

# Determining the distribution of word matches between Markovian sequences

Conrad Burden, Sylvain Forêt, \*Paul Leopardi

Mathematical Sciences Institute, Australian National University.

For presentation at CTAC, QUT, 2012.

Based on a presentation given by Conrad Burden at COMPSTAT Cyprus.

25 September 2012



AUSTRALIAN RESEARCH COUNCIL  
Centre of Excellence for Mathematics  
and Statistics of Complex Systems



# Acknowledgements

Sue Wilson (Australian National University, University of New South Wales).

Australian Research Council grant DP120101422.

## Definition of $D_2$

Given two sequences from a finite alphabet

$$A := (A_1, A_2, \dots, A_m) \text{ and } B := (B_1, B_2, \dots, B_n),$$

$D_2$  is the number of matches of words (including overlaps) of prespecified length  $k$  between two given sequences.

## Definition of $D_2$

Given two sequences from a finite alphabet

$$A := (A_1, A_2, \dots, A_m) \text{ and } B := (B_1, B_2, \dots, B_n),$$

$D_2$  is the number of matches of words (including overlaps) of prespecified length  $k$  between two given sequences.

Example: consider these two sequences and  $k = 7 \dots$

**A:** ATGCTTTGCTAGCGCTATGCTTTTCGCAAACATCAT

**B:** ATGCTTTTAAAACCGAGCTGGTCAGCGCTAAGCGCT

## Definition of $D_2$

Given two sequences from a finite alphabet

$$A := (A_1, A_2, \dots, A_m) \text{ and } B := (B_1, B_2, \dots, B_n),$$

$D_2$  is the number of matches of words (including overlaps) of prespecified length  $k$  between two given sequences.

Example: consider these two sequences and  $k = 7 \dots$

**A:** ATGCTTTGCTAGCGCTATGCTTTTCGCAAACATCAT

**B:** ATGCTTTTAAAACCGAGCTGGTCAGCGCTAAGCGCT

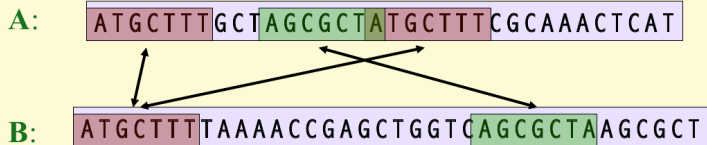
## Definition of $D_2$

Given two sequences from a finite alphabet

$$A := (A_1, A_2, \dots, A_m) \text{ and } B := (B_1, B_2, \dots, B_n),$$

$D_2$  is the number of matches of words (including overlaps) of prespecified length  $k$  between two given sequences.

Example: consider these two sequences and  $k = 7 \dots$



In this example, for  $k = 7$ ,  $D_2 = 3$ .

# Our previous results

(with Sue Wilson, Ruth Kantorovitz, and Junmei Jing)

We have exact formulas for  $\mathbf{E}(D_2)$  (simple) and  $\mathbf{Var}(D_2)$  (complicated), and accurate approximations for the distribution of  $D_2$  for the case of sequences composed of i.i.d. letters; i.e. each letter drawn independently from the same distribution.

# Periodic boundary conditions

To simplify the calculations (avoiding 'edge effects'), we imposed periodic boundary conditions:



ATGCTTTGCTAGCGCTATGCTTTCGCAAACATCAT



ATGCTTTTAAACCGAGCTGGTCAGCGCTAAGCGCT



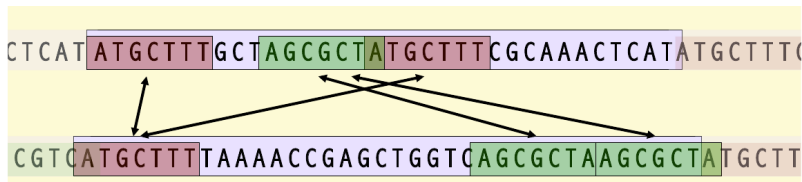
# Periodic boundary conditions

To simplify the calculations (avoiding 'edge effects'), we imposed periodic boundary conditions:

```
CTCAT ATGCTTTGCTAGCGCTATGCTTTCGCAAAC TCATA TGCTTT
CGTC ATGCTTTTAAACCGAGCTGGTCAGCGCTAAGCGCTATGCTT
```

# Periodic boundary conditions

To simplify the calculations (avoiding 'edge effects'), we imposed periodic boundary conditions:



Now, for  $k = 7$  we have  $D_2 = 4$ .

# Markovian sequences

The i.i.d. is not realistic. Real DNA sequences are more realistically modelled as Markovian (up to fifth order).

