdots w_m^r)} \right),$$

where $n(\cdot)$ denotes the number of occurrences inside $w^r$ and $n'(\cdot)$ denotes the number of occurrences inside $w^{r'}$. (When $r = r'$, the formula reduces to Equation (14).)

Note that, if one wants to study the total number of occurrences of a word family $\{w^r, r = 1, \ldots, q\}$, we have

$$\frac{1}{\sqrt{n}} \left( \sum_{r=1}^q N(w^r) - \sum_{r=1}^q \widehat{N}_m(w^r) \right) \xrightarrow{\mathcal{D}} \mathcal{N} \left( 0, \sum_{r,r'} \Sigma_m(w^r, w^{r'}) \right).$$

**Poisson and compound Poisson approximations for the joint distribution of the declumped counts and multiple word counts** We assume the model M1 since generalization to M$m$ follows the single pattern case. To give a bound on the error for a Poisson process approximation for overlapping

counts, Reinert and Schbath (1998) define the following quantities for all $r$ and $r'$ in $\{1,\ldots,q\}$, and for all $a \in \mathcal{A}$:

$$\Omega_r = \sum_{s=1}^{3\ell-\ell_r-2} \Pi^s,$$

$$\Omega_{r,r'} = \sum_{s=1}^{\ell_r+\ell_{r'}-2} \Pi^s,$$

$$M(w^r, w^{r'}) = \begin{cases} \displaystyle\sum_{p\in\mathcal{P}(w^r,w^{r'})} \frac{1}{\mu((w^{r'})^{(\ell_r-p)})} & \text{if } r \neq r', \\ 0 & \text{if } r = r', \end{cases}$$

$$T_1(w^r, w^{r'}) = (2n - \ell_r - \ell_{r'} + 2)\mu(w^r)\widetilde{\mu}(w^{r'})\left(\frac{\Omega_{r'}(w^{r'}_{\ell_{r'}}, w^r_1)}{\mu(w^r_1)} + M(w^{r'}, w^r)\right),$$

$$T_2(w^r, w^{r'}) = (n - \ell_r + 1)\left((\ell-1)(\widetilde{\mu}(w^r)\mu(w^{r'}) + \mu(w^r)\widetilde{\mu}(w^{r'})) + (6\ell-5)\widetilde{\mu}(w^r)\widetilde{\mu}(w^{r'})\right),$$

$$\begin{aligned} T_3(w^r, w^{r'}) = {} & (n - \ell_r + 1)\mu(w^r)\mu(w^{r'})\left(\frac{\Omega_{r,r'}(w^r_{\ell_r}, w^{r'}_1)}{\mu(w^{r'}_1)} + \frac{\Omega_{r,r'}(w^{r'}_{\ell_{r'}}, w^r_1)}{\mu(w^r_1)}\right) \\ & + \frac{(n - \ell_r + 1)(6\ell - 3\ell_r - 3\ell_{r'} + 2)}{\mu_{\min}}\widetilde{\mu}(w^r)\widetilde{\mu}(w^{r'}) \\ & + \frac{(n - \ell_r + 1)(\ell - 2)}{\mu_{\min}}\left(\mu(w^r)\widetilde{\mu}(w^{r'}) + \mu(w^{r'})\widetilde{\mu}(w^r)\right) \\ & + (n - \ell_r + 1)\mu(w^r)\mu(w^{r'})\left(M(w^r, w^{r'}) + M(w^{r'}, w^r)\right), \end{aligned}$$ (24)

$$\begin{aligned} \gamma_1(\ell_r, \ell, a) = {} & \sum_{x,y\in\mathcal{A}} \mu(x) \max_{b\in\mathcal{A}} \left| \frac{1}{\mu(b)} \sum_{(t,t')\neq(1,1)} \frac{\alpha_t^{2\ell-\ell_r}\alpha_{t'}^{2\ell-\ell_r}}{\alpha^\ell} Q_t(x,b)Q_{t'}(a,y) \right. \\ & \left. - \sum_{t=2}^{|\mathcal{A}|} \frac{\alpha_t^{4\ell-2}}{\alpha^\ell} Q_t(x,y) \right|, \end{aligned}$$

$$\begin{aligned} \gamma_2(\ell_r, \ell) = {} & \sum_{x,y\in\mathcal{A}} \mu(x) \max_{a,b\in\mathcal{A}} \left( \frac{1}{\mu(b)} \sum_{(t,t')\neq(1,1)} \left| \frac{\alpha_t^{2\ell-\ell_r}\alpha_{t'}^{2\ell-\ell_r}}{\alpha^\ell} Q_t(x,b)Q_{t'}(a,y) \right| \right. \\ & \left. + \sum_{t=2}^{|\mathcal{A}|} \left| \frac{\alpha_t^{5\ell-3}}{\alpha^\ell} Q_t(x,y) \right| \right). \end{aligned}$$

Here we choose as index set $I = \{1, 2, \ldots, q(n+1) - \sum_{s=1}^q \ell_s\}$; it can be written as the disjoint union $I = \bigcup_{r=1}^q I_r$ with

$$I_r = \left\{(r-1)(n+1) - \sum_{s=1}^{r-1} \ell_s + 1, \ldots, r(n+1) - \sum_{s=1}^r \ell_s\right\}.$$ (25)

We define $[i]$ by

$$[i] := i - (r-1)(n+1) + \sum_{s=1}^{r-1} \ell_s \quad \text{with } r = r(i) \text{ such that } i \in I_r.$$ (26)

To apply Theorem 13, the Bernoulli process $\underline{\widetilde{Y}} = (\widetilde{Y}_i)_{i\in I}$ and the Poisson process $\underline{Z} = (Z_i)_{i\in I}$ are given

32

by

$$\begin{aligned}
\widetilde{\mathbb{Y}}_i &= \widetilde{Y}_{[i]}(w^r), \\
Z_i &\sim \mathrm{Po}(\widetilde{\mu}(w^r)),
\end{aligned} \qquad (27)$$

where $r$ is such that $i \in I_r$. For $i \in I$, we choose the neighborhood $B_i := \{j \in I : |[j] - [i]| \leq 3\ell - 3\}$. Moreover, let $\ell = \max\{\ell_r, r = 1, \dots, q\}$ denote the maximal length of the words, and let $\ell_{\min}$ denote the minimal length of the words.

Then Reinert and Schbath (1998) prove the following results.

