# A comparison of computational strategies for multi-label prediction of protein subcellular localizations

Pingzhao Hu[1,2], Hui Jiang[2], Andrew Emili[1]

[1] Program in Proteomics and Bioinformatics, Banting and Best Department of Medical Research, University of Toronto, Toronto, Ontario, Canada

[2]Department of Computer Science, York University, Toronto, Ontario, Canada

Email: phu@cse.yorku.ca

*Abstract* **The subcellular localization of a protein is closely correlated with its functions. Although many machine learning algorithms have been developed and applied to predict protein compartments using different data sources, such as protein amino acid sequence and motif information, automatic prediction of subcellular localization remains a challenging problem. In this study, we compared three support vector machines (SVM)-based computational strategies for interpreting differential detection of proteins in isolated organellar compartments by high-throughput mass spectrometry. The main focus is on how to deal with multi-compartmental ambiguity in predicting protein subcellular localizations. When applied them to a global-scale proteomic study, their Area Under the Receiver Operating Characteristics Curves (ROC) for four major organellar compartments (cytosol, microsomes, mitochondria, and nucleus) are more than 0.75.**

## I. INTRODUCTION

Elucidation of gene function and protein regulatory mechanisms is a fundamental objective in human biology. It is well-known that determining the subcellular localization of a protein in a cell is a key to understanding its function and can facilitate biochemical experiments aimed at characterizing additional biological properties, such as purification. However, traditional experimental methods for examining subcellular localization are generally time-consuming and costly. Given the rapidly expanding plethora of uncharacterized proteins identified by the many ongoing genome sequencing projects, it is highly desirable to deduce or predict a protein's subcellular localization automatically [1, 2].

Currently, most of the automatic protein subcellular localization prediction methods fall into one of three categories [3]. The first one is prediction based on amino acid composition, as originally suggested by Nakashima and Nishikawa [4]. Different machine learning algorithms have been developed that make use amino acid composition information towards this end, including neural networks [5], support vector machines (SVM) [6], covariant discrimination [7] and augmented covariant discrimination methods [8], as

well as SVM incorporating quasi-sequence-order effects ([9]. The second major approach is prediction based on calculating a set of sequence-derived parameters and comparing these with a representation of a number of localization rules that have been collated from the literature. The most widely used algorithm in this category is the popular PSORT algorithm [10], which is a commonly-used bioinformatics tool. The key idea of this approach is to decide the presence of various sequence motifs that enable proteins to be localized to a certain compartment. Different types of prior knowledge are required for this determination, which are, actually, hard to get for uncharacterized proteins. The third category of prediction is the homology-based prediction [2, 11, 12], wherein the inferences are based on transference of knowledge from characterized to unknown homologous proteins.

One of main limitations in most of these studies is that their principle methods focus on mono-compartment prediction (that is, a protein is presumed to localize to a single organelle only). For example, Lu et al. constructed a parser to extract a simple ontological representation for proteins assigned to multiple compartments, without exploiting the information encoded by multi-localizations [2]. Park and Kanehisa did not included proteins annotated with two or more subcellular locations in their analysis [13].

As an alternate to sequence- or homology-based predictions, proteomic methods based on subcellular fractionation in combination with high-throughput protein mass spectrometry have emerged as a powerful alternative experimental platform for assessing subcellular localization directly. Indeed, substantive recent technical advances now make this the preferred approach for genome-wide protein identification and quantification with high sensitivity and accuracy [14]. Compared with previous sequence information-derived prediction methods, these newer proteomic profiling-based screening methods are also proving to be more effective for resolving ambiguous or difficult localization problems [15]. However, current procedures involving biochemical methods for subcellular fractionation

are still far from perfect, and artifacts due to cross-contamination can create misleading results.

An in-depth comparative proteomic analysis of the organelles of six representative mouse organs (adult brain, heart, kidney, liver, lung, and embryonic placenta) was recently carried out [16], in which computational and statistical procedures were used in combination with available conventional annotations and the observed proteomic profiles to create a high-quality reference map of the putative subcellular localization and tissue-selectivity of 4768 proteins. The present study is devoted to addressing the multi-compartment problem in subcellular localization prediction. The specific objective is to compare some computational strategies for multi-compartmental prediction of protein subcellular localizations based on support vector machine (SVM) methods, including one-class SVM [17], binary SVM [18] and multi-class SVM [19], which cannot handle multi-labeled prediction of subcellular localizations directly.