For first order:

$$\text{Prob}(A_{i+1} = u \mid A_i = v) = M_{u,v},$$
$$u, v \in \{A, C, G, T\}$$

where

$$0 \leq M_{u,v} \leq 1; \quad \sum_v M_{u,v} = 1.$$

# Markovian sequences

The i.i.d. is not realistic. Real DNA sequences are more realistically modelled as Markovian (up to fifth order).

For first order:

$$\text{Prob}(A_{i+1} = u \mid A_i = v) = M_{u,v},$$
$$u, v \in \{A, C, G, T\}$$

where

$$0 \leq M_{u,v} \leq 1; \quad \sum_v M_{u,v} = 1.$$

First problem:

How do we define periodic boundary conditions for a Markovian sequence?

# Markov chain with periodic boundary conditions

Define a Markov chain

$$\dots X_{n-1}, X_n, X_1, X_2, \dots, X_n, X_1, X_2, \dots$$

with periodic boundary conditions (PBCs) via the following algorithm:

1. Choose  $X_1$  from any distribution  $\pi(u)$ ,  $u \in \{1, \dots, d\}$ , where  $0 \leq \pi(u) \leq 1$ ;  $\sum_u \pi(u) = 1$ . Thus  $\Pr(X_1 = u) = \pi(u)$ .
2. Choose  $X_2, \dots, X_{n+1}$  via the Markov matrix  $M$ ,  $\Pr(X_{i+1} = v \mid X_i = u) = M_{u,v}$ ,  $i = 1, \dots, n$ .
3. If  $X_{n+1} = X_1$ , accept  $X_1, X_2, \dots, X_n$ , otherwise return to Step 1 and repeat the procedure.

# No privileged starting point

We further wish to restrict the definition to repeating Markov chains with no privileged starting point, by which we mean

$$\Pr(\mathbf{X} = \mathbf{x}) = \Pr(\mathbf{X} = (\mathbf{x}_{i+1} \dots \mathbf{x}_n, \mathbf{x}_1 \dots \mathbf{x}_i)),$$

for all  $i = 1, \dots, n - 1$ ,

where  $\mathbf{X} = (X_1 X_2 \dots X_n)$ .

# No privileged starting point

We further wish to restrict the definition to repeating Markov chains with no privileged starting point, by which we mean

$$\Pr(\mathbf{X} = \mathbf{x}) = \Pr(\mathbf{X} = (\mathbf{x}_{i+1} \dots \mathbf{x}_n, \mathbf{x}_1 \dots \mathbf{x}_i)),$$

for all  $i = 1, \dots, n - 1$ ,

where  $\mathbf{X} = (X_1 X_2 \dots X_n)$ .

## Theorem 1

*$\mathbf{X}$  has no privileged starting point if and only if  $\pi(\mathbf{u})$  is a uniform distribution:  $\pi(\mathbf{u}) = 1/d, \mathbf{u} = 1, \dots, d$ .*

# Probability of a specific sequence

## Corollary 2

*If  $X$  is a Markov chain with no privileged starting point, the probability of any given sequence  $x = (x_1 x_2 \dots x_n)$  is*

$$\Pr(X = x) = \frac{M_{x_1, x_2} M_{x_2, x_3} \dots M_{x_n, x_1}}{\text{tr}(M^n)}$$



# Mean of $D_2$

For two sequences  $A$  and  $B$  of length  $m$  and  $n$ , both generated using the matrix  $M$ , and word length  $k$ ,

$$\mathbf{E}(D_2) = \frac{mn \operatorname{tr} [(M^{m-k+1} \circ M^{n-k+1})(M \circ M)^{k-1}]}{\operatorname{tr}(M^m) \operatorname{tr}(M^n)},$$

where  $\circ$  indicates the Hadamard product of matrices

$$(P \circ Q)_{r,s} = P_{r,s} Q_{r,s}.$$

## Mean of $D_2$

Given two sequences

$A = (A_1, A_2, \dots, A_m)$  and  $B = (B_1, B_2, \dots, B_n)$ ,

define the word-match indicator

$$I_{i,j} = \begin{cases} 1 & \text{if } k\text{-word at position } i \text{ in } A \text{ matches} \\ & k\text{-word at position } j \text{ in } B, \\ 0 & \text{otherwise.} \end{cases}$$

