# Test on the Structure of Biological Sequences via Chaos Game Representation

Peggy Cénac

INRIA Domaine du Voluceau, B.P. 105, 78 153 Le Chesnay Cedex (France) Peggy.Cenac@inria.fr



1 Chaos Game Representation

#### 2 Testing the structure of a sequence

The CGR is both a graphical representation method of sequences and a storage tool. This iterative mapping technique was apparently for the first time applied to genomic sequences by Jeffrey [2]. From a given sequence –e.g. nucleotides in a DNA sequence or amino acids in a protein– one can define trajectories in a bounded set conserving all its statistical properties. **Each point of the CGR contains the whole history** of the sequence. One of the central goals of Cénac *et al.* [1] is to figure out whether the CGR provides more information than the classical methods based on word-counting.

#### **Definitions**

 $\succ U = u_1 \dots u_n$  a sequence of letters in a *d*-letter alphabet  $\mathcal{A}$ .

The Chaos Game Representation of U, on a bounded Borel set  $S \subset \mathbb{R}^q$  is a sequence  $\{X_0, \ldots, X_n\}$  defined by

$$\begin{cases} X_0 \in S \\ X_{i+1} = \theta \left( X_i + \ell_{u_{i+1}} \right) \stackrel{\text{\tiny def}}{=} T_{u_{i+1}}(X_i), \end{cases}$$

for  $0 < \theta < 1$ .

> Jeffrey's definition for DNA sequences,  $\mathcal{A} = \{A, C, G, T\}$ :

$$\begin{cases} S = [0, 1]^2, & \theta = \frac{1}{2}, \quad X_0 = (\frac{1}{2}, \frac{1}{2}) \\ \ell_A = (0, 0), & \ell_C = (0, 1), \quad \ell_G = (1, 1), \quad \ell_T = (1, 0) \end{cases}$$

 $\succ$  Counting points in  $Sw \stackrel{\text{\tiny def}}{=} \sum_{k=1}^{i} \theta^{i-k+1} \ell_{u_k} + \theta^i S \Leftrightarrow$  counting occurrences of the word w.

- $H_0$ : " $U = u_1 \dots u_N$  is an i.i.d. sequence"
- $H_m: "Uis a Markov chain of order m"$
- H : "U is a stationary ergodic sequence

## $\succ$ <u>Construction</u>

Denoting  $q_{\alpha}(d)$  the  $(1 - \alpha)$ -quantile of the chi-square distribution  $\chi^2(d)$ , and  $\hat{\pi}_n(E) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{E\}}(X_j)$  the empirical measure of  $\pi$ , for any partition  $\{B_1, \ldots, B_K\}$  of S, with K > 1, the following sets are reject region with asymptotic level  $\alpha$ , of a test of  $H_0$  against  $H \setminus H_0$  and respectively of  $H_m$  against  $H \setminus H_m$ 

$$\left\{\sum_{\substack{1\leq i\leq K\\v\in\mathcal{A}}}\frac{n\left(\hat{\pi}_n(Bv)-\hat{\pi}_n(B)\hat{\pi}_n(Sv)\right)^2}{\hat{\pi}_n(B)\hat{\pi}_n(Sv)} > q_{\alpha}\left[(d-1)(K-1)\right]\right\},$$
$$\sum_{\substack{wu\in\mathcal{A}^m\times\mathcal{A}\\1\leq i\leq k}}\frac{n\left(\hat{\pi}_n(Sw)\hat{\pi}_n(Bwu)-\hat{\pi}_n(Swu)\hat{\pi}_n(Bw)\right)^2}{\hat{\pi}_n(Sw)\hat{\pi}_n(Swu)\hat{\pi}_n(Bw)} > q_{\alpha}\left[d^m(d-1)(K-1)\right]\right\},$$

#### $\succ$ <u>Consistence</u>

Next, assume that  $H \setminus H_0$  (resp.  $H \setminus H_m$ ) holds, and let  $B \subset S$ ,  $w \in \mathcal{A}^d$  and  $v \in \mathcal{A}$  be such that

# $\pi(Bv) \neq \pi(B)\pi(Sv).$



**Figure** 1 : On the left, CGR of the 10 first nucleotides ATGCGAGTGT of the *E. Coli* threonine gene. Point number 3 corresponds to the first 3-letter word ATG. It is located in the corresponding quadrant. The second 3-letter word TGC corresponds to point 4 and so on. On the right, CGR of a sequence of length 80000 of *E. Coli*.

