

# Constitution Analysis

Israel T. Silva<sup>1,3,\*</sup>, and Rafael A. Rosales<sup>2</sup>

<sup>1</sup>Laboratory of Molecular Immunology, The Rockefeller University, 1230 York Avenue, New York, NY 10065

<sup>2</sup>Departamento de Computação e Matemática, Universidade de São Paulo. Av. Bandeirantes, 3900, Ribeirão Preto, CEP 14049-901, SP, Brazil

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXXX

## ABSTRACT

**Motivation:**

**Results:** <https://github.com/someone/rich>.

**Contact:** [someone@somewhere.world](mailto:someone@somewhere.world)

## 1 INTRODUCTION

The genome is targeted by a sophisticated and highly coordinated series of molecular events. Among these events, aberrant DNA methylation patterns in human malignancy [De et al. \(2013\)](#), somatic retrotransposition in human cancers [Lee et al. \(2012\)](#), AID-dependent chromosomal translocations (Klein, 2011) and HIV integration (Cohn, 2014), which arrives throughout DNA, are not randomly distributed but instead associated with chromosomal regions and contributes to disrupt the integrity of the genome and human disease.

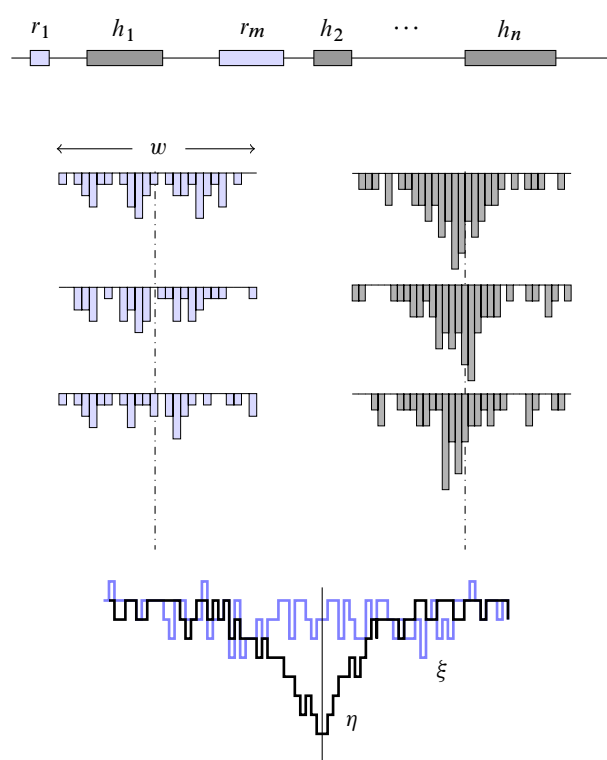
As result, these regions represents a genomic context in which are associate with multiple underlying mechanisms. The motif-based sequence analysis is the starting point to aim potential binding site of cis-regulatory elements associated. Nevertheless, the inherent low signal/noise ratio in sequence-based motif discovery is a limitation to detect a nucleic acid sequence pattern that has some biological significance. Moreover, these events may recognise a structural feature, rather than a specific sequence motif.

Others approaches were introduced to detect functional regions using methods for computing sequence complexities [Jin et al. \(2014\)](#); [Koslicki \(2011\)](#). In these methods, the complexity is measured by the entropy which evaluates the randomness of DNA sequence. In particular, topological entropy has been applied to compute the complexity of introns, exons and promoter regions. Due to the finite sample and high-dimensionality problems, efforts aimed to overcome these problems are put forward [Koslicki \(2011\)](#).<sup>1</sup>

Our work has some intersection whit the computation of ‘enrichment  $p$ -values’ considered in GO analysis. We may include the references [Huang et al. \(2009\)](#), [Rivals et al. \(2007\)](#) (just one!) or any other if you know a better alternative. We may like to mention the paper by [Bailey and Machanick \(2012\)](#) because it also considers a test for enrichment (although it is restricted to ChIP-seq peaks, and somehow different).

