

APPROACHES TO PREDICTION OF PROTEIN STRUCTURE: A REVIEW

Er. Amanpreet Kaur¹, Dr. Baljit Singh Khehra²

¹Er. Amanpreet Kaur, Student, Dept. of Computer Science & Engineering, Baba Banda Singh Bahadur Engineering College, Fatehgarh Sahib, Punjab, INDIA

²Dr. Baljit Singh Khehra, Professor and Head, Dept. of Computer Science & Engineering, Baba Banda Singh Bahadur Engineering College, Fatehgarh Sahib, Punjab, INDIA

Abstract - Data mining software is one of the most important analytical tool for analyzing data. It allows users to analyze data from many different dimensions, categorize it, and summarize the relationships identified. Protein structure prediction (PSP) is the most important and challenging problem in bioinformatics today. This is due to the fact that the biological function of the protein is determined by its structure. While there is a gap between the number of known protein structures and the number of known protein sequences, protein structure prediction aims at reducing this structure –sequence gap. Protein structure can be experimentally determined using either X-ray crystallography or Nuclear Magnetic Resonance (NMR). However, these empirical techniques are very time consuming, so various machine learning approaches have been developed for protein structure prediction like HMM, SVM and NN. In this paper we give a general introductory background to the area and a literature survey about the machine learning approaches. These approaches depends on the chemical and physical properties of the constituent amino acids. Not all machine learning algorithms have the same performance, so we represent the general success keys for any such algorithm.

Key Words: Data mining, Protein Structure Prediction (PSP), Hidden Markov Model, Support Vector Machines, Neural Networks.

1. INTRODUCTION

Data mining involves the use of practical data analysis tools to design previously exotic, valid patterns and relationships in huge data set. There are several applications for Machine Learning (ML), the most powerful of which is data mining. People are often prone to making mistakes during analysis or, possibly, when trying to establish relationships between multiple features. This makes it difficult for them to find solutions to specific problems. Machine learning can often be strongly applied to these problems, improving the designs of machines and the efficiency of the systems.

Classification is the most important technique to identify a specific character or group of them.

Different classification algorithms have been proposed by various researchers for classification of protein sequences. The Protein sequence consists of twenty different amino acids which are aligned in some specific sequences. Popular protein sequence classification techniques involve extraction of particular features from the sequences. These features depend on the structural and functional properties of amino acids. These features can be compared with their predefined values.

1.1 Proteins

Proteins are main building blocks of our life. Proteins form the basis of structures such as skin, hair, and tendon and they are responsible for catalyzing and regulating biochemical reactions, transporting molecules. The shape of protein is given by its amino acid sequence. There are 20 distinct types of amino acid and each amino acid is known by its side chain that determines the properties of amino acid [16].

Currently, protein plays a vital role in the research human body. The structural class, which is one of the important attribute of a protein plays an important role in both theoretical and experimental studies in protein science. In generally speaking, protein is the chief executor of important movement. On the one hand, the data of Protein sequence database has been growing very fast. On the other hand, the structure of the protein is comparably less identified. Protein tertiary structures have vital influence on the behaviour of the protein from long-term study. According to the definition by Levitt and Chothia, proteins are classified into the following four structural classes: (1) all-a class, which are essentially formed by helices and only includes small amount of strands, (2) all-b class, which are essentially formed by strands and only includes small amount of helices, (3) a/ b class, which includes both helices and mostly parallel strands, and (4) a+ b class, which includes both helices and mostly antiparallel strands. Prediction of tertiary structure,

however, still remains as an unsolved problem and various solution methods are urgently needed.

In its native environment, the chain of amino acids (or residues) of a protein folds into local secondary structures including alpha helices, beta strands, and non regular coils. The secondary structure is specified by a sequence classifying each amino acid into the corresponding secondary structure element (e.g., alpha, beta, or gamma). The secondary structure elements are further packed to form a tertiary structure depending on hydrophobic forces and side chain interactions, such as hydrogen bonding, between amino acids. The tertiary structure is described and coordinates of all the atoms of a protein or, in a more coarse description, by the coordinates of the backbone atoms. Finally, several related protein chains can interact or assemble together to form protein complexes. These protein complexes said to be the protein quaternary structure. The quaternary structure is described by the coordinates of all the atoms, or all the backbone atoms in a coarse version, associated with all the chains participating in the quaternary organization, given in the same frame of reference[4].

