

A Novel Predictor of Protein Subcellular Localization based on N-gram Features

Ho-Eun Park and Jinsuk Kim

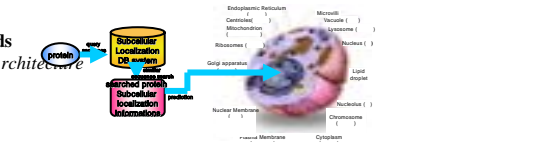
{hepark, jinsuk}@kisti.re.kr

Center for Computational Biology and Bioinformatics,
Korea Institute of Science and Technology Information

ProSLP prediction system for protein subcellular localization site(s), called as ProSLP. ProSLP classifies a protein sequence to one or more compartments based on the location of the top k sequences which have the highest weights against the input sequence. Currently ProSLP extracts n-grams as features of the sequence, computes scores of the potential localization site(s) using k -nearest neighbor (k NN) algorithm and finally presents the results and their associated scores. ProSLP is available through the World Wide Web at <http://proslp.kisti.re.kr>.

Introduction

As a result of large-scale genome sequencing projects, the number of genes and protein sequences of unknown functions are increasing tremendously. It is a time-consuming and costly job to identify the inherent functions of protein sequences. Given a protein sequence, how to determine its subcellular localization as an important clue to its functions is a problem vitally important to biologists and bioinformaticists (Chou and Elrod, 1999). To cooperate for a better understanding of the physiological function, proteins must be localized in the same cellular compartment. So subcellular localization is a key functional characteristic of proteins (Eisenhaber and Bork, 1998). Several automated subcellular localization prediction systems have been developed and made available online: TargetP (Nielsen et al., 2000), PSORT (Nakai and Horton, 1999), NNPSL (Nakai and Hubbard, 1998), SubLoc (Hua and Sun, 2001), MitoProt and MitoProt2 (Feng, 2002). In most of these approaches, features for a sequence are extracted in the form of amino acid composition and/or sorting signals located at the N-terminal amino acid sequence. And these features are further processed by various kinds of machine learning approaches such as neural networks, support vector machines, and k -nearest neighbor (k NN) classifiers, finally giving the subcellular localization site(s). Using k NN classification algorithm and n-gram features, a subcellular localization prediction service has been developed. This prediction system is called ProSLP (Protein Subcellular Localization Prediction) and now available via WWW at <http://proslp.kisti.re.kr>.



Feature Extraction

ProSLP adopts a simple feature extraction scheme, n-gram tokens, which is using intervals of fixed length n . For example, if the length is 5, a protein sequence "ACEDFIMMPAA" is segmented into "ACEDF", "CEDFI", "EDFIM", "DFIMP", "FIMPA". Current version of ProSLP system uses the interval length of 5 (n=5) as its feature extraction scheme. After such tokens are extracted and stored in the database, sequences can be searched with various information retrieval algorithms.

Similarity measure

The similarity between the query and target sequence is measured by vector space model regarding n-gram features as their representative vectors for the sequences. The similarity measure ($Sim(q,s)$) between query sequence q and target sequence s is defined as follows:

$$Sim(q,s) = \frac{1}{W_s} \cdot \sum_{t \in q \cap s} w_{s,t} w_{q,t}$$

with:

$$W_s = \log(1 + \sum_{t \in s} f_{s,t})$$

$$w_{s,t} = \log(f_{s,t} + 1) \cdot \log\left(\frac{N}{f_t}\right)$$

$$w_{q,t} = \log(f_{q,t} + 1) \cdot \log\left(\frac{N}{f_t}\right)$$

where $f_{s,t}$ is the frequency of n-gram token t in sequence s ; N is the total number of sequences in the data collection; f_t is the number of sequence where token t occurs more than or equal to once; $w_{s,t}$ is the weight of term t in sequence s ; W_s is the normalization factor for the length of target sequence s .

