



PROTEINS:
Structure, Function, and Bioinformatics

PSLDoc: Protein subcellular localization prediction based on gapped-dipeptides and probabilistic latent semantic analysis

Journal:	<i>PROTEINS: Structure, Function, and Bioinformatics</i>
Manuscript ID:	Prot-00301-2007.R2
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Chang, Jia-Ming; Academia Sinica, Institute of Information Science Su, Emily; Academia Sinica, Taiwan International Graduate Program Lo, Allan; Academia Sinica, Taiwan International Graduate Program Chiu, Hua-Sheng; Academia Sinica, Institute of Information Science Sung, Ting-Yi; Academia Sinica, Institute of Information Science Hsu, Wen-lian; Academia Sinica, Institute of Information Science
Key Words:	protein subcellular localization, document classification, vector space model, gapped-dipeptides, probabilistic latent semantic analysis, support vector machines



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

PSLDoc: Protein subcellular localization prediction based on gapped-dipeptides and probabilistic latent semantic analysis

Jia-Ming Chang¹, Emily Chia-Yu Su^{2,3}, Allan Lo^{2,4}, Hua-Sheng Chiu¹, Ting-Yi Sung¹, Wen-Lian Hsu¹

¹Bioinformatics Lab., Institute of Information Science, Academia Sinica, Taipei, Taiwan

²Bioinformatics Program, Taiwan International Graduate Program, Academia Sinica, Taipei, Taiwan

³Institute of Bioinformatics, National Chiao Tung University, Hsinchu, Taiwan

⁴Department of Life Sciences, National Tsing Hua University, Hsinchu, Taiwan

Email: {jmchang, cysu, allanlo, huasheng, tsung, hsu}@iis.sinica.edu.tw

Corresponding Author

Dr. Wen-Lian Hsu

Institute of Information Science

Academia Sinica

128 Academia Road, Section 2

Nankang, Taipei, Taiwan, ROC

Tel: +886-2-27883799 ext.1804

Tel: ++886-2-27824814

Email: hsu@iis.sinica.edu.tw

Keywords: protein subcellular localization; document classification; vector space model; gapped-dipeptides; probabilistic latent semantic analysis; support vector machines

ABSTRACT

Prediction of protein subcellular localization (PSL) is important for genome annotation, protein function prediction, and drug discovery. Many computational approaches for PSL prediction based on protein sequences have been proposed in recent years for Gram-negative bacteria. We present PSLDoc, a method based on gapped-dipeptides and probabilistic latent semantic analysis (PLSA) to solve this problem. A protein is considered as a term string composed by gapped-dipeptides, which are defined as any two residues separated by one or more positions. The weighting scheme of gapped-dipeptides is calculated according to a position specific score matrix, which includes sequence evolutionary information. Then, PLSA is applied for feature reduction, and reduced vectors are input to five one-versus-rest support vector machine classifiers. The localization site with the highest probability is assigned as the final prediction. It has been reported that there is a strong correlation between sequence homology and subcellular localization^{1,2}. To properly evaluate the performance of PSLDoc, a target protein can be classified into low- or high-homology data sets. PSLDoc's overall accuracy of low- and high-homology data sets reaches 86.84% and 98.21%, respectively, and it compares favorably with that of CELLO II². In addition, we set a confidence threshold in order to achieve a high precision at specified levels of recall rates. When the confidence threshold is set at 0.7, PSLDoc

1
2
3
4 achieves 97.89% in precision which is considerably better than that of PSORTb v.2.0³.

5
6
7 Our approach demonstrates that the specific feature representation for proteins can be
8
9
10 successfully applied to the prediction of protein subcellular localization and improves
11
12
13 prediction accuracy. Besides, because of the generality of the representation, our
14
15
16 method can be extended to eukaryotic proteomes in the future. The web server of
17
18
19 PSLDoc is publicly available at <http://bio-cluster.iis.sinica.edu.tw/~bioapp/PSLDoc/>.

20 21 22 23 24 25 26 **INTRODUCTION**

27 28 29 **Protein Subcellular Localization Prediction**

30 31 32 *Background*

33
34
35 Predicting protein subcellular localization (PSL) is key to elucidating many
36
37
38 biological problems, such as protein function prediction, genome annotation and drug
39
40
41 discovery. The task is to assign a protein to one or more localization sites
42
43
44 corresponding to the subcellular compartments based on its sequence. Recently, many
45
46
47 prediction methods for Gram-negative bacteria have been developed using different
48
49
50 computational techniques, including expert system⁴, *k*-nearest neighbors⁵, artificial
51
52
53 neural networks^{6,7}, support vector machines (SVM)^{2,8-13}, and Bayesian networks^{3,14,15}.
54
55
56 Among them, PSORTb v.2.0³ (updated from PSORTb v.1.1¹⁴) and CELLO II²
57
58
59 (updated from CELLO¹³) have been tested on a new Gram-negative bacteria data set¹⁶.
60

1
2
3
4 PSORTb v.1.1, released in 2003, integrates homology analyses, identification of
5
6
7
8 sorting signals and other motifs, and machine learning methods into an expert system
9
10 based on a Bayesian network to decide the final prediction. PSORTb v.2.0, released in
11
12
13
14 2005, uses SVM as the underlying machine learning model and takes frequent
15
16
17 subsequences occurring in proteins as input features. CELLO also uses SVM trained
18
19
20 by multiple feature vectors derived from n -peptide compositions. The updated
21
22
23 CELLO II is based on a two-level SVM system: the first-level SVM is comprised a
24
25
26 number of SVM classifiers using different feature vectors, and each classifier
27
28
29 generates a probability distribution of subcellular localization sites; the second-level
30
31
32 SVM is considered as a jury SVM that yields a final probability distribution based on
33
34
35 those generated in the previous stage and determines the final prediction as the site of
36
37
38 the highest probability. The authors of CELLO II also classify a query protein, whose
39
40
41 localization site is to be predicted, into low- or high-homology data sets depending on
42
43
44 its highest pairwise sequence identity with the training data set whether it is below or
45
46
47 above a similarity threshold of 30%. This classification of data is motivated by an
48
49
50 observation that sequence homology and subcellular localization have strong
51
52
53 correlation when the sequence identity is higher than 30%. Hence, they also propose a
54
55
56 hybrid method, called HYBRID², which uses the two-level SVM system for
57
58
59 low-homology proteins and a homology search method for high-homology proteins.
60

Document Classification Approach

In this paper, we formulate PSL prediction as a document classification problem. The document classification problem is to assign an electronic document to one or more categories, based on its contents. A protein sequence can be considered as the content of a document, and localization sites are considered as categories. To predict the localization site(s) of a protein is equivalent to predicting the category (e.g., sport, politics) of a document (e.g., a piece of news). This transformation is intuitive. Document classification methods have been successfully applied in many protein classification problems, such as protein function prediction¹⁷ and protein family classification¹⁸. King and Guda showed that using document classification techniques on the primary sequence can achieve good results on estimating subcellular proteomes of eukaryotes¹⁹.

Given a large number of documents, document classification is usually tackled by the following three steps. First, documents have to be transformed into feature vectors in which each distinct term corresponds to a feature. The value of a feature in a vector represents the weight of a term in a document. Another set of documents with known categories is used as a training set. Second, because of high-dimensional feature spaces, feature reduction is necessary before applying machine learning

1
2
3
4 methods, to improve generalization accuracy²⁰ and to avoid overfitting^{20,21}. The first
5
6
7 two steps could be considered as *feature representation*. Finally, these reduced feature
8
9
10 vectors are used to perform the category assignment automatically.
11

12
13
14 In this paper, we propose a specific feature representation embedded in a
15
16 prediction system called PSLDoc (Protein Subcellular Localization prediction based
17
18 on modified Document classification method), which uses SVM as the underlying
19
20 machine learning model. The design of PSLDoc's feature representation includes the
21
22 following tasks: (1) define the terms of a protein; (2) design a term weighting scheme;
23
24 and (3) apply a feature reduction and extraction method.
25
26
27
28
29
30

31
32 For a benchmark dataset of Gram-negative bacteria¹⁶, PSLDoc performs better
33
34 than HYBRID and PSORTb v.2.0. Our approach demonstrates that the specific feature
35
36 representation for proteins can be successfully applied to PSL prediction.
37
38
39
40
41
42
43

44 **A Baseline System Using TFIDF**

45
46
47 Before describing our method, we introduce a baseline system for performance
48
49 comparison that uses a traditional document classification method. Salton's vector
50
51 space model (VSM) is one of the most widely used methods for ad-hoc retrieval in
52
53 document classification²². Each document is represented by a feature vector (vector,
54
55 in short) composed of all terms in a collection of documents, where each entry (or
56
57
58
59
60

feature) of the vector corresponds to a term and its value is given by the weight of the term in the document²³. The similarity between two documents d and q , denoted by $sim(d,q)$, can be defined as the cosine of the angle between their vectors, called *cosine similarity*, as shown below:

$$sim(d,q) = \cos(\angle(\vec{d}, \vec{q})) = \frac{\vec{d} \cdot \vec{q}}{|\vec{d}| |\vec{q}|} \quad (1)$$

where \vec{d} denotes the vector for a document d . Given a collection of documents with known categories, we classify a document with unknown category (called *query document*) into the same category as the document whose cosine similarity with the query document is the largest. We refer to this prediction method as the *1-Nearest Neighboring* (1-NN) method based on cosine similarity. The advantage of the 1-NN method is that there is no training required as in general machine learning approach.

Weighting scheme, i.e., determining the weight of each entry in a vector, is crucial in document classification. In this baseline system, we use *term frequency–inverse document frequency* (TFIDF) as the weighting scheme. For a term t_i in a document d , a simple *term frequency* (TF) is the number of t_i 's occurrences in the document, denoted as n_i . However, to prevent a bias toward longer documents, term frequency $tf(t_i,d)$ is usually normalized as follows:

$$tf(t_i,d) = \frac{n_i}{\sum_k n_k}, \quad (2)$$

where the denominator is the number of occurrences of all terms. The term frequency $tf(t_i, d)$ gives a measure of the importance of the term t_i in the document d . The higher the term frequency, the more likely the term is a good description of the content of the document. In contrast, *inverse document frequency* (IDF) of t_i is a measure of the general importance of the term. A semantically important term will often occur several times in a document if it occurs at all. However, semantically unimportant terms are spread out homogeneously over all documents. A frequently used IDF for t_i , $idf(t_i)$, is defined as follows:

$$idf(t_i) = \log \frac{|D|}{|(d_i \supset t_i)|}, \quad (3)$$

where $|D|$ is the number of documents in the collection, and $|(d_i \supset t_i)|$ denotes the number of documents in which t_i appears. In the TFIDF scheme, the weight of the term t_i in a document d , $W(t_i, d)^a$, equals to $tf(t_i, d)$ multiplied by $idf(t_i)^{24}$. The values in a vector are normalized to (0~1] by dividing the maximum value in the vector.

METHODS

PSLDoc uses gapped-dipeptides²⁵ as the terms of a protein and calculates their weights according to a position specific score matrix (PSSM) instead of the TFIDF used in the baseline system. Probabilistic latent semantic analysis (PLSA) is used for

^a In this paper, we use the weights of the terms in a document and in a vector, denoted by $W(t_i, d)$ and $W(t_i, \vec{d})$, interchangeably.

1
2
3
4 feature reduction to improve learning efficiency and accuracy. The reduced feature
5
6
7
8 vectors are input to five one-versus-rest (1-v-r) SVM classifiers corresponding to five
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
localization sites. The probability estimated by a classifier can be considered as the
confidence level of a target protein belonging to the corresponding localization site.
The final prediction is determined to be the localization site whose corresponding
classifier outputs the largest confidence score.

