

This Provisional PDF corresponds to the article as it appeared upon acceptance. The fully-formatted PDF version will become available shortly after the date of publication, from the URL listed below.

## **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 doi:10.1186/1471-2105-6-174

Jiren Wang ([jiren@bii.a-star.edu.sg](mailto:jiren@bii.a-star.edu.sg))  
Wing-Kin Sung ([ksung@comp.nus.edu.sg](mailto:ksung@comp.nus.edu.sg))  
Arun Krishnan ([arun@bii.a-star.edu.sg](mailto:arun@bii.a-star.edu.sg))  
Kuo-Bin Li ([kuobin@bii.a-star.edu.sg](mailto:kuobin@bii.a-star.edu.sg))

**ISSN** 1471-2105

**Article type** Research article

**Submission date** 14 Feb 2005

**Acceptance date** 13 Jul 2005

**Publication date** 13 Jul 2005

**Article URL** <http://www.biomedcentral.com/1471-2105/6/174>

Like all articles in BMC journals, this peer-reviewed article was published immediately upon acceptance. It can be downloaded, printed and distributed freely for any purposes (see copyright notice below).

Articles in BMC journals are listed in PubMed and archived at PubMed Central.

For information about publishing your research in BMC journals or any BioMed Central journal, go to

<http://www.biomedcentral.com/info/authors/>

# **Protein subcellular localization prediction for Gram-negative bacteria using amino acid subalphabets and a combination of multiple support vector machines**

Jiren Wang<sup>1</sup>, Wing-Kin Sung<sup>2</sup>, Arun Krishnan<sup>1</sup>, Kuobin Li<sup>1</sup>

<sup>1</sup>Bioinformatics Institute, 30 Biopolis Street, #07-01 Matrix, Singapore 138671

<sup>2</sup>Department of Computer Science, National University of Singapore, 3 Science Drive 2, Singapore 117543

Email addresses:

JW: [jiren@bii.a-star.edu.sg](mailto:jiren@bii.a-star.edu.sg)

WS: [ksung@comp.nus.edu.sg](mailto:ksung@comp.nus.edu.sg)

AK: [arun@bii.a-star.edu.sg](mailto:arun@bii.a-star.edu.sg)

KL: [kuobin@bii.a-star.edu.sg](mailto:kuobin@bii.a-star.edu.sg)

# Abstract

## Background

Predicting the subcellular localization of proteins is an important order determining the function of proteins. Previous works focused on predicting protein localization in Gram-negative bacteria obtained good results. However, this method shows relatively low accuracies for the localization of extracellular proteins. This paper studies ways to improve the accuracy for predicting extracellular localization in Gram-negative bacteria.

## Results

We have developed a system for predicting the subcellular localization of proteins for Gram-negative bacteria based on a minimum number of combinations of multiple supportive features. The recall of the extracellular sites and overall recall of our predictor are 86.0% and 89.8%, respectively, in 5-fold cross-validation. To the best of our knowledge, these are the most accurate results for predicting subcellular localization in Gram-negative bacteria.

## Conclusions

Clustering amino acids into a few groups by the proposed greedy algorithm provides a new way to extract features from proteins. The dimensionality of the input vector of protein features. It was observed that a good amino acid grouping leads to an increase in prediction performance. Furthermore, a proper choice of a subset of complementary supportive features constructed by different features of proteins maximizes the prediction accuracy.

## Background

Subcellular localization is a key functional attribute of a protein. Since cellular functions are often localized in specific compartments, predicting the subcellular localization of unknown proteins may be useful to obtain useful information about their functions and to select proteins for further study. Moreover, studying the subcellular localization of proteins is also helpful in understanding disease mechanisms and for developing novel drugs.

As a result of large-scale genome sequencing efforts in recent years, protein data has accumulated in public databases at an increasing rate. Analyzing protein data to extract useful knowledge is thus essential for projects like automatic annotation. It is desirable to have an automated and reliable system for predicting subcellular localization of proteins from amino acid sequences.

Among the efforts [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21] have been made to predict protein subcellular localization. Most of these prediction methods are based on classification methods: on the one hand, the recognition of protein N-terminal sorting signals and their subsequent amino acid compositions [22].

Previous works have been focused on protein localization prediction of organelle-specific markers. There are five primary localizations: cytosol, nucleus, mitochondrion, chloroplast, and vacuole, which are the cytoplasm, the extracellular space, the inner membrane, the outer membrane, and the periplasm. P-SORT II [23] is the most widely used tool for

predicting multiple localizations of Gram-negative bacteria. This biological knowledge is represented by “if-then” rules for predicting protein localizations. Most of these rules were derived from experimental observations. However, the PSORT II does not consider the extracellular space site. Additionally, the overall recall for this dataset [24] only attains 60.9%.

Gardye et al. [24] presented PSORT-B, which improves prediction performance of PSORT II. PSORT-B combines information of the amino acid composition, similarity to proteins of known localization, presence of a signal peptide, transmembrane alpha-helices and motifs corresponding to specific localizations for a given protein sequence, through a probabilistic approach. It returns a list of five possible localizations with associated probabilities. It attains an overall recall of 74.8% for the dataset mentioned above.

Recently, Yue et al. [25] proposed a predictive system called CELLO for Gram-negative bacteria using support vector machines based on peptide compositions. They classified 20 amino acid side chain groups (charged, polar, aromatic and nonpolar) to reduce the dimensionality of the input vector. Forty SVM classifiers were used to predict the localizations. They overall recall was 88.9%. It was a significant improvement over the previous results of PSORT-B. However, the overall recall for extracellular proteins was still relatively low at 78.9%.