\*to whom correspondence should be addressed

<sup>1</sup> its: I would prefer to leave entropy out unless we have a really good point we want to make



**Fig. 1.** The method. Enhance: set notation equal to the one in main text

However, how exactly the pattern nucleotide composition could influence the selection of target site selection are not well understood. To further characterize at a genome-wide scale these regions, we introduce a new method to provide a quantitative measure of the differential spectra of  $k$ -mers (DNA ‘words’ of length  $k$ ) inside target DNA.

## 2 METHOD

Let  $\mathcal{A} = \{A, C, G, T\}$  and  $S \in \mathcal{A}^\ell$ , be a given specific string of length  $\ell \geq 1$ . In what follows, we describe a method to study the profile of  $S$  along a region of interest such as those defined by viral insertion or retrotranslocation hotspots. This provides the means

to assess the significance of a differently distributed profiles along two functionally defined regions. We specialise to viral insertion hotspots as described by Silva et al. (2014), but the scope is clearly not restricted to this particular application.

Let  $h = \{h_1, \dots, h_n\}$  be a set of viral insertion hotspots, namely a set of DNA segments characterized by having a substantially high density of viral insertion events. Let  $w$  be the length of the longest of such segments. The segments  $h_1, \dots, h_n$  are aligned with respect to their central base and then extended at both ends to have length  $w$ . Next, consider the partition of resulting set of segments into  $k$  evenly spaced bins of length  $\ell = w/k$ . Denote by  $h_{ij}$ ,  $1 \leq j \leq k$ , the  $j$ th bin of the  $i$ th segment. Consider now the set  $r = \{r_1, \dots, r_n\}$  of segments of width  $w$  that are either at the left or at the right of any one segment in  $h$ . Likewise, let  $r_{ij}$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq k$ , be the matrix formed by bins of length  $\ell$  that result by partitioning the elements of  $r$ . For any  $j = 1, \dots, k-1$  and  $i = 1, \dots, n$ , let  $\xi_{ij}$  and  $\eta_{ij}$  be the following Bernoulli random variables

$$\xi_{ij} = \begin{cases} 1, & \text{if } S \in h_{ij} \text{ and } S \notin h_{i,j+1} \\ 1, & \text{if } S \in h_{ij} \text{ and } S \in h_{i,j+1} \\ 0, & \text{otherwise} \end{cases}$$

$$\eta_{ij} = \begin{cases} 1, & \text{if } S \in r_{ij} \text{ and } S \notin r_{i,j+1} \\ 1, & \text{if } S \in r_{ij} \text{ and } S \in r_{i,j+1} \\ 0, & \text{otherwise} \end{cases}$$

Set  $\xi_{ik} = 1$  if  $S \notin h_{i,k-1}$  and  $S \in h_{i,k}$ , and  $\xi_{ik} = 0$  otherwise. Similarly define  $\eta_{ik}$  by using the information in  $r_{ik}$ . Finally, let

$$\xi_j = \sum_{i=1}^n \xi_{ij}, \quad \eta_j = \sum_{i=1}^n \eta_{ij}.$$

The variables  $\xi_j$  and  $\eta_j$ ,  $1 \leq j \leq k$ , count the number of times that the string  $S$  occurs along of a hotspot region and of a reference region respectively.

