

A Survey of Protein Classification Techniques using Machine Learning

Babasaheb .S. Satpute¹, Dr. Raghav Yadav²

^{1,2}Dept of Computer Science & I.T

^{1,2}SIET, SHUATS, Allahabad, India

Abstract- Protein classification or identification of family of protein is one of the demanding problems in bioinformatics and computational biology. Protein has immense importance in our life hence it becomes important to do impeccable classification of proteins and to also identify the family of the unknown protein. In this article we present a study report of the various classification techniques used in past to classify proteins. We also present a study of various features used for classification purpose.

Keywords- Bioinformatics, Computational biology, Protein classification, Machine Learning

I. INTRODUCTION

Bioinformatics [1,3], a contemporary science, is basically conceptualizing science as far as macromolecules and applying data innovation strategies to fathom and set up together the data related with these atoms. It primarily manages the utilization of PC and measurable strategies to the association of natural data. Bioinformatics has turned out as a one of the imperative research region in the ebb and flow time given that organic information is getting amassed at an exponential rate. This is as an outcome of the Human Genome Undertaking and other alike endeavors, alongside wonderful advancement of innovation for data stockpiling and access. In genome ventures, bioinformatics contain the improvement of strategies to seek information bases rapidly, to examine DNA succession data, and to foresee protein arrangement and structure from DNA grouping information. These have requested the advances of calculations which can mine helpful data from these information bases. As an answer for this predicament, various specialists have thought of strategies to break down and surmise the information, and determine thoughts in the DNA, RNA and protein information bases. One of the essential issues in this view is Ordering protein successions into superfamilies.

Proteins are fundamentally polymers formed from 20 dissimilar amino acids which can come about in random order. The quandary of protein superfamily classification can be properly defined as below [4, 5]. Given an unlabeled protein sequence P(formed from amino acids) and a set of known f

superfamilies $SF = SF1, SF2, \dots, SFf$, we should be able to decide with definite accuracy whether the protein sequence P belongs to one of the superfamilies SFi , $i = 1, \dots, f$ or not. At this point, a superfamily means a group of proteins which have certain similarity in structure and function. Alike protein sequences will nearly have similar biochemical functions. In this way, given an obscure protein, the primary errand is to group it into one of the known superfamilies. This will help in anticipating the protein work as well as structure of the obscure grouping; hence sparing, to an expansive degree, the costs brought about on costly organic (wet) explores in the research facility. Maybe, the most essential reasonable utilization of such learning is in tranquilize revelation. Assume we have acquired succession S from some sickness D and by our order technique we induce that S has a place with Fi . Keeping in mind the end goal to outline a medication for the infection D we may attempt a blend of existing medications for Fi .

II. RELATED WORK

Bandyopadhyay et.al proposed an productive method for grouping amino corrosive successions into various superfamilies. The proposed technique first concentrates 20 highlights from an arrangement of preparing successions. The separated highlights are with the end goal that they think about the probabilities of events of the amino acids in the distinctive places of the groupings. From that point, a hereditary fluffy grouping approach is utilized to naturally advance an arrangement of models speaking to each class. The normal for this bunching strategy is that it doesn't require from the earlier data about the quantity of groups, and is likewise ready to leave locally ideal setups. At last, the closest neighbor manage is utilized to group an obscure succession into a specific superfamily class, in light of its nearness to the models developed utilizing the hereditary fluffy bunching strategy. This outcomes in a huge change in the time required for ordering obscure arrangements. Results for three superfamilies, to be specific globin, trypsin and ras, show the viability of the proposed strategy concerning the situation where all the preparation successions are considered for arrangement utilizing a similar arrangement of highlights. Correlation with the notable procedure Impact likewise

demonstrates that the proposed strategy gives a huge change as far as the time required for grouping while at the same time giving equivalent characterization execution.