This paper studies ways to improve the accuracy of predicting extracellular localizations of Gram-negative bacteria. We explore a new way of extracting features from protein sequences for protein localization prediction by clustering 20 amino

acids into a few groups using a greedy algorithm. Our method for clustering amino acids considers the factors of both amino acids' physico-chemical properties and their contextual correlations. In contrast, the other method presented by Yu et al. classifies the amino acids into 4 groups (charged, polar, aromatic and nonpolar) based on physico-chemical properties of amino acids alone. Instead of simply combining multiple SVMs to give a better prediction, we propose a selection core function and a greedy algorithm to select a subset of SVMs to maximize their prediction accuracy.

Based on the proposed approaches, we have developed a system called P-CLASSIFIER for predicting the subcellular localization of Gram-negative bacteria by using a combination of multiple support vector machines. This has resulted in an improvement in the recall of extracellular proteins from 78.9% in CELLLO [25] (currently the best prediction system for Gram-negative bacteria) to 86.0% in P-CLASSIFIER. The overall recall of P-CLASSIFIER reaches 89.8%. To the best of our knowledge, this is the highest accuracy result for predicting protein subcellular localization in Gram-negative bacteria.

## Results

The dataset used in this study is from [24] and was extracted from SWISS-PROT release 40.29 [26]. It contains 1441 proteins of experimentally determined localization, where 1302 proteins are resident at single localizations and 139 proteins are resident at multiple localizations sites. Table 1 lists the number of protein sequences from different sites in the dataset.

The prediction performance of our prediction system is estimated from a 5-fold cross-validation where the given trainings samples are randomly partitioned into 5 mutually exclusive sets of a approximately equal size and a approximately equal class distribution.

In this experiment, the results are presented in terms of sequences in the dataset containing character "X". To avoid possible ambiguity of information, the sequences containing "X" in the training set, but excluded in the cross-validation training set, but included in the testing set are not shown.

Table 2 shows the prediction recall for single localization. The recall is calculated as  $TP_x / (TP_x + FN_x)$ , where  $TP_x$  and  $FN_x$  represent the number of samples correctly classified as X and the number of samples classified as not X that are actually X, respectively.

In the dataset, some proteins occur in different subcellular localizations. Since we are comparing our combined classifier P-CLASSIFIER with the P-SORT and CELLO classifiers, we followed the method of valuating the classifier for proteins resident in dual localizations, where we consider them as predicted correctly if one of their localizations is predicted correctly. Table 3 shows the prediction recall for dual localizations.

The Matthews correlation coefficient [27] is used to measure the predictive performance for five predicted classes. The Matthews correlation coefficient ( $MCC$ ) is defined by:

$$MCC = \frac{(TP_x)(TN_x) - (FP_x)(FN_x)}{\sqrt{(TP_x + FN_x)(TP_x + FP_x)(TN_x + FP_x)(TN_x + FN_x)}}$$

where  $TP_x$ ,  $TN_x$ ,  $FP_x$ , and  $FN_x$  are true positives, true negatives (the number of samples correctly predicted as not  $X$  that are actually not  $X$ ), false positives (the number of samples incorrectly predicted as  $X$  that are actually not  $X$ ), and false negatives of localization site  $X$ , respectively.  $MCC$  offers a comprehensive and robust measurement of the predictive performance of a classification model, considering both under- and over-predictions. The value of  $MCC$  equals 1 for a perfect prediction, and 0 for a completely random assignment.

Table 4 lists the performance comparisons among P-CLASSIFIER's (our system), PSORT-B's, and CELLO's [25] systems. As shown in Table 4, the values of  $MCC$  of all five systems are significantly higher than the quality of the values in CELLO's system, currently the best predicting system for Gram-negative bacteria. Moreover, we increase the recall of the extracellular site from 78.9% in CELLO to 86.0% in P-CLASSIFIER, a significant improvement of the extracellular site on the previous results. The overall recall of P-CLASSIFIER reaches 89.8%, which is better than previous results. To the best of our knowledge, these are the most accurate results for predicting Gram-negative bacterial localization.

## Discussion

To computationally analyse protein data, the representation of protein sequences is an important issue. A good input representation makes it easier for the SVM to identify underlying regularities and therefore is crucial to the success of SVM learning.

In this paper, we encode protein sequences by using the patterns of one amino acid, two adjacent amino acids, three adjacent amino acids, and four adjacent amino acids.



As there are 8000 and 16000 different patterns of orthetree and four adjacent amino acid sequences, clustering 20 amino acids into several groups provides a way to reduce the number of unique patterns in the SVM with very large number of features such as 160000 for all possible patterns of four adjacent amino acids. Since amino acids in proteins do not contribute to the function of proteins independently and functional patterns in proteins are embedded as sequence correlations, amino acids may not be grouped based on their physical-chemical properties alone [28]. For the prediction task, a good amino acid grouping leads to an improvement in prediction performance.

It is observed that the prediction results from SVMs constructed by different lengths of adjacent amino acid patterns, e.g. the patterns of a single amino acid and amino acid pairs, are complementary. That is, there are some cases where the prediction made by the SVM constructed by pairs of amino acids is correct while the prediction made by the SVM constructed by patterns of a single amino acid is incorrect, and vice versa. Therefore, combining complementary results provides a way to improve prediction accuracy. However, combining all complementary results together may not be a good choice. Therefore, we proposed to choose a subset of complementary support vectors to improve prediction accuracy.

After analyzing the prediction results, it is observed that there are some protein sequences that cannot be predicted correctly by any SVM or the combined classifier. It means that these protein sequences cannot be correctly classified by their

composition. This is the reason why the recall of some predictive items in Gram-negative bacteria cannot be further improved.