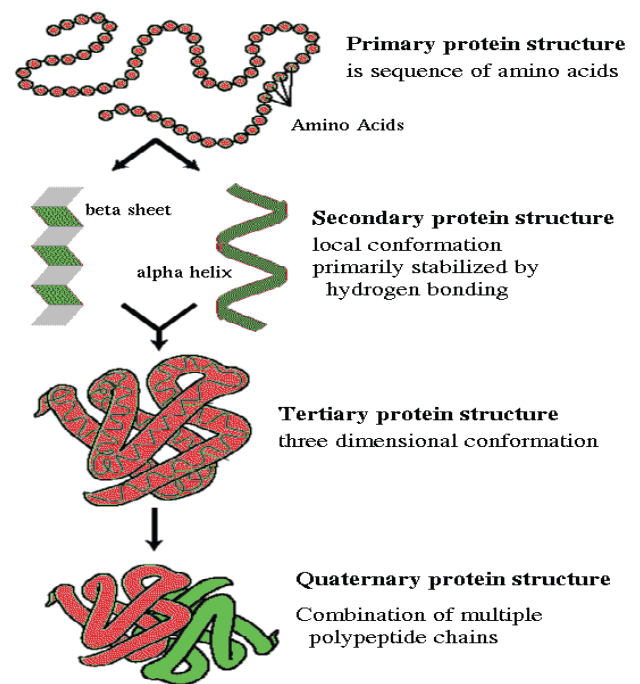


Fig-2: Levels of protein structure [13]

Primary Structure:

Chain of amino acid sequence is referring as primary structure. Every α -amino acid has of a backbone part that is present altogether in amino acid varieties, and a side chain that is distinctive to every variety of residue. Proline is an exception from this rule. The primary structure is held together by peptide bonds, that are created during the process of protein biosynthesis or translation. The primary structure of a protein is determined by the gene equivalent to the protein. A specific sequence of nucleotides in DNA is transcribed into mRNA, that is read by the ribosome in a process called translation. The sequence of a protein is exclusive to that protein, and defines the structure and function of the protein. The sequence of a protein may be determined by strategies such as Edman degradation or tandem mass spectrometry [13].

Secondary Structure:

The secondary structure consists of native folding regularities maintained by hydrogen bonds and is historically subdivided into 3 classes: alpha-helices (H), beta-sheets (E), and coil(C).Secondary structure contained localized and recurring fold of peptide chain, wherever 2 main regular structures are the α -helix and β -sheet. Hydrogen bond is answerable for secondary structure-helix may be considered the default state for secondary

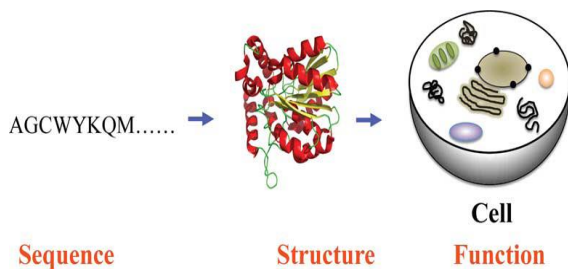


Fig-1: Protein sequence-structure-function relationship. A protein is a linear polypeptide chain composed of 20 different kinds of amino acids represented by a sequence of letters (left). It folds into a tertiary (3-D) structure (middle) composed of three kinds of local secondary structure elements (helix – red; betastrand– yellow; loop – green). The protein with its native 3-D structure can carry out several biological functions in the cell (right). [4]

1.2 Protein Structure

Proteins are large, organic molecules and are among the most vital components in the cells of living organisms. They are more diverse in structure and function than any other kind of molecule. It can act as Enzymes, antibodies, hormones, transport molecules, hair, skin, muscle, tendons, cartilage, claws, nails, horns, hooves, and feathers are all made of proteins. Protein structure has a basically four levels of category: Primary Structure, Secondary structure, Tertiary structure and Quaternary structure. Fig. 2 shows different levels of protein structure [16].

structure. It is most significant for higher understanding tertiary structure. It is extremely necessary because knowledge of secondary structure helps in the prediction of tertiary structure when structure discovery without sequence similarity within the datasets [13].