Gapped-dipeptides As the Terms of Proteins

When considering proteins as documents, many different types of terms have
been proposed, including single amino acid (AA)^{6,7,9,11,26-28} as a uni-gram descriptor,
and the general n -peptide¹³, i.e., peptides of length n without gaps. In particular, for n
 $= 2$, dipeptide (Dip) is a neighboring bi-gram descriptor. However, AA and Dip cannot
represent information between two gapped amino acids. The use of n -peptide to
capture long distance amino acid information will result in a high-dimensional vector
space. For example, the feature number of a vector is 3,200,000 ($= 20^5$), when n
equals five. “Gapped amino acid pair” was first proposed by Park and Kanehisa¹¹ for
protein representation. Later, Liang *et al.*²⁵ proposed a method based on a similar
encoding scheme, called amino acid-coupling patterns, to extract the information from
a protein sequence; the encoding scheme works well on distinguishing thermophilic

1
2
3
4 proteins. An amino acid-coupling pattern XdZ denotes the peptides of length $d + 2$
5
6
7 such that amino acids X and Z are separated by d amino acids, where d can be
8
9
10 negative depending on whether the position of X is closer to N-terminus or
11
12
13 C-terminus²⁵.
14

15
16 We adopt the same encoding scheme as in Liang *et al.* except with non-negative
17
18
19 d as the term of a protein sequence regardless whether the pattern appears near the
20
21
22 N-terminus or C-terminus. We call such amino acid-coupling pattern as
23
24
25 *gapped-dipeptides*. For example, the gapped-dipeptides for $d=0$ are dipeptide without
26
27
28 gaps (Dip's). Given a positive integer l as the upper bound of gapped distance, each
29
30
31 protein sequence is represented by a vector in the space of gapped-dipeptides with
32
33
34 each feature given by XdZ for $0 \leq d \leq l$. The length of vectors is the number of all
35
36
37 possible combinations of gapped-dipeptides, i.e., $(l+1) \times 20 \times 20$. For example, given
38
39
40
41 $l=10$, a protein is represented as a feature vector of 4,400 ($=11 \times 20 \times 20$) features.
42
43
44
45
46
47
48
49

50 Term Weighting - Position Specific Score Matrix Information

51 *Motivation*

52
53
54 Based on the finding in a previous work that sequence identity and subcellular
55
56
57 localizations of proteins have a strong correlation¹, Yu *et al.*² proposed a homology
58
59
60 search method for PSL prediction, which predicted the localization site of a query

1
2
3
4 protein by the most similar protein among the aligned protein sequences with known
5
6
7 localization sites generated by the global alignment program ALIGN²⁹. The authors
8
9
10 observed that, when the query protein and its most similar protein with known
11
12
13 localization site have sequence identity over 30%, the homology search method
14
15
16 performed very well with 97.7% accuracy. But the prediction performance dropped
17
18
19 significantly when the sequence identity is under 20%. In this case, it would be
20
21
22 difficult to predict the localization site of a query protein based on the sequence
23
24
25 identity or sequence information. To overcome this difficulty, we borrow the idea
26
27
28 from protein secondary structure prediction, in which homologous sequences are
29
30
31 usually removed from the testing and training data sets³⁰⁻³⁵.

32
33
34
35 Most of the prediction methods address the problem of **weak homology** by
36
37
38 utilizing sequence evolutionary information. One widely used representation of
39
40
41 evolutionary information is the PSSM generated by PSI-BLAST³⁶, which has been
42
43
44 used in PSIPRED³⁷, a very popular secondary structure prediction method.
45
46
47 PSI-BLAST finds remote homologues to a query protein from a chosen sequence
48
49
50 database (e.g., NCBI nr³⁸). Instead of TFIDF based on the sequence information, our
51
52
53 weighting scheme is based on PSSM.
54
55
56
57
58
59
60

Position specific score matrix

1
2
3
4 The PSSM of a sequence S of length n is represented by an $n \times 20$ matrix, in
5
6
7 which the n rows correspond to the amino acid sequence of S and the columns
8
9
10 correspond to the 20 distinct amino acids. Each row of a PSSM represents the
11
12
13 log-likelihood of the residue substitutions at the corresponding positions in S ³⁶. The
14
15
16 PSSM elements are normalized to the range from 0 to 1 using the following sigmoid
17
18
19 function³²:

$$f(x) = \frac{1}{1 + e^{-x}}, \quad (4)$$

22
23
24 where x is the original PSSM value. The higher the normalized value of the residue is,
25
26 the higher it is for the propensity of the residue in this position. In PSLDoc, the
27
28
29 PSI-BLAST's parameters were set to $j = 5$ (five iterations), $e = 10^{-2}$ (E -value < 0.01)
30
31
32 and the sequence database was NCBI nr which contains 3,747,820 sequences.
33
34
35
36
37
38
39
40

41 ***TFPSSM weighting scheme***

42
43
44 We design a term weighting scheme based on PSSM, denoted by TFPSSM as
45
46
47 follows. Given a protein sequence S of length n , any gapped-dipeptide XdZ of S has
48
49
50 PSSM entries corresponding to gapped-dipeptides $S(i)dS(i+d+1)$ for $1 \leq i \leq n-(d+1)$,
51
52
53 where $S(i)$ denotes the i th amino acid of S . For example, the PSSM (with original
54
55
56 value without normalization) of the sequence MPLDLYNTLT is shown in Figure 1.
57
58
59 From the sequence information, M2D only occurs once. However, in view of PSSM,
60

M2D may occur in the corresponding gapped-dipeptides obtained from the sequence, i.e., M2D, P2L, L2Y, D2N, L2T, Y2L, N2T. We define the weight of XdZ in S as

$$W(XdZ, S) = \sum_{1 \leq i \leq n-(d+1)} f(i, X) \times f(i+d+1, Z) \quad (5)$$

where $f(i, Y)$ denotes the normalized value of the PSSM entry at the i th row and the column corresponding to amino acid Y . In the above example, the weight of M2D based on PSSM is given by $f(1, M) \times f(4, D) + f(2, M) \times f(5, D) + \dots + f(7, M) \times f(10, D) = 0.99995 \times 0.04743 + 0.11920 \times 0.00247 + \dots + 0.00669 \times 0.26894$. It is unnecessary to incorporate IDF with term weighting based on PSSM because the term occurs in all documents based on PSSM.

(Figure1)

As mentioned before, each protein is represented by a vector, and each entry of the vector is given by TFPSSM of the corresponding gapped-dipeptides. Note that the values in each feature vector are normalized between (0~1] by dividing the maximum value in the vector.

Feature Reduction - Probabilistic Latent Semantic Analysis

Motivation

There are some limitations of the VSM for document classification. First, the vector space is high-dimensional²¹. Training and testing have to deal with the curse of

1
2
3
4 dimensionality. Second, document vectors are typically very sparse, i.e., most features
5
6
7 of a vector are zeros that are susceptible to noises³⁹, and cosine similarity could be
8
9
10 inaccurate. Finally, the inner product defining document similarity can only match
11
12
13 occurrences of the same terms. As a result, the vector representation does not capture
14
15
16 semantic relations between terms. Furthermore, this representation, which considers a
17
18
19 document as a bag of words, is unable to capture phrases and semantic/syntactic
20
21
22 regularities.
23

24
25
26 Hence, dimension reduction (feature reduction) is proposed for dealing with the
27
28
29 above limitations. The task of dimension reduction is to map similar terms to a similar
30
31
32 location in a low dimensional space called *latent semantic space*, which reflects
33
34
35 semantic associations. A frequently used dimension reduction technique is Latent
36
37
38 Semantic Analysis (LSA) (or called Latent Semantic Indexing in some papers), which
39
40
41 uses singular value decomposition (SVD) to do data mapping⁴⁰. The document
42
43
44 similarity based on the inner product is computed on the latent semantic space.
45
46
47 Experimentally, there are advantages of SVD over naive VSM. However, SVD still
48
49
50 has the following disadvantages⁶. First, the resulting dimensions might be difficult to
51
52
53 interpret. For instance, the size of a vector is reduced from three to two by LSA as
54
55
56 shown below;
57
58
59

$$\{(A0A), (A1A), (G0G)\} \rightarrow \{(1.3 * A0A + 0.2 * G0G), (A1A)\}$$

The value of the first reduced feature equals 1.3 multiplied by the value of the original first feature plus 0.2 multiplied by the value of the original third feature. This leads to results which might be justifiable on the mathematical level, but have no interpretable meaning in the original application. Second, the probabilistic model of LSA does not match observed data⁴¹. Third, the reconstruction may contain negative entries, which are inappropriate as a distance function for count vectors.

Probabilistic latent semantic analysis

Hofmann proposed probabilistic latent semantic analysis (PLSA) based on an aspect model to deal with those above disadvantages⁴¹. The aspect model is a latent variable model for co-occurrence data (i.e., documents and terms) that each observation is associated with an unobserved class variable $z \in Z = \{z_1, \dots, z_K\}$. The weight of the term w in a document d , $W(w, d)$, is considered as a joint probability $P(w, d)$ between w and d , which is modeled by z , a latent variable which can be loosely thought of as a topic or a reduced feature. Thus, the joint probability $P(w, d)$ based on PLSA model is

$$P(w, d) = P(d)P(w|d), P(w|d) = \sum_{z \in Z} P(w|z)P(z|d)^b \quad (6)$$

where $P(w|z)$ denotes the topic-conditional probability of a term conditioned on the

^bIt is assumed that the distribution of terms given a class is conditionally independent of the document, i.e., $P(w|z, d) = P(w|z)$.

unobserved topic, and $P(z|d)$ denotes a document-specific probability distribution over the latent variable space; that is, considering a vector \vec{d} in latent variable space, $P(z|d)$ denotes the weight of the latent variable z of the document d . Hence, a vector is mapped from the term space to latent space and its size is reduced from $|W|$ to $|Z|$.

PLSA model fitting (training)

A PLSA model is parameterized by $P(w|z)$ and $P(z|d)$ which are estimated by fitting $P(w, d)$ to a training corpus D with known $W(w, d)$. The fitting process is obtained by maximizing the log-likelihood function L given below⁴¹:

$$L = \sum_{w \in d} \sum_{d \in D} W(w, d) \log P(w, d) \tag{7}$$

The parameters of a PLSA model, $P(w|z)$ and $P(z|d)$, are estimated using the iterative Expectation-Maximization (EM) algorithm by maximizing the log-likelihood function L . $P(w|z)$ and $P(z|d)$ are initialized by random values in (0,1)-range. Then, the EM procedure iterates between the E-step and the M-step. In the E-step, the probability that a term w in a particular document d explained by the class corresponding to z , is estimated as

$$P(z | w, d) = \frac{P(z, w, d)}{P(w, d)} \tag{8}$$

$$P(z, w, d) = P(d)P(z | d)P(w | z)^c \tag{9}$$

^c This equation is derived from according to the Figure 1(a) of Hofmann³⁹.

Using Equations (6), (8) and (9), we can get

$$P(z | w, d) = \frac{P(d)P(z | d)P(w | z)}{P(d)\sum_{z'} P(w | z')P(z' | d)} = \frac{P(w | z)P(z | d)}{\sum_{z'} P(w | z')P(z' | d)} \quad (10)$$

In the M-step, we calculate

$$P(w | z) = \frac{\sum_d W(w, d)P(z | w, d)}{\sum_{w'} \sum_d W(w', d)P(z | w', d)}$$

$$P(z | d) = \frac{\sum_w W(w, d)P(z | w, d)}{\sum_{z'} \sum_w W(w, d)P(z' | w, d)} \quad (11)$$

where parameters $P(w|z)$ and $P(z|d)$ are re-estimated to maximize L .

PLSA model testing

After training, the estimated $P(w|z)$ parameters are used to estimate $P(z|q)$ for new (test) documents q through a *folding-in* process⁴¹. In the folding-in process, EM procedure runs in a similar manner to the training stage. The E-step is identical but the M-step keeps all the $P(w|z)$ constant and only re-calculates $P(z|q)$. Usually, a very small number of iterations of the EM algorithm are sufficient for folding-in process.

Feature reduction by PLSA

We apply PLSA not only for feature reduction but also for gapped-dipeptides semantic relation extraction. Vectors are mapped from the gapped-dipeptides space to the latent semantic space. This will lead to improvement in learning performance and

1
2
3
4 efficiency. Though it is not easy to determine an appropriate reduced feature size of
5
6
7
8 PLSA, it can be approximated by the reduced feature size of LSA. To determine the
9
10
11 reduced feature size of LSA, we calculate singular values of LSA and sort them in a
12
13
14 decreasing order. Then, the reduced feature size of LSA equals to n if the n -th largest
15
16
17 singular value is close to zero.

22 **The System Architecture of PSLDoc**

23
24
25
26 Prediction of PSL can be treated as a multi-class classification problem. For
27
28
29 multi-class classification, the 1-v-r SVM model has demonstrated a good
30
31
32 classification performance²⁸. For each class i , we construct a 1-v-r (C_i versus non- C_i)
33
34
35 binary classifier. PSLDoc consists of five 1-v-r SVM classifiers corresponding to five
36
37
38 localization sites in Gram-negative bacteria. Input features for all binary classifiers are
39
40
41 the same. The SVM program LIBSVM⁴² is used in PSLDoc, and it can generate
42
43
44 probability estimates that are used for determining the confidence levels of
45
46
47 classifications⁴³. For all classifiers, we use the Radial Basis Function kernel, and tune
48
49
50 the cost (c) and gamma (γ) parameters optimized by 10-fold cross-validation on the
51
52
53 training data set.
54

55
56
57 Given a protein, PSLDoc performs the following steps:

- 58
59
60 1. Use PSI-BLAST to generate PSSM of the protein.

2. Generate the feature vector of the protein, where each feature is defined as TFPSSM corresponding to a gapped-dipeptides.
3. Perform PLSA to generate a reduced feature vector, which will be input to each 1-v-r classifier.
4. Run five 1-v-r SVM classifiers.