## II. DATA SETS AND PREPROCESSING

In this global-scale mouse proteomic study, healthy adult brain, heart, kidney, liver, lung and embryonic placenta were excised from euthanized 6-8 week old ICR female mice. The tissues were gently disrupted and fractionated into the four major subcellular compartments (cytosol, microsomes, mitochondria, and nuclei) using differential ultracentrifugation. The proteins were identified by tandem mass spectrometry followed by database searches of the acquired spectra using the multidimensional protein identification technology [14]. The procedures for processing, searching and rigorously evaluating the proteomic expression profiles have been detailed by [20, 21]. A total of 4768 proteins were confidently identified in this analysis. Nearly half (2390) of the identified proteins lacked a previously assigned subcellular localization based on annotation obtained from the ExPASy Server. Protein relative abundance was estimated in the respective fractions based on the ratio of the cumulative number of spectra matching to any given protein in each sample [22].

In order to generate a suitable supervised learning approach for predicting the 2390 proteins with unknown subcellular localizations, we need to obtain a reference set of proteins with known subcellular localizations. For this, we obtained the annotations for 1558 proteins from the SWISS-PROT database (http://ca.expasy.org/sprot/). Additionally, we compiled a set of 820 proteins that had been independently identified in a single highly purified organelle in a previous proteomic study [15, 23-29]. Table 1 show a summary of the number of proteins in per monocompartmental and multicompartmental localizations. As we can see from the table, 580 of the 2378 annotated proteins have multicompartmental localizations.

TABLE I
THE NUMBER OF PROTEINS IN PER MONOCOMPARTMENTAL AND
MULTICOMPARTMENTAL LOCALIZATIONS

| Class Label | Subcellular localization | Number of proteins |
|---|---|---|
| 1 | Cytosol (Cyto) | $301^{[1]}$ ($741^{[2]}$) |
| 2 | Microsomes (Micro) | 640 (1052) |
| 3 | Mitochondria (Mito) | 253 (405) |
| 4 | Nucleus (Nuc) | 604 (821) |
| 5 | Cyto_Micro[3] | 219 |
| 6 | Cyto_Mito | 20 |
| 7 | Cyto_Nuc | 143 |
| 8 | Micro_Mito | 104 |
| 9 | Micro_Nuc | 31 |
| 10 | Mito_Nuc | 3 |
| 11 | Cyto_Micro_Mito | 20 |
| 12 | Cyto_Micro_Nuc | 35 |
| 13 | Cyto_Mito_Nuc | 2 |
| 14 | Micro_Mito_Nuc | 2 |
| 15 | Cyto_Micro_Mito_Nuc | 1 |
| 16 | Unknown | 2390 |

[1]The number of monocompartmental proteins in cytosol
[2]The number of monocompartmental and multicompartmental proteins in cytosol
[3]Multicompartment localization: Cytosol and Microsomes

## III. METHODS

As discussed above, most of the previously studies [4-9] focused on the monocompartmental prediction, which was formulated as a multi-class classification problem (see Figure 1(a)). Classes are mutually excusive by definition; that is, each case can be assigned to only one of several alternate classes. Essentially, protein subcellular localization prediction is a multicompartmental prediction problem, since some proteins may coexist in several different subcellular locations (see TABLE I). Therefore, we are indeed facing a multi-class multi-label classification problem (see Figure 2 (b)). So far, there are no effective computational procedures that can be used to treat this difficult multiplex (i.e., multi-label, multi-localization) problem [11].