Tertiary Structure:

Tertiary structure refers to 3-dimensional structure of a one protein molecule. It involves localized spatial interaction among primary structure parts, i.e. the amino acids. The alpha-helices and beta-sheets are folded into a compact ball. The folding is driven by the non-specific hydrophobic interactions, however the structure is stable only if when the elements of a protein domain are locked into place by specific tertiary interactions, like salt bridges, hydrogen bonds, and the tight packing of side chains and disulfide bonds. The disulfide bonds are extremely rare in cytoplasm proteins, since the cytoplasm is generally a reducing environment [13].

Quaternary Structure:

Quaternary structure is that the arrangement of multiple folded protein or coiling protein molecules in a multi-subunit complex. Several proteins are literally assemblies of more than one polypeptide chain, that within the context of the larger assemblage are called as protein subunits. Additionally to the tertiary structure of the subunits, multiple-subunit proteins possess a quaternary structure, that is the arrangement into which the subunits assemble. Enzymes composed of subunits with various functions are typically known as holoenzymes, in which some elements could also be called regulatory subunits and the functional core is called the catalytic subunit. Examples of proteins with quaternary structure include hemoglobin, DNA polymerase, and ion channels. Different assemblies referred to instead as multi-protein complexes also possess quaternary structure[13].

1.3 Protein Structure Prediction

Protein structure prediction is that the prediction of the 3-dimensional structure of a protein from its amino acid sequence therefore all activities of protein area unit depends upon its three dimensional structure. Structure prediction is essentially different from the inverse drawback of protein design. The 3-dimensional structure of a protein is decided by the network of covalent and non-covalent interactions. Though protein is built by the chemical process of 20 different amino acids into linear

chains, proteins form an incredible array of various tasks. A protein chain folds into a novel shape that is stabilized by non covalent interactions between regions within the linear sequence of amino acids. This spatial organization of a protein its shape in three dimensions could be a key to understanding its function. Only when a protein is in its correct three-dimensional structure, or conformation, is it ready to perform efficiently. A key idea in understanding how proteins work is that function is derived from 3-dimensional structure, and 3-dimensional structure is specified by amino acid sequences[16].

Protein structure prediction is that the prediction of the 3- dimensional structure of a protein from its amino acid sequence, i.e., the prediction of its secondary, tertiary, and quaternary structure from its primary structure. Structure prediction is basically different from the inverse problem

of protein design. In bioinformatics several prediction methods are available such as:

- Ab-initio, theoretical modelling, and conformation space search
- Homology modelling and threading

Primary and Secondary structure prediction:

Primary structure could be chain of 20 amino acid sequence, which is described as:

Protein Sequence: Input 1D
GRPRAINKHEQEQISRLLLEKGGHPRQQLAIF
↓
HCCCCCCEHECECCCCECHHCCCCCCCC

Protein Structure: Output 1D

Protein output 1D structure is getting using dictionary of secondary structure prediction (DSSP) methodology. Secondary structure prediction is that the classification of primary 1D structure in to three classes: Helix (H), Strand (E) and Coil(C).

Techniques used for Protein structure prediction are:

- Soft Computing Techniques like Artificial Neural networks.
- Probabilistic techniques like Hidden Markov Model.
- Evolutionary Computation like Genetic Algorithm.
- Statical techniques like SVM.
- Clustering algorithms etc. [14]

Bioinformatics techniques to protein secondary structure prediction largely depend upon the information out there in amino acid sequence. Evolutionary algorithms are like simple genetic algorithms (GA), messy GA, fast messy GA have addressed this problem. Support Vector Machine (SVM) represents a replacement approach to supervised pattern classification that has been with success applied to a large range of pattern recognition issues, as well as object recognition, speaker identification, gene function prediction with micro array expression profile, etc. In these cases, the performance of SVM either matches or is considerably higher than that of ancient machine learning approaches, as well as neural networks. However still SVMs are blackbox models. ANN is a good technique of protein structure prediction that relies on the sound theory of Back Propagation Algorithm. Protein secondary structure prediction has been satisfactorily performed by machine learning techniques like Artificial Neural Network and Support vector machines. Most secondary structure prediction programs target alpha helix and beta sheet structures and summarize all different structures within the random coil pseudo category. For the classification, ANN is employed as a binary classifier.