Since we are comparing our combined classifier P-CLASSIFIER with the P-SORTB and CELLLOC classifiers, we use the same datasets and methods. We did not check the sequence redundancy in these datasets. A slight level of sequence redundancy normally strongly affects prediction accuracy, removing those proteins sequences with high sequence identity (e.g. more than 40%) with each other in these datasets can avoid redundancy and bias.

Instead of giving full credit for dual-localized proteins if either of the sites is predicted correctly, we also evaluate the prediction performance by counting “half” correct when only one of the sites of dual-localized proteins is predicted correctly. Table 5 shows the prediction recalls. The full credit for dual-localized proteins is only given when both possible localizations of the protein are predicted with the optimal associated probabilities scores match the actual dual localizations of the protein. The corresponding overall recall for predicting dual localizations is only 67.3%. To properly deal with subcellular localizations for proteins residing in several different sites is a challenging problem. The paper [5] addressed the problem of subcellular localizations for proteins residing in several different sites.

There are three methods used for cross-validation test: the independent dataset test, n-fold cross-validation test, and the leave-one-out cross-validation test. Among these methods, the leave-one-out cross-validation test is the most rigorous and objective [29, 42]. However, the leave-one-out cross-validation test is very expensive

computationally and so far impractical for large datasets. The cross-validation test provides a bias-free estimation of the accuracy [30] and the reduced computational cost is considered as a more acceptable estimate of predictive performance for large datasets.

## Conclusions

This paper introduces a protein subcellular localization prediction method using amino acids alphabets and a combination of multiple support vector machines. The main contributions of our work include: (1) A new way of extracting features from protein sequences by clustering amino acids into a few groups using the proposed greedy algorithm to reduce the dimensionality of support vector machines. Our method or clustering amino acids considers not only the amino acids' physical-chemical properties but also the contextual correlations. (2) A selection core function and greedy algorithm are proposed to select a subset of candidates for support vector machines to maximize the cross-validation accuracy instead of simply combining multiple support vector machines together for prediction. (3) A web-based system has been developed for predicting proteins subcellular localization of Gram-negative bacteria. It allows people to submit multiple Gram-negative bacterial protein sequences for protein subcellular localization prediction. It is available at [43].

Clustering amino acids into a few groups by the proposed greedy algorithm provides a new way of extracting features to overcome adjacent amino acids from protein sequences and reduce the dimensionality of these features. Since amino acids

in practice, the contribution of each feature is not independent, it may not be a good idea to group a few features together. For the prediction task, a good feature grouping leads to an increase in prediction performance. Furthermore, properly choosing a subset of complementary support vectors can be constructed by different features of the support vectors to maximize the prediction accuracy.

## Methods

### Support vector machines

Support Vector Machines (SVMs) have been widely used in the analysis of biological data [32,33,34]. SVM is a relatively new family of learning methods and has some theoretical support from statistical learning theory [35,36]. SVM non-linearly maps the inputs space into a high-dimensional feature space, and seeks a hyperplane in this space that separates the positive samples from the negative ones with the largest possible margin. Instead of explicitly mapping the high-dimensional feature space, SVM usually works implicitly in the feature space by computing the corresponding kernel between the two objects.

Several parameters need to be set during the SVM training phase. These parameters include the regularization parameter, which controls the trade-off between good classification and large margin, the kernel type, and the kernel parameters. These parameters are set based on the cross-validation accuracy. The radial basis function (RBF) kernel is used for all our experiments and the software B SVM [44], a multi-class SVM [37], is used in this work.

## Protein features

The amino acid compositions of the full protein sequences are considered as global features, which represent the overall similarity among multiple protein sequences. In this paper, the global features are used as the SVM input to predict protein subcellular localization.

### a. W-gram protein encoding

Two types of features are considered in our work: W-gram and gapped2-gram. A W-gram is defined as a pattern of W ( $W \geq 1$ ) consecutive amino acid residues without any gaps and a gapped2-gram is defined as two amino acid residues with a specified number of gaps in a protein sequence. Here, a gapped2-gram is also referred to as a 2-gram. The main purpose of introducing the gapped encoding features or 2-grams is to increase the number of 2-gram feature candidates.

For each protein sequence  $P$  and each W-gram (or feature)  $F$ , let  $N(P, F)$  be the number of occurrences of  $F$  in the protein sequence  $P$ . Further, let  $T(P, W)$  be the total number of possible W-grams in  $P$ ,  $length(P)$  be the length of  $P$ , and  $G(F)$  be the specified number of gaps. We have  $T(P, W) = length(P) - W + 1 - G(F)$ , where  $G(F) = 0$  if  $W \neq 2$  and  $G(F) \geq 0$  if  $W = 2$ . The feature value  $U(P, F)$  with respect to the feature  $F$  and the sequence  $P$  is defined as  $N(P, F) / T(P, W)$ . For example, suppose  $P = \text{"LAEVLAAA"}$  and  $F = \text{"LA"}$  (without any gaps), then the feature value  $U(P, F)$  is  $2 / (8 - 2 + 1 - 0) = 0.2857$ , where  $F = \text{"LA"}$ ,  $N(P, F) = 2$ ,  $length(P) = 8$ ,  $W = 2$ ,  $G(F) = 0$ , and  $T(P, W) = 7$ . Intuitively,  $U(P, F)$  measures the proportion of occurrences of  $F$  among all possible W-grams in  $P$ . This measurement is length independent.

Int heW -grampr oteine ncodingm ethod,t het otal numberof di fferentpos sible featuresi s20<sup>w</sup>.