**Theorem 9** *Under assumption (A1) and with the notation (24) and (27), we have*

$$d_{TV}\left(\mathcal{L}(\underline{\widetilde{Y}}), \mathcal{L}(\underline{Z})\right)$$

$$\leq (n - \ell_{min} + 1)(6\ell - 5)\left(\sum_{r=1}^{q} \widetilde{\mu}(w^r)\right)^2 + \sum_{1 \leq r, r' \leq q} T_1(w^r, w^{r'})$$

$$+ |\alpha|^\ell \sum_{r=1}^{q} \gamma_1(\ell_r, \ell, w_{\ell_r}^r)(n - \ell_r + 1)\widetilde{\mu}(w^r).$$

**Corollary 3** *Let $(Z_r)_{r=1,\dots,m}$ be independent Poisson variables with $\mathbb{E}Z_r = (n - \ell_r + 1)\widetilde{\mu}(w^r)$. With the previous notation and under assumption (A1), we have*

$$d_{TV}\left(\mathcal{L}\big((\widetilde{N}(w^r))_{r=1,\dots,q}\big), \mathcal{L}\big((Z_r)_{r=1,\dots,q}\big)\right)$$

$$\leq (n - \ell_{min} + 1)(6\ell - 5)\left(\sum_{r=1}^{q} \widetilde{\mu}(w^r)\right)^2 + \sum_{1 \leq r, r' \leq q} T_1(w^r, w^{r'})$$

$$+ |\alpha|^\ell \sum_{r=1}^{q} \gamma_1(\ell_r, \ell, w_{\ell_r}^r)(n - \ell_r + 1)\widetilde{\mu}(w^r) + \sum_{r=1}^{q}(\ell_r - 1)\big(\mu(w^r) - \widetilde{\mu}(w^r)\big).$$

The proof is a direct application of Theorem 13, similar as in Section 5.

Moreover, in a similar way a compound Poisson approximation for the numbers of occurrences can be obtained. Choose as index set

$$I = \left\{1, 2, \dots, q(n+1) - \sum_{s=1}^{q} \ell_s\right\} \times \{1, 2, \dots\}.$$

To apply Theorem 13, the Bernoulli process $\underline{\widetilde{\mathbb{Y}}} = (\widetilde{\mathbb{Y}}_{i,k})_{(i,k) \in I}$ and the Poisson process $\underline{\mathbb{Z}} = (Z_{i,k})_{(i,k) \in I}$ are now defined as

$$\begin{aligned}
\widetilde{\mathbb{Y}}_{i,k} &= \widetilde{Y}_{[i],k}(w^r), \\
Z_{i,k} &\sim \mathrm{Po}(\widetilde{\mu}_k(w^r)),
\end{aligned}$$

where $r = r(i)$ is such that $i \in I_r$; $I_r$ and $[i]$ are given by (25) and (26). For $(i, k) \in I$, the neighborhood is still $B_{i,k} := \{(j, k') \in I : -(k' + 3)(\ell - 1) \leq [j] - [i] \leq (k + 3)(\ell - 1)\}$.
We make the following weak assumption on the overlap structure.

(A2)     $\forall r \neq r'$, $w^r$ is not a substring of any composed word in $\mathcal{C}_2(w^{r'})$.

**Theorem 10** *Under assumptions (A1), (A2) and with the notation (24), we have*

$$d_{TV}\left(\mathcal{L}(\underline{\widetilde{\mathbb{Y}}}), \mathcal{L}(\underline{\mathbb{Z}})\right) \leq \sum_{1 \leq r, r' \leq q} T_2(w^r, w^{r'}) + \sum_{1 \leq r, r' \leq q} T_3(w^r, w^{r'}) + |\alpha|^\ell \sum_{r=1}^{q} \gamma_2(\ell_r, \ell)(n - \ell_r + 1)\widetilde{\mu}(w^r).$$

33

Moreover the following corollary is easily obtained.

**Corollary 4** *Let $(Z_k)_{k \geq 1}$ be independent Poisson variables with expectation $\mathbb{E}Z_k = \sum_{r=1}^{q}(n - \ell_r + 1)\widetilde{\mu}_k(w^r)$; CP denotes the compound Poisson distribution of $\sum_{k \geq 1} kZ_k$. Under assumptions (A1), (A2) and with the notation (24), we have*

$$
d_{TV}\left(\mathcal{L}\left(\sum_{r=1}^{q} N(w^r)\right), \text{CP}\right) \leq \sum_{1 \leq r, r' \leq q} T_2(w^r, w^{r'}) + \sum_{1 \leq r, r' \leq q} T_3(w^r, w^{r'})
$$

$$
+ |\alpha|^\ell \sum_{r=1}^{q} \gamma_2(\ell_r, \ell)(n - \ell_r + 1)\widetilde{\mu}(w^r)
$$

$$
+ 2 \sum_{r=1}^{q} (\ell_r - 1)\big(\widetilde{\mu}(w^r) - \mu(w^r)\big).
$$

Erhardsson (1997) derived a compound Poisson approximation for the special case of counting only the total number of overlapping occurrences of words from the word family. As he uses a Markov chain coupling rather than the local approach, his bound on the error has a slightly different flavor.

**Poisson approximation for the renewal count distribution** Related results for renewal counts are available. For a Poisson approximation, the problem can be reduced to declumped counts, like in the case of a single word. As in Tanushev (1996) we consider non-overlapping occurrences in competition with each other. For example, in the sequence `CGTATATTAAAAATATTAGA`, the set of words `TAT`, `TTA` and `AA` has renewal occurrences of `TAT` at position 3 and 14, of `TTA` at position 7, and of `AA` at positions 10 and 12. The occurrences of `TAT` at position 5, of `TTA` at position 16, and of `AA` at positions 9 and 11 are not counted because they overlap with some already counted words.

Let

$$
\mathbb{I}_i^c(w^r) = \mathbb{I}\{\text{a competing renewal of } w^r \text{ starts at position } i \text{ in } X_1 \cdots X_n\},
$$

and let

$$
R_n^c(w^r) = \sum_{i=1}^{n-\ell_r+1} \mathbb{I}_i^c(w^r)
$$

be the number of competing renewals of $w^r$ in the sequence $X_1 X_2 \cdots X_n$. For a Poisson process approximation (and, following from that, a Poisson approximation for the counts), we want to assess $\mathbb{P}(\mathbb{I}_i^c(w^r) \neq \widetilde{Y}_i(w^r))$. First consider $\mathbb{P}(\mathbb{I}_i^c(w^r) = 1, \widetilde{Y}_i(w^r) = 0)$. Note that, from (22), for $i > \ell_r$, to have $\mathbb{I}_i^c(w^r) = 1, \widetilde{Y}_i(w^r) = 0$, there must be an occurrence of $w^r$ at position $i$, and this occurrence cannot be the start of a clump of $w^r$, so that there must be an overlapping occurrence of $w^r$ at some position $j = i - \ell_r + 1, \ldots, i - 1$. Moreover, this occurrence cannot be a competing renewal, so there must be another word $w^{r'}$ overlapping this occurrence. Hence we may bound

$$
\mathbb{P}(\mathbb{I}_i^c(w^r) = 1, \widetilde{Y}_i(w^r) = 0) \leq \mu^2(w^r) \sum_{p \in \mathcal{P}(w^r)} \frac{1}{\mu((w^r)^{(\ell_r - p)})} \sum_{r'=1}^{q} \mu(w^{r'}) M(w^{r'}, w^r).
$$

For $i \leq \ell_r$ the above bound is still valid (the probability is even smaller since not always would there be enough space for these clumps to occur). Secondly, consider $\mathbb{P}(\mathbb{I}_i^c(w^r) = 0, \widetilde{Y}_i(w^r) = 1)$. For $\mathbb{I}_i^c(w^r) = 0, \widetilde{Y}_i(w^r) = 1$ to occur, there must be an occurrence of $w^r$ at position $i$, overlapped by an occurrence of a different word $w^{r'}$, so that we may bound

$$
\mathbb{P}(\mathbb{I}_i^c(w^r) = 0, \widetilde{Y}_i(w^r) = 1) \leq \mu(w^r) \sum_{r'=1}^{q} \mu(w^{r'}) M(w^{r'}, w^r).
$$