### A. Computational Strategies for multi-label prediction of protein subcellular localizations

As a first step towards resolving this multiplex problem, here we compared three computational strategies as follows:
**Strategy one:** It is called "cross-training" [30], a variant of "One vs. All" algorithm [31], in which each class is
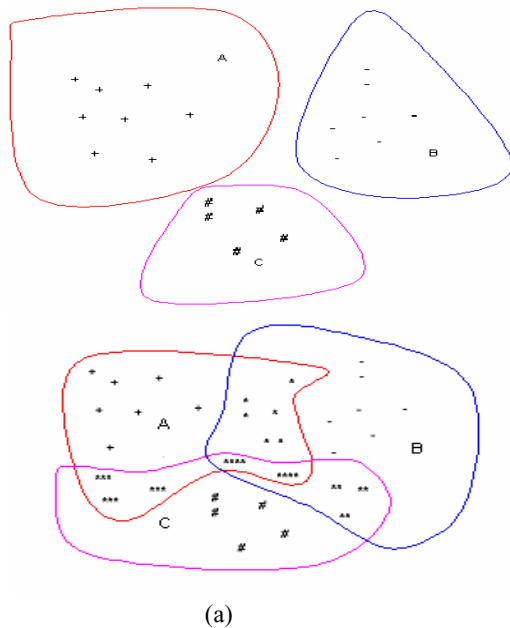
(a)

Fig 1. Considering three classes (e.g. subcellular locations) A, B and C, we denote proteins that belong to A, B, C are circled by red, green and purple curves, respectively. Figure 1(a) is the classic classification problem while Figure 1 (b) is a multi-label classification problem occurred in our study. Proteins uniquely belonging to classes A, B and C are denoted by "+", "-" and "#", respectively. Proteins belonging to both A and B classes are denoted by "*", both B and C classes are denoted by "**" and both A and C classes are denoted by "***". Proteins belonging to classes A, B and C simultaneously are denoted by "****".

compared to all others. After the binary classification problems have been solved, the resulting set of binary classifiers was combined in some way. The cross-training algorithm has been applied with some success as a means of rationalizing pattern recognition as applied to multi-label semantic scene classification. In our implementation of this approach, we used the multi-labeled proteins as positive examples for each of the four associated localization classes during training. For example, if a protein was annotated as both nuclear and mitochondrial, it was considered as a positive example during training of both the nuclear and mitochondrial classes, but never as a negative example of either category. Therefore, the number of positive training proteins for cytosol, microsomes, mitochondria and nuclei is 741, 1052, 405 and 821 respectively. After this processing, we can train a binary-class SVM for each subcellular localization.

**Strategy two:** It is called "super-class training". The algorithm works as: a new class is defined for the proteins in each combination of compartments and a model is built for it. For example, in Figure 1 (b), we can define a new class for the proteins belonging to both A and B classes (denoted by "*"), both B and C classes (denoted by "**"), both A and C

classes (denoted by "***") and classes A, B and C simultaneously (denoted by "****"). The algorithm is also called MODEL-n (n stands for "new" class) by Boutell et al. [30], but they did not test and evaluate the performance of this method. TABLE II shows the number of training proteins used for training in each monocompartmental subcellular localization and multicompartmental subcellular localization. Since some of these combinations have few proteins, we just kept the localizations with at least 20 proteins. After this processing, we can train a multi-class SVM for all selected subcellular localizations.

TABLE II
THE NUMBER OF PROTEINS IN PER MONOCOMPARTMENTAL AND MULTICOMPARTMENTAL LOCALIZATIONS USED FOR SUPER-TRAINING

| Index | Label | Subcellular localization | # training | # prediction |
|---|---|---|---|---|
| A | 1 | Cytosol (Cyto) | 301 | |
| B | 2 | Microsomes (Micro) | 640 | |
| C | 3 | Mitochondria (Mito) | 253 | |
| D | 4 | Nucleus (Nuc) | 604 | 2390 |
| AB | 5 | Cyto_Micro | 219 | |
| AC | 6 | Cyto_Mito | 20 | |
| AD | 7 | Cyto_Nuc | 143 | |
| BC | 8 | Micro_Mito | 104 | |
| BD | 9 | Micro_Nuc | 31 | |
| ABC | 10 | Cyto_Micro_Mito | 20 | |
| ABD | 11 | Cyto_Micro_Nuc | 35 | |

**Strategy three:** It is called one-class SVM [17]. Differentiation of members between known classes (i.e., subcellular localizations) is achieved by a data domain description. This is done by estimating a binary function that is positive where most of the data are located and negative elsewhere. A hyperplane with the largest possible margin is chosen to separate the training data from where the novel data are assumed to be. Once the data domain description is known, the problem is reduced to a classification task where only one class exists. The number of training proteins for cytosol, microsomes, mitochondria and nuclei is 741, 1052, 405 and 821 respectively. The algorithm was detailed by Manevitz and Yousef [17].