#### b. Amino acid subalphabets

Iti sdi fficultt ot raint he SVMw ithve ryl argenu mberof f eatures ucha s 8000f or3 - gram.T or educedi mensionalitiy,one w ayi st oc lassifyt he20 aminoa cidsi ntos mall numberof groupsba sed ont heirph ysical-chemicalpr operties.A llm embersi nt he samegr oupc anbe r epresentedb yone s ymbol.T hem erged aminoa cida lphabeta s fewert han20s ymbolsa ndi sc alledt hea minoa cids ubalphabet,w hichc anbe us edt o re-encodet heo riginalpr oteins equences.T her e-encodedpr oteins equencesha ve fewerf eatures. Fore xample,i ft henum berof s ymbolsi na na lphabeta sr educedf rom 20t o6,t henum berof 3- gramf eaturesi sr educed f rom4000( 20× 20 ×20 )t o216( 6 ×6× 6) .R educingt hen umberof f eatures t o a m anageables izef orS VMsc anhe lpt o improvet hepr edictivep erformance.

Thispa pers uggestsopt imizingt heg roupingb yu singt hepr oposed greedya lgorithm, whichc onsiderst he factorsof bot ht hea minoa cids' ph ysical-chemicalpr opertiesa nd theirc ontextualc orrelations,i nsteadof us ingt he groupingb asedont heirp hysical-chemicalpr opertiesa lone.N otet hatt herea rea n exponentialnum berof w ayst o group the20a minoa cids.F ore xample,t herea re580606 446a nd45232115901w ayst o divide20a minoa cidsi nto3a nd4gr oups,r espectively.T henum berof s ubalphabets withm groups( 1 ≤m ≤ 20)f ort hepr oteina lphabets izeof 20,  $N(m)$ c anbe calculated byt hef ormula[ 28]be low.

$$N(m) = \begin{cases} 1 & i = fm = 1 \\ \frac{m^{20}}{m!} - \sum_{k=1}^{m-1} \frac{N(m-k)}{k!} & \text{if } 1 < m \leq 20 \end{cases}$$

We learn the local optimal grouping based on a greedy algorithm using the SVM classification algorithm to evaluate the fitness of each candidate subalphabet, where the criteria for evaluation is the 5-fold cross-validation accuracy.

### c. Search for amino acid subalphabets

This section presents our greedy algorithm for finding a good grouping for the amino acids. Given a particular subalphabet encoding schema  $S$ , supposing  $N_g$  and  $T_c$  are the predefined number of groups and threshold of cross-validation accuracy, respectively. Further, we assume the parameters of a SVM to evaluate the fitness of a candidate subalphabet are given. These SVM parameters can be set either by the values suggested by the SVM software or by the tuning result of the SVM, which is constructed from features re-encoded by grouping 20 amino acids based on their physical-chemical properties, according to the criteria of cross-validation accuracy. For a particular subalphabet encoding schema  $S$ , let the grouping score  $h(S)$  be the cross-validation accuracy when prediction is done by a SVM using  $W$ -gram and the subalphabet scheme  $S$ .  $h(S)$  can be used to measure the goodness of the grouping  $S$ .

Table 6 shows an example of clustering 20 amino acids into 4 groups for the 4-gram protein encoding method using the proposed greedy algorithm. The initial node with 4-group assignment is set to  $\{(A, G, I, L, M, P, V), (C, N, Q, S, T), (D, E, K, H, R), (F, W, Y)\}$ , which is based on the physical-chemical property of amino acids. The

process for searching for an amino acid subalphabet is iterated until it reaches a local maximal grouping score at 79.0285%, where the final four groups are {(I, L, M, V), (N, S, T), (C, D, E, H, K, Q, R, Y), (A, F, G, P, W)}. Note that some group members in the classified result have the same physical-chemical property of amino acids. For example, the amino acids A, F, G and W in the fourth group (A, F, G, P, W) are all hydrophobic. In particular, the amino acids F and W are aromatic while amino acids A and G are tiny. Further, the hydrophilicity scale indices of A, G, P, and W have approximately the same values in the amino acid index database [38], which suggests that the hydrophilicity of amino acids may be an important property in classifying the 20 amino acids.

The proposed greedy algorithm to search for amino acid subalphabets is described in Table 7. The greedy local search [39] has been used for learning the subalphabets. In the search tree [39], every node represents an amino acid subalphabet encoding schema. The child nodes of a node are subalphabets encoding schemata, which are generated by moving every group member to each other group if the number of members in this group is greater than one.

This algorithm is composed of the following four steps. First, the 20 amino acids are initially divided into  $N_g$  groups either randomly with approximately the same size or based on some physical-chemical properties of the 20 amino acids. Amino acids in the same group are denoted by one symbol in a subalphabet. Suppose the current subalphabet encodings chemically represented by current node, its grouping score is calculated where the groupings core is the cross-validation accuracy when prediction is done by a SVM using W-gram and the subalphabet scheme.



Second, a new child node is generated from the current node. If there is only one member in the group, then the member is annotated with the group. Otherwise, the total number of groups will be less than  $N_g$ . Therefore, at most  $20 \times (N_g - 1)$  possible child nodes are generated. The minimum number of amino acids is 20 and the maximum number of amino acids is  $(N_g - 1)$  for the groups. If the highest group is chosen, then the child node is selected as the group. If the highest group is chosen, then the child node will become the current node.

Third, the above process is repeated for the child node. If the group is chosen among all child nodes, then the group is chosen as the child node. If the group is chosen, then the group is chosen as the current node.

Fourth, if the group is chosen, then the current node is selected as the group. If the group is chosen, then the current node is selected as the group. If the group is chosen, then the current node is selected as the group.