Ulavappa B. Angadi and M. Venkatesulu proposed a strategy where they have built a database and comparability lattice utilizing P-values got from an all-against-all Impact [8] run and prepared the system with the ART2 unsupervised learning calculation utilizing the columns of the similitude grid as information vectors, empowering the prepared system to order the proteins from 0.82 to 0.97 f-measure exactness. The execution of ART2 [7] has been contrasted and that of otherworldly grouping, Arbitrary backwoods, SVM, and HHpred. The performance of ART2 is superior to the others aside from HHpred. The results and performance of HHpred are superior to anything ART2 and the aggregate of mistakes is littler than that of alternate techniques assessed.

The objective of the said system is to order a given set of groupings at the level of SCOP [6] level and relegate protein areas to the SCOP superfamily in view of Impact E-esteem. Impact E-esteem is outstanding and acknowledged by scientists.

An imperative concern to apply neural systems to protein grouping is the issue of protein encoding successions as information esteems to the neural systems. To be sure, a succession may not be the best portrayal as an info incentive to neural systems. Great info portrayal makes it less demanding for the neural system to perceive the fundamental regularities, and is urgent to the accomplishment of neural system learning. Here, they have proposed Impact P-esteem as an information esteem, got from E-esteem by running Impact all-versus-all with the sigmoid capacity. They prepared the system with the ART2 unsupervised learning calculation. Fig. 1 represents the general engineering of their strategy and its application. Spotted line rectangle in Fig. 1 demonstrates ART2 preparing and testing process. Run of the mill engineering of ART2 is outfitted in Fig. 2. It includes standardization, commotion concealment utilizing clamor concealment parameter θ , reset system on watchfulness parameter ρ , refreshing of weights in the learning procedure utilizing learning rate α and input examples, and bunching.

There are three subsystems in the proposed method. To begin with subsystem is the development of a database and similitude network of P-values got from Impact all versus-all run. Second subsystem is the preparation of ART2 unsupervised neural system with the columns of similitude grid as info vectors. Third subsystem includes figuring Impact P-esteem vector for obscure question sequence(s) with Impact inquiry sequence(s)- versus-database run, and ordering the

successions into SCOP superfamilies or appointing the grouping to a suitable SCOP superfamily utilizing the trained network.

