

# Protein Structure Prediction Using Robust Principal Component Analysis and Support Vector Machine

Nur Aini Zakaria<sup>1</sup>, Zuraini Ali Shah\*<sup>2</sup>

Department of Software Engineering, Faculty of Computing, Universiti Teknologi Malaysia,  
81310 Johor Bahru, Johor, Malaysia

<sup>1</sup>ainizakaria@gmail.com, <sup>2</sup>aszuraini@utm.my

## Abstract

*Existence of bioinformatics is to increase the further understanding of biological process. Proteins structure is one of the major challenges in structural bioinformatics. With former knowledge of the structure, the quality of secondary structure, prediction of tertiary structure, and prediction function of amino acid from its sequence increase significantly. Recently, the gap between sequence known and structure known proteins had increase dramatically. So it is compulsory to understand on proteins structure to overcome this problem so further functional analysis could be easier. The research applying RPCA algorithm to extract the essential features from the original high-dimensional input vectors. Then the process followed by experimenting SVM with RBF kernel. The proposed method obtains accuracy by 84.41% for training dataset and 89.09% for testing dataset. The result then compared with the same method but PCA was applied as the feature extraction. The prediction assessment is conducted by analyzing the accuracy and number of principal component selected. It shows that combination of RPCA and SVM produce a high quality classification of protein structure.*

**Keywords:** Protein Structure, Robust Principal Component analysis, Support Vector Machine

## 1.0 Introduction

The functional and structural annotation of protein domain is one of the important roles in bioinformatics. In this context, protein structure information plays an important information key of their structural part also the features related to the biological function (Sahu et al., 2009) such as prediction of DNA binding site, implementation of a heuristic approach to find tertiary structure, reduction of conformation search space and also characterizing the folding type of a protein or its domain. Li, et al. (2012) state that the exponentially growth of newly discovered protein sequences by different scientific community caused a large gap between the number of sequence-known and the number of structure-known proteins. Hence, there exist critical challenges to develop automated method for fast and accurate determination of

the structures of proteins in order to reduce gap. Therefore, there is a compulsory to implement reliable and effective computational methods for identifying the structural class of newly discovered protein based on their primary sequences. Biological data is a big and massive dataset. In order to avoid facing difficulties in handling a large amount of data like redundancy, large memory amount requirement and high computation power consumption, feature extraction has been used by many researchers to reduce amount of resource required.

The aim in this study is to compare and evaluate the performance of the commonly used supervised and unsupervised machine learning tools specifically for protein structural class prediction. Therefore the purposes of this research are to implement Robust Principal Component Analysis (RPCA) to determine the number of principal component Support Vector Machine (SVM) for protein structure classification. Lastly, to evaluate the performance of RPCA and SVM based on accuracy.

## 2.0 Background of The Study

For the past decades, a lot of research was conducted to predict protein structure due to its important role in bioinformatics. Protein secondary structure prediction was proposed by Faraggi, et al. (2012), Y-F Huang and S-Y Chen (2013), and Misconai et al., (2015). Faraggi, et al., (2012) using a multi-step neural network algorithm (SPINEX) that improves the previous method, SPINE and Cheol Jeong, J., Lin, X., & Chen, X. W. (2013) proposed position-specific scoring matrix (PSSM) together with physiochemical features as based for prediction by SVM. While Misconai, András, et al., (2015) developed  $\beta$ -structure selection (BeStSel) which involving the twist of  $\beta$ -structures. Singh, Lavneet, Girija Chetty, and Dharmendra Sharma (2012), and Sułkowska, Joanna I., et al., (2012) predict protein structure based on their folding information. Singh, Lavneet, Girija Chetty, and Dharmendra Sharma (2012) applies extreme learning machine (ELM) algorithm and SVM along with PCA and LDA as feature selection. Sułkowska, Joanna I., et al., (2012) proposed a hybrid method, direct coupling analysis (DCA)-fold that incorporates DCA contact with local information. Hopf, Thomas A., et al., (2012) shows genomic sequencing can predict three-dimensional protein structure through the development of EVfold\_membrane algorithm. A mix of protein structure network (PSN) and elastic network model (ENM), PSN-ENM proposed by Raimondi, Francesco et al., (2013) for structural communication prediction in biomolecular systems. Braun, Tatjana, Julia Koehler Leman, and Oliver F. Lange, (2015) present RASREC a structure prediction using evolutionary information with a resolution adapted structural recombination approach of Rosetta. A wide range of studies using machine learning-based methods for protein structure prediction has been developed by the past researchers to solve the case in order to facilitate an automated, high throughput assignment. So far, among all the methods developed, the Support Vector Machine (SVM) machine learning is the most used since it was considered as one of the best computer algorithm in this field for classification. It is also easy to apply and flexible. Feature extraction is the best method to extract valuable information of a dataset besides reduces the data dimension. PCA is a good method and widely used method for feature extraction because it emphasizes variation and bring out strong patterns in a dataset so the data easily explored and visualized. However, PCA is unable to detect and ignore outliers. Therefore, RPCA is the best solution for feature extraction since it has better extraction quality and work better in grossly corrupted dataset observation.