In the training stage of PSLDoc, to train PLSA model with different topic sizes and the SVM classifiers, proteins with known localization sites are used to estimate $P(w|z)$ and $P(z|d)$, and the reduced vectors are used to determine the c and γ parameters of the RBF kernel of each classifier. In the testing stage of PSLDoc, Step 3 of PSLDoc performs PLSA folding-in process on trained $P(w|z)$. Step 4 of PSLDoc is performed on the trained SVM classifiers. The localization site of the protein is predicted as the class with the highest probability ($prob_i$: the confidence of the query protein predicted as class i ; $0 \leq prob_i \leq 1$) generated from the five 1-v-r classifiers. The system architecture of PSLDoc is shown in Figure 2.

(Figure2)

Data Sets

To evaluate the performance of PSLDoc, we utilize a benchmark data set of proteins from Gram-negative bacteria with single localization that have been used in previous works^{3,13}. It consists of 1444 proteins with experimentally determined

1
2
3
4 localizations, referred to as PS1444¹⁶. Table I lists the distribution of localization sites
5
6
7 of the data set.
8

9
10 (Table I)
11

12
13 To analyze the performance of PSLDoc under the effect of sequence homology
14 information, we further classify each protein in PS1444 into two data sets, the high- or
15
16 low-homology data sets based on whether or not the protein's highest sequence
17
18 identity of all-against-all alignment by ClustalW is greater than an identity threshold
19
20 of 30%. The high-homology data set, referred to as PSHigh783, consists of 783
21
22 proteins and the low-homology set, referred to as PSLow661, consists of 661 proteins.
23
24
25

26
27 The three data sets are available at
28
29
30
31
32 <http://bio-cluster.iis.sinica.edu.tw/~bioapp/PSLDoc/DataSet.htm>.
33
34
35

36 37 38 **Evaluation Measures** 39

40
41 To evaluate the performance of our method, we follow the same measures used
42
43 in previous works^{4,12-14} for comparison with other approaches. These measures
44
45 include accuracy (*Acc*), precision, recall, Matthew's correlation coefficient (*MCC*)⁴⁴
46
47 for five localization sites, and the overall accuracy defined in Eq. (12), (13), (14), (15)
48
49 and (16) below:
50
51
52
53
54

$$55 \quad Acc_i = TP_i / N_i \quad (12)$$

$$56 \quad Precision_i = TP_i / (TP_i + FP_i) \quad (13)$$

$$Recall_i = TP_i / (TP_i + FN_i) \quad (14)$$

$$MCC_i = \frac{(TP_i)(TN_i) - (FP_i)(FN_i)}{\sqrt{(TP_i + FN_i)(TP_i + FP_i)(TN_i + FP_i)(TN_i + FN_i)}} \quad (15)$$

$$Acc = \sum_{i=1}^l TP_i / \sum_{i=1}^l N_i, \quad (16)$$

where $l = 5$ is the number of localization sites, and TP_i , TN_i , FP_i , FN_i , and N_i are the number of true positives, true negatives, false positives, false negatives, and proteins in localization site i , respectively. MCC considers both under- and over-predictions, and takes range from -1 to 1 , where $MCC = 1$ indicates a perfect prediction; $MCC = 0$ indicates a completely random assignment; and $MCC = -1$ indicates a perfectly reverse correlation. The Acc_i is the same as $Recall_i$ because N_i equals to the sum of TP_i and FN_i . We will use Acc_i or $Recall_i$ interchangeably in the experiments depending on which method is compared.

Five Simple PSL Prediction Methods

To evaluate the benefit of each step in our document classification method, we propose two simple prediction methods: 1NN_TFIDF and 1NN_TFPSSM, which consist of different parts of PSLDoc. To further analyze the effect of the PSSM information generated from databases of different sizes, we propose two methods based on PSI-BLAST: 1NN_PSI-BLAST_{ps} and 1NN_PSI-BLAST_{nr}. In addition, we

1
2
3
4 also construct a homology search method, 1NN_ClustalW, which is similar to Yu *et*.
5
6
7
8 *al.*'s for comparison with PSLDoc.

13 ***1NN_TFIDF***

16 1NN_TFIDF solely incorporates protein encoding scheme, the gapped-dipeptides
17
18 of PSLDoc. The remaining steps are the same as the baseline system. That is, terms
19
20 are weighted according to the TFIDF weighting scheme, and a query protein is
21
22
23
24
25
26 predicted by 1-NN method based on cosine similarity.

32 ***1NN_TFPSSM***

35 1NN_TFPSSM incorporates two parts of PSLDoc, the gapped-dipeptide
36
37 encoding scheme and the TFPSSM weighting scheme. It predicts a query protein
38
39
40
41 using 1-NN method based on cosine similarity.

44 ***1NN_PSI-BLAST_{ps}***

47 1NN_PSI-BLAST_{ps} performs two PSI-BLAST searches, one of which for
48
49 generating a PSSM and the other for searching the most similar protein using the
50
51 PSSM generated in the previous step. First, for each query protein, PSI-BLAST
52
53
54 search is performed against the training data and its parameters are the same as those
55
56
57
58
59
60 in PSLDoc. Then, 1NN_PSI-BLAST_{ps} performs a one-run PSI-BLAST search (i.e., j

1
2
3
4 = 1)^d against the training data using the obtained PSSM^e. Finally, the localization site
5
6
7 of the protein with the highest e-value is assigned as the predicted localization for the
8
9
10 query protein. In a five-fold cross-validation, the PSSM information used in
11
12
13 1NN_PSI-BLAST_{ps} is generated from a small database which consists of
14
15
16 approximately 1,155 (=1,444×4/5) sequences from PS1444.
17
18

19 20 21 22 23 ***1NN_PSI-BLAST_{nr}*** 24

25
26 Although 1NN_PSI-BLAST_{ps} utilizes the PSSM information, the source
27
28 database used is not as large as that of 1NN_TFIDF and PSLDoc. For fair comparison
29
30 with 1NN_TFIDF and PSLDoc, we construct 1NN_PSI-BLAST_{nr} which uses PSSM
31
32 generated from the NCBI nr database. The only difference between
33
34 1NN_PSI-BLAST_{nr} and 1NN_PSI-BLAST_{ps} is the size of the databases searched in
35
36 the first step, and the remaining steps are all the same, including the generation of
37
38 PSSM, followed by performing a second PSI-BLAST search, and lastly, the
39
40 prediction of the localization site of the query protein.
41
42
43
44
45
46
47
48
49
50

51 52 53 54 ***1NN_ClustalW*** 55

56
57 1NN_ClustalW differs from Yu *et al.*'s method only in the pairwise sequence
58

59 ^d The parameters of e-value are ignored because we want to find the most similar protein instead of
60 constructing a PSSM.

^e Please refer to the last example on blastpgp's document for how to save a PSSM and perform
PSI-BLAST search from the PSSM (<http://biowulf.nih.gov/apps/blast/doc/blastpgp.html>).

1
2
3
4 alignment algorithm used, i.e., ClustalW in the former and ALIGN in the latter. For a
5
6
7 query protein, we calculate its pairwise sequence identities with the remaining
8
9
10 proteins by performing 1-against-others pairwise sequence alignment. Then, the
11
12
13 localization site of the query protein is predicted by the 1-NN method based on
14
15
16 pairwise sequence identity, that is, its localization site is assigned as that of the protein
17
18
19 whose pairwise sequence identity is highest.
20
21
22
23
24
25

26 **Experiment Design**

27
28
29 We conduct the following experiments to evaluate the benefit of each step in our
30
31
32 document classification model where the gapped distance upper bound, l , ranges from
33
34
35 3 to 15. We follow the same validation procedures for the performance measurement
36
37
38 as those of the other approaches^{2,3}. All experiments are carried out in five-fold
39
40
41 cross-validation, that is, the data is equally divided into five parts. In each run, four
42
43
44 folds are used for training and the remaining fold is used for testing. All reported
45
46
47 results are average over the five folds. We have conducted the following six
48
49
50 experiments:
51
52
53
54
55
56

57 ***Experiment 1: Comparison between 1NN_TFIDF and 1NN_TFPSSM on the***
58
59 ***PS1444, PSHigh783, and PSLow661 data sets***
60

1
2
3
4 The purpose of this experiment is to evaluate the benefit of using the TFPSSM
5
6
7 weighting scheme because the simple 1NN prediction method can reflect the relation
8
9
10 between performance and weighting schemes avoiding the effect of the prediction
11
12
13 algorithm. The distribution of benefit among 1444 protein sequences is further
14
15
16 analyzed by comparing their performance on PSHigh783 and PSLow661.
17
18
19
20
21
22

23 ***Experiment 2: Comparison among 1NN_TFPSSM, 1NN_ClustalW,***
24 ***1NN_PSI-BLAST_{ps}, and 1NN_PSI-BLAST_{nr} on the PSHigh783 and***
25 ***PSLow661 data sets***
26
27
28
29
30
31

32 To compare the effect of utilizing PSSM, we compare the performance of
33
34 1NN_TFPSSM, 1NN_ClustalW, 1NN_PSI-BLAST_{ps}, and 1NN_PSI-BLAST_{nr}.
35
36 1NN_ClustalW is based on a pairwise sequence alignment in which no PSSM
37
38 information is incorporated. We further analyze the relationship between the effect of
39
40 PSSM and the size of databases used in the construction of PSSM. Compared with
41
42 1NN_PSI-BLAST_{ps}, both 1NN_TFPSSM and 1NN_PSI-BLAST_{nr} incorporate a
43
44 larger database for PSSM construction. Finally, the comparison between
45
46 1NN_TFPSSM and 1NN_PSI-BLAST_{nr} serves to highlight the benefit of
47
48 gapped-dipeptide encoding scheme.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 ***Experiment 3: Comparison between PSLDoc and PSLDoc-PLSA on the***
5
6
7 ***PS1444 data set***
8
9

10 PSLDoc-PLSA represents PSLDoc without PLSA, which simply applies SVM on
11 the original feature vectors. The overall accuracies of PSLDoc and PSLDoc-PLSA are
12 compared in order to evaluate the benefit of PLSA feature reduction for SVM
13 learning.
14
15
16
17
18
19
20
21
22
23
24
25

26 ***Experiment 4: Comparison among PSLDoc, 1NN_TFPSSM, and***
27
28 ***1NN_ClustalW on the PSHigh783 and PSLow661 data sets***
29
30
31

32 Using the PSHigh783 data set, we can verify whether PSLDoc can replace
33 1NN_ClustalW. Using PSLow661, we can investigate whether PSLDoc can improve
34 1NN_TFPSSM by applying PLSA and SVM classification. Hence, we could
35 determine whether PSLDoc is suitable for both high- and low-homology data sets.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52

53 ***Experiment 5: Comparison among PSLDoc, HYBRID and PSORTb v.2.0 on***
54
55 ***the PS1444 data set***
56
57
58
59
60

We compare the performance of PSLDoc, HYBRID, and PSORTb v.2.0. Besides,
we also assess the performance of PSLDoc using a three-way data split procedure⁴⁵,
which is commonly used in machine learning to prevent overestimation of the

1
2
3
4 performance. The data set is randomly divided into three disjoint sets, i.e., a training
5
6
7 set for classifier learning, a validation set for feature selection and parameter tuning,
8
9
10 and a test set for performance evaluation. Hence, for each run in the original five-fold
11
12
13 cross-validation, we divide the training data set into four distinct sets: three for
14
15
16 training, one for validation. Then, we select the gapped distance upper bound and
17
18
19 PLSA reduced feature size based on the validation set instead of the test set. Then
20
21
22 PSLDoc performance is evaluated under the selected parameters in the original
23
24
25 five-fold cross-validation.
26
27
28
29
30
31

32 ***Experiment 6: PSLDoc under different Prediction Thresholds versus***
33
34
35 ***PSORTb v.2.0 on the PS1444 data set***
36
37

38 The precision and recall of PSLDoc is evaluated under different prediction
39
40
41 thresholds to compare with PSORTb v.2.0.
42
43
44
45
46
47
48
49

50 **RESULTS AND DISCUSSION**

51 **Experimental Results**

52
53 ***Experiment 1: The benefit of using the TFPSSM weighting scheme***
54
55

56
57 The overall accuracy of 1NN_TFIDF and 1NN_TFPSSM for each gapped
58
59
60 distance are shown in Figure 3. The highest overall accuracy of 1NN_TFPSSM is

1
2
3
4 89.47% when l equals 4, 5, and 13 and it is considerably higher than the best

5
6
7 1NN_TFIDF score 74.38% when l equals to 4. Therefore, adopting the TFPSSM

8
9
10 weighting scheme significantly improves the performance of 1NN_TFIDF.