The amino acid sequences are divided into two parts: one part is used for the subalphabet and the other part is used for the evaluation of the performance of the subalphabet.

The greedy algorithm is applied to reduce the number of  $W$ -gram features. In particular, for 3-gram, we classify the 20 amino acids into 6, 7, and 8 groups. For 4-gram, we classify the 20 amino acids into 4 groups. The number of features is  $m^W$ , where  $m$  is the number of groups and  $W$  is the number of protein peptides in  $W$ -gram.

encoding methods. For example, the number of features is  $6 \times 6 \times 6 = 216$  for 6 groups in 3-gram encoding method.

### Multiple SVMs

Due to the difficulty of the multi-class classification, it may not be easy to obtain a single SVM that can return high accuracies for the subcellular localization prediction.

Therefore, multiple SVMs are trained from different features and their results are combined using voting.

Currently most of the existing protein subcellular localization prediction systems using SVM only use the features generated from 1-gram or 2-gram protein encoding methods. For example, the extracted features of amino acid compositions [2] and features of amino acid pair and gap amino acid pair compositions [40] are considered as the features generated from the 1-gram and 2-gram encoding methods, respectively.

As many functional patterns are embedded as sequence correlations, it is expected that more information will be included by combining classifiers constructed from features generated by 1-gram, 2-gram, 3-gram, and 4-gram protein encoding methods, instead of just using the classifiers constructed from 1-gram and 2-gram encoding methods since more adjacent amino acid residues will be considered.

In this paper, the following four types of features are extracted from protein sequences. The first type is the 1-gram encoding feature, which includes amino acid compositions and the partitioned amino acid compositions, which are protein

sequence is partitioned into  $P$  parts with approximately equal length. The total number of these features is  $20 \times P$ . In this work,  $P$  is set from 2 to 6. The second one is 2-gram encoding feature, which includes amino acid pairs and gapped amino acid pair compositions, where the number of features is  $400(20 \times 20)$  and the number of gaps is set from 1 to 2. The purpose of introducing the gapped encoding features only for 2-gram is to increase the number of 2-gram feature candidates. The third one is the 3-gram encoding feature, where the 20 amino acids are divided into 6, 7, and 8 groups whose number of features are  $216(6 \times 6 \times 6)$ ,  $343(7 \times 7 \times 7)$ , and  $512(8 \times 8 \times 8)$ , respectively. The last one is the 4-gram encoding method, where the 20 amino acids are divided into 4 groups, whose number of features is  $256(4 \times 4 \times 4 \times 4)$ .

### Feature selection

We applied the wrapper approach [41] in this work to eliminate irrelevant features and select the features subset for SVM classification and use 5-fold cross-validation accuracy as the criteria for evaluation.

Let  $SVM_a$  and  $SVM_b$  be the SVM classifiers using all features and features selected by the wrapper approach, respectively. Although the prediction accuracy of  $SVM_b$  is improved, the prediction results from  $SVM_a$  and  $SVM_b$  are different. There are some cases where the prediction made by  $SVM_a$  is correct while the prediction made by  $SVM_b$  is not correct, and vice versa. Therefore, both  $SVM_a$  and  $SVM_b$  can be considered as candidates to build the final combined classifier.

### SVM subset selection

Different SVMs provided different predictions. One new way to combine their predictions is by voting. That is, each prediction is assigned to a class with the most votes. For cases where more than one class gets the most votes, we assign these cases to the prediction results by one of the SVMs, which gets them the most correct predictions for all these cases.

Suppose  $S$  is a set of  $n$  patterns,  $N$  is the number of candidate SVMs,  $M = \{SVM_1, SVM_2, \dots, SVM_N\}$  is the set of candidate SVMs defined previously,  $V_1(S, M)$  is the number of correct predictions classified by  $M$  with one class corresponding to the most votes, and  $V_2(S, M)$  is the number of the correct predictions by the assigned SVM with more than one class corresponding to the most votes. The selection score function  $V(S, M)$  is defined as  $V_1(S, M) + V_2(S, M)$  and is used to select a subset of all candidate SVMs to form a combined classifier, which maximizes the cross-validation accuracy. The proposed greedy algorithm to select a subset of  $M$  is described in Table 8.

This greedy algorithm consists of the following steps. First, set  $M = \{SVM_1, SVM_2, \dots, SVM_N\}$ ,  $Score_{max} = V(S, M)$ ,  $Set_{max} = M$ , and  $i = N - 1$ . Second, for every member  $SVM_r \in M$  ( $1 \leq r \leq N$ ), remove  $SVM_r$  from  $M$  and calculate the value of its corresponding selection score function  $V(S, M - \{SVM_r\})$  ( $1 \leq r \leq N$ ). Suppose for some  $SVM_j$  ( $1 \leq j \leq N$ ),  $V(S, M - \{SVM_j\})$  is equal to  $V_{max}$ , the maximal value of all  $V(S, M - \{SVM_r\})$  ( $1 \leq r \leq N$ ), then update the following:  $M = M - \{SVM_j\}$ ,  $Score_{max} = V_{max}$ ,  $Set_{max} = M$ , and  $i = i - 1$ . The process of removing some  $SVM_p$  ( $1 \leq p \leq N$ ) will continue until  $i = 1$ , that is, only one SVM is left. Then  $Set_{max}$  is selected to be the combined classifier.

We can use the prediction results of four-fifth training proteins sequence as electa subset of S VMs and use the prediction result of the rest of one-fifth training protein sequence to evaluate the performance of the S VMs subsets election.