Performance measures

To evaluate the prediction performance, we used the standard definition of precision (p) and recall (r), F_1 and break-even point (BeP).

$$p = \frac{\text{locations relevant and retrieved}}{\text{locations retrieved}} = \frac{tp}{tp + fp}$$

$$r = \frac{\text{locations relevant and retrieved}}{\text{locations relevant}} = \frac{tp}{tp + fn}$$

$$F_1 = \frac{2pr}{p+r} \quad (\text{A special point of } F_1 \text{ measure where } p = r \text{ is called BeP})$$

Results

Prediction effectiveness

For k -nearest neighbor (k NN) classifiers, one of the most crucial steps to improve categorization performance is to determine suitable k for given data collection (Yang, 1999). Using two data collections, k value for the k NN classifier is experimentally optimized by obtaining the k value which shows the best prediction effectiveness. Figure 1 shows the prediction effectiveness plotted against varying k value using a) data collection 1 and b) data collection 2.

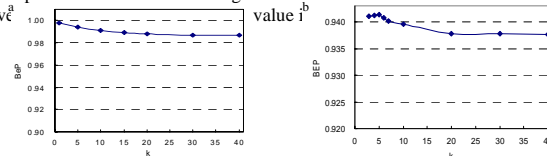


Figure 1. Prediction effectiveness plotted against varying k value using a) data collection 1 and b) data collection 2.

Data collection 1 consists of training set (86,720 sequence entries) and test set (10,840 sequences) and data collection 2 consists of training set (43,240 sequence entries) and test set (8,645 sequences).

Prediction accuracy of ProSLP is within the range of 93% to 98% and recall for all k points.

Prediction times

Since ProSLP data collection is one or two-magnitudes larger than other prediction systems, one of our major focuses on this work is to reduce the prediction time, which are plotted against query sequence lengths in Figure 2. These experiments were carried out on a Pentium-III dual processor system running the Linux operating system. 48 query sequences were randomly selected and classification time for each query was measured. Average length of query sequence is 262 amino acids and average search time for each query is shown in Figure 2.

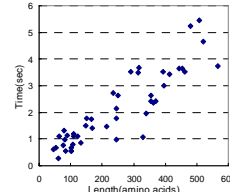


Figure 2. Prediction times plotted against the amino acid lengths of query protein sequences. We have used 48 query sequences, of which average length is 262 amino acids ranging 53 to 600.

Discussion

We showed that large-scale prediction system for subcellular localization possible with a k NN classifier and n-gram features of the protein sequence. ProSLP, a novel subcellular location prediction system, gives increased prediction accuracy. In addition, its prediction time is practical even on a single computer architecture. We hope that ProSLP can help biologists and bioinformaticists study various biological problems related to proteins.

References

- Chou,K.C. and Elrod,D.W. (1999) Protein subcellular location prediction. *Protein Eng.* 12:107-108.
- Emanuelsson,O., Nielsen,H., Brunak,S., and Hejine,G.V. (2000) Predicting the subcellular localization of Proteins Based on their N-terminal Amino Acid sequence. *J. Mol. Biol.* 300:1005-1016.
- Feng,Z.P. (2002) An overview on predicting the subcellular location of protein. *Int. J. Bioinformatics* 2:291-303.
- Gardy,J.L., Spencer,C., Wang,K., Ester,M., Tusnady,G.E., Simon,I., Hua,S., Lambert,C., Nakai,K., and Brinkman,F.S.L. (2003) PSORT-B: improving protein subcellular localization prediction for Gram-negative bacteria. *Nucleic Acids Res.* 31:3613-3617.
- Hua,S. and Sun,Z. (2001) Support vector machine approach for protein subcellular localization prediction. *Bioinformatics*, 17:721-728.
- Nakai,K. and Horton,P. (1999) PSORT: a program for detecting the sorting signals and predicting their subcellular localization. *Trends Biochem. Sci.*, 24:34-35.
- van Rijsbergen,C. (1979) *Information Retrieval* Butterworths, London, 1979.
- Yang,Y. (1999) An Evaluation of Statistical Approaches to Text Categorization. *Information Retrieval*, 1:67-88.