11
12
13
14 (Figure3)

15
16 The performance of 1NN_TFIDF and 1NN_TFPSSM in the high- and

17
18 low-homology data sets is shown in Table II. 1NN_TFPSSM dramatically improves

19
20 the performance of 1NN_TFIDF by about 26% in overall accuracy on PSLow661.

21
22 Hence, the incorporation of PSSM information into the weighting scheme is useful for

23
24 improving performance due to insufficient sequence information in the low-homology

25
26 data set.

27
28
29
30
31
32
33
34
35 (Table II)

36
37
38
39
40
41 ***Experiment 2: The effect of incorporating PSSM information and***

42
43
44 ***gapped-dipeptide encoding scheme***

45
46

47 Table III shows the performance of 1NN_TFPSSM, 1NN_ClustalW,

48
49 1NN_PSI-BLAST_{ps}, and 1NN_PSI-BLAST_{nr} on the PSHigh783 and PSLow661 data

50
51 sets. The overall accuracy on the PSHigh783 data set is very similar for all methods.

52
53 However, for the PSLow661 data set, 1NN_ClustalW, 1NN_PSI-BLAST_{ps}, and

54
55
56
57
58
59 1NN_PSI-BLAST_{nr} attain 42.97%, 57.94% and 66.57%, respectively, in overall

60

1
2
3
4 accuracy. This result reveals that better performance can be achieved when a larger
5
6
7 database is used in constructing PSSM. This also lends support to our assumption that
8
9
10 incorporating more information into PSSM is more effective for the prediction of
11
12
13 proteins with low sequence identity to the training set. Most notably, 1NN_TFPSSM
14
15
16 outperforms 1NN_PSI-BLAST_{nr} by 12.86% in overall accuracy. This suggests that the
17
18
19 incorporation of PSSM based on gapped-dipeptide encoding scheme significantly
20
21
22 improves the predictive performance, especially for proteins of low sequence identity.
23
24

25
26 (Table III)
27
28
29
30
31

32 ***Experiment 3: The benefit of PLSA feature reduction***

33 ***Determine the reduced size of PLSA***

34
35
36 The size of PLSA is determined by LSA singular values. Figure 4 shows the
37
38
39 singular values in decreasing order on different gapped distances upper bound data
40
41
42 sets.
43
44
45
46

47 (Figure4)
48
49

50
51 The 40-th largest singular value is close to zero in Figure 4, but in the inset the
52
53
54 160-th largest singular value is close to zero. Hence, the reduced feature size of PLSA
55
56
57 is set to 40, 80 and 160. However, we do not test larger PLSA reduced size or
58
59
60 one-by-one PLSA reduced size in consideration of the training efficiency and

1
2
3
4 avoidance of data overfitting.
5
6

7 For one PLSA reduced size, the training and testing procedures of PSLDoc take
8
9
10 1.5 hours and about 2~3 minutes for all gapped distances, respectively. However,
11
12
13 PSLDoc_{PLSA} takes about 180 and 1.4 hours in training and testing, respectively.
14
15
16 Figure 5 shows the performance of PSLDoc_{PLSA} and PSLDoc, where PSLDoc_{F x}
17
18 denotes PSLDoc with PLSA reduced size x .
19
20

21
22 (Figure5)
23
24

25
26 The highest overall accuracy among all gapped distances of PSLDoc_{F40},
27
28 PSLDoc_{F80}, and PSLDoc_{F160} is 92.31%, 93.01%, and 92.52%, respectively,
29
30 which is 0.83%, 1.52%, and 1.04% better than that of PSLDoc_{PLSA}. Using PLSA not
31
32 only improves learning efficiency but also performance. In the following experiments,
33
34
35 PSLDoc takes the gapped distance 13 and PLSA at reduced size 80.
36
37
38
39
40
41
42
43

44 ***Experiment 4: The benefit of SVM and PLSA feature reduction***

45
46

47 Table IV shows the performance of PSLDoc, 1NN_TFPSSM and 1NN_ClustalW
48
49 on PSHigh783 and PSLow661. The overall accuracy of 1NN_ClustalW on
50
51 PSHigh783 (97.32%) is very similar to that of Yu *et. al.*'s (97.7%). 1NN_TFPSSM
52
53 and PSLDoc perform better than 1NN_ClustalW on PSHigh783. On the other hand,
54
55 PSLDoc improves 1NN_TFPSSM on PSLow661 by 7.41% due to the non-linear
56
57
58
59
60

1
2
3
4 SVM classification and PLSA feature reduction and extraction. This shows that
5
6
7 PSLDoc is suitable for both the high- and low-homology data sets.
8
9

10 (Table IV)
11

12 ***Experiment 5: Comparison of PSLDoc, HYBRID and PSORTb v.2.0***

13
14
15
16
17 Table V shows the performance of PSLDoc, HYBRID, and PSORTb v2.0 on
18
19
20 PS1444. PSLDoc achieves the best performance of 93.01%, better than HYBRID of
21
22
23 91.6% and PSORTb of 82.6%.
24

25
26 (Table V)
27
28
29
30
31

32 ***Experiment 6: PSLDoc under different Prediction Thresholds versus*** 33 34 35 ***PSORTb v.2.0 on the PS1444 data set***

36 37 38 ***Prediction confidence***

39
40
41 The probability estimated by LIBSVM is used for determining the confidence
42
43
44 levels of classifications. The class with the largest probability is chosen as the final
45
46
47 predicted class. The confidence of the final predicted class, *prediction confidence*³²,
48
49
50 could be defined as the value of the largest probability minus the second largest
51
52
53 probability. Figure 6 shows the relationship between accuracy and prediction
54
55
56 confidence. For proteins with prediction confidence in the range [0.9-1], the
57
58
59 prediction accuracy is near 100% (99.12%).
60

1
2
3
4 (Figure6)
5
6
7
8
9

10 ***Prediction threshold*** 11

12
13 Gardy *et al.* suggested that when a prediction system is unable to generate a
14 confident prediction, the program outputs a result of “*Unknown*” because biologists
15 usually prefer correct predictions (high precision) over prediction coverage (recall)³.
16
17 To provide prediction results with higher precision, we determine a *prediction*
18 *threshold* to filter out prediction results with low confidence. That is, the SVM
19 classifier predicts results only when the prediction confidence is above the threshold,
20 otherwise the SVM classifier will output “*Unknown*”^{3,14}. The recall and precision for
21 each prediction threshold are shown in Figure 7.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

38 (Figure7)
39
40
41
42
43

44 Table VI shows the performance of PSLDoc under different prediction thresholds.
45
46 Setting the prediction threshold to 0.7, PSLDoc achieves slightly better recall than
47 PSORTb v.2.0 (83.66% versus 82.6%), whereas the precision of PSLDoc is better
48 than PSORTb v.2.0 (97.89% versus 95.8%). In addition, when the prediction
49 threshold is set to 0.3, PSLDoc achieves comparable precision to PSORTb v.2.0
50 (95.77% vs. 95.8%), and PSLDoc’s recall is much better than that of PSORTb v.2.0
51
52
53
54
55
56
57
58
59
60

(89.27% vs. 82.6%).

(Table VI)

Discussion

In PLSA, we associate proteins and gapped-dipeptides with topics. Through analyzing the trained PLSA model with $P(w|z)$ and $P(z|d)$ for gapped-dipeptide w , topic z and protein d , gapped-dipeptide signatures in proteins with different localization sites are discovered for the PS1444 data set. Some of these signatures have been reported in the literature as motifs critical for stability or localization. We also discuss the problem of polysemy and solve it through the PLSA model.

Gapped-dipeptide signatures for Gram-negative bacteria localization sites

In Figure 8, we show the distribution of topic versus protein as visualized by $P(z|d)$ for topic $z \in Z$ and protein $d \in D$. In the figure, the size of topic ($|Z|$) is set to 80 according to the conclusion from Experiment 2.

(Figure8)

To find site-topic preference, we then cluster proteins according to their localization sites for examining preferred topics for each localization site. The *site-topic preference* of the topic z for a localization site l is calculated by averaging

1
2
3
4 $P(z|d)$, where d (a protein) belongs to l class (i.e., d has localization site l .) The
5
6
7 site-topic preference over topics per localization site is shown in Figure 9. We can
8
9
10 observe from the figure that topics can be divided into five groups such that each
11
12
13 group “prefers” a specific localization site.
14

15
16
17 (Figure9)

18
19
20 We say a topic z *prefers* a localization size l , if the corresponding site-topic
21
22 preference is the largest of all localization sites. For some topics preferring PP and EC
23
24 classes, the difference of the site-topic preference between their own preferring site
25
26 and other sites are not obvious in Figure 9. This also reflects the relative poor
27
28 performance of PSLDoc in PP and EC classes.
29
30
31
32
33

34
35 The distribution of topic versus gapped-dipeptide is visualized by $P(w|z)$ for
36
37 gapped-dipeptide $w \in W$ and topic $z \in Z$ as shown in Figure 10. In the figure, the size
38
39 of gapped-dipeptides ($|W|$) is set to 5,600 ($=14 \times 20 \times 20$) following the conclusion of
40
41
42
43
44 Experiment 2.
45

46
47 (Figure10)

48
49
50 To list gapped-dipeptides of interest, we select ten preferred topics for each
51
52 localization site according to *site-preference confidence*, which is defined as the
53
54 largest site-topic preference minus the second largest site-topic preference. For each
55
56
57 topic, five most frequent gapped-dipeptides are selected. We list the
58
59
60

1
2
3
4 gapped-dipeptides signatures of ten preferred topics corresponding to each of the
5
6
7 localization sites in Table VII.

8
9
10 (Table VII)

11 12 *Gapped-dipeptide signatures reflecting motifs relevant to protein localization sites*

13
14
15
16
17 Interestingly, some of the signatures in Table VII found by PSLDoc have been
18
19 reported in the literature as motifs critical for stability or localization. One example is
20
21 observed in the integral membrane (IM) proteins, in which helix-helix interactions are
22
23 stabilized by aromatic residues⁴⁶. Specifically, the aromatic motif (WXXW or W2W)
24
25 is involved in the dimerization of transmembrane (TM) domains by π - π interactions⁴⁶.
26
27 Remarkably, one preferred topic predicted for the IM class includes this motif (W2W)
28
29 among other signatures of aromatic residues. Another example is found in the outer
30
31 membrane (OM) class, where the C-terminal signature sequence is recognized by the
32
33 assembly factor, OMP85, regulating the insertion and integration of OM proteins in
34
35 the outer membrane of gram-negative bacteria⁴⁷. The C-terminal signature sequence
36
37 contains a Phe (F) at the C-terminal position, preceded by a strong preference for a
38
39 basic amino acid (K, R)⁴⁷. One of the preferred topics indeed contains this motif
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54 (R0F.)

55
56
57 The above findings demonstrate the sensitivity of PSLDoc for capturing
58
59 gapped-dipeptide signatures relevant to localization sites. Thus, the predicted
60

1
2
3
4 signatures can provide important clues for further studies of uncharacterized sequence
5
6
7 motifs related to protein localization.
8
9

10 *Comparison of gapped-dipeptide signature encoding and amino acid composition*

11
12
13
14 Figure 11 shows the amino acid compositions of single residues and
15
16
17 gapped-dipeptide signatures for each localization site, respectively. It is observed that
18
19
20 the distributions of 20 amino acids calculated from single residues and
21
22
23 gapped-dipeptide signatures are quite different. The distribution from single residues
24
25
26 [Fig. 11(A)] has no clear separation for some amino acids but the distribution from
27
28
29 gapped-dipeptide signatures [Fig. 11(B)] has a clear separation among five classes.
30

31
32 (Figure 11)
33

34
35 From Fig. 11(A) and (B), it is observed that for some amino acids, general amino acid
36
37
38 composition bias have an effect on the gapped-dipeptide signatures (e.g., CP: E; IM: I,
39
40
41 L; PP: P, K; OM: Y; EC: G, N). That is, amino acids having high composition in a
42
43
44 localization site tend to also have high composition in gapped dipeptide signatures of
45
46
47 the localization site. For example, there are relatively high proportions for Ile and Leu
48
49
50 in both single residue and gapped-dipeptide signature compositions in IM proteins.
51
52
53 However, many amino acids have high compositions in at least two localization sites.
54
55
56 Therefore, it is difficult to predict localization site based on single residue
57
58
59 compositions. From the amino acid composition of gapped-dipeptide signatures, we
60

observe a clear separation among different localizations for several amino acids, which are indistinguishable at the single residue level (i.e., A, M, V, Q, S, H, W). Specifically, Met, Val, and Trp have similar proportions across all five localizations in single residue composition. The small differences in single amino acid composition for these residues are amplified by examining the gapped-dipeptide signature compositions and thus, they can be used for predicting localization site in a discriminative manner. We further analyze the correlation between single amino acid and gapped-dipeptide signature compositions by the Pearson correlation coefficient whose definition for a series of n measurements of variables X and Y is as follows:

$$r_{xy} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}} \quad (17)$$

The Pearson correlation coefficient (r) between the two compositions (single residues vs. gapped dipeptide signatures) for CP, IM, PP, OM, EC and all localization sites are 0.29, 0.50, 0.41, 0.07, 0.50 and 0.36, respectively. The correlation for all localization sites is medium (in range 0.30 to 0.49)⁴⁸.

