

Chapter 54

Open–Source Essential Protein Prediction Model by Integrating Chi–Square and Support Vector Machine

S. R. Mani Sekhar

REVA University, India & Ramaiah Institute of Technology, India

Siddesh G. M.

Ramaiah Institute of Technology, India

Sunilkumar S. Manvi

REVA University, India

ABSTRACT

Identification and analysis of protein play a vital role in drug design and disease prediction. There are several open-source applications that have been developed for identifying essential proteins which are based on biological or topological features. These techniques infer the possibility of proteins to be essential by using the network topology and feature selection, which can ignore some of the features to reduce the complexity and, subsequently, results in less accuracy. In the paper, the authors have used selenium driver to scrap the dataset. Later, the authors integrated the chi-square method with support vector machine for the prediction of essential proteins in baker yeast. Here, chi-square is a test of dissimilarity used for altering the record, and afterward, the support vector machine is used to classify the test dataset. The results show that the proposed model Chi-SVM model achieves an accuracy of 99.56%, whereas BC and CC achieved an accuracy of 84.0% and 86.0%. Finally, the proposed model is validated using Statistical performance measures such as PPA, NPA, SA, and STA.

DOI: 10.4018/979-8-3693-3026-5.ch054

1. INTRODUCTION

In the current computational domain, proteomics act as a sub-branch of bioinformatics, which studies the structure and functions of proteins. It also helps in identifying the procedure and encoding order used in the analysis of the large scale data set. Free and open-source software (FOSS) allows the developer to use, change, and distribute software without significant limitations. However, it also provides a complete software package to users freely. Now a day's requirement of open source software in bioinformatics and chem-informatics etc is increasing rapidly. cPath (Cerami, Bader, Gross, & Sander, 2006) an open-source software contains a web interface for gathering, storing, and querying bio-informatics data. It also supports the build-in mapping service and helps to link to external resources (Breitwieser & Colinge, 2013). Developed a tool to analysis proteins localization and processing of Tandem Mass Tags peptides. In addition it also supports statistical computation with user friendly reports generation. For study of protein function and sequence (Deng, Yuan, Huang, & Wang, 2013) created a SFAPS module by using spectrum function. This function helps in converting numerical sequences into the characteristic frequency of the protein interface. Subsequently helps in analyze sequence of proteins. ATGme (Daniel et al., 2015) designed an open-source web-based tool for optimizing the protein sequence by classifies organisms into rare and very high rare codons. After that, help the user to generate an individual custom optimized DNA sequence.

Proteins are the best adaptable macromolecules in living creatures and also play important roles in biological procedures. They are made of individual or several chains of amino acid residues. An amino acid is made up of an organic composite that has an essential carbon atom, called an alpha carbon ($C\alpha$). The bond is consists of four valences with a link to the carboxylic group, hydrogen atom, an amino group, and an adjacent chain. The amino acids are the strong structured chunks of protein. This sequence of amino acids in a given polymer helps in identifying protein structure and protein function. It can also be called as a polypeptide, though a small chain of amino acids. This chained bond starts with the amino cluster and ends with the carboxyl cluster. According to Bioinformatics researcher DNA, RNA, Protein, and other molecules do not work separately. They can be linked together to perform a specific biological action. Interaction is a process when two molecules linked to perform a certain function. These type of interaction are mainly divided into four groups by considering their molecule types as shown below

- *Protein-protein interactions*: collaboration between proteins to derive biological procedures.
- *Gene regulatory interactions*: interaction of genetic data to normalize protein expression level.
- *Metabolic interactions*: create support between enzyme proteins and helps in changing a substrate cell into product cell through numerous catalysis responses.
- *RNA-DNA interactions*: provide collaboration among RNA-RNA or RNA-DNA connections and also plays a vital role in critical diseases.

PPI performance numerous functions in a different cell, include metabolic pathways, DNA, replication, transduction, and transcription etc. In a PPI network, a vertical position signifies the proteins, and the linking line shows the communication link between proteins. This network can be noticed as a graph. Thus the closely connected protein clusters can be computed efficiently.

Essential Proteins are made up of a set of minimal genome, which can support cell survival with basic requirements. A living creature cannot breed and live without them (Kamath et al., 2003; Wang, X Peng, W Peng, & Wu, 2014). Proteins are mostly classified into two types Essential proteins and non-essential

15 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/open-source-essential-protein-prediction/342572

Related Content

Estimation of Fractal Dimension in Different Color Model

Sumitra Kisan, Sarojananda Mishra, Ajay Chawda and Sanjay Nayak (2018). *International Journal of Knowledge Discovery in Bioinformatics* (pp. 75-93).

www.irma-international.org/article/estimation-of-fractal-dimension-in-different-color-model/202365

A Two-Layer Learning Architecture for Multi-Class Protein Folds Classification

Ruofei Wang and Xieping Gao (2013). *Bioinformatics: Concepts, Methodologies, Tools, and Applications* (pp. 786-797).

www.irma-international.org/chapter/two-layer-learning-architecture-multi/76094

Computer Aided Tissue Engineering from Modeling to Manufacturing

Mohammad Haghighpanahi, Mohammad Nikkhoo and Habib Allah Peirovi (2010). *Biocomputation and Biomedical Informatics: Case Studies and Applications* (pp. 75-88).

www.irma-international.org/chapter/computer-aided-tissue-engineering-modeling/39604

The Avatar as a Self-Representation Model for Expressive and Intelligent Driven Visualizations in Immersive Virtual Worlds: A Background to Understand Online Identity Formation, Selfhood, and Virtual Interactions

Colina Demirdjian and Hripsime Demirdjian (2020). *International Journal of Applied Research in Bioinformatics* (pp. 1-9).

www.irma-international.org/article/the-avatar-as-a-self-representation-model-for-expressive-and-intelligent-driven-visualizations-in-immersive-virtual-worlds/261865

Using Biomedical Terminological Resources for Information Retrieval

Piotr Pezik, Antonio Jimeno Yepes and Dietrich Rebholz-Schuhmann (2009). *Information Retrieval in Biomedicine: Natural Language Processing for Knowledge Integration* (pp. 58-77).

www.irma-international.org/chapter/using-biomedical-terminological-resources-information/23055