In this work, 15 S VMs are selected and combined to form the final classifier. Table 9 shows the encoding methods of input vectors of the fifteen elected S VMs. Rows 12, 13, and 14 represent three different merge alphabets, which are  $\{(A, F, G, P, W), (C, D, E, H, K, Q, R, Y), (N, S, T), (I, L, M, V)\}$ ,  $\{(A, C, M, P, V), (F, I, L, W), (D, E, H, Q, R), (G, K, N, S, T, Y)\}$ , and  $\{(A, G, P, Q, Y), (C, D, E, H, K, M, R), (N, S, T), (F, I, L, M, V)\}$ , respectively. Rows 4, 7, and 15 represent the same encoding methods as the rows 3, 6, and 14 but with features election.

We have conducted some experiments on constructing S VMs by using 5-gram encoding method. Preliminary experimental results show that the cross-validation accuracies predicted by S VMs constructed by 3-gram, 4-gram, and 5-gram encoding methods are not satisfactory when the number of groups is less than 6, 4, and 4, respectively. When we increase the number of groups to 4 for 5-gram, the time required to train the corresponding S VMs and calculate the 5-fold cross-validation accuracy is relatively low at the number of features is eaches 1024 ( $4 \times 4 \times 4 \times 4$ ). Therefore, only 1-gram, 2-gram, 3-gram, and 4-gram encoding methods are considered in this paper. Furthermore, the 20 amino acids are classified into 6, 7, and 8 groups for 3-gram and 4-gram encoding methods, respectively.

Since herea ret oom any zeroe lementsi nt hee ncodingr esults,2- gram,3- gram,a nd 4-grampr otein's encodingm ethodsa renot a ppliedt ot hosec asesw heret hepr otein sequences arep artitionedi ntoP ( P> 1) pa rtsw itha pproximatelys amel ength.

## Authors' contributions

JWde velopedt hem ethods,bui ltt hes ystem andd raftedt hem anuscript.W S,A Ka nd KLp articipatedi ns ystemde sign,pr ovidedv aluablec omments,a ndhe lpedt odr aft them anuscript.

## Acknowledgements

Thea uthorsw ouldl iket ot hantk hea nonymousr eviewers whosec ommentsha ve helpedus i mprovet hem anuscript.

## References

1. Emanuelsson,O .,Nielsen,H .,Brunak,S .,vonHeijne,G .: **Predicting subcellular localization of proteins based on their N-terminal amino acid sequence.** *JM olB iol.*,2000,300( 4):1005 -1016
2. Hua,S .a ndSun,Z .: **Support vector machine approach for protein subcellular localization prediction.** *Bioinformatics*,2001,17( 8):721 -728
3. Horton,P .a ndNakai,K . **Better Prediction of Protein Cellular Localization Sites with the k-Nearest Neighbors Classifier.** *Intellig. Syst. M ol. B iol.* 1997,5: 147- 152
4. Nakashima,H . andNishikawa,K .: **Discrimination of Intracellular and Extracellular Proteins Using Amino Acid Composition and Residue-pair Frequencies.** *JM olB iol.*,1994,238( 1):54 -61
5. Cai,Y .D.a ndChou,K .C.: **Predicting 22 protein localizations in budding yeast.** *Biochemical and Biophysical Research Communications*,2004,323: 425-428
6. Chou,K .C.: **Prediction of protein cellular attributes using pseudo-amino acid composition.** *PROTEINS: Structure, Function, and Genetics*,2001, 43: 246-255( Erratum: ibid.,2001,44: 60)
7. Chou,K .C.a ndCai,Y .D.: **An ew hybrid approach predicts subcellular localization of proteins by incorporating Gene Ontology.** *Biochemical and Biophysical Research Communications*,2003,311: 743 -747
8. Chou,K .C.a ndCai,Y .D.: **Prediction and classification of protein subcellular localization: sequence-order effects and pseudoamino acid**