34

Again, for $i \leq \ell_r$ the above bound remains valid. Thus we have

$$\mathbb{P}(\underline{\mathbb{I}}^c(w^r) \neq \underline{\widetilde{Y}}(w^r)) \leq (n - \ell_r + 1)\mu(w^r) \sum_{r'=1}^{q} \mu(w^{r'})M(w^{r'}, w^r)$$

$$\left(1 + \mu(w^r) \sum_{p \in \mathcal{P}(w^r)} \frac{1}{\mu((w^r)^{(\ell_r - p)})}\right).$$

Hence

$$\mathbb{P}(\underline{\mathbb{I}}^c \neq \underline{\widetilde{Y}}) \leq \sum_{r=1}^{q}(n - \ell_r + 1)\mu(w^r) \sum_{r'=1}^{q} \mu(w^{r'})M(w^{r'}, w^r)$$

$$\left(1 + \mu(w^r) \sum_{p \in \mathcal{P}(w^r)} \frac{1}{\mu((w^r)^{(\ell_r - p)})}\right).$$

Thus we obtain as a corollary of Theorem 9

**Corollary 5** *Under assumption (A1) and with the notation (24) and (27), we have*

$$d_{TV}\left(\mathcal{L}(\underline{\mathbb{I}}^c), \mathcal{L}(\underline{Z})\right) \leq (n - \ell_{min} + 1)(6\ell - 5)\left(\sum_{r=1}^{q} \widetilde{\mu}(w^r)\right)^2 + \sum_{1 \leq r, r' \leq q} T_1(w^r, w^{r'})$$

$$+ |\alpha|^{\ell} \sum_{r=1}^{q} \gamma_1(\ell_r, \ell, w_{\ell_r}^r)(n - \ell_r + 1)\widetilde{\mu}(w^r)$$

$$+ \sum_{r=1}^{q}(n - \ell_r + 1)\mu(w^r) \sum_{r'=1}^{q} \mu(w^{r'})M(w^{r'}, w^r)$$

$$\left(1 + \mu(w^r) \sum_{p \in \mathcal{P}(w^r)} \frac{1}{\mu((w^r)^{(\ell_r - p)})}\right).$$

Note that the order of the approximation is the same as in Theorem 9; the additional error terms are comparable to $T_1$ and $T_2$, respectively. A Poisson approximation for the competing renewal counts follows immediately.

**Gaussian approximation for the joint distribution of competing renewal counts**    A multivariate normal approximation has been obtained by Tanushev (1996). The main problem here is to specify the covariance structure. To state the result, quite a bit of notation is needed. For a matrix $A$ denote its transposed matrix by $A^T$, and, if $A$ is a square matrix, $Diag(A)$ represents the vector of the diagonal elements of $A$. Define the probabilities of ending a word for $1 \leq j \leq \ell_r - 1$ as

$$\mathbb{P}_r(j) = \mathbb{P}(\text{ collect final } j \text{ letters of } w^r | \text{ start with correct } \ell_r - j \text{ initial}$$
$$\text{letters of } w^r)$$

$$= \frac{\mu(w^r)}{\mu((w^r)^{(\ell_r - j)})}.$$

Then the overlap-matching polynomials are defined as

$$Q_{r,r'}(z) = \sum_{p \in \mathcal{P}(w^r, w^{r'}) \cup \{0\}} z^p \mathbb{P}_{r'}(p).$$

Define the $q \times q$ matrix

$$\Delta(z) = (Q_{r,r'}(z))_{r,r'=1,\ldots,q}$$

35

and

$$\Lambda(z) = (\Delta^{-1})(z))^T$$
$$\Lambda = \Lambda(1).$$

Furthermore denote

$$\widetilde{K}_r(z) = z^{\ell_r - 1} \mathbb{P}_r(\ell_r - 1)$$

and the vector

$$\widetilde{K}(z) = (\widetilde{K}_1(z), \ldots, \widetilde{K}_q(z))^T.$$

Denote by

$$Diag(\widetilde{K}(z))$$

the $q \times q$ diagonal matrix with the components of $\widetilde{K}(z)$ as diagonal elements. Put

$$\widetilde{K} = \widetilde{K}(1).$$

Moreover put $K_r = \mu(w_1^r)\mathbb{P}_r(\ell_r - 1)$ and define the vector

$$K = (K_1, \ldots, K_q)^T.$$

Put

$$H(z) = \frac{d}{dz}\Lambda(z)$$
$$H = H(1).$$

Define the vector

$$L = (\ell_1 K_1, \ldots, \ell_q K_q).$$

and the matrix

$$\widetilde{Z} = Z_{[\psi]},$$

where $I$ is the identity matrix, $Z$ is defined in (21), and for a matrix $A$ the matrix $A_{[\psi]}$ is the $q \times q$ matrix whose $(r, r')$ entry is the element of $A$ at the row corresponding to the last letter $w_{\ell_r}^r$ of the word $w^r$, and at the columns corresponding to the first letter $w_1^{r'}$ of $w^{r'}$. Define the variance-covariance matrix

$$C = \frac{1}{2}\left(\Lambda K(\Lambda K - 2HK - 2\Lambda L)^T + (\Lambda K - 2HK - 2\Lambda L)(\Lambda K)^T\right)$$
$$+ Diag(\Lambda K)\widetilde{Z}Diag(\widetilde{K})\Lambda^T + \Lambda Diag(\widetilde{K})\widetilde{Z}^T Diag(\Lambda K) + Diag(\Lambda K).$$

Define the mean $\mu_R^c(w^r)$ by

$$(\mu_R^c(w^1), \ldots, \mu_R^c(w^q))^T = \Lambda K.$$

Denote weak convergence by $\overset{w}{\Rightarrow}$. Now we have all the ingredients to state the normal approximation obtained by Tanushev (1996).

**Theorem 11** *Under Assumption (A1) we have*

$$\left(\frac{R_n^c(w^r) - n\mu_R^c(w^r)}{\sqrt{n}}\right)_{r=1,\ldots,q} \overset{w}{\Rightarrow} \mathcal{N}(0, C).$$

In the case of a single pattern, this theorem reduces to Theorem 8.

# 8    Sequencing by Hybridization

As a slightly more involved example of how statistics and probability on words are applied in DNA sequence analysis, we describe a problem related to sequencing by hybridization; see Arratia *et al.* (1996) for more details. Sequencing by hybridization is a tool to determine a DNA sequence from the unordered list of all $\ell$-tuples contained in this sequence; typical numbers for $\ell$ are $\ell = 8, 10, 12$. It is based on the fact that DNA nucleotides bind or hybridize with each other: A and T hybridize, and A and G hybridize. For example, the sequence TGTGTGAGTG hybridizes with ACACACTCAC. In a sequencing chip, all $4^\ell$ possible oligonucleotides ("probes") of length $\ell$ are attached to the surface of a substrate, each fragment at a distinct location.