Then

$$D_2 = \sum_{i=1}^m \sum_{j=1}^n I_{i,j}$$

and

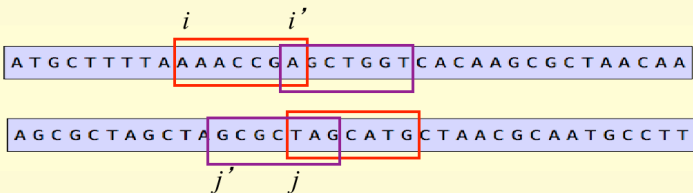
$$\mathbf{E}(D_2) = \sum_{i=1}^m \sum_{j=1}^n \mathbf{E}(I_{i,j}) = \sum_{i=1}^m \sum_{j=1}^n \Pr(I_{i,j} = 1).$$

## Variance of $D_2$

The variance of  $D_2$  is much harder but can be done, at least for Markov order 1:

$$\begin{aligned}\text{Var}(D_2) &= \text{Var}\left(\sum_{i,j} I_{i,j}\right) = \mathbf{E}\left(\left(\sum_{i,j} I_{i,j}\right)^2\right) - \left(\mathbf{E}\left(\sum_{i,j} I_{i,j}\right)\right)^2 \\ &= \left(\sum_{i,j,i',j'} \mathbf{E}(I_{i,j}, I_{i',j'})\right) - \mathbf{E}(D_2)^2.\end{aligned}$$

The difficult part is  $\mathbf{E}(I_{i,j}, I_{i',j'})$ , the probability of word matches like this:



# Variance of $D_2$

The formula for  $\text{Var}(D_2)$  with periodic boundary conditions and Markov order 1 is complicated ...

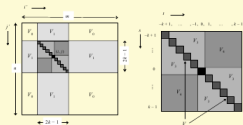


Figure 2: Contributions to  $\text{Var}(D_2)$  via the sum in Eq. (24). The left-hand diagram shows the  $(i, j)$  plane for a fixed value of  $(i, j)$ , shown as the black square. The right-hand diagram is an expanded view of the 'accident' region  $-k+1 \leq k, j \leq k-1$ , where  $i = i' + a$  and  $j = j' + b$  in FBCs.

and  $j'$  in sequence  $Y$ ,

$$\begin{aligned} E(D_2^2) &= \sum_{i', j'} \sum_{i, j} E(I_{ij} I_{i'j'}) \\ &= \sum_{i', j'} \sum_{i, j} \text{Prob}(I_{ij} = 1, I_{i'j'} = 1) \\ &= V_1 + V_2 + V_3 + V_4, \end{aligned} \quad (24)$$

The partitioning reflects the degree of overlap between words in each of the two sequences, and is illustrated in Fig. 2. We assume  $a, n \geq 2k$ , which will almost certainly be the case in any biological application.

We will write a Hadamard product of  $q$  factors,  $M \circ \dots \circ M$ , using the shorthand notation  $M^{\circ q}$ . With this notation, the contributions to the variance are:

$$V_1 = \frac{\text{tr}(M^{\circ 2})}{\text{tr}(M^{\circ 2}) \text{tr}(M^{\circ 2})} \times \sum_{i=1}^{n-2k+1} \sum_{j=1}^{n-2k+1} \left[ \text{tr}(M^{i+1} \circ M^{i+1}) \text{tr}(M \circ M)^{i-1} \times (M^{n-2k+i+1} \circ M^{n-2k+i+1}) \text{tr}(M \circ M)^{i-1} \right], \quad (25)$$

$$\begin{aligned} V_2 &= \frac{\text{tr}(M^{\circ 2})}{\text{tr}(M^{\circ 2}) \text{tr}(M^{\circ 2})} \times \left\{ \sum_{i=1}^{n-2k} \left[ \text{tr} \left( (M \circ M \circ M)^{i-1} \circ (M^{i+1})^i \right) \times (M^{n-k+i} \circ M^{n-k+i+1}) \right] \right. \\ &\quad + 2 \sum_{i=1}^{k-1} \left[ \text{tr} \left( (M \circ M)^i \right) \times \left. \left[ \text{tr} \left( (M \circ M \circ M)^{k-i-1} \circ (M^{i+1})^i \right) \times (M \circ M)^i \left( M^{n-k+i+1} \circ M^{n-2k+i+1} \right) \right] \right] \right\} \\ &\quad + \text{the same with } n \text{ and } n \text{ interchanged}, \end{aligned} \quad (26)$$

$$V_3 = \frac{\text{tr}(M^{\circ 2})}{\text{tr}(M^{\circ 2}) \text{tr}(M^{\circ 2})} \times \left\{ \text{tr} \left( M^{n-k+1} \circ M^{n-k+1} \right) \text{tr} \left( M \circ M \right)^{n-1} \right. \\ \left. + 2 \sum_{i=1}^{k-1} \left[ \text{tr} \left( M^{n-k+i+1} \circ M^{n-k+i+1} \right) \times \text{tr} \left( M \circ M \right)^{i-1} \right] \right\}, \quad (27)$$

