

ProSLP : Penta-gram

{hepark,jinsuk}@kisti.re.kr

122 , 305-600

ProSLP	(penta-gram)	k -nearest neighbor(kNN)
	가	가
	가	가
(Test set)	가	(Training set)
	가	(Test Collection)
break-even point	(Precision),	(Recall)
93%	가	ProSLP
98%	가	break-even point
k		ProSLP

<http://proslp.kisti.re.kr>

Abstract

A new prediction system for protein subcellular localization site(s), called ProSLP, has been developed. The ProSLP system classifies a protein sequence to one or more subcellular compartments based on the location of top k sequences which show the highest similarity against the input sequence. Currently the ProSLP extracts penta-grams as features of the sequences, computes scores of the potential localization site(s) using k -nearest neighbor (kNN) algorithm and finally presents the predicted site(s) and their associated scores. For two test collection tested so far, the ProSLP shows precision-recall break-even points of more than 0.93 and 0.98, respectively. The ProSLP service is available at <http://proslp.kisti.re.kr>.

1.

가 ,
 가 .
 가
 가 .
 (Chou and Elrod, 1999).

(Eisenhaber and Bork, 1998).
 가
 , , 3
 가
 (Gardy et al., 2003).

가 가
 가
 가
 TargetP (Emanuelsson *et al.*, 2000),
 PSORT (Nakai and Horton, 1999), NNPSL (Reinhardt and Hubbard, 1998),
 MitoProt/Predotar (Feng, 2002), PLOC (Park and Kanehisa, 2003)
 가 . 가
 . PLOC, NNPSL, SubLoc

(Signal Peptide)

가

(Text categorization)

(Feature)

(Effectiveness)

(Sebastiani, 2002).

kNN

(Yang, 1994; Sebastiani, 2002).

(Feature extraction)

n- (gram)

kNN n- (gram)

ProSLP(Protein Subcellular Localization Prediction

)

<http://proslp.kisti.re.kr>

가

2.

2.1

ProSLP
(Linux)

2GB

가

-III

KRISTAL - 2000¹

k-nearest neighbor(kNN)

n- (gram)

n- (gram)

(Inverted File)

KRISTAL - 2000

2.2 (Feature Extraction)

(Word)

(Phrase)

21

ProSLP

n

n- (gram)

가

n

5

¹ <http://giis.kisti.re.kr>

KRISTAL - 2000

"ACDEFLERR" "ACDEF", "CDEFL", "DEFLE",
 "EFLER", "FLERR"
 (Indexing)

ProSLP n- (gram)
 5 (penta-gram)

2.3 (Similariy Measures)

가 가 k , k 가
 (Target) (penta-gram) (Query)
 (Vector space model)

q s

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

with:

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

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

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

$f_{s,t}$ s (penta-gram) t ; N
 ; f_t (penta-gram) t가 ;
 $w_{s,t}$ s t 가 ; W_s s

k , $K = \{s_1, s_2, s_3, \dots, s_k\}$ 가
 q (c) 가 w_c

$$w_c = \sum_{i=1}^k sim(q, s_i) \cdot (c \wedge C_i)$$

C_i s_i K가 가
 가 , (Threshold) 가 가
 가

2.4 (Data Set)

kNN

(Eisenhaber and Bork, 1998)

CELL_LOC SWISS-PROT 'Access
number' CELL_LOC SWISS-
PROT (Bairoch and APweiler, 2000) release 39.0
(1).
1 97,560 , 13
(1).
10,798 (Test Set),
(Training Set)

1. CELL_LOC

1

13

Symbol	Description
intracellular	Proteins inside the cell
cytoplasmic	Cytoplasm
nuclear	Nucleus
endop_r_golgi	Endoplasmic Reticulum or Gogi Apparatus
mitochondrion	Mitochondrial proteins
chloroplast	Chloroplast
DNAb	DNA-binding proteins
RNAb	RNA-binding proteins
extracellular	Proteins secreted to outside the cell
Membrane	Membrane-floating proteins
Transmem	Trans-membrane proteins
Viral	Proteins of virus origin
hypothetical	Still undefined

SWISS-PROT

release 40.0

(2).