To use an SBH chip, the single-stranded target DNA is amplified, labeled by a fluorescent, and exposed to the sequencing chip. The probes on the chip will hybridize to a copy of the single-stranded target DNA if the substring complementary to the probe exists in the target. These probes are then detected with a spectroscopic detector. For example, if $\ell = 4$, the sequence TGTGTGAGTG will hybridize to the probes ACAC, ACTC, CACA, CACT, CTCA and TCAC.

As chips can be washed and used again, and due to automatization, this method is not only fast but also inexpensive. Originally developed to sequence DNA, as to date a major use has been to compare different DNA strings in order to detect mutations. In particular, SBH chips are employed in the analysis of HIV blood samples, to decide whether a virus in the sample is a known form of HIV or a new mutation. Also such arrays are used in gene expression studies.

There are still technical difficulties in producing an error-free chip; moreover the SBH image may be difficult to read. However, even if these sources of errors are eliminated, a main drawback of the SBH procedure is that more than one sequence may produce the same SBH data. For example, if $\ell = 4$, the sequence ACACTCACAC will hybridize to the same probes as the sequence ACACACTCAC.

To control this error resulting from non-unique recoverability, we are interested in an estimate for the probability that a sequence is uniquely recoverable, that is, the sequence is unambiguous. This probability will depend on the probe length $\ell$, on the length $n$ of the target sequence, and on the frequencies of the different nucleotides, A, C, G and T, in the sequence. Furthermore we need to bound the error made in estimating the probability of unique recoverability in order to make assertions about the reliability of the chip.

As a simplification, we assume that we not only know the set of all $\ell$-tuples in the sequence but also their multiplicity (but not the order in which they occur). This multiset is called the $\ell$-spectrum of the sequence. In the sequel, unique recoverability is understood to mean unique recoverability of a sequence from its $\ell$-spectrum.

Pevzner (1989) characterizes unique recoverability from the $\ell$-spectrum using the de Bruijn-graph (see van Lint and Wilson (1992)) whose vertices are the $(\ell - 1)$-tuples in the sequence. Two vertices $v$ and $w$ are joined by a directed edge from $v$ to $w$ if the $\ell$-spectrum contains an $\ell$-tuple for which the first $(\ell - 1)$ nucleotides coincide with $v$ and the last $(\ell - 1)$ nucleotides coincide with $w$. Pevzner (1989) showed that a sequence is uniquely recoverable from its $\ell$-spectrum if and only if there is a unique (Eulerian) path connecting all the vertices. Ukkonen (1992) conjectured and Pevzner (1995) proved that there are exactly three structures that prevent unique recoverability:

**1. Rotation.** The sequence starts and ends with the same $(\ell - 1)$-tuple. In this case, the de Bruijn-graph is a cycle, and any vertex could be chosen as starting point.

**2. Transposition with a three-way repeat.** If an $(\ell - 1)$-tuple occurs three times in the sequence, then the de Bruijn-graph has two loops at this vertex, and the order in which these loops are passed is not fixed.

**3. Transposition with two interleaved pairs of repeats.** There are two "interleaved" pairs of $(\ell - 1)$-tuple repeats, i.e. in the de Bruijn-graph there are two vertices $x$ and $y$ connected by a path of the form $\ldots c \ldots y \ldots x \ldots y \ldots$, where we described a path connecting all the vertices by listing the vertices in the order they are used in the path. This implies that there are two ways of going from $x$ to $y$ in the graph.

**Example 1** *The sequence ACACACTCAC possesses as 4-spectrum the multiset {ACAC, ACAC, CACA, CACT, ACTC, CTCA, TCAC}. The competing sequence ACACTCACAC has the same 4-spectrum. The de Bruijn-graph for the sequence  ACACACTCAC has as vertices ACA, CAC, ACT, CTC and TCA. There are two directed edges from ACA to CAC, and one directed edge each from CAC to ACA, from CAC to ACT, from ACT to CTC, from CTC to TCA, and from TCA to CAC. The competing sequence ACACTCACAC has the same de Bruijn-graph.*

*For the sequence ACACACTCAC, a path connecting all vertices is*

$$ACA, CAC, ACA, (CAC, ACT, CTC, TCA), CAC.$$

*The alternate path*

$$ACA, (CAC, ACT, CTC, TCA), CAC, ACA, CAC$$

*also connecting all the vertices, corresponds to the sequence, ACACTCACAC, with the same 4-spectrum.*

Thus unique recoverability can be described in terms of possibly overlapping repeats of $(\ell - 1)$-tuples within a single sequence. We model DNA as a random sequence $X_1 X_2 \ldots X_n$, where $X_1, X_2, \ldots$ are independent and identically distributed over the alphabet $\{A, C, G, T\}$. For a sequence to be uniquely recoverable, the event of an $(\ell - 1)$-tuple repeat should be rare. This implies that we consider the occurrence of $(\ell - 1)$-tuples under a Poisson regime. (Note that we are interested in the configuration in which the repeats occur; hence we need a Poisson process approximation for the process of repeats rather than a Poisson approximation for the number of repeats.) If repeats are rare, then three-way repeats are negligible, and so is the probability that a sequence starts and ends with the same $(\ell - 1)$-tuple. After bounding these probabilities, we thus restrict our attention to interleaved pairs of repeats. Under the Poisson regime, if there are $k$ pairs of repeats, then the occurrences of these repeats are discrete uniform. Additional randomization makes the position of the repeats continuously uniform, so that all orderings of these pairs will be approximately equally likely. This allows the application of a combinatorial argument using Catalan numbers to obtain that the number of interleaved pairs of repeats, if $k$ repeats are present, is approximately $2^k/(k + 1)!$. If $\lambda$ is the expected number of repeats of $\ell$-tuples in a single sequence, we hence get, for the probability $P_\ell$ that $X_1 X_2 \ldots X_n$ is uniquely recoverable from its $\ell$-spectrum,

$$P_\ell \quad \approx \quad e^{-\lambda} \sum_{k \geq 0} \frac{(2\lambda)^k}{k!(k + 1)!}.$$

The Chen-Stein method for Poisson approximation (Theorem 13) provides explicit bounds for the error terms in this approximation, as follows.

In the sequence $X_1 \ldots X_n$ of independent identically distributed letters, let $p$ be the probability that two random letters match. We write $t$ for $\ell - 1$, as we are interested in $(\ell - 1)$-repeats. Again we have to declump: Define $Y_{i,i} = 0$ for all $i$, and

$$Y_{i,j} = \begin{cases} \mathbb{I}\{X_1 \cdots X_t = X_{j+1} \cdots X_{j+t}\} & \text{if } i = 0 \\ (1 - \mathbb{I}\{X_i = X_j\})\mathbb{I}\{X_{i+1} \cdots X_{i+t} = X_{j+1} \cdots X_{j+t}\} & \text{otherwise.} \end{cases}$$

Thus $Y_{i,j} = 1$ if and only if there is a leftmost repeat starting after $i$ and $j$. Put $I = \{(i, j), 1 \leq i, j \leq n - \ell + 1\}$. A careful analysis (see Arratia *et al.* (1996)) yields that the process $\underline{Y} = (Y_\alpha)_{\alpha \in I}$ is sufficient to decide whether a sequence is uniquely recoverable from its $\ell$-spectrum (although $\underline{Y}$ contains strictly less information than the process of indicators of occurrences).