*B. Evaluating Performance of the Learning Strategies*

We evaluated the performance of our machine learning strategies using the standard method of stratified 10-fold cross-validation [32]. In this procedure, we randomly divided the training set associated with each subcellular compartment $m$ $(m=1,…, k)$ into 10 sub-groups ($G_{m1},...,G_{m10}$), keeping the number of proteins in the localization class approximately the same across each training category. For "cross-training" and one-class SVM algorithms, we constructed 10 different classifiers ($C_{m1}$, $C_{m2}$, …, $C_{m10}$) for each of the four subcellular compartments , where $C_{mi}$ use all of the training

proteins from all of the groups except $G_{mi}$ $(m=1,..., 4)$. Proteins in group $G_{mi}$ were used for testing classifier $C_{mi}$. For "super-class training" algorithm, we constructed 10 different classifiers ($C_1$, $C_2$, ..., $C_{10}$) for all 11 subcellular compartments, where $C_i$ use all of the training proteins from all of the groups except ($G_{1i},...,G_{ki}$)$(k=11)$. Facing such a multicompartmental localization prediction problem, different performance metrics have been applied. For example, Lu et al. used four standard statistics: specificity, precision, sensitivity and recall (the last two are identical) [2]. Scott et al. applied precision and sensitivity [3]. Here for each of the three strategies, we used the following measures to assess performance:

**Total Accuracy (TA)**: the rate of the total of correct predictions (true positives) $T_m$ in each subcellular localization $m$ compared to all predictions N, that is $TA = \frac{\sum_{M=1}^{K} T_m}{N}$, where $K$ is the total number of subcellular localizations

**F-measure:** It is equal to 2RP/(R+P), where P (Precision) is the portion of true positives with respect all predicted positive for a given location, that is P=TP/(TP+FP), and R (Recall/ Sensitivity) is the portion of true positive with respect the sum of true positive and false negatives for a given location, that is R=TP/(TP+FN).TP, FP, TN, FN denote the total number of true positives, false positives, true negatives and false negatives, respectively.

**Area Under the Receiver Operating Characteristics (ROC) Curves:** ROC curves and their Area Under the Curves (AUC) can also be used to evaluate the power of different classifiers for predicting protein subcellular localization. ROC curves have been used to depict the pattern of sensitivity and specificity observed when the performance of a classifier is evaluated at different thresholds [33]. Since the prediction confidence (probability) from trained classifiers varies between 0 and 1, we created 100 thresholds of equal interval across the range of prediction confidence. For each of the 100 thresholds, we calculated classifier specificity, sensitivity based on the cross-validation results.

## IV. RESULTS

We applied the discussed SVM-based strategies with radial basis kernel function to our proteome-wide mouse data. TABLE III shows the 10-fold cross validation results for training data using cross-training method and one class SVMs. The 10-fold cross validation results for training data using super-class training method are shown in the TABLE IV. Overall, the classification performances of classifiers built for mitochondria and nucleus are much better than cytosol and microsomes. The former two reached F-measure values more than 65% in all three strategies while the latter two have only 43-67% F-measure values.

TABLE III
10-FOLD CROSS-VALIDATION PERFORMANCE ON TRAINING
DATA USING CROSS-TRAINING AND ONE-CLASS SVM

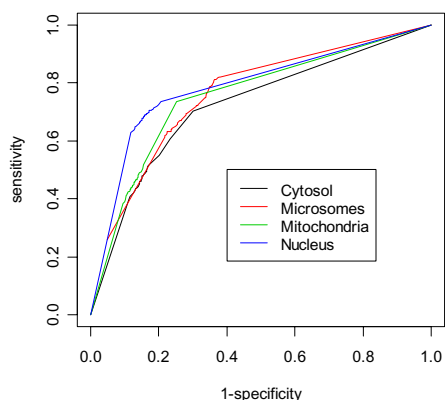| Subcellular localization | Cross-Training | | One Class SVM | |
| --- | --- | --- | --- | --- |
| | F-measure (%) | Total Accuracy (%) | F-measure (%) | Total Accuracy (%) |
| Cytosol | 43.6 | | 54.9 | 59.6 |
| Microsomes | 51.6 | 57.4 | 62.7 | |
| Mitochondria | 65.8 | | 68.6 | |
| Nucleus | 68.5 | | 71.6 | |