1.4 Necessities of PSSP(Protein Secondary Structure Prediction)

PSSP is receiving significance in the recent area of research due to the following:

- Being the difficulty of structural bioinformatics, protein secondary structure prediction can give prediction and analysis of macromolecules which are the basis of an organism.
- Protein secondary structure prediction(PSSP) give structure function relationship. That is which particular protein structure is responsible for which particular function would be known by PSSP. So by changing the protein's structure or by synthesizing new proteins, functions could be added or removed or required functions could be attained.
- Structure of the viral proteins can be examined by PSSP and this examination of the structures of the viral proteins provides the way to design drugs for specific viruses.
- PSSP lowers the sequence structure gap. The sequence structure gap can be best defined by giving the example of large scale sequencing projects like Human Genome Project. In these

type of projects, protein sequences are produced at a very rapid speed which results in a huge gap between the number of known protein structures (>150,000) and the no. of known protein sequences (>4,000). This gap is called sequence structure gap and PSSP can successfully minimize this gap. Experimental approaches are not capable of structure determination of few proteins like membrane proteins. So the prediction of protein structure using computational tool is of good interest.[2]

Sequence-Structure Gap and the Need for Structure Prediction

With the advent of recombinant DNA technology it has become possible to signify the amino acid sequences of proteins quite fast. However, signifying the 3- dimensional structure of proteins is a time taking task and hence there exists a large gap between the number of proteins of known amino acid sequence and that of known structures. This is known as the sequence-structure gap. As the knowledge of the 3-D structure of a protein is very important to understand its function, it is essential to develop techniques to predict the structure of a protein from its amino acid sequence.

1.5 Fold Recognition

Proteins fold due to hydrophobic effect, electrostatic forces, Vander Waals interaction and Hydrogen bonding. Protein threading, also called fold recognition, is a methodology of protein modelling (i.e. computational protein structure prediction) which is used to model those proteins those have the same fold as proteins of known structures, but do not have similar proteins with known structure. Protein folding is the method by which a protein assumes its 3D structure. All protein molecules are endowed with a primary structure having the polypeptide chain. Fold recognition need a criterion to identify the best template for single target sequence. The protein fold-recognition method to structure prediction aims to identify the known structural framework that accommodates the target protein sequence in the best way. Typically, a fold-recognition program comprises four components:

- (1) The representation of the template structures (usually corresponding to proteins from the Protein Data Bank database),

(2) The evaluation of the compatibility between the target sequence and a template fold,

(3) The algorithm to compute the optimal alignment between the target sequence and the template structure, and the method the ranking is computed and the statistical significance is estimated [16]

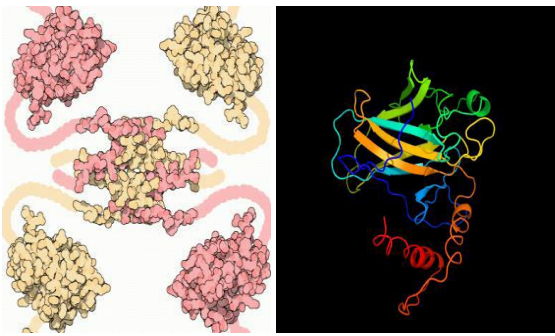


Fig. 3.1- Structure of protein P53 **Fig. 3.2-** Folded structure of protein P53 [16]