For a Poisson process approximation, we first identify the expected number $\lambda$ of leftmost repeats. If $\alpha = (i, j)$ does not have self-overlap, that is, if $j - i > t$, then

$$\mathbb{E}(Y_\alpha) = \begin{cases} p^t & \text{if } i = 0 \\ (1 - p)p^t & \text{otherwise.} \end{cases}$$

Hence the expected number $\lambda^*$ of repeats without self-overlap is

$$\lambda^* = \binom{n - 2t}{2}(1 - p)p^t + (n - 2t)p^t.$$

If $\alpha$ does have self-overlap, then, in order to have a leftmost repeat at $\alpha$, for indices in the overlapping set, two matches are required, and for indices in the non-overlapping set, one match is required. Let $d = j - i$; then $E(Y_\alpha)$ depends on the decomposition of $t + d$ into a quotient $q$ of $d$ and a remainder $r$ (such that $t + d = qd + r$): if $p_q$ is the probability that $q$ random letters match, then

$$\mathbb{E}(Y_\alpha) = \begin{cases} p_{q+1}^r p_q^{d-r} & \text{if } i = 0 \\ (p_q - p_{q+1})^r p_q^{d-r} & \text{otherwise.} \end{cases}$$

If $\lambda^*$ is bounded away from 0 and infinity, which corresponds to having $t = 2log_{1/p}(n) + c$ for some constant $c$, then it can be seen that

$$\lambda \approx \frac{n^2}{2}(1-p)p^t.$$

Under the regime that $\lambda$ is bounded away from 0 and infinity, here is a short version of the general result in Arratia $et\ al.$ (1996). Let $\xi_* = \max_a(\mathbb{P}(X_i = a))$ be the probability of the most likely letter.

**Theorem 12** *Let $\underline{Z} \equiv (Z_\alpha)_{\alpha \in I}$ be a process with independent Poisson distributed coordinates $Y_\alpha$, with $\mathbb{E}Z_\alpha = \mathbb{E}Y_\alpha, \alpha \in I$. Then*

$$d_{TV}(\underline{Y}, \underline{Z}) \leq b(n, t) \sim \left\{ \begin{array}{ll} 16\lambda^2 \frac{t}{n} & \textit{in the uniform case} \\ n\xi_*^t & \textit{in the nonuniform case.} \end{array} \right.$$

In Arratia $et\ al.$ (1996), a more general result is derived for general alphabets, and explicit bounds are obtained. These bounds can be used to approximate the probability of unique recoverability. Recently Arratia $et\ al.$ (1999) have obtained results on the number of possible reconstructions for a given sequence (when the reconstruction is not unique).

# 9 Some probabilistic and statistical tools

## 9.1 The Chen-Stein method

The Chen-Stein method is a powerful tool for deriving Poisson approximations and compound Poisson approximations in terms of bounds on the total variation distance. For any two random processes $\underline{Y}$ and $\underline{Z}$ with values in the same space $E$, the total variation distance between their probability distributions is defined by

$$\begin{aligned} d_{\mathrm{TV}}(\mathcal{L}(\underline{Y}), \mathcal{L}(\underline{Z})) &= \sup_{B \subset E, \text{measurable}} |\mathbb{P}(\underline{Y} \in B) - \mathbb{P}(\underline{Z} \in B)| \\ &= \sup_{h: E \to [0,1], \text{measurable}} |\mathbb{E}h(\underline{Y}) - \mathbb{E}h(\underline{Z})|. \end{aligned}$$

First published by Chen (1975) as the Poisson analog to Stein's method for normal approximations (Stein (1972)), it has found widespread application; word counts being just one of them. A friendly exposition is found in Arratia $et\ al.$ (1989) and a description with many examples can be found in Arratia $et\ al.$ (1990) and Barbour $et\ al.$ (1992a). The key theorem for word counts in stationary Markov chains is Theorem 1 in Arratia $et\ al.$ (1990) with an improved bound by Barbour $et\ al.$ (1992a) (Theorem 1.A and Theorem 10.A), giving the following theorem.

**Theorem 13** *Let $I$ be an index set. For each $\alpha \in I$, let $Y_\alpha$ be a Bernoulli random variable with $p_\alpha = \mathbb{P}(Y_\alpha = 1) > 0$. Suppose that, for each $\alpha \in I$, we have chosen $B_\alpha \subset I$ with $\alpha \in B_\alpha$. Let $Z_\alpha, \alpha \in I$, be independent Poisson variables with mean $p_\alpha$. The total variation distance between the dependent Bernoulli process $\underline{Y} = (Y_\alpha, \alpha \in I)$ and the Poisson process $\underline{Z} = (Z_\alpha, \alpha \in I)$ satisfies*

$$d_{TV}(\mathcal{L}(\underline{Y}), \mathcal{L}(\underline{Z})) \leq b_1 + b_2 + b_3,$$

*where*

$$b_1 = \sum_{\alpha \in I} \sum_{\beta \in B_\alpha} p_\alpha p_\beta \tag{28}$$

$$b_2 = \sum_{\alpha \in I} \sum_{\beta \in B_\alpha, \beta \neq \alpha} \mathbb{E}(Y_\alpha Y_\beta) \tag{29}$$

$$b_3 = \sum_{\alpha \in I} \mathbb{E}|\mathbb{E}\{Y_\alpha - p_\alpha | \sigma(Y_\beta, \beta \notin B_\alpha)\}|. \tag{30}$$

*Moreover, if $W = \sum_{\alpha \in I} Y_\alpha$ and $\lambda = \sum_{\alpha \in I} p_\alpha < \infty$, then*

$$d_{TV}(\mathcal{L}(W), \mathrm{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda}(b_1 + b_2) + \min\left(1, \sqrt{\frac{1}{\lambda e}}\right) b_3.$$

Note that $b_3 = 0$ if $Y_\alpha$ is independent of $\sigma(Y_\beta, \beta \notin B_\alpha)$. We think of $B_\alpha$ as a neighborhood of strong dependence of $Y_\alpha$.

One consequence of this theorem is that for any indicator of an event, i.e. for any measurable functional $h$ from $E$ to $[0,1]$, there is an error bound of the form $|\mathbb{E}h(\underline{Y}) - \mathbb{E}h(\underline{Z})| \le d_{\mathrm{TV}}(\mathcal{L}(\underline{Y}), \mathcal{L}(\underline{Z}))$. Thus, if $T(\underline{Y})$ is a test statistic then, for all $t \in \mathbb{R}$,

$$|\mathbb{P}\left(T(\underline{Y}) \ge t\right) - \mathbb{P}\left(T(\underline{Z}) \ge t\right)| \le b_1 + b_2 + b_3\,,$$

which can be used to construct confidence intervals and to find p-values for tests based on this statistic.