$$V_4 = \frac{2\text{tr}(M^{\circ 2})}{\text{tr}(M^{\circ 2}) \text{tr}(M^{\circ 2})} \times \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} \left[ \text{tr} \left( (M \circ M)^i \circ (M \circ M)^j \right) \times (M^{n-k+i+1} \circ M^{n-k+i+1} + M^{n-k+i+1} \circ M^{n-k+i+1}) \right], \quad (28)$$

where

$$Q = \begin{cases} (M^{2k-2i-1})^{i-1} \circ (M^{2k-2i})^{j-i-1} \circ (M^{2k-2i})^j & \text{if } i > j, \\ (M^{2k-2i-1})^{j-1} \circ (M^{2k-2i})^i & \text{if } i = j, \\ (M^{2k-2i-1})^{i-1} \circ (M^{2k-2i})^i & \text{if } i < j, \end{cases} \quad (29)$$

and

$$v = \frac{k-i}{r+1}, \quad \rho = (k-i) \bmod (r-i), \quad (30)$$

Finally,

$$V_4 = \frac{2\text{tr}(M^{\circ 2})}{\text{tr}(M^{\circ 2}) \text{tr}(M^{\circ 2})} \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} v U, \quad (31)$$

where

$$U = \begin{cases} (M^{2k-2i-1})^{j-1} \circ (M^{n-k+i+1})^j & \text{if } i < j, \\ (M^{2k-2i-1})^{i-1} \circ (M^{n-k+i+1})^i & \text{if } i = j, \\ (M^{2k-2i-1})^{i-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^{j-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^{i-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^{j-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^{j-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^{i-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^{j-1} \circ (M^{n-k+i+1})^i & \text{if } i < j, \\ (M^{2k-2i-1})^i & \text{if } i < j, \end{cases} \quad (32)$$

... but is easily evaluated.

# Verification by simulation

1. For a given order 1 Markov matrix, generate 10,000 random pairs of Markovian sequences with periodic boundary conditions (R scripts).
2. Obtain the value of  $D_2$  for each pair (SAFT program, written in C).
3. Compare empirical cumulative distribution function of  $D_2$  with that of Normal and Pólya-Aeppli (compound Poisson) distributions using theoretical  $E(D_2)$  and  $\text{Var}(D_2)$  (R scripts).

# Generation of sequences

Sequences are generated by using the algorithm shown previously.

Seemingly more efficient algorithms, such as continuing with  $X_{n+2}, X_{n+3}, \dots$  until  $X_k$  matches  $X_{n+k}$  can yield different distributions.

## Calculating $D_2$

The SAFT program was written to compare a given sequence against a database of sequences, by calculating  $D_2$  for each pair.

The program was adapted to take two lists of sequences  $\mathcal{A}, \mathcal{B}$ , and determine  $D_2$  for each corresponding pair  $\mathcal{A}_k, \mathcal{B}_k$ .

# Comparing against known distributions

The R statistical system was used to

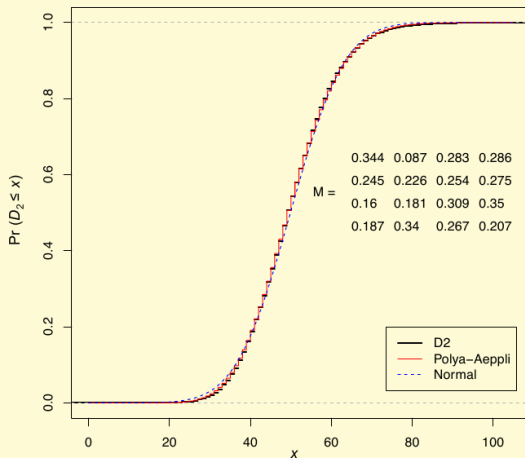
- ▶ Use the `t.test` function in R to compare the empirical and theoretical mean values.
- ▶ Use the `qchisq` function in R to produce a confidence interval for the variance, and test the variance using this interval.
- ▶ Use the `ks.test` function in R to compare the empirical distribution of  $D_2$  against the Normal and the Pólya-Aepli distributions, using the theoretical mean and variance. The R function `pPolyaAepli` was written for this purpose.
- ▶ Plot cumulative distributions.



# Results for one random Markov matrix

Randomly chosen Markov matrix  $M$

$n = 100$   $k = 4$



# Results for one random Markov matrix