In summary, the gapped-dipeptide signatures predicted by PSLDoc can (1) successfully capture the compositional bias inherent at the single residue level; and (2) better resolve ambiguity in discriminating amino acid compositions for each localization site.

The physicochemical preference of gapped-dipeptide signatures

To further analyze the physicochemical preference of gapped-dipeptide signatures, each amino acid is classified into one of the four groups: non-polar (AIGLMV), polar (CNPQST), charged (DEHKR), and aromatic (FYW). Figure 12 shows the grouped amino acid compositions of single residues and gapped-dipeptide signatures for each localization class. The grouped amino acid composition of single residues for each localization site has very similar preferences, but different preferences are observed for gapped-dipeptide signature composition. For example, in Fig. 12(A), IM, PP, OM, and EC have similar distribution, but in Fig. 12(B), each localization has distinct distribution of grouped amino acid composition. This also lends support to the second point in the previous section, that gapped-dipeptide signature can better resolve ambiguity in discriminating amino acid compositions for each localization. Furthermore, our analysis shows that the amino acid compositions of the predicted gapped-dipeptide signatures exhibit some over-represented patterns for a particular compartment.

(Figure12)

Gapped-dipeptide signatures predicted for CP, IM, and EC classes have distinct preferences for different groups of amino acids, possibly reflecting the physico-chemical constraints imposed by the environment of a subcellular

1
2
3
4 compartment. In particular, the signatures predicted for IM has a high percentage of
5
6
7 non-polar amino acids (60%) and no charged (0%) amino acids. This can be explained
8
9
10 in terms of the physico-chemical properties of the lipid bilayer, in which non-polar
11
12
13 amino acids are favored in the transmembrane domains of IM proteins⁴⁹. In contrast,
14
15
16 charged amino acids are disfavored due to the penalty incurred in energy terms during
17
18
19 the assembly of IM proteins⁵⁰. CP and EC classes are found to contain a high
20
21
22 percentage of charged and polar amino acids, respectively. The role of charged amino
23
24
25 acids in the cytoplasm is probably related to pH homeostasis in which they act as
26
27
28 buffers, whereas secreted proteins in the EC classes may require more polar amino
29
30
31 acids for promoting interactions in the solvent environment⁵¹.
32
33
34

35 Although gapped-dipeptide signatures are found, PSLDoc performs training and
36
37
38 testing procedures solely based on the topics of the PLSA model. In addition,
39
40
41 Hofmann⁴¹ also noted that PLSA can capture the semantic meaning of words, in our
42
43
44 case, the gapped-dipeptides. This part will be discussed in the following section.
45
46
47
48
49

50 ***PSLDoc's capability to solve the polysemy of Gapped-dipeptides***

51
52

53 In document classification, a word with two different meanings is called
54
55
56 polyseme (e.g., 'bank' means (i) an organization that provides various financial
57
58
59 services or (ii) the side of a river). Hofmann mentioned PLSA could deal with
60

1
2
3
4 polysemy and gave an example about the word “segment” ((i) an image region or (ii)
5
6
7 a phonetic segment)^{41,52}. Such a word w would have a high probability in two
8
9
10 different topics. The hidden topic variable, $P(w|z)$, associated with each word
11
12
13 occurrence in a particular document is used to determine which particular topic w is
14
15
16 assigned to, depending on the context of the document. Sivic *et al.*⁵³ applied PLSA to
17
18
19 images and discussed the polysemy on images. We discuss the polysemy effect on
20
21
22 gapped-dipeptides.
23
24
25

26 A gapped-dipeptide may prefer two localization sites, e.g., “A6A” prefer CP and
27
28
29 PP in Table VII. It is sometimes difficult to determine the localization site of a protein
30
31
32 based on the weight of a polysemous gapped-dipeptide. PLSA can be used to remedy
33
34
35 the polysemy effect of a gapped-dipeptide by associating the gapped-dipeptide with
36
37
38 different topics. For example, “A6A” is among the top five frequent dipeptides of
39
40
41 Topic 73 in CP and Topic 6 in PP that their probabilities $P(w|z)$ are sorted in a
42
43
44 decreasing order as shown in Table VIII.
45
46
47

48 (Table VIII)
49
50

51 For example, two proteins from PS1444 data set, chemotaxis protein cheZ and
52
53
54 Endoglucanase B^f, contain subsequences of the polysemous gapped-dipeptide “A6A.”
55
56
57 They are in different classes, CP class and PP class, respectively, and some of their
58
59
60

^f chemotaxis protein cheZ and Endoglucanase B are 44-th and 680-th proteins in PS1444, respectively. We use d_{44} and d_{680} to denote them for ease.

relevant information is listed in Table IX. Using the original vector space, the two proteins have $P\{w = \text{"A6A"}, d_{44}\} = 0.7001$ and $P\{w = \text{"A6A"}, d_{680}\} = 0.651$, which differ slightly, and thus it is difficult to distinguish them. However, using the posterior probabilities of Topic 73 and Topic 6, given the different occurrences of "A6A" based on the PLSA reduced vector space can distinguish the two proteins and determine their classes. That is, since $(P\{z_{73}|w = \text{"A6A"}, d_{44}\}, P\{z_6|w = \text{"A6A"}, d_{44}\}) = (0.0794, 0.0)$ and $(P\{z_{73}|w = \text{"A6A"}, d_{680}\}, P\{z_6|w = \text{"A6A"}, d_{680}\}) = (0.0, 0.0596)$ and Topics 73 and 6 are associated with different classes, the proteins d_{44} and d_{680} can be distinguished to be in CP and PP classes. This example demonstrates PLSA's capability to remedy the polysemy effect of gapped-dipeptides.

(Table IX)

CONCLUSION

We present a new PSL prediction method, PSLDoc, based on gapped-dipeptides and PLSA and demonstrate that it is suitable for proteins of a wide range of sequence homologies. PSLDoc extracts features from gapped-dipeptides of various distances, where evolutionary information from the PSSM is utilized to determine the weighting of each gapped-dipeptides such that its performance is comparable to the homology search method in the high-homology data set. These features are further reduced by PLSA and incorporated as input vectors for SVM classifiers. PSLDoc performs very

1
2
3
4 well in low-homology data set with overall accuracy of 86.84%. It can also achieve
5
6
7 very high precision by using a flexible prediction threshold. Experiments show
8
9
10 PSLDoc performs better than some of the current methods in overall accuracy by
11
12
13 1.51%. Because of the generality of this method, it can be extended to other species or
14
15
16 multiple localization sites in the future. Through analyzing the amino acid
17
18
19 composition of gapped-dipeptide signatures, there is a relationship between the amino
20
21
22 acid group and localization sites. For future work, we will incorporate the amino acid
23
24
25 groups with gapped-dipeptides to design a new representation of terms for predicting
26
27
28 protein subcellular localization.
29
30
31
32
33
34

35 **ACKNOWLEDGEMENTS**

36
37
38 We thank Ching-Tai Chen and Han-Kuen Liang for helpful discussion. The
39
40
41 research was supported in part by the thematic program of Academia Sinica under
42
43
44 grant AS94B003 and AS95ASIA02.
45
46
47
48
49

50 **REFERENCES**

- 51 1. Nair R, Rost B. Sequence conserved for subcellular localization. *Protein Sci*
52 2002;11(12):2836-2847.
- 53 2. Yu CS, Chen YC, Lu CH, Hwang JK. Prediction of protein subcellular
54 localization. *Proteins* 2006;64(3):643-651.
- 55 3. Gardy JL, Laird MR, Chen F, Rey S, Walsh CJ, Ester M, Brinkman FSL.
56 PSORTb v.2.0: Expanded prediction of bacterial protein subcellular
57 localization and insights gained from comparative proteome analysis.
58 *Bioinformatics* 2005;21(5):617-623.
59
60

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
4. Nakai K, Kanehisa M. Expert system for predicting protein localization sites in gram-negative bacteria. *Proteins* 1991;11(2):95-110.
 5. Nakai K, Horton P. PSORT: a program for detecting sorting signals in proteins and predicting their subcellular localization. *Trends Biochem Sci* 1999;24(1):34-35.
 6. Nair R, Rost B. Better prediction of sub-cellular localization by combining evolutionary and structural information. *Proteins-Structure Function and Genetics* 2003;53(4):917-930.
 7. Reinhardt A, Hubbard T. Using neural networks for prediction of the subcellular location of proteins. *Nucleic Acids Research* 1998;26(9):2230-2236.
 8. Bhasin M, Garg A, Raghava GPS. PSLpred: prediction of subcellular localization of bacterial proteins. *Bioinformatics* 2005;21(10):2522-2524.
 9. Hua SJ, Sun ZR. Support vector machine approach for protein subcellular localization prediction. *Bioinformatics* 2001;17(8):721-728.
 10. Nair R, Rost B. Mimicking cellular sorting improves prediction of subcellular localization. *J Mol Biol* 2005;348(1):85-100.
 11. Park KJ, Kanehisa M. Prediction of protein subcellular locations by support vector machines using compositions of amino acids and amino acid pairs. *Bioinformatics* 2003;19(13):1656-1663.
 12. Wang J, Sung WK, Krishnan A, Li KB. Protein subcellular localization prediction for Gram-negative bacteria using amino acid subalphabets and a combination of multiple support vector machines. *Bmc Bioinformatics* 2005;6:174.
 13. Yu CS, Lin CJ, Hwang JK. Predicting subcellular localization of proteins for Gram-negative bacteria by support vector machines based on n-peptide compositions. *Protein Sci* 2004;13(5):1402-1406.
 14. Gardy JL, Spencer C, Wang K, Ester M, Tusnady GE, Simon I, Hua S, deFays K, Lambert C, Nakai K, Brinkman FSL. PSORT-B: improving protein subcellular localization prediction for Gram-negative bacteria. *Nucleic Acids Research* 2003;31(13):3613-3617.
 15. Lu Z, Szafron D, Greiner R, Lu P, Wishart DS, Poulin B, Anvik J, Macdonell C, Eisner R. Predicting subcellular localization of proteins using machine-learned classifiers. *Bioinformatics* 2004;20(4):547-556.
 16. Rey S, Acab M, Gardy JL, Laird MR, DeFays K, Lambert C, Brinkman FSL. PSORTdb: a protein subcellular localization database for bacteria. *Nucleic Acids Research* 2005;33:D164-D168.
 17. Costa EP, Lorena AC, Carvalho AeCPLF, Freitas AA, Holden N. Comparing

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- Several Approaches for Hierarchical Classification of Proteins with Decision Trees. 2007.
18. Cheng BYM, Carbonell JG, Klein-Seetharaman J. Protein classification based on text document classification techniques. *Proteins* 2005;58(4):955-970.
 19. King BR, Guda C. ngLOC: an n-gram-based Bayesian method for estimating the subcellular proteomes of eukaryotes. *Genome Biology* 2007;8(5):-.
 20. Valdes-Perez RE, Pereira F, Pericliev V. Concise, intelligible, and approximate profiling of multiple classes. *Int J Hum-Comput St* 2000;53(3):411-436.
 21. Namburu SM, Tu H, Luo J, Pattipati KR. Experiments on Supervised Learning Algorithms for Text Categorization. 2005. p 1-8.
 22. Manning CD, Schütze H. Foundations of statistical natural language processing. Cambridge, Mass.: MIT Press; 1999. xxxvii, 680 p. p.
 23. Salton G, Wong A, Yang CS. Vector-Space Model for Automatic Indexing. *Communications of the Acm* 1975;18(11):613-620.
 24. Salton G, Buckley C. Term-Weighting Approaches in Automatic Text Retrieval. *Inform Process Manag* 1988;24(5):513-523.
 25. Liang HK, Huang CM, Ko MT, Hwang JK. Amino acid coupling patterns in thermophilic proteins. *Proteins* 2005;59(1):58-63.
 26. Cedano J, Aloy P, PerezPons JA, Querol E. Relation between amino acid composition and cellular location of proteins. *J Mol Biol* 1997;266(3):594-600.
 27. Chou KC, Elrod DW. Protein subcellular location prediction. *Protein Engineering* 1999;12(2):107-118.
 28. Garg A, Bhasin M, Raghava GPS. Support vector machine-based method for subcellular localization of human proteins using amino acid compositions, their order, and similarity search. *Journal of Biological Chemistry* 2005;280(15):14427-14432.
 29. Myers EW, Miller W. Optimal Alignments in Linear-Space. *Comput Appl Biosci* 1988;4(1):11-17.
 30. Cuff JA, Barton GJ. Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. *Proteins-Structure Function and Genetics* 1999;34(4):508-519.
 31. Hua SJ, Sun ZR. A novel method of protein secondary structure prediction with high segment overlap measure: Support vector machine approach. *J Mol Biol* 2001;308(2):397-407.
 32. Jones DT. Protein secondary structure prediction based on position-specific scoring matrices. *J Mol Biol* 1999;292(2):195-202.