TABLE IV
10-FOLD CROSS-VALIDATION PERFORMANCE ON TRAINING
DATA USING SUPER-CLASS TRAINING

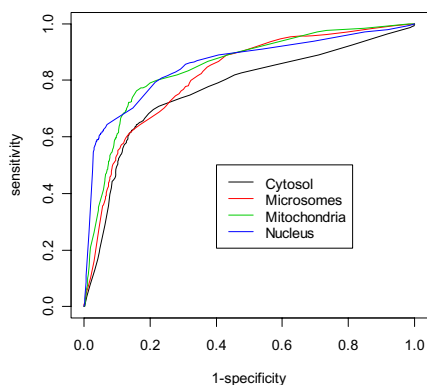| Subcellular localization | F-measure (%) | Total Accuracy (%) |
| --- | --- | --- |
| Cytosol (Cyto) | 61.6 | |
| Microsomes (Micro) | 66.8 | |
| Mitochondria (Mito) | 69.5 | |
| Nucleus (Nuc) | 74.1 | |
| Cyto_Micro | 55.2 | 63.8 |
| Cyto_Mito | 57.5 | |
| Cyto_Nuc | 45.5 | |
| Micro_Mito | 40.8 | |
| Micro_Nuc | 47.3 | |
| Cyto_Micro_Mito | NA[1] | |
| Cyto_Micro_Nuc | NA | |

[1]The denominator is zero in the formula to calculate precision

For the three learning strategies, super-class training is better than cross-training and one-class SVM, since its total accuracy is 63.8%, which is higher than that of cross-training (57.4%) and one-class SVM (59.6%). The results based on super-class training algorithm show that the performance of classifiers learned on monocompartmental localizations is better than that of classifiers learned on multicompartmental localizations. The classifiers trained on three-compartmental localizations are not learnable in this data set. When we closely look the predicted results for the proteins which have multicompartmental localizations, we found that many of these proteins are partial-correctly predicted. For example, protein A is in cytosol and microsomes, but the classifier based on super-class training can only assign it to cytosol or microsomes, not both.
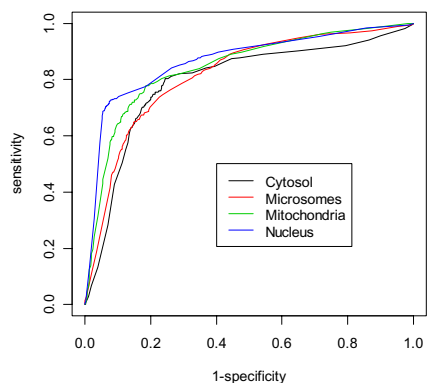
Figure 2 shows the ROC curves of the four subcellular compartments based on the specificity and sensitivity calculated from the 100 thresholds of predicted confidence for the three learning strategies. For the super-class training, we just show the results of four base-subcellular compartments, since other super-classes have no good performance (see TABLE IV). The AUC values of these ROC curves shown in Figure 2 are listed in TABLE V. Overall, the performance of super-class training approach is still a little better than those of cross-training method and one-class SVM. Most of the

(a) Cross-training



(b) One-class SVM



(c) Super-class Training

Fig. 2. ROC plots for the performance of the classifiers using three learning strategies

AUC values based on the super-class training algorithm are larger than 80%. The classifiers for mitochondria nucleus have the best performance in the three algorithms.

