

Alignment-free comparison of biological sequences

Conrad Burden, Sylvain Forêt, *Paul Leopardi

Mathematical Sciences Institute, Australian National University.

For presentation at ANZIAM, Newcastle, 2013.

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

4 February 2013



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

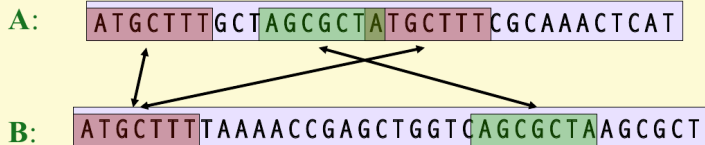
Definition of D_2

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 $k = 7$, $D_2 = 3$.

Markovian sequences

Real DNA sequences are modelled as Markovian.

For first order:

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

where

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

Periodic boundary conditions

To simplify the calculations of theoretical mean and variance (avoiding 'edge effects'), we impose periodic boundary conditions:



ATGCTTTGCTAGCGCTATGCTTTCGCAAACATCAT



ATGCTTTTAAAACCGAGCTGGTCAGCGCTAAGCGCT

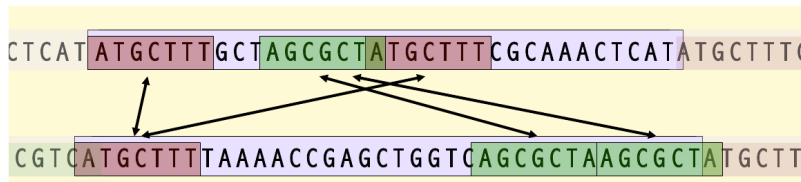
Periodic boundary conditions

To simplify the calculations of theoretical mean and variance (avoiding 'edge effects'), we impose periodic boundary conditions:

```
CTCAT ATGCTTTGCTAGCGCTATGCTTTCGCAAAC TCATA TGCTTT
CGTC ATGCTTTTAAAACCGAGCTGGTCAGCGCTAAGCGCTA TGCTT
```

Periodic boundary conditions

To simplify the calculations of theoretical mean and variance (avoiding 'edge effects'), we impose periodic boundary conditions:



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

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

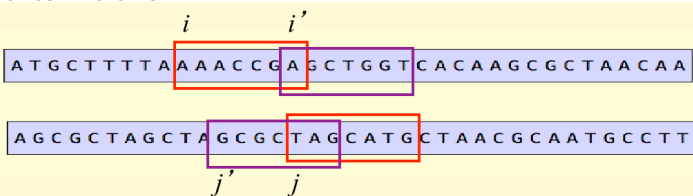
$$\mathbf{E}(D_2) = \sum_{i=1}^m \sum_{j=1}^n \mathbf{E}(I_{i,j}) = \sum_{i=1}^m \sum_{j=1}^n \Pr(I_{i,j} = 1).$$

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 . . .

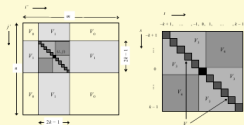


Figure 2: Contributions to $\text{Var}(D_2)$ via the sum in Eq. (24). The left-hand diagram shows the (i, j) plane for a fixed value of (i, j) , shown as the black square. The right-hand diagram is an expanded view of the 'accident' region $-k+1 \leq k \leq k-1$, where $i = i' + a$ and $j = j' + b$ in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

and j' in sequence \mathbf{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:

$$V_1 = \frac{\text{tr}(M^{\circ 2k})}{\text{tr}(M^{\circ n}) \text{tr}(M^{\circ n})} \times \sum_{a=0}^{n-2k} \sum_{b=0}^{n-2k} \left[\text{tr}(M^{a+1} \circ M^{b+1}) \text{tr}(M \circ M)^{k-1} \times (M^{n-2k-a-1} \circ M^{n-2k-b-1}) \text{tr}(M \circ M)^{k-1} \right]. \quad (25)$$

$$\begin{aligned} V_2 &= \frac{\text{tr}(M^{\circ 2k})}{\text{tr}(M^{\circ n}) \text{tr}(M^{\circ n})} \times \left\{ \sum_{a=0}^{n-2k} \left[\text{tr} \left((M \circ M \circ M)^{k-1} \circ (M^{a+1})^2 \right) \times (M^{n-k-1} \circ M^{n-2k-a-1}) \right] \right. \\ &\quad + 2 \sum_{a=0}^{k-1} \left[\text{tr} \left(M \circ M \right)^k \times \left. \left[(M \circ M \circ M)^{k-1-a} \text{tr}(M^{a+1})^2 \right] \times (M \circ M)^{k-1-a} \circ M^{n-2k-a-1} \right] \right\} \\ &\quad + \text{the same with } n \text{ and } n \text{ interchanged.} \end{aligned} \quad (26)$$