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
33. Lin HN, Chang JM, Wu KP, Sung TY, Hsu WL. HYPROSP II - A knowledge-based hybrid method for protein secondary structure prediction based on local prediction confidence. *Bioinformatics* 2005;21(15):3227-3233.
34. Rost B, Sander C. Prediction of Protein Secondary Structure at Better Than 70-Percent Accuracy. *J Mol Biol* 1993;232(2):584-599.
35. Wu KP, Lin HN, Chang JM, Sung TY, Hsu WL. HYPROSP: a hybrid protein secondary structure prediction algorithm - a knowledge-based approach. *Nucleic Acids Research* 2004;32(17):5059-5065.
36. Altschul SF, Madden TL, Schaffer AA, Zhang JH, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* 1997;25(17):3389-3402.
37. McGuffin LJ, Bryson K, Jones DT. The PSIPRED protein structure prediction server. *Bioinformatics* 2000;16(4):404-405.
38. Wheeler DL, Chappay C, Lash AE, Leipe DD, Madden TL, Schuler GD, Tatusova TA, Rapp BA. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Research* 2000;28(1):10-14.
39. Kumar CA, Gupta A, Batool M, Trehan S. Latent Semantic Indexing-Based Intelligent Information Retrieval System for Digital Libraries. *Journal of Computing and Information Technology* 2006;14:191-196.
40. Deerwester S, Dumais ST, Furnas GW, Landauer TK, Harshman R. Indexing by Latent Semantic Analysis. *J Am Soc Inform Sci* 1990;41(6):391-407.
41. Hofmann T. Unsupervised learning by probabilistic latent semantic analysis. *Machine Learning* 2001;42(1-2):177-196.
42. Chang C-C, Lin C-J. LIBSVM : a library for support vector machines; 2001.
43. Wu TF, Lin CJ, Weng RC. Probability estimates for multi-class classification by pairwise coupling. *Journal of Machine Learning Research* 2004;5:975-1005.
44. Matthews BW. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta* 1975;405(2):442-451.
45. Ritchie MD, White BC, Parker JS, Hahn LW, Moore JH. Optimization of neural network architecture using genetic programming improves detection and modeling of gene-gene interactions in studies of human diseases. *Bmc Bioinformatics* 2003;4:28.
46. Sal-Man N, Gerber D, Bloch I, Shai Y. Specificity in transmembrane helix-helix interactions mediated by aromatic residues. *Journal of Biological Chemistry* 2007;282(27):19753-19761.
47. Robert V, Volokhina EB, Senf F, Bos MP, Van Gelder P, Tommassen J. Assembly factor Omp85 recognizes its outer membrane protein substrates by a

- 1
2
3
4 species-specific C-terminal motif. *Plos Biology* 2006;4(11):1984-1995.
- 5 48. Cohen J. *Statistical power analysis for the behavioral sciences*. Hillsdale, N.J.:
6 L. Erlbaum Associates; 1988. xxi, 567 p.
- 7
8 49. Ulmschneider MB, Sansom MSP, Di Nola A. Properties of integral membrane
9 protein structures: Derivation of an implicit membrane potential. *Proteins*
10 2005;59(2):252-265.
- 11
12 50. Hessa T, Kim H, Bihlmaier K, Lundin C, Boekel J, Andersson H, Nilsson I,
13 White SH, von Heijne G. Recognition of transmembrane helices by the
14 endoplasmic reticulum translocon. *Nature* 2005;433(7024):377-381.
- 15
16 51. Booth IR. Regulation of Cytoplasmic Ph in Bacteria. *Microbiological Reviews*
17 1985;49(4):359-378.
- 18
19 52. Hofmann T. *Probabilistic Latent Semantic Analysis* 1999; Stockholm.
- 20
21 53. Sivic J, Russell BC, Efros AA, Zisserman A, Freeman WT. *Discovering*
22 *Object Categories in Image Collections*. 2005.
- 23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table I. Number of proteins in different localization sites.

Localization sites	No.
Cytoplasmic (CP)	278
Inner membrane (IM)	309
Periplasmic (PP)	276
Outer membrane (OM)	391
Extracellular (EC)	190
All sites	1,444

For Peer Review

Table II. The comparison of 1NN_TFIDF and 1NN_TFPSSM on the PSHigh783 and PSLow661 data sets.

Loc. Sites	PSHigh783				PSLow661			
	1NN_TFPSSM		1NN_TFIDF		1NN_TFPSSM		1NN_TFIDF	
	<i>Acc.</i>	<i>MCC</i>	<i>Acc.</i>	<i>MCC</i>	<i>Acc.</i>	<i>MCC</i>	<i>Acc.</i>	<i>MCC</i>
CP	94.20	0.96	71.01	0.74	83.25	0.77	41.15	0.36
IM	99.31	0.99	98.62	0.89	82.93	0.82	84.15	0.48
PP	95.86	0.94	86.21	0.89	74.05	0.63	38.17	0.46
OM	99.66	0.99	95.88	0.95	85	0.82	66.00	0.48
EC	96.99	0.96	92.48	0.91	57.89	0.51	28.07	0.26
Overall	97.96	-	91.83	-	79.43	-	53.86	-

Table III. Comparison of 1NN_TFPSSM, 1NN_ClustalW, 1NN_PSI-BLAST_{ps} and 1NN_PSI-BLAST_{nr} for the PSHigh783 and PSLow661 data sets

PSHigh783									
Loc.	1NN_TFPSSM		1NN_ClustalW		1NN_PSI-BLAST _{ps}		1NN_PSI-BLAST _{nr}		
Sites	<i>Acc.(%)</i>	<i>MCC</i>	<i>Acc.(%)</i>	<i>MCC</i>	<i>Acc.(%)</i>	<i>MCC</i>	<i>Acc.(%)</i>	<i>MCC</i>	
CP	94.20	0.96	89.86	0.90	88.41	0.92	86.96	0.90	
IM	99.31	0.99	98.62	0.97	99.31	0.98	99.31	0.98	
PP	95.86	0.94	93.79	0.93	93.79	0.93	92.41	0.91	
OM	99.66	0.99	99.66	0.99	99.66	0.99	99.66	0.99	
EC	96.99	0.96	98.50	0.98	98.50	0.98	98.50	0.98	
Overall	97.96	-	97.32	-	97.32	-	96.93	-	
PSLow661									
Loc.	1NN_TFPSSM		1NN_ClustalW		1NN_PSI-BLAST _{ps}		1NN_PSI-BLAST _{nr}		
Sites	<i>Acc.(%)</i>	<i>MCC</i>	<i>Acc.(%)</i>	<i>MCC</i>	<i>Acc.(%)</i>	<i>MCC</i>	<i>Acc.(%)</i>	<i>MCC</i>	
CP	83.25	0.77	39.23	0.23	36.84	0.40	55.50	0.53	
IM	82.93	0.82	46.95	0.33	68.29	0.57	75.00	0.66	
PP	74.05	0.63	41.98	0.44	59.54	0.51	64.12	0.54	
OM	85.00	0.82	45.00	0.47	87.00	0.57	87.00	0.66	
EC	57.89	0.51	43.86	0.10	50.88	0.37	52.63	0.45	
Overall	79.43	-	42.97	-	57.94	-	66.57	-	

Table IV. Comparison of PSLDoc, 1NN_TFPSSM, and 1NN_ClustalW for the PSHigh783 and PSLow661 data sets.

Loc. Sites	PSHigh783						PSLow661					
	PSLDoc		1NN_TFPSSM		1NN_ClustalW		PSLDoc		1NN_TFPSSM		1NN_ClustalW	
	<i>Acc. (%)</i>	<i>MCC</i>	<i>Acc. (%)</i>	<i>MCC</i>	<i>Acc. (%)</i>	<i>MCC</i>	<i>Acc. (%)</i>	<i>MCC</i>	<i>Acc. (%)</i>	<i>MCC</i>	<i>Acc. (%)</i>	<i>MCC</i>
CP	95.65	0.96	94.2	0.96	91.3	0.89	94.74	0.88	83.25	0.77	39.23	0.23
IM	99.31	0.99	99.31	0.99	97.93	0.97	87.80	0.88	82.93	0.82	46.95	0.33
PP	95.17	0.94	95.86	0.94	93.1	0.93	82.44	0.78	74.05	0.63	41.98	0.44
OM	99.66	0.99	99.66	0.99	99.66	0.99	84.00	0.84	85.00	0.82	45.00	0.47
EC	98.5	0.98	96.99	0.96	99.25	0.99	70.18	0.65	57.89	0.51	43.86	0.10
Overall	98.21	-	97.96	-	97.32	-	86.84	-	79.43	-	42.97	-

For Peer Review

Table V. Comparison of PSLDoc, HYBRID and PSORTb v.2.0 on the PS1444 data sets. The PSLDoc performance of incorporating a three-way data split procedure is indicated in the parentheses.

	PSLDoc		HYBRID		PSORTb v.2.0			
	Loc. Sites	Acc.(%)	MCC	Acc.(%)	MCC	Acc.(%)	MCC	
CP	94.96	(94.24)	0.91	(0.91)	95.00	0.89	70.10	0.77
IM	93.20	(93.53)	0.94	(0.94)	90.60	0.92	92.60	0.92
PP	89.13	(89.13)	0.87	(0.85)	88.80	0.84	69.20	0.78
OM	95.65	(95.14)	0.95	(0.94)	95.10	0.93	94.90	0.95
EC	90.00	(87.37)	0.87	(0.86)	85.30	0.87	78.90	0.86
Overall	93.01	(92.45)	-		91.60	-	82.60	-

Table VI. Comparison of PSLDoc under the prediction threshold 0.7, PSLDoc under the prediction threshold 0.3 and PSORTbv.2.0.

Loc.	PSLDoc_PreThr=0.7					PSLDoc_PreThr=0.3					PSORTb v.2.0				
Sites	TP	FP	FN	Pre.	Rec.	TP	FP	FN	Pre.	Rec.	TP	FP	FN	Pre.	Rec.
CP	216	6	62	97.30	77.70	243	13	35	94.92	87.41	195	15	83	92.86	70.14
IM	273	3	36	98.91	88.35	285	6	24	97.94	92.23	286	14	23	95.33	92.56
PP	202	8	74	96.19	73.19	226	17	50	93.00	81.88	191	9	85	95.50	69.20
OM	366	2	25	99.46	93.61	372	6	19	98.41	95.14	371	10	20	97.38	94.88
EC	151	7	39	95.57	79.47	163	15	27	91.57	85.79	150	4	40	97.40	78.95
Total	1208	26	236	97.89	83.66	1289	57	155	95.77	89.27	1193	52	251	95.82	82.62

For Peer Review

Table VII. Gapped-dipeptide signatures for each Gram-negative bacteria localization site.

Site	Gapped-dipeptide signatures		
CP	E0E, K1I, K5V, K1V, D0E; R3R, R6R, R2R, R0R, R9R; H3H, H1H, H7H, H13H, H10H; E4E, K6E, E6E, E3E, E0E	L1H, L5H, L3H, H4L, H0L; A6A, A13A, A7A, A10A, A11A; H1M, H2M, H11M, M0H, H0M;	A12C, A9C, A13C, A5C, A7C; I0E, R6I, I3R, I3K, R6V; A4E, E1E, A2E, V4E, A9E;
IM	I2I, I3I, I0I, L0I, I0F; V2I, V2V, V3I, V3V, I0V; W3W, W0W, W2W, W6W, W4W; F10P, F8P, F12P, F3P, F13P	L7L, L4L, L10L, L3L, L6L; T2F, T6F, F3F, T4F, T8F; Y12L, Y1L, Y11L, L0Y, L1L;	M3M, M2M, M0M, M8M, M6M; A1A, A7L, A4A, A1C, A11L; M2T, M3T, M10T, M4T, M0L;
PP	A1A, A2A, A0A, A3A, M4A; D0D, Q0D, D3D, D3Q, D11D; A3A, A7A, A1P, A6R, A10R; A10A, A11A, A6A, A12A, A3A	M0H, W1Q, W1H, W1K, W5Q; W0E, E4W, W11E, E0W, W13E; P3N, N4P, N3P, N5P, N0P;	P1E, P0E, E0P, P0K, E1P; K3K, K0K, K2K, K1K, K7K; H6G, G3M, H7D, G11H, H11G;
OM	T1R, R3T, R1T, T5R, P0P; Q6Q, Q1Q, Q3Q, Q13Q, Q4Q; N1Q, N1N, Q1Q, N12N, Q11V; Y1Y, Y0Y, Y5Y, Y4Y, Y12Y	R0F, R4F, Y13R, R6F, R2F; S0F, A3F, F0S, R9F, F7F; W2N, N2W, N0W, D2W, N13W;	N4N, N0N, N10N, N7N, F1N; G0G, A0G, A1G, G1A, G3A; Q5R, R1Q, Q1R, Q3R, R2Q;
EC	S6S, S2S, T11T, S13S, T6S; N10N, N9N, N13N, N11N, N12N; Q2N, N1Q, Q1Q, N3Q, Q7Q; N0N, N12V, N4V, V12N, N9V	G8G, G0G, G7G, G9G, G6G; N1N, N3N, N4N, N11N, N1T; K1S, S6S, S5S, S11M, S0S;	T1T, T3T, T5T, T9T, T10T; I5Y, Y12S, Y3S, Y9S, Y6I; S3G, G3G, G4S, G3S, G2G;

1
2
3
4 Table VIII. "A6A" is among the top five frequent dipeptides of Topic 73 and Topic 6,
5 where gapped-dipeptides are arranged to the decreasing order of $P(wlz)$.
6
7

Topic 73	Topic 6
<u>A6A</u>	A10A
A13A	A11A
A7A	<u>A6A</u>
A10A	A12A
A11A	A3A

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Table IX. Examples of two proteins in PS1444 having the gapped-dipeptide "A6A".