### 3.0 Methodology

Firstly, the current issues of protein structure prediction are investigated followed by collecting research materials such as journals, articles, conference paper and others. The data preprocessing conducted to gain higher and better prediction success rate and system performance. It also help to minimizing error in preparation be validated by machine learning algorithm. Datasets by Ding, Chris HQ, and Inna Dubchak (2012) filtered to remove unnecessary values and information. Research continues by applying Principal component analysis (PCA) and RPCA (Croux, Christophe, and Anne Ruiz-Gazen, 2005)) algorithm to extract the essential features from the original high-dimensional input vectors. The process continued by experimenting SVM with RBF kernel using the reduced and normalized features by PCA and RPCA. The final phase is the prediction assessment of the application of RPCA and SVM by the comparison of recognition ratio compared between different methods and methods used by previous researcher. Performance testing of this research by comparing classification result of protein by overall accuracy that expressed in equation 1.

$$\frac{\text{correctly recognize protein} = \text{correctly recognize number of query protein}}{\text{total number of protein}} \quad (1)$$

### 4.0 Result

The experiment was conducted by using three approaches in order to analyse the performance of RPCA and SVM. In order to gives a clear view on performance of RPCA, the method was compared with the PCA (the basic of RPCA) and SVM. In order to select the components that contain >60% of variance, the number of PC selected are different accordingly. Table 1 shows that number of selected PC in training dataset is lower compared to testing dataset. Table 2 shows the accuracy percentage of tested approach divided by training and testing datasets.

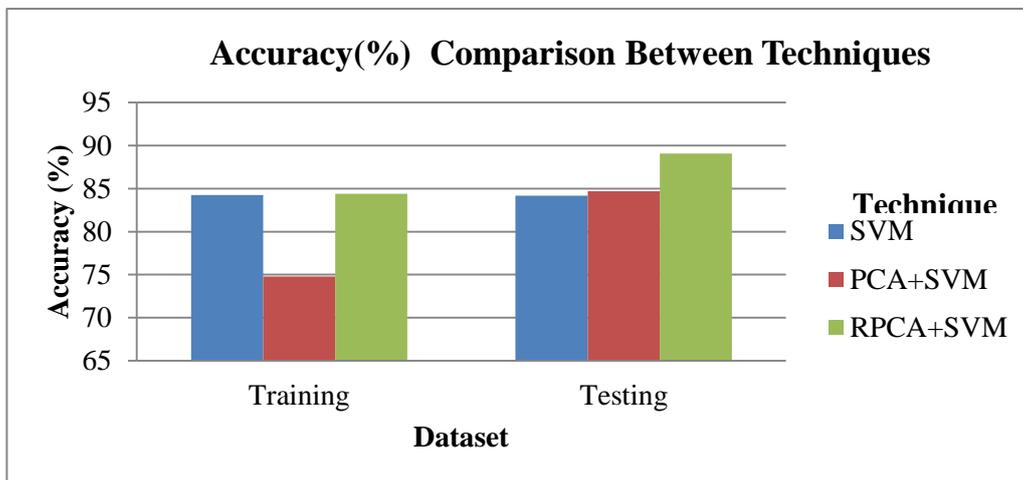
**Table 1** Number of PC selected for classification

Feature extraction	Number of PC selected for classification	
	Training	Testing
PCA	2	3
RPCA	2	4

**Table 2** Comparison of SVM, PCA+SVM and RPCA+SVM

Technique	Training Dataset Accuracy (%)	Testing Dataset Accuracy (%)
SVM	84.25	84.16
PCA+SVM	74.79	84.68
RPCA+SVM	84.41	89.09

## 5.0 Discussion



**Figure 1** Accuracy comparison between techniques.

Based on this analysis, it can be assumed that differences in data characteristics will influence the number of sufficient PCs required in both PCA and RPCA approaches. The number of PCs required for the training dataset is less than for the testing dataset since the size of the training dataset is larger, so it may contain higher information and better interpretation of features compared to the testing dataset.

From the results in Figure 1, it can be seen that the non-extracted features technique (only SVM) achieves a high percentage of accuracy (84.25% and 84.16%). However, the result is doubtful since models built on extracted features may be of higher quality, as the data is described by fewer, more meaningful attributes. Results obtained by combining PCA and SVM are 74.79% on the training dataset and 84.68% on the testing dataset. The accuracy on both datasets is quite high but still lower than the combination of RPCA and SVM technique (84.41% on the training dataset and 89.09% on the testing dataset). The gap seems to be higher in the training dataset, possibly due to the larger number of outliers. RPCA appears to perform best since this method is not influenced much by outliers and its ability to detect exact fit situations.

Table 3 shows the comparison of accuracy percentage of PCA and RPCA combination with SVM. Even according by number of component, the RPCA method always seems to lead in terms of accuracy. This proves the effectiveness of RPCA approach. Table 3 also shows the increasing pattern of the accuracy for both datasets. It can be assume that higher number of PC contain much more data information lead to higher accuracy.