Note that this method can also be used to prove compound Poisson approximations. For multivariate compound Poisson approximations it is very convenient. For univariate compound Poisson approximations, better bounds are at hand, see Roos (1993), Barbour and Utev (1998), Barbour $et\ al.$ (1992b).

## 9.2 Stein's method for normal approximations

Stein's method for the normal approximation was first published by Stein (1972). Recently Rinott and Rotar (1996) applied it to obtain multivariate normal approximations with a bound on the error in the distance of suprema over convex sets.

Let $\mathcal{H}$ denote the class of indicator functions of convex sets in $\mathbf{R}^d$. Let $Y_j$ be random vectors taking values in $\mathbf{R}^d$, and let $W = \sum_{j=1}^{n} Y_j$ the vector of sums. Assume there is a constant $B$ such that $|Y_j| := \sum_{i=1}^{d} |Y_{(j,i)}| \le B$.

**Theorem 14** Let $\mathcal{S}_i$ and $\mathcal{N}_i$ be subsets of $\{1, \ldots, n\}$, such that $i \in \mathcal{S}_i \subset \mathcal{N}_i, i = 1, \ldots, n$. Assume that there exist constants $D_1 \le D_2$ such that

$$\max\{|\mathcal{S}_i|, i = 1, \ldots, n\} \le D_1$$

and

$$\max\{|\mathcal{N}_i|, i = 1, \ldots, n\} \le D_2,$$

where for sets $|\cdot|$ denotes cardinality.

Then, for $d = 1$ there exists a universal constant $c$ such that

$$\sup_{h \in \mathcal{H}} |\mathbb{E}h(W) - \Phi h| \le c\{aD_2 B + n(a + \sqrt{\mathbb{E}W^2})D_1 D_2 B^3 + \chi_1 + \chi_2 + \chi_3\}.$$

For $d \ge 1$ there exists a constant $c$ depending only on the dimension $d$ such that

$$
\begin{aligned}
\sup_{h \in \mathcal{H}} |\mathbb{E}h(W) - \Phi h| \quad \le \quad & c\{aD_2 B + naD_1 D_2 B^3(|\log B| + \log n) \\
& + \chi_1 + (|\log B| + \log n)(\chi_2 + \chi_3)\}.
\end{aligned}
$$

Here $a = 2\sqrt{d}$ and

$$
\begin{aligned}
\chi_1 \quad &= \quad \sum_{j=1}^{n} \mathbb{E}|\mathbb{E}(Y_j \mid \sum_{k \notin \mathcal{S}_j} Y_k)| \\
\chi_2 \quad &= \quad \sum_{j=1}^{n} \mathbb{E}|\mathbb{E}(Y_j\,(\sum_{k \in \mathcal{S}_|} Y_k)^T) - \mathbb{E}(Y_j\,(\sum_{k \in \mathcal{S}_|} Y_k)^T| \sum_{l \notin \mathcal{N}_j} Y_l)|) \\
\chi_3 \quad &= \quad |I - \sum_{j=1}^{n} \mathbb{E}(Y_j\,(\sum_{k \in \mathcal{S}_|} Y_k)^T)|.
\end{aligned}
$$

## 9.3 Moment-generating function

Here is a short outline on moment-generating functions; see, e.g., Rice (1995). The *moment-generating function $M$* of a random variable $X$ is defined as

$$\Phi_X(t) = \mathbb{E}(e^{tX}).$$

So, if $X$ has a discrete distribution $p$, we have that

$$\Phi_X(t) = \sum_x e^{tx} p(x).$$

If the moment-generating function exists for all $t$ in an open interval containing zero, it uniquely determines the probability distribution.

In particular, under regularity conditions the moments of a random variable can be obtained via the moment-generating function using differentiation. Namely, if $\Phi_X(t)$ is finite, we have

$$\Phi'_X(t) = \frac{d}{dt}\mathbb{E}(e^{tX}) = \mathbb{E}(Xe^{tX}).$$

Thus

$$\Phi'_X(0) = \mathbb{E}X$$

if both sides of the equation exist. Similarly, differentiating $r$ times we obtain

$$\Phi_X^{(r)}(0) = \mathbb{E}(X^r).$$

A special case is when the moment-generating function $\Phi_X(t)$ is rational, that is, when $\Phi_X(t)$ can be written as

$$\Phi_X(t) = \frac{p_0 + p_1 t + \ldots + p_r t^r}{q_0 + q_1 t + \ldots + q_s t^s} = \sum_d f(d) t^d,$$

for some $r, s$ and coefficients $p_1, \ldots, p_r, q_1, \ldots, q_s$. By normalization we may assume $q_0 = 1$. Then

$$p_0 + p_1 t + \ldots + p_r t^r = \sum_d f(d) t^d (1 + q_1 t + \ldots + q_s t^s).$$

Identification of the coefficients of $t^i$ on both sides yields

$$p_i = \sum_{d=0}^{i} f(d) q_{i-d} \text{ for } i \leq r$$

$$0 = \sum_{d=0}^{i} f(d) q_{i-d} \text{ for } i > r.$$

This gives a recurrence formula for the coefficients $f(d)$; we have

$$f(0) = p_0$$

$$f(d) = p_d - \sum_{i=1}^{\min(d,s)} f(d-i) q_i, \quad d \geq 1$$

where $p_d = 0$ for $d > r$.

## 9.4   The $\delta$-method

In general, the $\delta$-method, or *propagation of error*, is a linear approximation (Taylor expansion) of a nonlinear function of random variables. Here we are particularly interested in the validity of a normal approximation for functions of random vectors (see, e.g., Rice (1995)). The following theorem can be found on p.313 in Waterman (1995).

**Theorem 15** *Let $\underline{X}_n = (X_{n1}, X_{n2}, \ldots, X_{nk})$ be a sequence of random vectors satisfying*

$$b_n(\underline{X}_n - \underline{\mu}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma)$$

41

with $b_n \to \infty$. The vector valued function $\underline{g}(\underline{x}) = (g_1(\underline{x}), \ldots, g_\ell(\underline{x}))$ has real valued $g_i(\underline{x})$ with non-zero differential

$$\frac{\partial g_i}{\partial g_{\underline{x}}} = \left( \frac{\partial g_i}{\partial g_{x_1}}, \ldots, \frac{\partial g_i}{\partial g_{x_k}} \right).$$

Define $\mathbf{D} = (d_{i,j})$ where $d_{i,j} = \frac{\partial g_i}{\partial g_{x_j}}(\underline{\mu})$. Then

$$b_n(g(\underline{X}_n) - g(\underline{\mu})) \xrightarrow{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{D\Sigma D^T}).$$

## 9.5 A Large Deviation Principle

Assume $X_1 \cdots X_n$ is an irreducible Markov chain on a finite alphabet $\mathcal{A}$ with transition probabilities $\pi(a, b)$, $a, b \in \mathcal{A}$. The following theorem for Markov chain can be found on p.78 in Bucklew (1990).