2. LITERATURE SURVEY

W. Dianhui, L. K. Wung *et al.*[17] presents a modular neural classifier for protein sequences with improved classification criteria. The intelligent classification techniques described in this paper aims to enhance the performance of single neural classifiers based on a centralized information structure in terms of recognition rate, generalization and reliability. F. Mhamdi *et al.* [5] presents the classification of proteins by basing on its primary structures. The sequence of proteins can be collected in a file. The application of text mining technique is proposed for extracting the features. An algorithm is also developed which extracts all the n-grams existing in the file of data and produced a learning file. D. H. Shing, Y. Y. Chi [4] implement a genetic algorithm to cluster the training set before a prediction model is built. Using position specific scoring matrix (PSSM) as part of the input, the hybrid method achieves good performances on sets of 513 non redundant protein sequences and 294 partially redundant sequences. The results also show that clustering achieves the goal of data preprocessing differently on redundant and non-redundant sets, and it seems almost preferable to cluster the data before prediction is preformed. C. Jianlin, N. T. Allison *et al.* [2] reviews the development and application of hidden Markov models, neural networks, support vector machines, Bayesian methods, and clustering methods in 1-D, 2-D, 3-D, and 4-D protein structure predictions. P. Sun and J. Zhang [12] presents a prediction method of protein

contact on the basis of information granules and RBF neural network have been brought forward. This method improved the encoding approach of protein structure data and classifier performance to enhance the predicting accuracy of protein contact. V. Swati, S. K. Bithin *et al.* [16] discusses three types of neural networks such as feed forward neural network, probabilistic neural network and radial basis function neural network. The main objective of the paper is to build up an efficient classifier using neural networks. The measures used to estimate the performance of the classifier are Precision, Sensitivity and Specificity. N. Mathuriya *et al.*[10] observed that the K-means clustering algorithm is not very much suitable for the problem and the back propagation neural network has the high performance. The artificial neural network (ANN) is the technique of data mining that is different from traditional techniques. It is the nonlinear auto-fit dynamic system made of various cells with simulating the construction of biology neural systems. S. Saha *et al.*[15] presents a review with three different classification models such as neural network model, fuzzy ARTMAP model and Rough set classifier model. This is followed by a new technique for classifying protein sequences. The proposed model is typically implemented with an own designed tool and tries to reduce the computational overheads encountered by earlier approaches and increase the accuracy of classification.

K. Yasaman, F. Mahmood *et al.*[7] discuss the most important algorithms like evolutionary algorithms, particle swarm or ant colony optimization and computational methods introduced in this field and their challenges. B. Wenzheng, C. Yiming *et al.*[1] purposed forward novel approach for predicting the tertiary structure of protein and construct an Error Correcting Output Codes(ECOC) classification model on the basis of Particle swarm optimization(PSO) and neural network(NN). Three feature extraction methods, which are Amino Acid Composition, Amino Acid Frequency and Hydrophobic Amino Acid Combination, respectively, are employed to extract the features of protein sequences. To evaluate the efficiency of the proposed method we choose a benchmark protein sequence dataset (640 dataset) as the test data set. The final results show that our method is efficient for protein structure prediction. R. N. Chandrayani, K. Manali [13] presents the results of protein p53. P. Mayuri, S. Hitesh [11] uses model based (i.e., supervised learning) approach for protein secondary structure prediction and our objective is to enhance the prediction of 2D protein structure problem using advance

machine learning techniques like, linear and non-linear support vector machine with different kernel functions. M. Vidyasagar,[9] review some of the many challenging problems in computational biology that are amenable to treatment using a systems approach. Specific problems discussed include string alignment and protein structure prediction. M. Sonal, P. Yadunath *et al.*[8] explored the machine learning classification models with six physical and chemical properties to classify the root mean square deviation (RMSD) of the protein structure in absence of its true native state and each protein structure lies between 0A° to 6A° RMSD space. Physical and chemical properties used in this paper are total surface area, Euclidean distance, total empirical energy, secondary structure penalty, residue length, and pair number. Artificial bee colony algorithm is used to determine the feature importance. To measure the robustness of the best classification model, K-fold cross validation is used. R. S. Prashant, S. Harish *et al.*[14] explore nine machine learning methods with six physicochemical properties to predict the RMSD (Root Mean Square Deviation), TM-score (Template Modelling) and GDT TS-score (Global Distance Test) of modelled protein structure in the absence of its true native state. Physicochemical properties namely total surface area, euclidean distance, total empirical energy, secondary structure penalty, sequence length and pair numbers are used. The K-fold cross validation is used to measure the robustness of the best predictive method.