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

$$V_4 = \frac{2\text{tr}(M^{\circ 2k})}{\text{tr}(M^{\circ n}) \text{tr}(M^{\circ n})} \times \sum_{a=0}^{k-1} \sum_{b=0}^{k-1} \left[\text{tr} \left((M \circ M)^a \text{tr}(M \circ M)^b \times (M^{n-k+1-a} \circ M^{n-k+1-b} \circ M^{n-k+1-a} \circ M^{n-k+1-b}) \right) \right], \quad (28)$$

where

$$\rho = \begin{cases} (M^{2k-2b})^{a-1} \text{tr}(M^{2k-2b})^{a-1} & \text{if } \rho > 0, \\ (M^{2k-2b})^{a-1} \text{tr}(M^{2k-2b})^a & \text{if } \rho = 0, \end{cases} \quad (29)$$

and

$$\psi = \frac{k-a}{r+1}, \quad \rho = (k-x) \bmod (r-x). \quad (30)$$

Finally,

$$V_4 = \frac{2\text{tr}(M^{\circ 2k})}{\text{tr}(M^{\circ n}) \text{tr}(M^{\circ n})} \sum_{a=0}^{k-1} U_a. \quad (31)$$

where

$$U_a = \begin{cases} (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^{a-1} M^{2k-2b} \times \left\{ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \right\} M^{2k-2b} & \text{if } \zeta = 0, \\ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \times \left\{ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \right\} \times (M^{2k-2b})^{a-1} & \text{if } 0 < \zeta \leq a, \\ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \times \left\{ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \right\} \times (M^{2k-2b})^{a-1} & \text{if } a < \zeta \leq r, \\ \text{(as above with } n \text{ and } n \text{ interchanged and } r \text{ and } r \text{ interchanged)} & \text{if } r < \zeta \leq r, \\ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \times \left\{ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \right\} \times (M^{2k-2b})^{a-1} & \text{if } r < \zeta \leq r, \\ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \times \left\{ (M^{2k-2b})^{a-1} \text{tr}(M^{n-k+1})^a \right\} \times (M^{2k-2b})^{a-1} & \text{if } r < \zeta \leq r, \end{cases} \quad (32)$$

. . . but is easily evaluated.

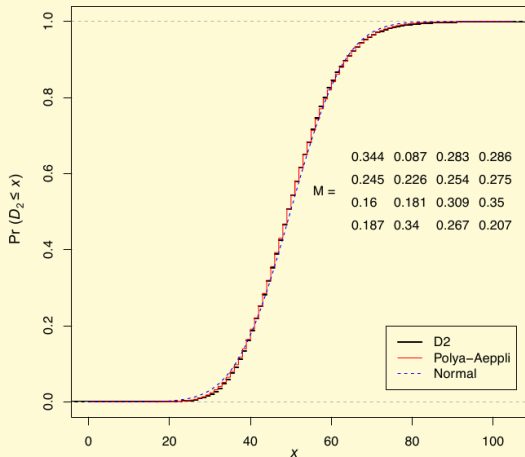
Verification by simulation

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).