**Theorem 16 (Miller)** *Let $f$ be a function mapping $\mathcal{A}$ into $\mathbb{R}$. Then, $n^{-1} \sum_{i=1}^{n} f(X_i)$ obeys a large deviation principle with rate function $I$ defined below: for every closed subset $F \subset \mathbb{R}$ and every open subset $O \subset \mathbb{R}$,*

$$\limsup_{n \to +\infty} \frac{1}{n} \log \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^{n} f(X_i) \in F \right) \leq - \inf_{x \in F} I(x),$$

$$\liminf_{n \to +\infty} \frac{1}{n} \log \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^{n} f(X_i) \in O \right) \geq - \inf_{x \in O} I(x).$$

*The rate function $I$ is positive, convex, uniquely equal to zero at $x = \mathbb{E}f(X_1)$ and given by*

$$I(x) = \sup_{\theta}(\theta x - \log \lambda(\theta))$$

*where $\lambda(\theta)$ is the largest eigenvalue of the matrix $\left( e^{\theta f(b)} \pi(a, b) \right)_{a,b \in \mathcal{A}}$.*

## 9.6 A CLT for martingales

The following theorem can be found in Dacunha-Castelle and Duflo (1983) p.80.

**Theorem 17** *Let $(\xi_{n,i})_{i=1,\ldots,n}$ be a triangular array of $d$-dimensional random vectors such that $\mathbb{E}\|\xi_{n,i}\|_2^2 < \infty$, and $V$ be a positive $d \times d$ matrix. Put $\mathcal{F}_{n,i} = \sigma(\xi_{n,1}, \ldots, \xi_{n,i})$; $\mathbb{E}(\xi_{n,i} \mid \mathcal{F}_{n,i-1})$ denotes the conditional expectation vector of $\xi_{n,i}$ and $Cov(\xi_{n,i} \mid \mathcal{F}_{n,i-1})$ denotes the conditional covariance matrix of $\xi_{n,i}$. If as $n \to \infty$*

*(i)* $\sum_{i=1}^{n} \mathbb{E}(\xi_{n,i} \mid \mathcal{F}_{n,i-1}) \xrightarrow{\mathbb{P}} 0$,

*(ii)* $\sum_{i=1}^{n} Cov(\xi_{n,i} \mid \mathcal{F}_{n,i-1}) \to V$,

*(iii)* $\forall \varepsilon > 0$, $\sum_{i=1}^{n} \mathbb{P}(|\xi_{n,i}| > \varepsilon \mid \mathcal{F}_{n,i-1}) \xrightarrow{\mathbb{P}} 0$,

*then $\sum_{i=1}^{n} \xi_{n,i} \xrightarrow{\mathcal{D}} \mathcal{N}(0, V)$.*

# References

Apostolico, A., Bock, M.E., and Xuyan, X.. 1998. Annotated statistical indices for sequence analysis. *Proceedings of Compression and Complexity of SEQUENCES 97*, 215–229. IEEE Computer Society Press.

Arratia, R., Bollobas, B., Coppersmith, C. and Sorkin, G. 1999. Euler circuits and DNA sequencing by hybridization. Preprint.

Arratia, R., Goldstein, L. and Gordon, L. 1989. Two moments suffice for Poisson approximations: the Chen-Stein method. *Ann. Prob.*, **17**, 9–25.

Arratia, R., Goldstein, L. and Gordon, L. 1990. Poisson approximation and the Chen-Stein method. *Statistical Science*, **5**, 403–434.

Arratia, R., Martin, D., Reinert, G., and Waterman, M.S. 1996. Poisson Approximation for long repeats in a random sequence with Application to sequencing by hybridization. *J. Comp. Biol.*, **3**, 425–463.

Barbour, A. D. and Xia, A. 1999. Poisson Perturbations. Preprint.

Barbour, A. D., Holst, L. and Janson, S. 1992a. *Poisson approximation.* Oxford - University Press.

Barbour, A. D., Chen, L. H. Y. and Loh, W.-L. 1992b. Compound Poisson approximation for nonnegative random variables via Stein's method. *Ann. Prob.*, **20**, 1843–1866.

Barbour, A. D. and Utev, U. 1998. Solving the Stein equation in compound Poisson approximation. *Adv. Appl. Prob.*, **30**, 449–475.

Biaudet, V., El Karoui, M., and Gruss, A. 1998. Codon usage can explain GT-rich islands surrounding Chi sites on the *Escherichia coli* genome. *Mol. Microbiol.*, **29**, 666–669.

Biggins, J., and Cannings, C. 1987. Markov renewal processes, counters and repeated sequences in Markov chains. *Ann. Appl. Prob.*, 19, 521–545.

Blom, G. and Thorburn, D. 1982. How many random digits are required until given sequences are obtained? *J. Appl. Prob.*, **19**, 518–531.

Breen, S., Waterman, M.S., and Zhang, N. 1985. Renewal theory for several patterns. *J. Appl. Prob.*, **22**, 228–234.

Brendel, V., Beckmann, J. S. and Trifonov, E. N. 1986. Linguistics of nucleotide sequences : Morphology and comparison of vocabularies. *J. Biomol. Struct. Dynamics*, **4**, 11-21.

Bucklew, J. A. 1990. *Large Deviation Techniques in Decision, Simulation, and Estimation.* Wiley.

Chedin, F., Noirot, P., Biaudet, V. and Ehrlich, S. 1998. A five-nucleotide sequence protects DNA from exonucleolytic degradation by AddAB, the RecBCD analogue of *Bacillus subtilis*. *Mol. Microbiol.*, **31**, 1369–1377.

Chen, L. H. Y. 1975. Poisson approximation for dependent trials. *Ann. Prob.*, **3**, 534–545.

Chryssaphinou, O. and Papastavridis, S. 1988a. A limit theorem for the number of non-overlapping occurrences of a pattern in a sequence of independent trials. *J. Appl. Prob.*, **25**, 428–431.

Chryssaphinou, O. and Papastavridis, S. 1988b. A limit theorem on the number of overlapping appearances of a pattern in a sequence of independent trials. *Prob. Theory Rel. Fields*, **79**, 129–143.

Chung, K.L. 1974. *A Course in Probability Theory.* $2^{nd}$ ed. Academic Press.

Churchill, G.A. 1989. Stochastic models for heterogeneous DNA sequences. *Bull. Math. Biol.*, **51**, 79–94.

Cowan, R. 1991. Expected frequencies of DNA patterns using Whittle's formula. *J. Appl. Prob.*, **28**, 886–892.

Dacunha-Castelle, D. and Duflo, M. 1983. *Probabilités et Statistiques 2.Problèmes à Temps Mobile* Masson.

Deheuvels, P., Devroye, L. and Lynch, J. 1986. Exact convergence rate in the limit theorems of Erdös-Rnyi and Shepp. *Ann. Prob.*, **14**, 209–223.

Dembo, A. and Karlin, S. 1992. Poisson Approximations for $r$-scan processes. *Ann. Appl. Prob.*, **2**, 329–357.

Durbin, R., Eddy, S., Krogh, A. and Mitchison, G. 1998. *Biological Sequence Analysis.* Cambridge University Press.