12

, SWISS-PROT

CC (Field)

("-!-

SUBCELLULAR LOCATION")

(Park and Kanehisa, 2003).

2 51,885 , 8,645
43,240 .

2. 2 SWISS-PROT

12 (Park and Kanehisa, 2003).

Symbol	Keywords
Chloroplast	Chloroplast
Cytoplasmic	Cytoplasmic
Cytoskeleton	Cytoskeleton, Filament, Microtubule
endoplasmic reticulum	Endoplasmic Reticulum
Extracellular	Extracellular, Secreted
golgi apparatus	Golgi
Lysosomal	Lysosomal
Mitochondrial	Mitochondrial
Nuclear	Nuclear
Peroxisomal	Peroxisomal, Microsomes, Gloxysomal, Glycosomal
plasma membrane	Integral membrane
Vacuolar	Vacuolar, Vacuole

가

ProSLP

가

kNN

k

1 86,720

가

10,840

2 43,240

8,645

2.5 (Performance Measure)

(Precision)

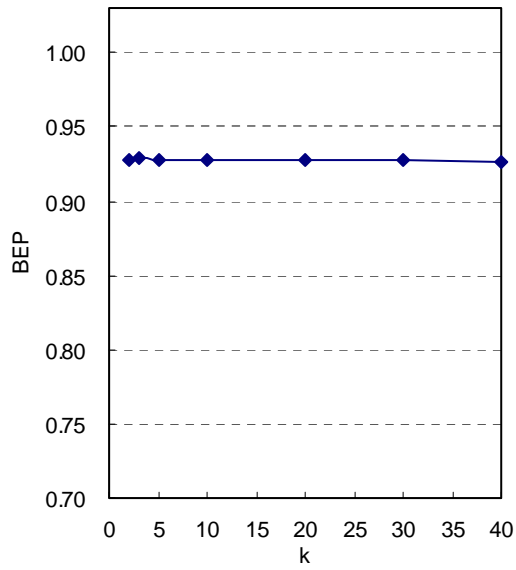
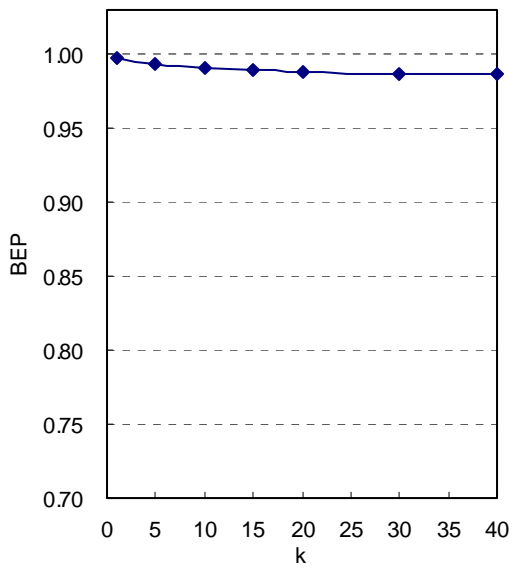
(Recall)

(p)

(r)

1, BeP 3 4
 1 k 가 k=1 가
 BeP 0.987 0.998

2 k BeP 가



1. 1(a) 2(b) ProSLP

ProSLP 60%~85% 93%

ProSLP 86,720 43,240
 ProSLP

가

ProSLP 가 k 1 3 10
 가 k 10 가

k=1 가 k

3. 1 k
1 86,720 10,840 가

k	Precision	Recall	BeP ^a
1	0.997	0.999	0.998 ^b
5	0.994	0.994	0.994
10	0.991	0.991	0.991
15	0.989	0.989	0.989
20	0.988	0.988	0.988
30	0.987	0.987	0.987
40	0.987	0.987	0.987