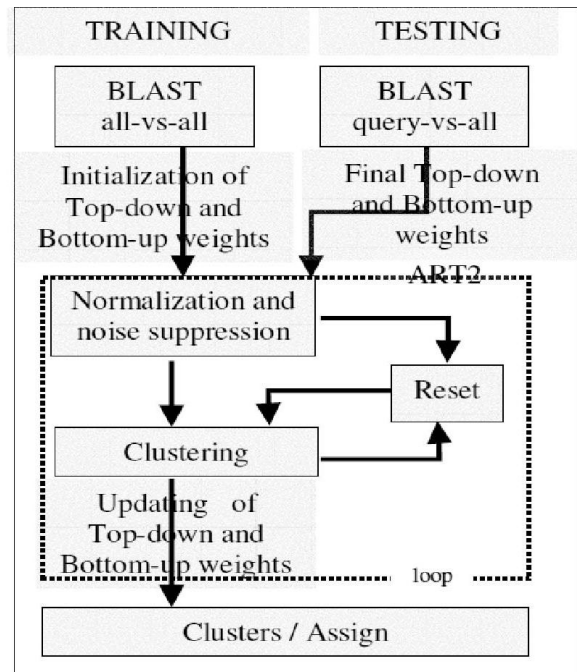


Figure1. Architecture of ART 2 based system for protein classification

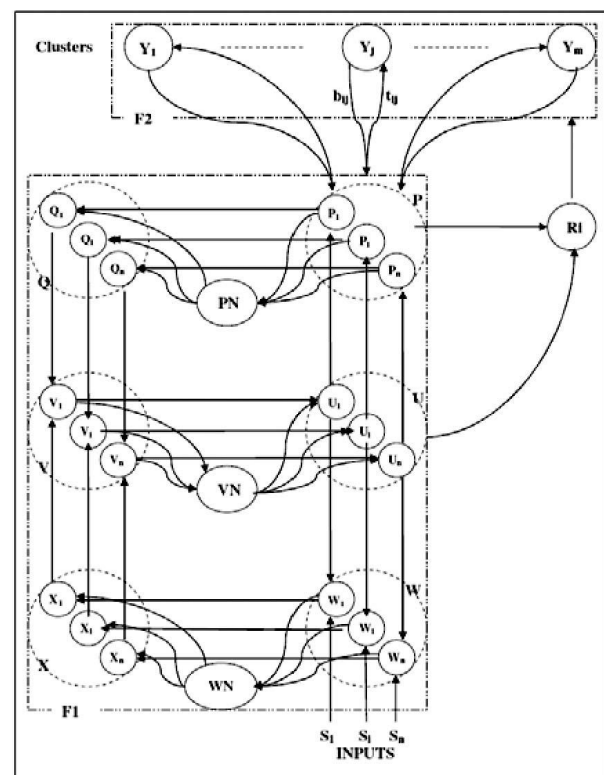


Figure2. Classic architecture of ART2 neural network

Swati Vipsita executed Probabilistic Neural Network (PNN) for protein superfamily characterization issue. The arrangement assignment sorts out proteins into their

superfamilies and aides in rectify forecast of structure and capacity of newfound proteins. The two principle ventures for any example order issue are highlight choice and highlight extraction. The bi-gram hashing capacity is utilized which concentrates and tallies the events of bi-gram designs from long strings of amino corrosive successions. The bi-gram technique maps successions of various length into input vectors of same length, yet the significant disadvantage of this strategy is that, the span of the information include vector has a tendency to be substantial. Choice of ideal number of highlights remains a basic issue for any example order issue. Foremost Component Analysis(PCA), a capable factual procedure, is utilized to diminish the measurement of the expansive information vector without much loss of data and consequently distinguishing design in information of high measurement. Customary PCA makes a straight change where as Kernel PCA (KPCA) is utilized when information are appropriated nonlinearly. Numerical recreations have demonstrated that for protein information conveyed non-straightly, KPCA outflanks PCA regarding precision, affectability and specificity.

Dianhui Wang et.al actualized Extreme Learning Machine (ELM) and used to group protein arrangements with ten classes of super-families downloaded from an open area database. A similar report on framework execution is directed amongst ELM and the primary ordinary neural system classifier - Backpropagation Neural Networks. Results demonstrate that ELM needs up to four requests of greatness less preparing time contrasted with BP Network. The arrangement precision of ELM is likewise higher than that of BP organize. There are no control parameters like ceasing criteria, learning rate, learning epochs etc. available in ELM for given system design to be physically tuned and could be actualized effectively.

BRIEF OF THE EXTREME LEARNING MACHINE

It is realized that generally every one of the feedforward systems parameters should be tuned and along these lines there be present the reliance between various layers of parameters (weights and predispositions). For past decades slope plunge based strategies have for the most part been utilized as a part of different learning calculations of feedforward neural systems. Be that as it may, plainly inclination plunge based learning strategies are for the most part ease back because of uncalled for learning steps or may effectively unite to nearby minima. Also, such learning calculations require numerous iterative learning steps keeping in mind the end goal to acquire better learning execution.

Dissimilar to well known usage, for example, Back-Propagation (BP), Huang, et al[9], [10], [11] have as of late developed another learning calculation which was named as Extreme Learning Machine (ELM) for Single-shrouded Layer Feedforward neural Networks (SLFNs) that are either added substance neurons or bit based plans. For added substance neurons based SLFNs someone may haphazardly pick and after that fix the information weights (connecting the information layer to the shrouded layer) and the concealed neurons' inclinations and scientifically decide the yield weights (connecting the concealed layer to the yield layer) of SLFNs[9]. Information weights are the weights of the associations between input neurons and shrouded neurons and yield weights are the weights of the associations between concealed neurons and yield neurons. After the information weights and the shrouded layer predispositions are picked self-assertively, SLFNs can be essentially considered as a direct framework and the yield weights (connecting the concealed layer to the yield layer) of SLFNs can be logically decided through straightforward summed up backwards operation of the concealed layer yield networks. Rather than tuning the focuses and effect widths of RBF bits we can just haphazardly pick and fix these portion parameters and scientifically compute the yield weights of RBF networks[10], [11], for Radial Basis Function (RBF) part based SLFNs. After the info weights and the concealed layer inclinations are picked discretionarily, SLFNs can be basically considered as a direct framework and the yield weights (connecting the shrouded layer to the yield layer) of SLFNs can be diagnostically decided through straightforward summed up reverse operation of the concealed layer yield networks. As examined by Huang, et al[9], [10], [11], ELM has a tendency to have great speculation execution and can be actualized effectively. Not at all like other tuning/change techniques which may not be reasonable for non-differential initiation capacities nor keep the upsetting issues, for example, halting criteria, learning rate, learning epoches, and nearby minima, the ELM calculation can stay away from these challenges extremely well.