TABLE V
AREA UNDER THE CURVES (AUC) FOR THE FOUR SUBCELLULAR LOCALIZATIONS

| Subcellular localization | Cross-Training (%) | One Class SVM (%) | Super-class Training (%) |
|---|---|---|---|
| Cytosol | 72.3 | 76.8 | 79.5 |
| Microsomes | 76.7 | 81.6 | 81.6 |
| Mitochondria | 75.5 | 84.8 | 84.4 |
| Nucleus | 78.7 | 85.6 | 86.7 |

V. CONCLUSION

In this study, we have evaluated three SVM-based computational strategies for predicting protein subcellular localizations on a proteome-wide scale. The approaches address some of the key problems associated with predicting multiple organellar compartments given proteins of uncertain association. The first extension is based on the classic "One vs All" method, called as "cross-training"; the second extension is based on the classic multi-class SVM method, named as "super-class training"; the last one is based on one-class SVM.

Using large-scale proteomics data as a building block, the first algorithm achieved a total accuracy 57.4%, the second algorithm reached a total accuracy 63.8% and the last one have the total accuracy 59.68% over the four major cellular compartments.

Some reasons may explain why the cross-training algorithm has not shown good performance. One of these is probably due to the unbalance between the size of positive samples and that of the negative samples in the training data. A possible solution to this is that we may need to subsampling the large size of negative samples. The worse performance of the one-class SVM algorithm is probably because the approach tries to catch the distribution of a given class while the distributions of other classes have not been seen in the training. For super-class training algorithm, when the number of monocompartments is very large, it will also have some deficiencies since we may generate a lot of trivial super-classes, which may not learnable as shown in this study, such as super-classes: Cyto_Micro_Mito and Cyto_Micro_Nuc.

REFERENCES

[1] Cai, Y.D., Chou, K.C. (2004). Predicting subcellular localization of proteins in a hybridization space. *Bioinformatics,* 20:1151-1156.

[2] Lu, Z., Szafron, D., Greiner, R., Lu, P., Wishart, D.S., Poulin, B., Anvik, J., Macdonell, C., Eisner, R. (2004). Predicting subcellular localization of proteins using machine-learned classifiers. *Bioinformatics,* 20:547-556.

[3] Scott, M.S., Thomas, D.Y., Hallett, M.T. (2004). Predicting subcellular localization via protein motif co-occurrence. *Genome Res.*, 14:1957-1966.

[4] Nakashima, H., Nishikawa, K. (1994). Discrimination of intracellular and extracellular proteins using amino acid composition and residue-pair frequencies. *J. Mol. Biol.,* 238:54-61.

[5] Reinhardt, A., Hubbard, T. (1998). Using neural networks for prediction of the subcellular location of proteins, *Nucleic Acids Res.,* 26:2230-2236.

[6] Hua, S., Sun, Z. (2001). Support vector machine approach for protein subcellular localization prediction. *Bioinformatics*, 17:721–728.

[7] Chou, K. C., Elrod, D.W. (1998). Using discriminant function for prediction of subcellular location of prokaryotic proteins. *Biochem. Biophys. Res Commu,* 252:63-68.

[8] Chou, K. C. (2000). Prediction of protein subcellular locations by incorporating quasi-sequence-order effect. *Biochem. Biophys. Res. Commun*, 278:477–483.

[9] Cai, Y. D., Liu, X. J., Xu, X. B., Chou, K. C. (2002). Support vector machines for prediction of protein subcellular location by incorporating quasi-sequence-order effect. *J. Cell. Biochem,* 84:343–348.

[10] Nakai, K., Kanehisa, M. (1992). A knowledge base for predicting protein localization sites in eukaryotic cells. *Genomics,* 14:897-911.

[11] Chou, K.C., Cai, Y.D. (2005). Predicting protein localization in budding yeast. *Bioinformatics,* 21:994-950.

[12] Mott, R., Schultz, J., Bork, P., Ponting, C.P. (2002). Predicting protein cellular localization using a domain projection method. *Genome Res.,* 12:1168-1174.

[13] Park, J.K., Kanehisa, M. (2003). Prediction of protein subcellular localizations by support vector machines using compositions of amino acids and amino acid pairs. *Bioinformatics*, 19:1656-1663.

[14] Yates, J. R. (2004). Mass spectral analysis in proteomics. *Annu Rev Biophys Biomol Struc.,* 33:297-316.

[15] Schirmer, E. C., Florens, L., Guan, T., Yates, J. R., Gerace, L. (2005). Identification of novel integral membrane proteins of the nuclear envelope with potential disease links using subtractive proteomics. *Novartis Found Symp,* 264:63-76; discussion 76-80, 227-230.