^a Micro-averaged break-even points

^b Micro-averaged F₁ measure

4. 2 k
2 43,240 8,645 가

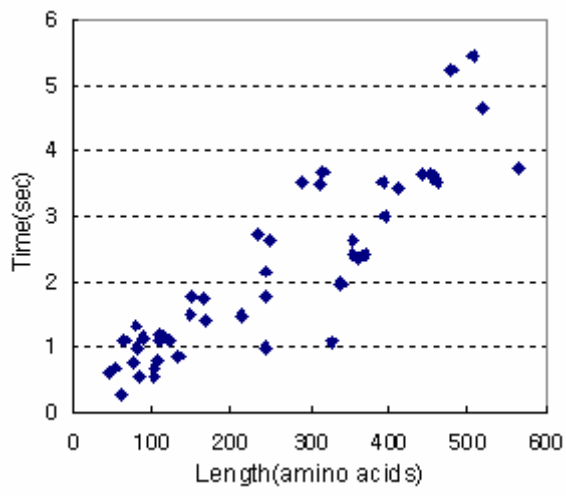
k	Precision	Recall	BeP ^a
1	0.9292	0.9302	0.9297 ^b
2	0.9282	0.9282	0.9282
3	0.9292	0.9292	0.9292
5	0.9282	0.9282	0.9282
10	0.9279	0.9279	0.9279
20	0.9272	0.9272	0.9272
30	0.9271	0.9271	0.9271
40	0.9269	0.9269	0.9269

^a Micro-averaged break-even points

^b Micro-averaged F₁ measure

ProSLP 가 2

ProSLP (Penta-gram) /-5+1 가 (Feature Selection) - ProSLP (penta-gram) n- (gram)



2. 48 53~600 262

4.

가

n- (gram)

kNN

가

ProSLP 가 가

ProSLP

<http://proslp.kisti.re.kr>

ProSLP

ProSLP

가

ProSLP

5.

- [1] Bairoch, A. and Apweiler, R. (2000) The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res.*, **28**, 45-48
- [2] Chou, K. C. and Elrod, D. W. (1999) Protein subcellular locations prediction. *Protein Engineering*, **12**, 107-118
- [3] Eisenhaber, F. and Bork, P. (1998) Evaluation of human-readable annotation in biomolecular sequence database with biological rule libraries. *bioinformatics*, **15**, 528-535
- [4] Emanuelsson, O., Nielsen, H., Brunak, S., and Heijne, G. V. (2000) Predicting Subcellular Localization of Proteins Based on their N-terminal Amino Acid Sequence. *J. Mol. Biol.*, **300**, 1005-1016
- [5] Feng, Z. P. (2002) An overview on predicting the subcellular location of protein. *In silico Biol.*, **2**, 291-303
- [6] Gardy, J. L., Spencer, C., Wang, K., Ester, M., Tusnady, G. E., Simon, I., Hua, S., deFays, K., 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
- [7] Park, K. J. and Kanehisa, M. (2003) Prediction of protein subcellular locations by support vector machines using compositions of amino acids and amino acid pairs. *Bioinformatics*, **19**, 1656-1663
- [8] Nakai, k. and Horton, P. (1999) PSORT: a program for detecting the sorting signals of proteins and predicting their subcellular localization. *Trends Biochem. Sci.*, **24**, 34-35
- [9] Reinhardt, A. and Hubbard, T. (1998) Using neural networks for prediction of the subcellular location of proteins. *Nucleic Acids Res.*, **26**, 2230-2236
- [10] Sebastiani, F. (2002) Machine Learning in Automated Text categorization. *ACM computing Surveys*, **34**, 1-47

van Rijsbergen, C. (1979) Information Retrieval. Butterworths, London, 1979

[11] Yuan, Z. (1999) Prediction of protein subcellular locations using Markov chain models. *FEBS Letters*, **451**, 23-26

[12] Yang, Y. (1994) Expert Network: Effective and Efficient Learning from Human Decisions in Text Categorization and Retrieval. *Proceeding of the 17th Annual International ACM/SIGIR conferences on Research and Development in Information Retrieval*, 13-22

[13] Yang, Y (1999) An Evaluation of Statistical Approaches to Text Categorization. *Journal of Information Retrieval*, **1**, 67-88