- composition.** *Journal of Cellular Biochemistry*, 2003, 90: 1250-1260  
(Addendum, *in situ*, 2004, 91(5):1085)
9. Chou, K.C. and Cai, Y.D.: **Predicting subcellular localization of proteins by hybridizing functional domain composition and pseudo-amino acid composition.** *Journal of Cellular Biochemistry*, 2004, 91: 1197-1203
  10. Chou, K.C. and Cai, Y.D.: **Prediction of protein subcellular localizations by GO-FunD-PseAAC predictor.** *Biochemical and Biophysical Research Communications*, 2004, 323: 1236-1239.
  11. Feng, Z.P.: **Prediction of the subcellular localization of prokaryotic proteins based on an evolutionary representation of the amino acid composition.** *Biopolymers*, 2001, 58: 491-499
  12. Feng, Z.P. and Zhang, C.T.: **Prediction of membrane protein types based on the hydrophobic index of amino acids.** *Journal of Protein Chemistry*, 2000, 19: 269-275
  13. Feng, Z.P. and Zhang, C.T.: **Prediction of the subcellular localization of prokaryotic proteins based on the hydrophobicity index of amino acids.** *Int J Biol Macromol*, 2001, 28: 255-261
  14. Pan, Y.X., Zhang, Z.Z., Guo, Z.M., Feng, G.Y., Huang, Z.D. and He, L.: **Application of pseudo amino acid composition for predicting protein subcellular localization: stochastic signal processing approach.** *Journal of Protein Chemistry*, 2003, 22: 395-402
  15. Wang, M., Yang, J., Liu, G.P., Xu, Z.J. and Chou, K.C.: **Weighted-support vector machines for predicting membrane protein types based on pseudo amino acid composition.** *Protein Engineering, Design, and Selection*, 2004, 17: 509-516
  16. Wang, M., Yang, J., Xu, Z.J. and Chou, K.C.: **SLE for predicting membrane protein types.** *Journal of Theoretical Biology*, 2004, 232: 7-15
  17. Xiao, X., Shao, S., Ding, Y., Huang, Z., Huang, Y. and Chou, K.C.: **Using complexity measure for topological prediction of protein subcellular localization.** *Amino Acids*, 2005, 28(1): 57-61
  18. Yuan, Z.: **Prediction of protein subcellular localizations using Markov chain models.** *FEBS Letters*, 1999, 451: 23-26
  19. Zhou, G.P.: **An intriguing controversy over protein structural class prediction.** *Journal of Protein Chemistry*, 1998, 17: 729-738
  20. Zhou, G.P. & Assa-Munt, N.: **Some insights into protein structural class prediction.** *PROTEINS: Structure, Function, and Genetics*, 2001, 44: 57-59
  21. Zhou, G.P. and Doctor, K.: **Subcellular localization prediction of apoptosis proteins.** *PROTEINS: Structure, Function, and Genetics*, 2003, 50: 44-48
  22. Nakai, K.: **Protein sorting signals and prediction of subcellular localization.** *Adv. Protein Chem.*, 2000, 54: 277-344
  23. Nakai, K. and Kanehisa, M.: **Expert system for predicting protein localizations in Gram-negative bacteria.** *Proteins*, 1991, 11(2): 95-110
  24. Jennifer L. Gardy, Cory Spencer, Kevin Wang, Martin Ester, Gabriel T. Uszady, Istvan Simon, Susan Hua, Katalin Fays, Christophe Lambert, Kent Nakai and Fiona L. Brinkman. **PSORT-B: improving protein subcellular localization prediction for Gram-negative bacteria.** *Nucleic Acid Research*, 2003, 31: 3613-3617.
  25. Yu, C.-S., Lin, C.-J., and Hwang, J.-K.: **Predicting subcellular localization of proteins for Gram-negative bacteria by support vector machines based on peptide compositions.** *Protein Science*, 2004, 13(5): 1402-1406

26. Bairoch, A. and Apweiler, R. : **The SWISS-PROT protein sequence database and the TrEMBL database in 2000.** *Nucleic Acids Research*. 2000, 28: 45- 48
27. Matthews, B. W., **Comparison of the predicted and observed secondary structure of T4 phage lysozyme.** *Biochim. Biophys. Acta*, 1975, 405( 2): 442-451
28. Andersen, C. A. F. and Brunak, S. **Representation of Protein-Sequence Information by a Minimal Alphabet.** *AIM Magazine*. 2004, 25( 1): 97-104
29. Mardia, K. V., Kent, J. T., Bibby, J. M.: **Multivariate Analysis.** London: Academic Press, 1979, 322-381
30. Stone, M. : **Cross-validation of the assessment of statistical predictions.** *Journal of the Royal Statistical Society* 1974, 36: 111-147.
31. Kohavi, R. : **Wrapper for performance enhancement and bias reduction graphs.** PhD thesis, Stanford University 1995.
32. Brown, M. P., Grundy, W. N., Lin, D., Cristianini, N., Sugnet, C. W., Furey, T. S., Ares, M., and Haussler, D. : **Knowledge-based analysis of microarray gene expression data by support vector machines.** *PNAS*. 2000, 97: 262-267
33. Lee, Y. and Lee, C.-K. : **Classification of multiple classes by multicategory support vector machines using gene expression data.** *Bioinformatics*, 2003, 19 : 1132-1139
34. Ward, J. J., McGuffin, L. J., Buxton, B. F., and Jones, D. T. : **Secondary structure prediction with support vector machines.** *Bioinformatics*, 2003, 19: 1650 -1655
35. Vapnik, V. : **The Nature of Statistical Learning Theory.** Springer-Verlag, New York, 1995
36. Vapnik, V. , **Statistical learning theory,** John-Wiley, New York, 1998
37. Hsu, C.-W. and Lin, C.-J., **As implied by composition method for support vector machines.** *Machine Learning*, 2002, 46: 291- 314.
38. Kawashima, S. and Kanehisa, M. : **A index: amino acid index database.** *Nucleic Acids Research*. 2000, 28: 374
39. Russell, S. J. and Norvig, P.: **Artificial Intelligence: a modern approach,** Prentice Hall. 2003
40. Park, K.-J. and Kanehisa, M. : **Prediction of protein subcellular localizations by support vector machines using composition of amino acids and amino acid pairs.** *Bioinformatics*, 2003, 19( 13): 1656 -1663
41. Kohavi, R. and John, G. H.: **Wrapper for feature subsets selection.** *Artificial intelligence*. 1997, 97( 1-2): 273- 324
42. Chou, K. C. & Zhang, C. T. . **Review: Prediction of protein structural classes.** *Critical Reviews in Biochemistry and Molecular Biology*, 1995, 30, 275-349
43. **Protein subcellular localization prediction for Gram-negative bacteria** [<http://protein.bio.nyu.edu/localization/gram-negative/>].
44. **BSVM** [<http://www.csie.ntu.edu.tw/~cjlin/bsvm/index.html>].