#### Stochastic properties of the CGR

- > U is supposed to be a stationary ergodic sequence.
- $> (X_n)_{n \ge 0}$  is a Markov chain of order 1, and converges almost surely to a random vector X with distribution  $\pi$ .
- > When U is i.i.d. and uniformly distributed,  $\pi$  is the Lebesgue measure on S. Whenever S is not uniformly distributed,  $\pi$  is continuous, singular with respect to the Lebesgue measure. The law of large number holds, and the empirical measures converge.

#### Characterization of Structure

➤ For any word  $w = u_1 \dots u_i$  and for any set  $B \subset S$ ,  $Bw \stackrel{\text{\tiny def}}{=} T_{u_i} \circ \dots \circ T_{u_1}(B)$ .

**Proposition 1.1.** The stationary random sequence U is

respectively

## $\pi(Sw)\pi(Bwu)\neq\pi(Bw)\pi(Swu).$

If B is one of the set forming the partition, then the test is asymptotically consistent.

For this test, Reinert *et al.* [3] propose to make use of Pearson statistics. In the special case when the partition is  $\{Su, u \in A\}$ , the tests are asymptotically the same.

# **3** Numerical experiments



**Figure** 2 : The 4 different partitions of the square  $[0, 1]^2$  arbitrary chosen for the test.

We did generate 1000 sequences of length n of Markov chains of order m with random transition matrix, for various values of n.

der	n	Pearson	$B^{(1)}$	$B^{(2)}$	$B^{(3)}$	order	n	Pearson	$B^{(2)}$	
0	100	4.2%	3.6%	3.4%	2.5%	1	500	5.2%	5.4%	
	500	6.1%	4.8%	4.8%	3.8%		1000	4.1%	5.2%	
	1000	5.0%	5.6%	6.2%	5.3%		10000	5.3%	6.4%	
	10000	6.5%	4.8%	5.1%	7.7%	2	100	60.8%	25.8%	1
1	100	86.4%	12.9%	51.1%	28.9%		500	100%	98.4%	i
	500	100%	54.2%	98.7%	94.5%		1000	100%	99.9%	i
	1000	100%	70.9%	99.9%	99.0%		10000	100%	100%	
	10000	100%	97.6%	100%	100%	4	500	22.8%	20.2%	1
5	1000	8.6%	6.8%	8.6 %	8.4 %		1000	51.0%	44.3%	i
	10000	54.6%	28.7%	55.6%	85.3%		5000	99.9%	99.6%	i
	80000	99.4%	84.5%	99.6%	100%		10000	100%	100%	
2 mixed	500	5.8%	16.5%	49.9%	76.8%	3 mixed	l 500	5.4%	29.6%	1
	1000	7.0%	26.9%	73.7%	95.1%		1000	8.1%	55.6%	ł
	10000	7.3%	73.2%	99.8%	100%		5000	8.3%	96.4%	ł
5 mixed	80000	5.8%	29.7%	76.7%	85.8%		10000	6.3%	99.0%	i

*an i.i.d. sequence if and only if* 

 $\pi(Bu) = \pi(B)\pi(Su), \quad \forall u \in \mathcal{A}, \ \forall B \subset S.$ 

 $\mathscr{A} a \text{ Markov chain of order } m \text{ if and only if } \forall B \subset S, \forall w \in \mathcal{A}^m, \forall u \in \mathcal{A},$ 

$$\frac{\pi(Bwu)}{\pi(Bw)} = \frac{\pi(Swu)}{\pi(Sw)}, \quad \forall B \subset S, \quad \forall w \in \mathcal{A}^m, \quad \forall u \in \mathcal{A}$$

- In particular the ratio  $\frac{\pi(Bwu)}{\pi(Bw)}$  does not depend on B.
- $\succ$  Characterization of independence and of Markov chains
- $\blacktriangleright$  construction of a test.
- ▶ genomic signature (see Cénac *et al.* [1])

**Table** 1 : The left (resp. right) hand side of the following Table shows the fraction of cases when  $H_0$  (resp.  $H_1$ ) is rejected.

#### $\succ$ The choice of the partition is crucial.

- $\succ$  The reject of long dependence Markov chains increases with the size of the partition.
- ➤ In the special case of markov chains of order-*m* given by the aggregation of *m* indepedent markov chains of order 1, the CGR-based test behaves pretty well. This illustrates the strength of the CGR: it does not impose any constraint on the input sequence besides stationarity.

## References

P. Cénac, G. Fayolle, and J.M. Lasgouttes. Dynamical systems in the analysis of biological sequences. Technical Report 5351, INRIA, october 2004.
H.J. Jeffrey. Chaos Game Representation of gene structure. *Nucleic Acid. Res*, 18:2163–2170, 1990.
G. Reinert, S. Schbath, and M.S. Waterman. Probabilistic and statistical properties of words: An overview. *Journal of Computational Biology*, 7(1/2):1–46, 2000.