The basic question we like to address is whether the distribution profile of the string  $S$  is significantly different along a typical hotspot region and a reference region. This may be assessed by considering the following  $2 \times k$  contingency table

$$\begin{bmatrix} \xi_1 & \xi_2 & \dots & \xi_k \\ \eta_1 & \eta_2 & \dots & \eta_k \end{bmatrix},$$

obtained by merging the two vectors of counts  $\xi_j$  and  $\eta_j$ . Provided the number of counts in each of the cells of this table is sufficiently large, the significance of a differential profile can be determined by using Pearson's  $\chi^2$  statistic, which is distributed according to a  $\chi^2$  density with  $k-1$  degrees of freedom. Other alternatives for the large sample case exist, see for instance Read and Cressie (1988), but we do not pursue this further here. It is well known that this procedure can give a poor approximation when several cells present low counts (smaller than 10). This may be the case in the current setting when analysing the profile distribution of

longer strings with  $\ell \geq 10$  or even smaller but rarely occurring strings. In these situations the significance for a differential profile is more appropriately determined by using Fisher's exact test, see for instance Agresti (2012). The computations necessary to derive a  $p$ -value are not feasible because of the large number of contingency tables that have to be considered as a reference set when  $k$  is large. The significance may however be approximately computed by considering a permutation test using the method in Patefield (1981). We found that R's implementation via `fisher.test` takes only few seconds for relatively large tables with  $k = 1000$ .

We provide examples for the two scenarios just mentioned by considering strings formed by a single base and strings defined by longer motifs with  $\ell = 15$ . The former provides an example where the  $\chi^2$  statistic is appropriate and the latter one that is amenable to the analysis with Fisher's exact test.

### 3 RESULTS

Put the plots and the  $p$ -values.

### 4 DISCUSSION

Mention that the results are surprising and important from the perspective of the virus insertion problem. Then talk very briefly about the scope of this method: what kind of problems can we consider.

### REFERENCES

- Agresti, A. (2012). *Categorical Data Analysis*. Wiley Series in Probability and Statistics. Wiley-Interscience, 3rd edition.
- Bailey, T. L. and Machanick, P. (2012). Inferring direct DNA binding from ChIP-seq. *Nucl. Acids Res.*, **40**(17), 1–10.
- De, S., Shaknovich, R., Riester, M., Elemento, O., Geng, H., Kormaksson, M., Jiang, Y., Woolcock, B., Johnson, N., Polo, J. M., Cerchietti, L., Gascoyne, R. D., Melnick, A., and Michor, F. (2013). Aberration in DNA methylation in B-cell lymphomas has a complex origin and increases with disease severity. *PLoS Genet.*, **9**(1), e1003137.
- Huang, D. W., Sherman, B. T., and Lempicki, R. A. (2009). Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. *Nucl. Acids Res.*, **37**(1), 1–13.
- Jin, S., Tan, R., Jiang, Q., Xu, L., Peng, J., and Wang, Y. (2014). A generalized topological entropy for analyzing the complexity of DNA sequences. *PLoS ONE*, **9**(2), e88519.
- Koslicki, D. (2011). Topological entropy of DNA sequences. *Bioinformatics*, **27**(8), 1061–1067.
- Lee, E., Iskow, R., Yang, L., Gokcumen, O., Haseley, P., Luquette, L. J., Lohr, J. G., Harris, C. C., Ding, L., Wilson, R. K., Wheeler, D. A., Gibbs, R. A., Kucherlapati, R., Lee, C., Kharchenko, P. V., and Park, P. J. (2012). Landscape of somatic retrotransposition in human cancers. *Science*, **337**(6097), 967–971.
- Patefield, W. M. (1981). Algorithm AS 159: An efficient method of generating random  $R \times C$  tables with given row and column totals. *J. Appl. Stat.*, **30**(1), 91–97.
- Read, T. R. C. and Cressie, N. A. C. (1988). *Goodness-of-Fit Statistics for Discrete Multivariate Data*. Springer Series in Statistics. Springer Verlag, New York.
- Rivals, I., Personnaz, L., Taing, L., and Potier, M. C. (2007). Enrichment or depletion of a GO category within a class of genes: which test? *Bioinformatics*, **23**(4), 401–407.
- Silva, I. T., Rosales, R. A., Holanda, A. J., Nussenzweig, M., and Jankovic, M. (2014). Identification of chromosomal translocation hotspots via scan statistics. *Bioinformatics*, **30**(18), 2551–2558.