Protein Name	$P\{w_j = \text{"A6A"}, d_i\}$	$(P\{z_{73} d_i, w_j = \text{"A6A"}\}, P\{z_{61} d_i, w_j = \text{"A6A"}\})$	"A6A"	Localization Site
116294	0.7001	(0.0794, 0.0)	AIAEAAEA(40*) ASQPHQDA(75)	CP
121816	0.651	(0.0, 0.0596)	APGDPGSA(362) AQWGVNSA(409) AQYGGFLA(420)	PP

*The number in brackets denotes the starting position of the gapped-dipeptide "A6A."

For Peer Review

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
1 M	-3	-3	-4	-5	-3	-3	-4	-5	-4	0	1	-3	10	-2	-5	-4	-3	-4	-3	-1
2 P	2	-3	-3	-1	-3	-1	-1	-1	-4	-2	-4	-2	-2	-5	4	2	4	-5	-4	-3
3 L	-4	-5	-6	-6	-4	-3	-5	-6	-5	3	5	-5	4	0	-5	-5	-3	-4	-3	2
4 D	-2	5	-1	-3	-4	2	-1	-4	2	-5	-3	5	-2	-2	-4	-2	0	-1	0	-3
5 L	-4	-5	-6	-6	-4	-5	-6	-6	-4	4	4	-5	0	1	-5	-5	-3	-4	-3	3
6 Y	-4	-3	-3	-5	-5	-3	-4	-5	4	-4	-3	-3	-2	4	-5	-3	-2	2	8	-4
7 N	-4	-3	8	4	-6	-3	-2	-3	-2	-6	-6	-3	-5	-6	-4	-1	-3	-7	-5	-6
8 T	-2	-3	-1	-3	-1	-3	-3	-4	-3	-4	-4	-1	-4	-4	-4	4	6	-5	-4	-2
9 L	0	-1	-5	-5	-4	-3	-4	-4	-3	-1	5	-3	3	0	-4	-3	-3	-3	-2	-1
10 T	-1	-3	-1	-1	-4	-2	-3	-2	-1	-4	-3	-1	-3	-4	-4	3	6	-5	-4	-3

Fig. 1. PSSM of the sequence MPLDLYNTL, where each entry is the original value without normalization.

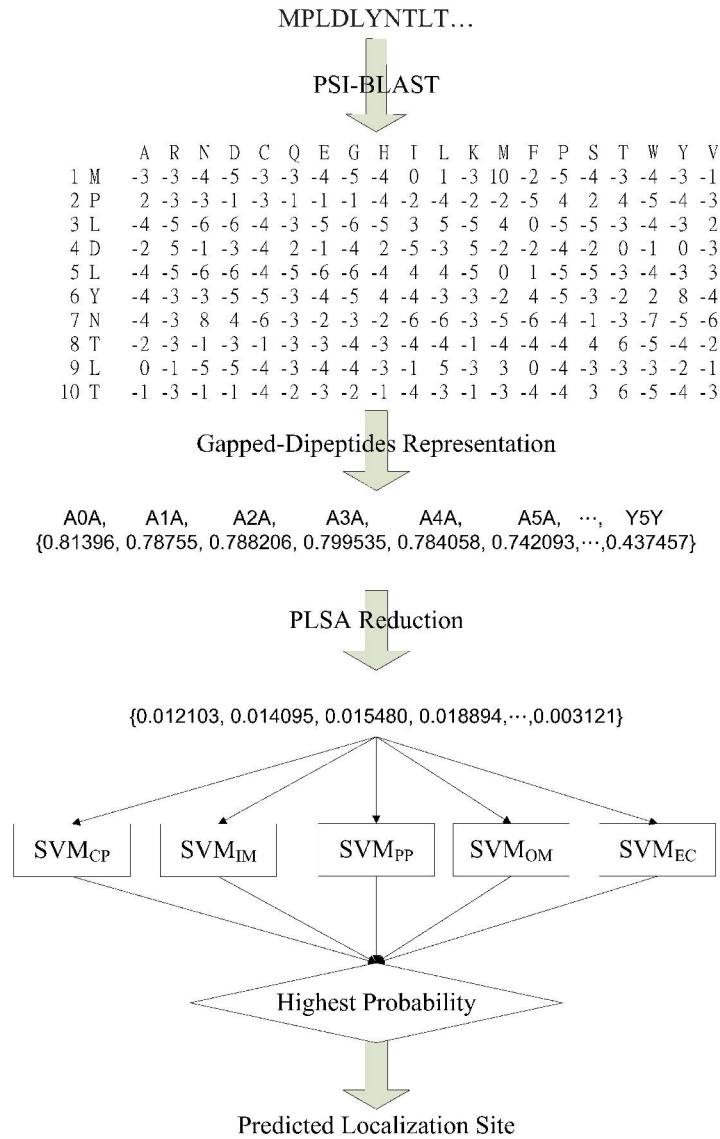


Fig. 2. System architecture of PSLDoc based on 1-v-r SVM models using reduced/transformed feature vectors.

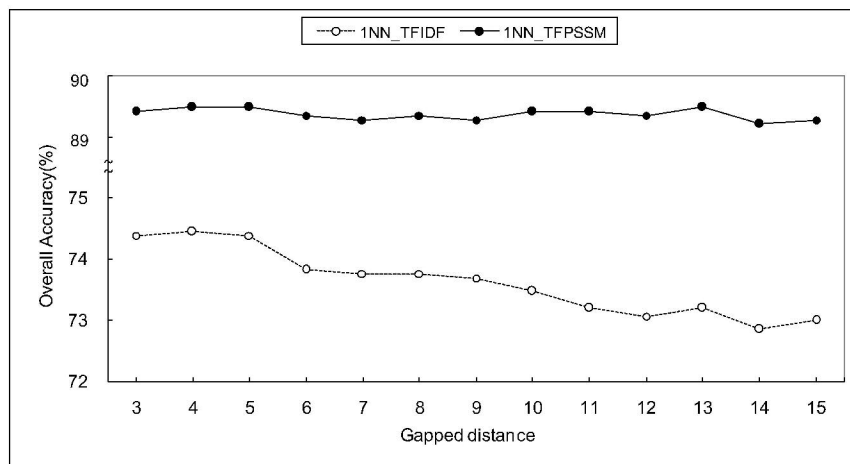


Fig. 3. Overall accuracy of 1NN_TFIDF and 1NN_TFPSSM with respect to maximum allowed gapped distances on the PS1444 data set.

review

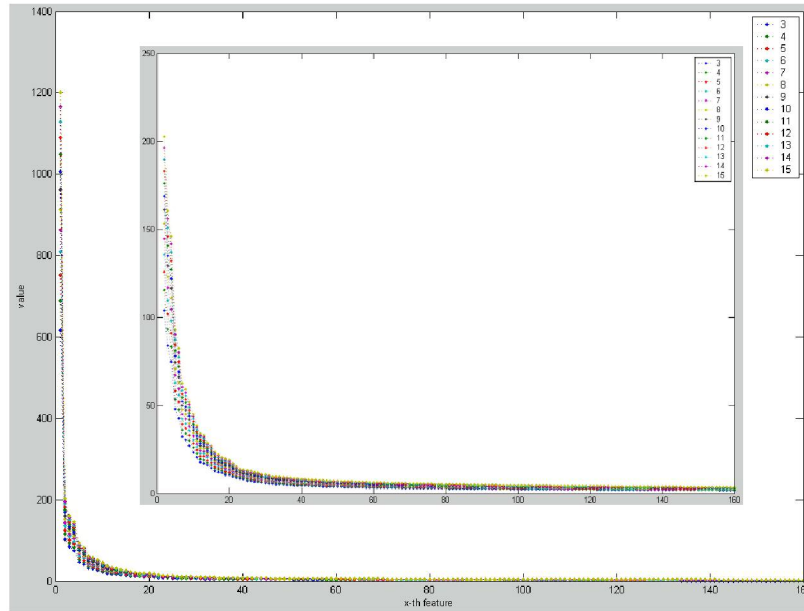


Fig. 4. Singular values in decreasing order of each gapped distance. The inset shows singular values without 1-th largest one for detailed representation.

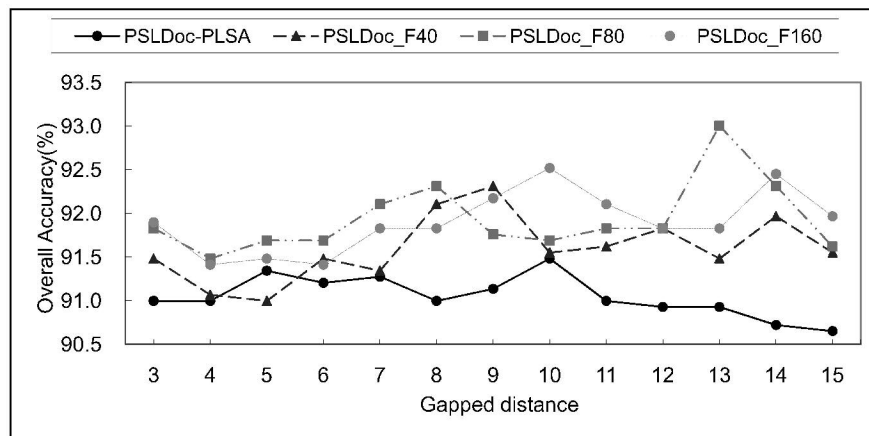


Fig. 5. Overall accuracy of PSLDoc_F40, PSLDoc_F80, PSLDoc_F160 and PSLDoc-PLSA with respect to gapped distance on the PS1444 dataset.

review

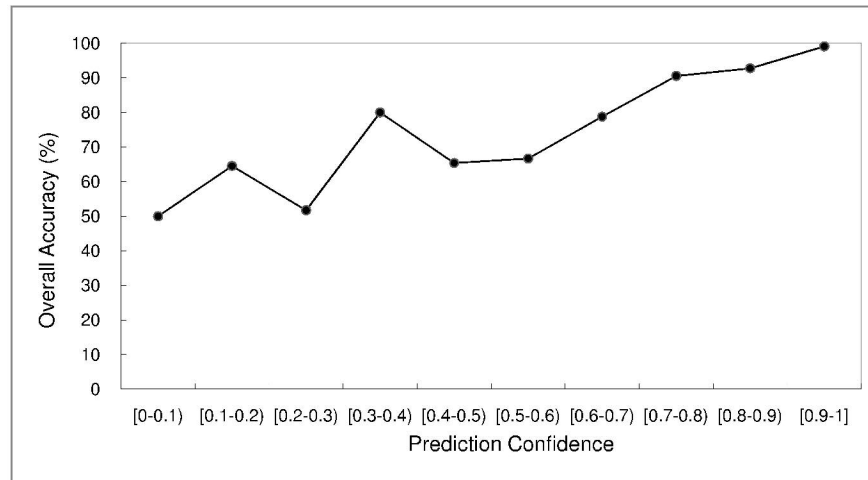


Fig. 6. Overall accuracy of PSLDoc with respect to prediction confidence. [x,y) represents the prediction confidence is more than x but under y.