Eichelsbacher, P. and Roos, M. 1999. Compound Poisson Approximation for dissociated random variables via Stein's method. *To appear in Combinatorics, Probability and Computing.*

Erdös, P. and Rényi, A. 1970. On a new law of large numbers. *J. Analyse Math.*, **23**, 103–111.

Erhardsson, T. 1997. *Compound Poisson Ann. Prob. proximation for Markov chains.* Ph.D. thesis, Royal Institute of Technology, Stockholm.

Fu, J.C. 1993. Poisson convergence in reliability of a large linearly connected system as related to coin tossing. *Statistica Sinica*, **3**, 261–275.

Gentleman, J. 1994. The distribution of the frequency of subsequences in alphabetic sequences, as exemplified by deoxyribonucleic acid.*Appl. Statist.*, **43**, 404–414.

Gentleman, J., and Mullin, R. 1989. The distribution of the frequency of occurrence of nucleotide subsequences, based on their overlap capability. *Biometrics*, **45**, 35–52.

Geske, M. X., Godbole, A. P., Schaffner, A. A., Skolnick, A. M. and Wallstrom, G. L. 1995. Compound Poisson approximations for word patterns under Markovian hypotheses. *J. Appl. Prob.*, **32**, 877–892.

Godbole, A.P. 1991. Poisson Approximations for runs and patterns of rare events. *Adv. Appl. Prob.*, **23**, 851–865.

Godbole, A.P., and Schaffner, A.A. 1993. Improved Poisson approximations for word patterns. *Adv. Appl. Prob.*, **25**, 334–347.

Guibas, L. J. and Odlyzko, A. M. 1980. Long repetitive patterns in random sequences. *Z. Wahrscheinlichkeithsth.*, **53**, 241–262.

Gusfield, D. 1997. *Algorithms on Strings, Trees, and Sequences.* Cambridge University Press.

Hirano, K. and Aki, S. 1993. On number of occurrences of sucess runs of specified length in a two-state Markov chain. *Statistica Sinica*, **3**, 313–320.

Karlin, S. and Burge, C. and Campbell, A. M. 1992. Statistical analyses of counts and distributions of restriction sites in DNA sequences. *Nucl. Acids Res.*, **20**, 1363–1370.

Karlin, S., and Macken, C. 1991. Assessment of inhomogeneities in an *E. coli* physical map. *Nucl. Acids Res.*. **19**, 4241–4246.

Karlin, S. and Taylor, H. M. 1975. *A first course in stochastic processes.* $2^{nd}$ ed. Academic Press.

Kleffe, J., and Borodovsky, M. 1992. First and second moment of counts of words in random texts generated by Markov chains. *Comp. Applic. Biosci.*, **8**, 433-441.

Kleffe, J., and Langbecker, U. 1990. . Exact computation of pattern probabilities in random sequences generated by Markov chains. *Comp. Applic. Biosci.*, **6**, 347–353.

Leung, M. Y. Marsh, G. M. and Speed, T. P. 1996. Over and underrepresentation of short DNA words in Herpesvirus genomes. *J. Comp. Biol.*, **3**, 345–360.

Li, S. Y. 1980. A martingale approach to the study of occurrence of sequence patterns in repeated experiments. *Ann. Prob.*, **8**, 1171-1176.

Lothaire, M. 1983. *Combinatorics on words.* Addison-Wesley.

Lundstrom, R. 1990. *Stochastic models and statistical methods for DNA sequence data.* Ph.D. Thesis, Department of Mathematics, University of Utah.

Muri, F. 1998. "Modelling Bacterial Genomes using Hidden Markov Models. In *Compstat'98 Proceedings in Computational Statistics* (eds Payne, R. and Green, P.J.), 89–100. Physica-Verlag, Heildelberg.

Pevzner. P.A. 1989. l-tuple DNA Sequencing: Computer Analysis. *J. Biomol. Struct. Dynamics*, **7**, 63–73.

Pevzner, P.A. 1995. DNA physical mapping and alternating Eulerian cycles in colored graphs. *Algorithmica*, **13**, 77–105.

Prum, B., Rodolphe, F. and Turckheim, É. de 1995. Finding words with unexpected frequencies in DNA sequences. *J. R. Statist. Soc. B*, **57,** 205–220.

Rajarshi, M.B. 1974. Success runs in a two-state Markov chain. *J. Appl. Prob.*, **11**, 190–192.

Régnier, M. 1998. A unified approach to word statistics", In *Proceedings of the Second Annual International Conference on Computational Molecular Biology, RECOMB*, 207–231.

Régnier, M. 1998. A unified approach to word occurrence probabilities", To appear in Discrete Applied Mathematics, Special Issue on Computational Biology.

Reinert, G., and Schbath, S. 1998. Compound Poisson and Poisson process approximations for occurrences of multiple words in Markov chains. *J. Comp. Biol.*, **5**, 223–253.

Rice, J. 1995. *Mathematical Statistics and Data Analysis.* Duxbury Press.

Rinott, Y., and Rotar, V. 1996. On coupling constructions and rates in the CLT for dependent summands with applications to the antivoter model and weighted $U$-statistics. *J. Multivariate Analysis*, **56**, 333–350.

Robin, S., and Daudin, J.-J. 1999. Exact distribution of word occurrences in a random sequence of letters. *J. Appl. Prob.*, **36**, 179–193.

Rocha, E.P.C. and Viari, A. and Danchin, A. 1998. Oligonucleotide bias in *Bacillus subtilis*: general trends and taxonomic comparisons. *Nucl. Acids Res.*, **26**, 2971–2980.

Roos, M. 1993. *Stein-Chen method for compound Poisson approximation.* Ph.D. thesis, University of Zurich.

Roos, M., and Stark, D. 1999. Poisson approximation for mixing sequences using coupling. Preprint.

Schbath, S. 1995a. Compound Poisson approximation of word counts in DNA sequences. *ESAIM: Probability and Statistics.* **1** 1–16.

Schbath, S. 1995b. *Étude asymptotique du nombre d'occurrences d'un mot dans une chaîne de Markov et application á la recherche de mots de fréquence exceptionnelle dans les séquences d'ADN.* PhD thesis, Université René Descartes, Paris V.

Senoussi, R. 1990. Statistique asymptotique presque-sûre de modèles statistiques convexes, *Ann. Inst. Henri Poincaré*, **26**, 19–44.

Sourice, S., Biaudet, V., El Karoui, M., Ehrlich, S.D., and Gruss, A. 1998. Identification of the Chi site of *Haemophilus influenzae* as several sequences related to the *Escherichia coli* Chi site. *Mol. Microbiol.*, **27**, 1021–1029.

Stein, C. 1972. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proc. Sixth Berkeley Symp. Math. Statist. Probab.*, Vol.2, 583–602. Univ. California Press, Berkeley.

Tanushev, M. 1996. Central limit theorem for several patterns in a Markov chain sequence of letters. Preprint.

Ukkonen, E. 1992. Approximate string-matching with q-grams and maximal matches. *Theoret. Comp. Science,*, **92**, 191–211.

Van Lint, J.H., and Wilson, R.M. 1992. *A Course in Combinatorics.* Cambridge University Press.

Waterman, M. S. (995. *Introduction to Computational Biology.* Chapman & Hall.