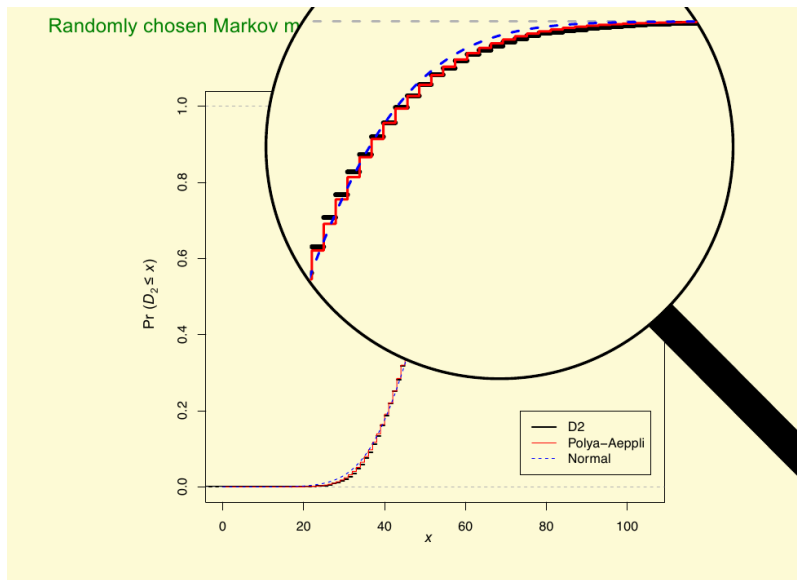
Results for a random Markov matrix

Randomly chosen Markov matrix M

$n = 100$ $k = 4$



Results for a random Markov matrix



DNA is messy

Real DNA is messy

messy with repeats of different length and complexity,
and contains unknown regions.

- ▶ The **Ensembl** database marks unknown regions and masks repetitive regions including *tandem repeats*.
- ▶ The **tantan** program masks simple repeats.

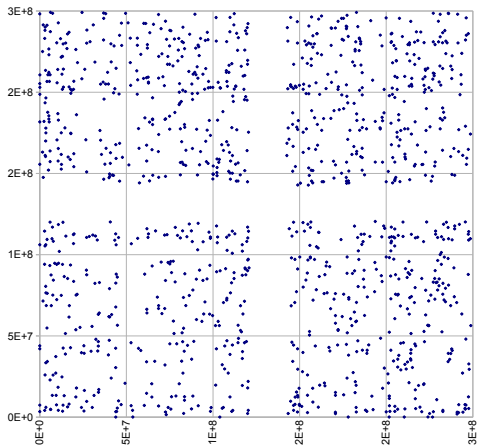
(Ensembl: Wellcome Trust Sanger Institute and European Bioinformatics Institute, 2012; tantan: Firth, 2011)

To compare D_2 from DNA with Markov models:

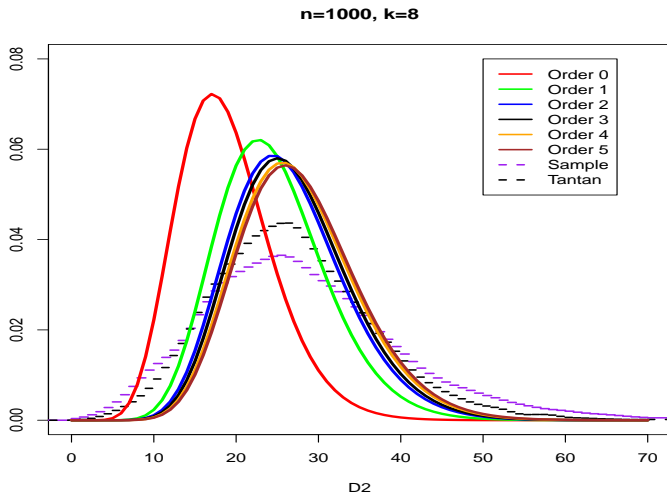
1. Obtain and mask a DNA sequence, yielding a series of unmasked regions;
2. For a fixed length n , produce a random sample of 10,000 pairs of sequences from the regions, using word length k to yield sequences with periodic boundary conditions;
3. Use **SAFT** program to calculate the D_2 value of each pair;
4. Given Markov order ω , compute the Markov matrix M from the DNA regions;
5. Given Markov matrix M , word length k and sequence length n , compute the theoretical mean and variance;
6. Compare the empirical distribution of D_2 values with a Gamma distribution using the theoretical mean and variance.

Human Chromosome 1: sample pairs

Note the gap around the centromere.



Human Chromosome 1: D_2 vs Markov models



What's next?

1. Compare chromosomes with their *stationary k -mer spectrum* to look for regions of interest.
Have determined the theoretical mean and variance for this case. The formula for the variance is much simpler than the variance for D_2 between two sequences.
Need to modify the SAFT program to compare a stationary k -mer spectrum to a database of sequences.
2. Scale up the SAFT program to work quickly with large databases.
This includes testing parallel code on NCI clusters during 2013.
3. Release the SAFT program as open source software.
Anticipated some time in 2013.