Singh, Lavneet, Girija Chetty, and Dharmendra Sharma (2012) apply the same dataset (feature vector described by Ding, Chris HQ, and Inna Dubchak, 2001) to predict protein structure using PCA and LDA based in Extreme Learning Machine (ELM). According to the Table 4, it can be seen that proposed method used in this research shows promising results in term of accuracy obtained compare to the proposed method proposed by Singh, Lavneet, Girija Chetty, and Dharmendra Sharma (2012). This shows that feature extraction using RPCA and classification using SVM is an efficient method for protein structure prediction. It also shows that method proposed by Singh, Lavneet, Girija Chetty, and Dharmendra Sharma (2012) has drawbacks in due to the outliers and low ability in detection of exact fit situation.

**Table 3** Comparison of PCA+SVM and RPCA+SVM based on number of component

Number of Principal Component (PC)	Accuracy (%)		Accuracy (%)	
	For Training Dataset		For Testing Dataset	
	PCA +SVM	RPCA+SVM	PCA+SVM	RPCA+SVM
1	51.91	80.60	55.84	77.92
2	74.79	84.41	84.68	84.94
3	82.75	86.90	87.53	88.05

**Table 4** Accuracy Comparison between Method

Method	Accuracy (%)
LDA-ELM	77.67
PCA-ELM	82.45
RPCA-SVM	89.09

## 6.0 Conclusion

This research focus is on protein structural classification. Protein Structure classification is important for identification of protein function. As the protein structure classification is a first and key step in protein structure prediction, it becomes an increasingly challenging task. Recently, the exponentially increase of sequence data protein cause the increasing of the requirements for reliable and effective computational method for protein structure classification. Protein structure classification is very important in bioinformatics field. Proposed feature extraction method, Robust Principal Component Analysis (RPCA) combines with Support Vector Machine (SVM) shows that data with extracted features can obtain higher accuracy (84.41% for training dataset and 89.09% for testing dataset). It also shows that RPCA works well with highly corrupted data especially dataset with outliers.

## References

- Braun, T., Leman, J. K., & Lange, O. F. (2015). Combining evolutionary information and an iterative sampling strategy for accurate protein structure prediction. *PLoS Comput Biol*, 11(12), e1004661.
- Cheol Jeong, Jong, Xiaotong Lin, and Xue-Wen Chen. "On position-specific scoring matrix for protein function prediction." *IEEE/ACM transactions on computational biology and bioinformatics* 8.2 (2011): 308-315.
- Croux, C., & Ruiz-Gazen, A. (2005). High breakdown estimators for principal components: the projection-pursuit approach revisited. *Journal of Multivariate Analysis*, 95(1), 206- 226.
- Ding, C. H., & Dubchak, I. (2001). Multi-class protein fold recognition using support vector machines and neural networks. *Bioinformatics*, 17(4), 349-358.
- Faraggi, Eshel, et al. "SPINE X: improving protein secondary structure prediction by multistep learning coupled with prediction of solvent accessible surface area and backbone torsion angles." *Journal of computational chemistry* 33.3 (2012): 259-267.
- Hopf, T. A., Colwell, L. J., Sheridan, R., Rost, B., Sander, C., & Marks, D. S. (2012). Three dimensional structures of membrane proteins from genomic sequencing. *Cell*, 149(7),1607-1621.
- Li L, Cui X, Yu S, Zhang Y, Luo Z, Yang H, et al. (2014) "PSSP-RFE: Accurate Prediction of Protein structure by Recursive Feature Extraction from PSI-BLAST Profile, Physical Chemical Property and Functional Annotations." *PLoS ONE* 9(3): e92863. doi:10.1371/journal.pone.0092863,
- Micsonai, A., Wien, F., Kernya, L., Lee, Y. H., Goto, Y., Réfrégiers, M., & Kardos, J. (2015). Accurate secondary structure prediction and fold recognition for circular dichroism spectroscopy. *Proceedings of the National Academy of Sciences*, 112(24),E3095-E3103.
- Raimondi, F., Felling, A., Seeber, M., Mariani, S., & Fanelli, F. (2013). A mixed protein structure network and elastic network model approach to predict the structural communication in biomolecular systems: the PDZ2 domain from tyrosine phosphatase 1E as a case study. *Journal of Chemical Theory and Computation*, 9(5), 2504-2518.
- Singh, L., Chetty, G., & Sharma, D. (2012, July). A hybrid approach to increase the performance of protein folding recognition using support vector machines. In *International Workshop on Machine Learning and Data Mining in Pattern Recognition* (pp. 660-668). Springer Berlin Heidelberg.

Sahu, S. S., Panda, G., Nanda, S. J., & Mishra, S. K. (2009). Protein Structural Class Prediction Using Differential Evolution.

Sułkowska, J. I., Morcos, F., Weigt, M., Hwa, T., & Onuchic, J. N. (2012). Genomics-aided structure prediction. *Proceedings of the National Academy of Sciences*, 109(26), 10340-10345.