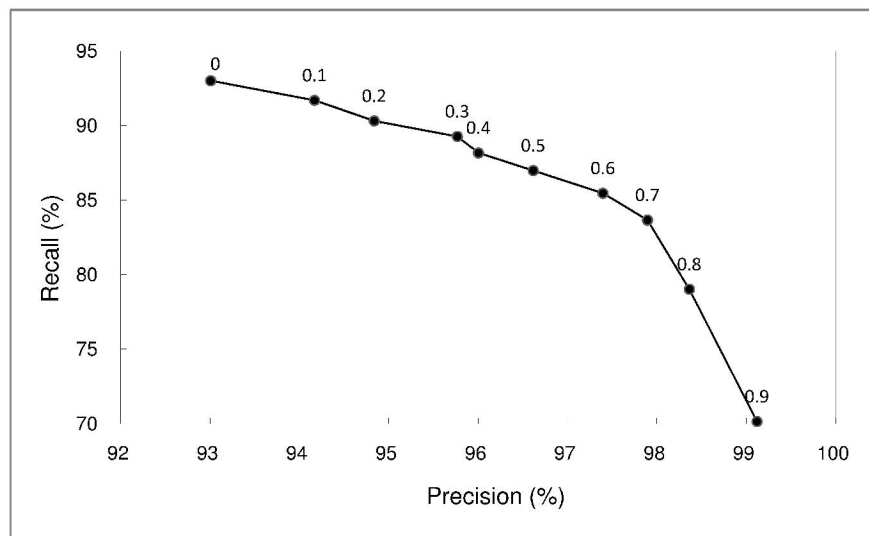


Fig. 7. Overall accuracy of PSLDoc with respect to prediction confidence. The value above the point denotes the corresponding prediction threshold.

review

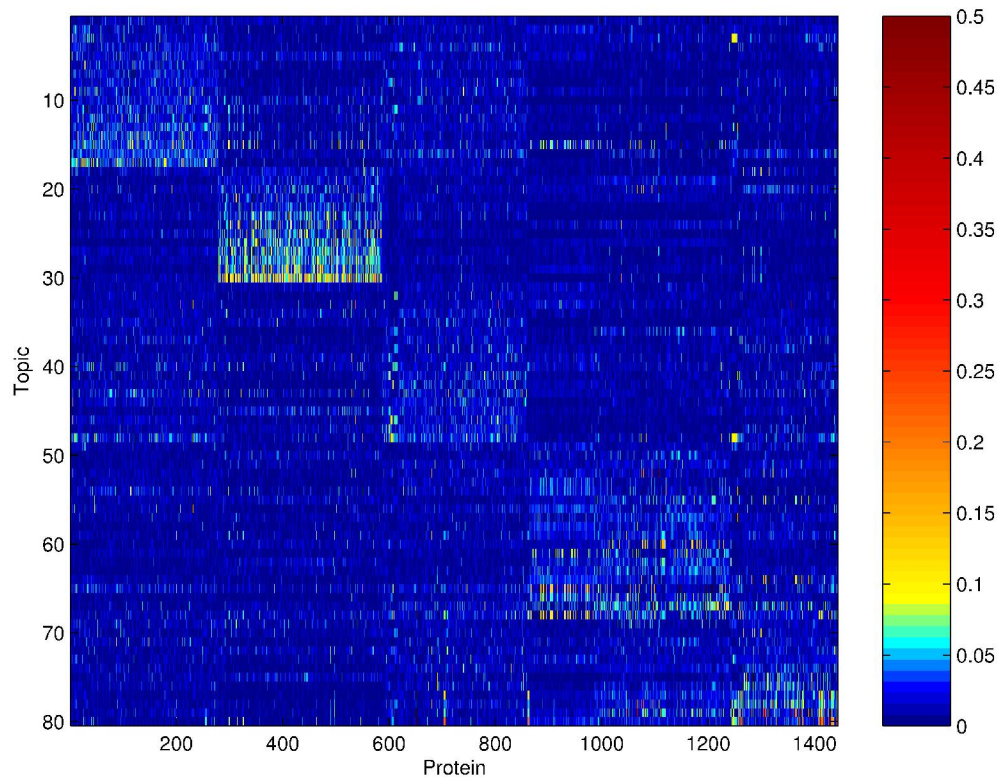


Fig. 8. Distribution of topic versus protein plotted as an image with its colormap^d, where the topics are sorted such that topics "preferring" (to be explained in the third paragraph) the same localization site are grouped together.

175x137mm (600 x 600 DPI)

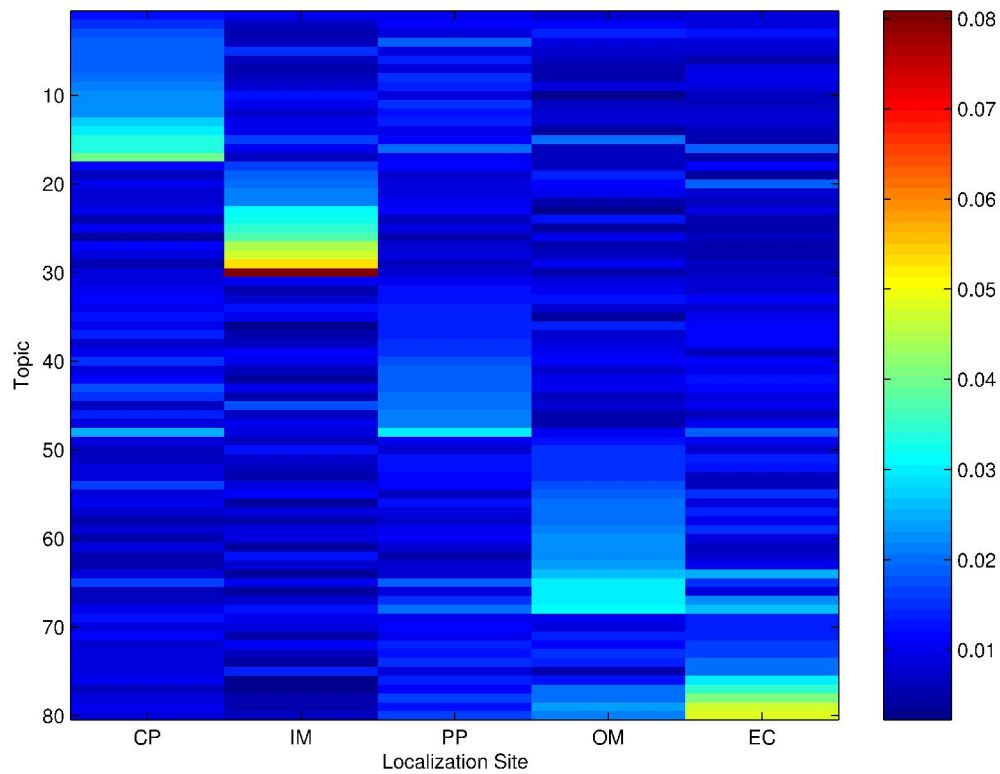


Fig. 9. Distribution of site-topic preference versus localization site. The 80 topics are divided into five groups of 17, 13, 18, 20 and 12 topics that prefer CP, IM, PP, OM and EC, respectively.

175x136mm (600 x 600 DPI)

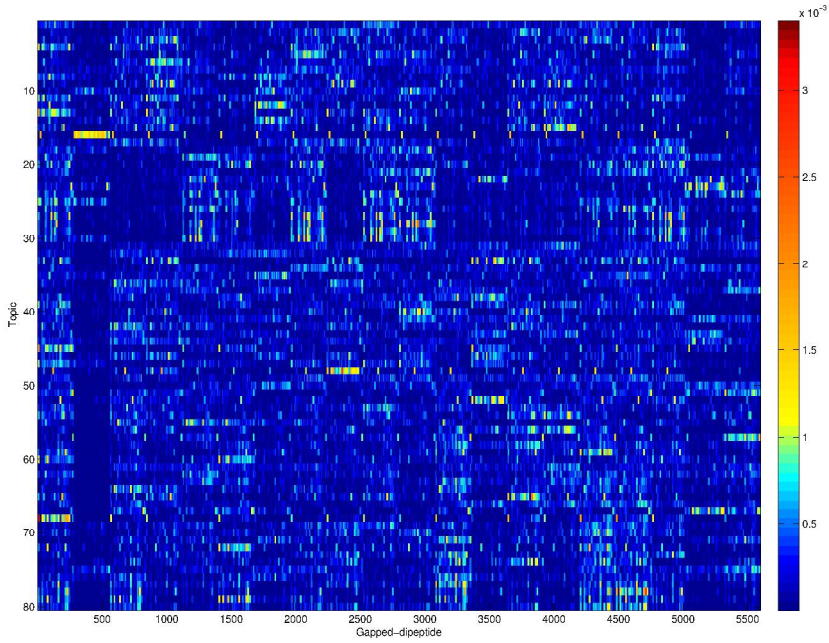


Fig. 10. Distribution of topic versus gapped-dipeptide.
338x245mm (600 x 600 DPI)

Review

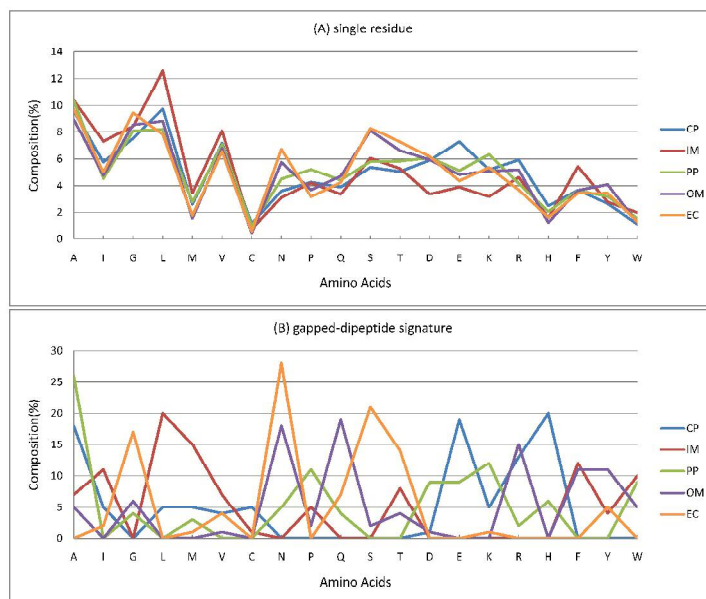


Fig. 11. The amino acid compositions of single residues (A) and selected gapped-dipeptide signatures (B) in different localization sites.
297x209mm (600 x 600 DPI)

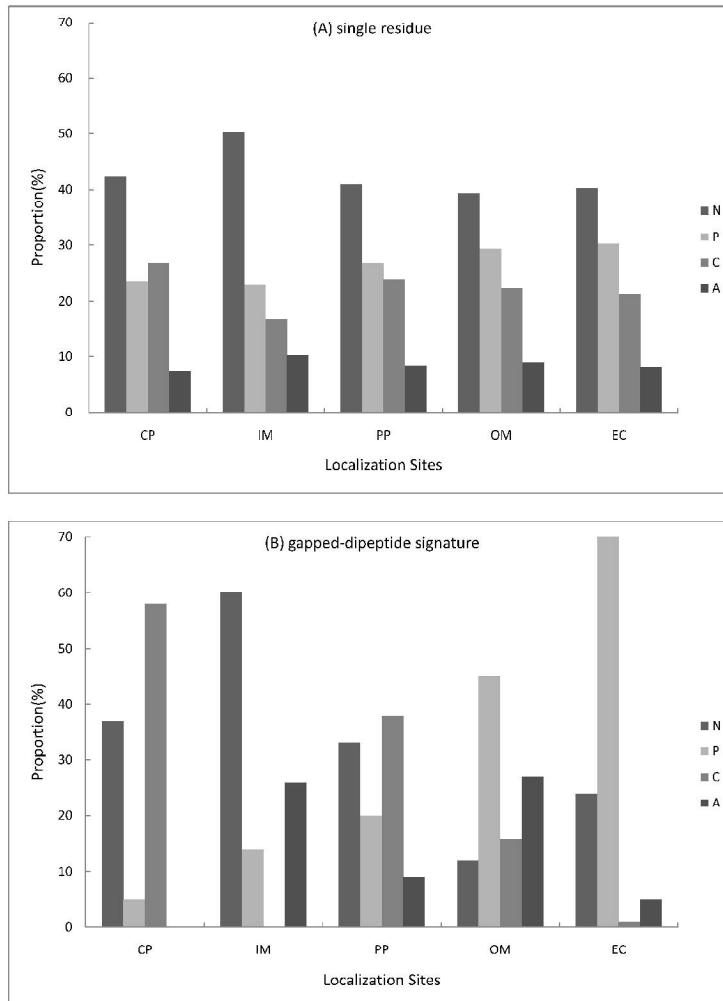


Fig. 12. The amino acid compositions of single residues (A) and predicted gapped-dipeptide signatures (B) for each protein class distinguished by the localization site. Localization sites: CP, IM, PP, OM, and EC. Amino acid groups: N (non-polar: AIGLMV), P (polar: CNPQST), C (charged: DEHKR), and A (aromatic: FYW).
209x297mm (600 x 600 DPI)