## Tables

**Table 1 - Number of protein sequences in different sites**

Localizations ites	No.
cytoplasmic	248
inner membrane	268
periplasmic	244
outmembrane	352
extracellular	190
cytoplasmic/ i nnerm embrane	14
membrane/ pe riplasmic	49
outerm embrane/ e xtracellular	76
Alls ites	1441

**Table 2 - Prediction recall for a single localization.**

Localization	Recall( $TP_x/(TP_x+FN_x)$ )
Cytoplasmic	94.8%( 235/ 248)
Extracellular	83.2%( 158/ 190)
Innermembrane	88.1%( 236/ 268)
Outermembrane	93.2%( 328/ 352)
Periplasmic	86.9%( 212/ 244)
Overallr ecall	89.8%( 1169/1302)

**Table 3 - Prediction recall for dual localizations.**

Localization	Recall ( $TP_x/(TP_x+FN_x)$ )
Cytoplasmic/ i nnermembrane	92.9%( 13/14)
Outermembrane/ e xtracellular	98.9%( 75/76)
Periplasmic/ i nnermembrane	75.5%( 37/49)
Overallr ecall	89.9%( 125/139)

**Table 4 - Performance comparisons among P-CLASSIFIER's, PSORT-B's, and CELLO's methods.**

Localization	P-CLASSIFIER		CELLO		PSORT-B	
	Recall	<i>MCC</i>	Recall	<i>MCC</i>	Recall	<i>MCC</i>
Cytoplasmic	94.6%	0.85	90.7%	0.85	69.4%	0.79
Extracellular	86.0%	0.89	78.9%	0.82	70.0%	0.79
Innermembrane	87.1%	0.92	88.4%	0.92	78.7%	0.85
Outermembrane	93.6%	0.90	94.6%	0.90	90.3%	0.93
Periplasmic	85.9%	0.81	86.9%	0.80	57.6%	0.69
Overallr ecall	89.8%	-	88.9%	-	74.8%	-

**Table 5 - Prediction recall for dual localizations when “half” predictions are only counted as half correct.**

Localization	Recall ( $TP_x/(TP_x+FN_x)$ )
Cytoplasmic/ innermembrane	75.0%( 10.5/14)
Outermembrane/ extracellular	84.2%( 64/76)
Periplasmic/ innermembrane	38.8%( 19/49)
Overall recall	67.3%( 93.5/139)

**Table 6 - An example of clustering 20 amino acids into 4 groups.**

Searchings tates	Cross-validation accuracy	Actions
(A,G , I, L,M ,P ,V ) (C,N ,Q ,S ,T ) (D,E ,H ,K ,R ) (F,W ,Y )	71.2413%	Move‘ G’f rom group 1t og roup4
(A, I, L,M ,P ,V ) (C,N ,Q ,S ,T ) (D,E ,H ,K ,R ) (F,G ,W ,Y )	74.0941%	Move‘ A’f rom group 1t og roup4
(I, L,M ,P ,V ) (C,N ,Q ,S ,T ) (D,E ,H ,K ,R ) (A,F ,G ,W ,Y )	75.9445%	Move‘ P’f rom group 1t og roup4
(I, L,M ,V ) (C,N ,Q ,S ,T ) (D,E ,H ,K ,R ) (A,F ,G ,P ,W ,Y )	77.5636%	Move‘ C’f rom group 2t og roup3
(I, L,M ,V ) (N,Q ,S ,T ) (C,D ,E ,H ,K ,R ) (A,F ,G ,P ,W ,Y )	78.4888%	Move‘ Q’f rom group 2t og roup3
(I, L,M ,V ) (N,S ,T ) (C,D ,E ,H ,K ,Q ,R ) (A,F ,G ,P ,W ,Y )	78.9514%	Move‘ Y’f rom group 4t og roup3
(I, L,M ,V ) (N,S ,T ) (C,D ,E ,H ,K ,Q ,R ,Y ) (A,F ,G ,P ,W )	79.0285%	Reachl ocalm aximal groupings core and stop.

**Table 7 - Algorithm for amino acid subalphabets searching**

1	current_node ← the initial group assignment by dividing the 20 amino acids into $N_g$ groups.
2	REPEAT
3	best_node ← current_node
4	REPEAT
5	current_node ← generate a new node from the current node.
6	if $h(\text{best\_node}) < h(\text{current\_node})$ then
7	best_node ← current_node
8	UNTIL $h(\text{best\_node}) < h(\text{current\_node})$
9	IF $h(\text{current\_node}) < T_c$ THEN
10	current_node ← randomly re-generate initial group assignment
11	ENDIF
12	UNTIL $h(\text{current\_node}) \geq T_c$

**Table 8 - Algorithm for SVM subset selection**

1	Let $M = \{SVM_1, SVM_2, \dots, SVM_N\}$ be the set of candidate SVMs
2	Let $S_{core\_max} = V(S, M)$ and $set\_max = M$
3	FOR $i = 1$ to $N$
4	$V_{max} = \max\{V(S, M - \{SVM_r\}) \mid SVM_r \in M, 1 \leq r \leq N\}$
5	IF $V(S, M - \{SVM_j\}) = V_{max}$ ( $1 \leq j \leq N$ ) THEN
6	$set\_max = M - \{SVM_j\}$
7	ENDIF
8	IF $V_{max} \geq S_{core\_max}$ THEN
9	$S_{core\_max} = V_{max}$
10	$set\_max = M$
11	ENDIF
12	ENDFOR

**Table 9 - The encoding methods of input vectors in the fifteen selected SVMs.**

No.	Encoding methods of input vectors
1	1-gram without partitioned parts
2	1-gram without partitioned parts
3	1-gram without partitioned parts
4	1-gram without partitioned parts (apply features election on 0.3)
5	1-gram without partitioned parts
6	2-gram without gaps
7	2-gram without gaps (apply features election on 0.6)
8	2-gram with one gap
9	3-gram with merged groups
10	3-gram with merged groups
11	3-gram with merged groups
12	4-gram with merged groups
13	4-gram with merged groups
14	4-gram with merged groups
15	4-gram with merged groups (apply features election on 0.14)