[16] Kislinger, T., Cox, B., Kannan, A., Chung, C., Hu, P., Ignatchenko, A., Scott, M.S., Gramolini, A., Morris, Q., Hughes, T., Rossant, J., Frey, B., Emili, A. (2006) Global survey of organ and organelle protein expression in mouse: combined proteomic and transcriptomic profiling. *Cell,* 125:173-186.

[17] Manevitz, L.M., Yousef, M. (2002). One-class svms for document classification. *Journal of Machine Learning Research*, 2:139-154.

[18] Vapnik, V. (1998). Statistical learning theory. New York: Wiley.

[19] Weston, J., Watkins, C. (1999). Support vector machines for multiclass pattern recognition. *In Proceedings of the Seventh European Symposium on Artificial Neural Networks*. M. Verleysen, Ed., Brussels, Belgium: D-Facto Public, pp219-224.

[20] Kislinger, T., Emili, A. (2003). Going global: protein expression profiling using shotgun mass spectrometry. *Curr Opin Mol Ther,* 5:285-293.

[21] Kislinger, T., Rahman, K., Radulovic, D., Cox, B., Rossant, J., Emili, A. (2003). PRISM, a Generic Large Scale Proteomic Investigation Strategy for Mammals. *Mol Cell Proteomics*, 2:96-106.

[22] Liu, H., Sadygov, R. G., Yates, J. R. (2004). A model for random sampling and estimation of relative protein abundance in shotgun proteomics. *Anal Chem.* 76:4193-4201.

[23] Andersen, J. S., Lam, Y. W., Leung, A. K., Ong, S. E., Lyon, C. E., Lamond, A. I., Mann, M. (2005). Nucleolar proteome dynamics, *Nature,* 433:77-83.

[24] Beausoleil, S. A., Jedrychowski, M., Schwartz, D., Elias, J. E., Villen, J., Li, J., Cohn, M. A., Cantley, L. C., Gygi, S. P. (2004). Large-scale characterization of HeLa cell nuclear phosphoproteins. *Proc Natl Acad Sci U S A,* 101:12130-12135.

[25] Krapfenbauer, K., Fountoulakis, M., Lubec, G. (2003). A rat brain protein expression map including cytosolic and enriched mitochondrial and microsomal fractions. *Electrophoresis*, 24:1847-1870.

[26] Mootha, V. K., Bunkenborg, J., Olsen, J. V., Hjerrild, M., Wisniewski, J. R., Stahl, E., Bolouri, M. S., Ray, H. N., Sihag, S., Kamal, M., et al. (2003). Integrated analysis of protein composition, tissue diversity, and gene regulation in mouse mitochondria. *Cell*, 115:629-640.

[27] Nielsen, P. A., Olsen, J. V., Podtelejnikov, A. V., Andersen, J. R., Mann, M., Wisniewski, J. R. (2005). Proteomic mapping of brain plasma membrane proteins. *Mol Cell Proteomics,* 4:402-408.

[28] Wu, C. C., MacCoss, M. J., Howell, K. E., Yates, J. R. (2003). A method for the comprehensive proteomic analysis of membrane proteins. *Nat Biotechnol,* 21:532-538.

[29] Wu, C. C., MacCoss, M. J., Mardones, G., Finnigan, C., Mogelsvang, S., Yates, J. R., Howell, K. E. (2004). Organellar proteomics reveals Golgi arginine dimethylation. *Mol Biol Cell,* 15:2907-2919.

[30] Boutell, M., Shen, X., Luo, J., Brown, C. (2004). Learning multi-label semantic scene classification. *Pattern Recognition*, 37:1757-1771.

[31] Yeang, C.H., Ramaswamy, S., Tamayo, P., Mukherjee, S., Rifkin, R.M., Angelo, M., Reich, M., Lander, E., Mesirov, J., Golub, T. (2001). Molecular classification of multiple tumor types. *Bioinformatics*. 17 suppl., S316-S322.

[32] Mitchell, T.M. (1997). Machine Learning. McGraw-Hill, N.Y.

[33] Bradley, A.P. (1997). The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition*, 30:1145-1159.