ENTROPIC PROFILER Efficient whole genome analysis using information theory and statistical concepts

Francisco Fernandes^a, Ana T. Freitas^a, Jonas S. Almeida^c, Susana Vinga^{a,b}

a INESC-ID/IST, Portugal ^b FCM/UNL, Portugal ^c MDAnderson Cancer Center, Houston, USA

Email: atf@inesc-id.pt and svinga@kdbio.inesc-id.pt

Abstract

ISMB ECCB 2009

Entropic Profiles (EP) are local information plots that indicate overall conservation of motifs in genomes. They are based on Information Theory concepts, in particular to the Renyi entropy of biological sequences. The present tool implementation, based on new data structures and algorithmic simplifications, allows to process whole genomes in few minutes. ENTROPIC PROFILER is freely available as a web interface and downloadable source code at http://kdbio.inesc-id.pt/software/ep/ and is fully described in [1].

Introduction

In a recent paper [2], the authors presented the concept of Entropic Profiles (EP), a new method to extract and classify relevant and statistically significant segments of DNA sequence. The study of these motifs is very relevant because under or overrepresentation segments are often associated with significant biological meaning. In this work EP plots express the relative abundance of corresponding motifs for each position and allow estimating local relevant scales. Its calculation is based on the continuous Rényi quadratic entropy [3], using the Parzen window estimation method [4] applied to the Chaos Game Representation (CGR) of a sequence [5].



- Entropic profiles (EP) provide useful local information about global features of DNA. The current implementation exhibits excellent performance for sequences up to 2Gbp.
- Tests on whole genomes corroborate the strengths of this approach to detect biologically aningful DNA segments, related with the detection of local scales and suffix/motifs ove
- Further improvements include: search for motifs in IUPAC code and analyze both forward and reverse strands and development of a web-based application with a simple Representational State Transfer (REST) I/O. As a consequence, this functionality will be also easily accessible as a URL POST which enables using the application as a web service.

[1] Fernandes et al. (2009) BMC Research Notes 2:72 [2] Vinga S, Almeida JS (2007) BMC Bioinformatics 8:393. [3] Rényi A (1961) University of California Press 547-561. [4] Parzen E (1962) The Annals of Math Stat 33:1065-1076. [5] Jeffrey HJ (1990) Nucleic Acids Res 18:2163-2170. Acknowledgments

The authors acknowledge financial support by projects DynaMo (PTDC/EEA-ACR/ 69530/2006, FCT) and ARN (PTDC/EIA/67722/2006, FCT).