Jong Cheol Jeong et.al created Position-Specific Scoring Matrix for the prediction of function of protein. They propose new highlights separated from protein arrangement and which are based on machine learning strategies for computational capacity expectation. These highlights are gotten from a position-particular scoring network, which has demonstrated incredible potential in different bininformatics issues. They assess those highlights utilizing four distinct classifiers and information related to yeast protein. Their trial comes about demonstrate that highlights got from the position-particular scoring framework are suitable for programmed work comment.

METHOD DESCRIPTION

Position-Specific Scoring Matrix (PSSM)

Scientists first introduced PSSM for identifying remotely correlated proteins. It was produced from a gathering of groupings beforehand adjusted by auxiliary or succession closeness [12].

Place-particular repeated BLAST (PSI-BLAST) [1], [13] is the most usually utilized program, that analyzes PSSM profiles for distinguishing distantly connected similar proteins or DNA. The first PSSM presented by Gribskov et al. [12] comprises of the accompanying segments:

1. Site, which shows the successively expanded file of every amino corrosive buildups in an arrangement after different grouping arrangement;
2. Test, that is a gathering of run of the mill successions of practically allied proteins those have been adjusted by arrangement or auxiliary comparability;
3. Profile, that is a network comprising of 20 segments relating to twenty amino acids; and
4. Agreement, that is a succession of amino corrosive deposits which are highly like all the arrangement buildups of tests at every specific place. It is produced by choosing the most noteworthy score in profile at every specific place.

PSSM was first presented for distinguishing remotely related proteins. It was created from a gathering of successions beforehand adjusted by auxiliary or arrangement comparability [12].

Place-particular repeated BLAST (PSI-BLAST), [13] is the widely utilized computer program, that thinks about PSSM profiles for recognizing distantly related alike proteins or DNA. The first PSSM presented by Gribskov et al. [12] comprises of the accompanying segments:

1. Place that demonstrates the successively expanded record of every amino corrosive buildups in an arrangement after various grouping arrangement;
2. Test, which is a gathering of run of the mill groupings of practically related proteins that have been adjusted by succession or basic similitude;

3. Profile, that is a network comprising of 20 sections relating to 20 different amino acids; and
4. Agreement, that is a grouping of amino corrosive deposits which are the majorly like all the arrangement buildups of tests at every place. It is produced by choosing the most noteworthy score in the profile at every place.

A PSSM for a question protein is $N \times 20$ lattice $P = \{P_{ij} : i = 1 \text{ to } N \text{ and } j = 1 \text{ to } 20\}$, here N is the length of the protein grouping. It relegates a score P_{ij} for the j th amino corrosive in the i th position of the question succession with a vast esteem showing a very moderated place and a little esteem demonstrating a feebly saved place [12]. While the developments of PSSM are marginally unique in relation to one application to another, the standards are particularly the same.

PSSM-Derived Features

Out of few difficulties in methodologies based on machine learning for protein work comment is to make enlightening highlights. In here, they can't change over the PSSMs straightforwardly to highlight vectors as proteins have diverse quantities of amino acids, that will prompt distinctive sizes of highlight vectors. They deal with this issue by creating highlights which are arrived at the midpoint of over a nearby area, as they examine next. Point to note is that every one of the highlights are gotten from PSSMs which are removed from protein groupings.

