

Chapter 10

Biological Big Data Analysis and Visualization: A Survey

Vignesh U
VIT University, India

Parvathi R
VIT University, India

ABSTRACT

The chapter deals with the big data in biology. The largest collection of biological data maintenance paves the way for big data analytics and big data mining due to its inefficiency in finding noisy and voluminous data from normal database management systems. This provides the domains such as bioinformatics, image informatics, clinical informatics, public health informatics, etc. for big data analytics to achieve better results with higher efficiency and accuracy in clustering, classification and association mining. The complexity measures of the health care data leads to EHR (Evidence-based Healthcare) technology for maintenance. EHR includes major challenges such as patient details in structured and unstructured format, medical image data mining, genome analysis and patient communications analysis through sensors – biomarkers, etc. The big biological data have many complications in their data management and maintenance especially after completing the latest genome sequencing technology, next generation sequencing which provides large data in zettabyte size.

DOI: 10.4018/978-1-5225-2805-0.ch010

INTRODUCTION

The chapter was initiated by requirement of higher and efficient methodologies to analyze big data in a faster manner. The deficiency has motivated us to investigate the problems in an existing technology and frame a feasible model for this big data analysis. On the other hand, there is a considerable interest in the development of new techniques using dynamic programming algorithms to work faster for bioinformatics methods. High throughput sequencing workflow systems provide easy and cost reduced perspective to genome sequencing with timely detection of functions, accurate and fast solutions for big data in bioinformatics. The table 1 shows the detailed view of the different workflow systems that can support high throughput sequencing technologies which includes a big data incorporated in it for analysis.

Bioinformatics is an interdisciplinary area that deals with the biology, computer and statistics. It involves the major aspects of genomics and proteomics with the genome sequencing, which are very sensitive in nature as representing the individual letter for a single nucleotide in case of DNA sequencing. Since 1970, the biological databases are digitized and their sensitivity factors with efficiency are maintained in a perfect manner but due to the vast amount of increasing data the maintenance aspect and extraction of information from gene expression becomes so complex,

Table 1. High Throughput Sequencing Workflow Systems

Name	Illumina	Solid	Requirements	GUI	CLI	Online	Cloud
Ergatis	yes	yes	Linux, MAC OS X, Windows	yes	no	yes	Yes
Galaxy	yes	yes	Linux, MAC OS X	yes	no	yes	yes
Genboree Workbench	yes	yes	Linux, MAC OS X, Windows	yes	no	yes	Yes
GenePattern	yes	yes	Linux, MAC OS X, Windows	yes	no	yes	No
GeneProf	yes	yes	Linux (it is not tested on Others yet)	yes	no	yes	No
Kepler (bioKepler)	yes	yes	Linux, MAC OS X, Windows; > 1 GB RAM, 2 GHz CPU	yes	no	no	No
KNIME	yes	-	Linux, MAC OS X, Windows	yes	yes	no	Yes
LONI Pipeline	yes	yes	Linux, MAC OS X, Windows	yes	yes	no	No
Moa	yes	yes	Linux	yes	yes	no	No
Tavaxy	yes	yes	Linux	yes	no	yes	Yes
Taverna	yes	yes	Linux, MAC OS X, Windows	yes	yes	no	yes
Yabi	-	-	Linux	yes	yes	yes	yes

14 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/biological-big-data-analysis-and-visualization/185986

Related Content

Mining Free Text for Structure

Vladimir A. Kulyukinand Robin Burke (2003). *Data Mining: Opportunities and Challenges* (pp. 278-300).

www.irma-international.org/chapter/mining-free-text-structure/7605

Managing Late Measurements in Data Warehouses

Matteo Golfarelliand Stefano Rizzi (2007). *International Journal of Data Warehousing and Mining* (pp. 51-67).

www.irma-international.org/article/managing-late-measurements-data-warehouses/1793

Conceptual Model and Design of Semantic Trajectory Data Warehouse

Michael Mireku Kwakye (2020). *International Journal of Data Warehousing and Mining* (pp. 108-131).

www.irma-international.org/article/conceptual-model-and-design-of-semantic-trajectory-data-warehouse/256165

Selection of High Quality Rules in Associative Classification

Silvia Chiusanoand Paolo Garza (2009). *Post-Mining of Association Rules: Techniques for Effective Knowledge Extraction* (pp. 173-198).

www.irma-international.org/chapter/selection-high-quality-rules-associative/8443

Load Balancing in Cloud Computing: Challenges and Management Techniques

Pradeep Kumar Tiwari, Geeta Rani, Tarun Jain, Ankit Mundraand Rohit Kumar Gupta (2020). *Critical Approaches to Information Retrieval Research* (pp. 294-316).

www.irma-international.org/chapter/load-balancing-in-cloud-computing/237652