They remove 3 distinctive capabilities present in PSSMs. The main set (include set number 1) depends on the found the middle value of PSSM profiles over hinders, every valu with 5 percent of a succession. In this way, a protein arrangement, paying little heed to its span, is separated into 20 pieces and each square comprises of twenty highlights (got from the twenty segments in PSSMs). The method of reasoning at the back this method is the buildup preservation propensities in a similar area family are comparative, and the areas of spaces in a similar family are firmly identified with the span of the arrangement. Numerically, for the j th obstruct in include set #1, the element F_j is a 1×20 highlight vector, which is produced by utilizing the accompanying conditions:

$$F_j = \frac{1}{B_j} \sum_{i=1}^{B_j} P_i^{(j)}, \quad (1)$$

where B_j is the extent of the j th square, that is five percent of the span of an arrangement and $P_i^{(j)}$ is a 1×20 vector removed from the profile of the PSSM at the i th place

in the j th piece. For every arrangement, there are a sum of 20 squares; in this manner, the last element is a 1×400 vector.

The second set (include set number 2) is persuaded by a perception which is not quite the same as the one utilized for highlight set number 1. Rather than thinking about the areas of spaces in a succession, they concentrate on the areas with comparable protection rates, as spaces of a similar family can have few other comparable preservation inclinations than areas of various families. Illuminated by this perception, they amass area families in view of their preservation scores in the profiles of PSSM. Point to note here is, this is managed without the knowledge that the particular area districts in a protein. The thought is like the test idea utilized as a part of microarray innovations, where tests are utilized to recognize qualities. For the comfort, they call it buildup testing strategy. In this program, every test is an amino corrosive, which compares to a specific segment in the profiles of the PSSM. For every test, they normal the PSSM scores of each and every amino acids in the related segment with a PSSM esteem more noteworthy than zero in the arrangement, which prompts a 1×20 highlight vector. By and by, for the 20 tests, the last element for every protein succession is a 1×400 vector.

For the third set (include set number 3), they mull over the physicochemical characteristics of examined deposits acquired from highlight set number two. Nonetheless, rather than utilizing unique protein arrangements, they utilize agreement successions acquired from the PSSM profiles. The utilization of agreement successions depends on the understanding that every accord grouping gives less one-sided data in the midst of homologous proteins. For every accord succession, they utilize 9 physicochemical characteristics, those can be sorted into two gatherings, for example, normal and thickness gatherings: isoelectric point, hydrophobicity, and wreckage scale are found the middle value of, and nonpolar, polar, hydrophilic, hydrophobic, positive, and negative charge deposits are utilized for computing densities. The extent of definite element vectors in this gathering is 1×180 .

PATRICK et. al proposed UPSEC: An Algorithm for Classifying Unaligned Protein Sequences into Functional Families.

To group proteins into practical families in view of their essential groupings, prevalent calculations, for example, the k -NN-, HMM-, and SVM-based calculations are regularly utilized. For a significant number of these calculations to play out their assignments, protein successions should be legitimately adjusted first. Since the arrangement procedure can be mistake inclined, protein grouping may not be

performed precisely. To enhance characterization precision, they propose a calculation, called the Unaligned Protein SEquence Classifier (UPSEC), which can play out its undertakings without arrangement. UPSEC makes utilization of a probabilistic measure to recognize buildups that are valuable for arrangement in both positive and negative preparing tests, and can deal with multi-class order with a solitary classifier and a solitary go through the preparation information. UPSEC has been tried with genuine protein informational collections. Test comes about demonstrate that UPSEC can adequately arrange unaligned protein groupings into their relating practical families, and the examples it finds amid the preparation procedure can be organically important.

PROPOSED ALGORITHM

Given an arrangement of N unaligned protein groupings, $S_1 \dots \dots S_i \dots \dots S_N$, of differing length, $L_1 \dots \dots L_i \dots \dots L_N$,

Separately. After the transformation of each arrangement into various subsequences, a three-advance information mining strategy is performed. In Step 1, UPSEC decides whether deposits are helpful for order by deciding whether the affiliation designs amongst buildups and a specific class mark are intriguing (in preparing set). Provided that this is true, in Step 2, the heaviness of each found example will be resolved. At last, in Step 3, these examples will be utilized for the order of successions not initially in the database (in testing set).

III. CONCLUSION

Building a proficient insightful protein classification framework for successfully looking, an extensive natural database is thought to be a noteworthy accomplishment in bioinformatics. In this paper they reviewed various protein classification systems. The classifiers are implemented by using Artificial neural network, machine learning methods like support vector machines, decision tree classifiers, Naive Bayes classifiers etc. Unsupervised classifiers are also successfully used for classification purpose. Various features of proteins based on the amino acid residues and structure have been used successfully. The major challenge lies in getting more accurate classification.

REFERENCES

- [1] Baxevanis, F.B.F. Ouellette (Eds.), Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins, Wiley, New York, 1998

- [2] P. Jain, J.M. Garibaldi, and J.D. Hirst, "Supervised Machine Learning Algorithms for Protein Structure Classification," *Computational Biology and Chemistry*, vol. 33, no. 3, pp. 216-223, 2009.
- [3] N.M. Luscombe, D. Greenbaum, M. Gerstein, What is bioinformatics: a proposed definition and overview of the field, *Methods Inform. Med.* 40 (2001) 346–358
- [4] J.T.L.Wang, Q.C. Ma, D. Shasha, C.H.Wu, Application of neural networks to biological data mining: a case study in protein sequence classification, in: *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2000)*, Boston, MA, 2000, pp. 305–309.
- [5] J.T.L.Wang, Q.C. Ma, D. Shasha, C.H.Wu, New techniques for extracting features from protein sequences, *IBM Systems J.* 40 (2) (2001) 426–441 (Special Issue on Deep Comput. Life Sci.).
- [6] A.G. Murzin, S.E. Brenner, T. Hubbard, and C. Chothia, "SCOP: A Structural Classification of Proteins Database for the Investigation of Sequences and Structures," *J. Molecular Biology*, vol. 247, no. 4, pp. 536-540, 1995.
- [7] G.A. Carpenter and S. Grossberg, "ART 2: Self-Organization of Stable Category Recognition Codes for Analog Input Patterns," *Applied Optics*, vol. 26, no. 23, pp. 4919-4930, 1987.
- [8] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, and D.J. Lipman, "Basic Local Alignment Search Tool," *J. Molecular Biology*, vol. 215, no. 3, pp. 403-410, 1990.
- [9] G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, "Extreme learning machine: A new learning scheme of feedforward neural networks," in *Proceedings of International Joint Conference on Neural Networks (IJCNN2004)*, (Budapest, Hungary), 25-29 July, 2004.
- [10] G.-B. Huang and C.-K. Siew, "Extreme learning machine: RBF network case," in *Proceedings of the Eighth International Conference on Control, Automation, Robotics and Vision (ICARCV 2004)*, (Kunming, China), 6-9 Dec, 2004.
- [11] G.-B. Huang and C.-K. Siew, "Extreme learning machine with randomly assigned RBF kernels," *International Journal of Information Technology*, vol. 11, no. 1, 2005.
- [12] M. Gribskov et al., "Profile Analysis: Detection of Distantly Related Proteins," *Proc. Nat'l Academy of Sciences USA*, vol. 84, pp. 4355-4358, July 1987.
- [13] S.F. Altschul and E.V. Koonin, "Iterated Profile Searches with PSIBLAST— A Tool for Discovery in Protein Databases," *Trends Biochemical Sciences*, vol. 23, pp. 444-447, Nov. 1998.
- [14] P. Domingos and M. Pazzani, "ON the Optimality of the Simple Bayesian Classifier under Zero-One Loss," *Machine Learning*, vol. 29, pp. 103-130, 1997.
- [15] V.N. Vapnik, "An Overview of Statistical Learning Theory," *IEEE Trans Neural Networks*, vol. 10, pp. 988-999, Sept. 1999.
- [16] L. Breiman et al., *Classification and Regression Trees*. Chapman and Hall/CRC 1984.
- [17] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, pp. 5-32, 2001.