C. Nandini *et al.*[3] explains several techniques used by different researches for the classification of proteins and also provides an overview of different protein sequence classification methods. From the vast data we have to derive the hidden knowledge so that it is used in wide range of areas to design drug, to identify diseases, and in classification of protein sequence etc. I. S. Mohammad, A. Hakimeh [6] proposed approach to reduce predicted RMSD Error than the actual amount for RMSD and calculate mean absolute error (MAE), through feed forward neural network, adaptive neuro fuzzy method. ANFIS is achieved better and more accurate results. W. Bo,L. Yongkui *et al.*[18] summarise some of the recent studies adopting this SVM learning machine for prediction structure prediction are the one which used frequent profiles with evolutionary information. W. Jian and L. Jian-Ping [24] discussed about neural network, an improvement scheme that iterative matrix replace secondary derivative has been developed by introduced quasi-Newton algorithm. Profile code based on probability has been used and comparison of window width and

learning training has been completed. The experiment results indicate that the prediction for secondary structures of protein obtain a very good effect based on neural network and quasi-Newton algorithm. W. L. George, P. Marius *et al.*[25] describe a large scale application of a back-propagation neural network to the analysis, classification and prediction of protein secondary and tertiary structure from the sequence information alone. W. J. Barry [23] presented the tutorial. this tutorial begins with a short history of neural network research, and a review of chemical applications. The bulk, however, is devoted to providing a clear and detailed introduction to the theory behind backpropagation neural networks, along with a discussion of practical issues facing developers.

Z. Zhen, J. Nan [27] proposed a new technique based on radial basis function neural networks for prediction of protein secondary structure. To make the technique comparable to other secondary structure prediction methods, they used the benchmark evaluation data set of 126 protein chains in this paper. They also analyzed how to use evolutionary information to increase the prediction accuracy. The paper discussed the influence of data selection and structure design on the performance of the networks. The results show that this method is feasible and effective. W. Leyi and Z. Quan [26] reviews some machine learning methods. They conduct a detail survey of recent computational methods, especially machine learning-based methods, for protein fold recognition. This review is anticipated to assist researchers in their pursuit to systematically understand the computational recognition of protein folds. P. Rojalina, D. Nilamadhab *et al.* [14] investigates the protein secondary structure prediction problem by ancient learning techniques such as Artificial Neural Network where Back propagation algorithm is used for learning. It measures the efficiency and accuracy of the machine learning methods through Mean Square Error. B. Hemashree, S.K. Kandarpa [2] discuss about the ANN approach for protein structure prediction. The Artificial Neural Network (ANN) technique for prediction of protein secondary structure is the most successful one among all the techniques used. In this method, ANNs are trained to make them capable of performing recognition of amino acid patterns in known secondary structure units and these patterns are used to differentiate between the different types of secondary structures. This work is related to the prediction of secondary structure of proteins employing artificial neural network though it is

restricted initially to three structures only. A. Shivani, B. Arushi et al.[1] discuss multilayer feed forward artificial neural network. A tool used for the secondary structure prediction of proteins from the amino acid sequence is multilayer feed forward artificial neural network with back propagation. This approach is a machine learning methodology in which the network is trained using the recognized data sets for which the widely used benchmark is Protein Data Bank (PDB), maintained by Research Collaboratory for Structural Bioinformatics (RCSB). The algorithm used for the classification is Define Secondary Structure Prediction (DSSP) that classifies the sequences in the 3-level subclasses: Helix (H), Sheet (E) and Coil (C). The objective is to get the maximum predictive accuracy with the minimalized error.

3. APPROACHES USED IN PSP

3.1 Markov Chains and Hidden Markov Model(HMM)

In this section we briefly review some standard material on Markov chains. Then the discussion is extended to so-called hidden Markov models (HMM's). Hidden Markov models (HMM's) are used to distinct the coding regions of a Prokaryote gene from the non-coding regions, and also to classify a protein into one of a small number of previously classified protein families.

Markov Chains

Suppose $X := \{s_1, \dots, s_n\}$ is a finite set. A stochastic process

$\{X_t\}_{t \geq 0}$ assuming values in X is known to be a Markov chain if

$$\Pr\{X_t | X_{t-1}, X_{t-2}, \dots\} = \Pr\{X_t | X_{t-1}\}. \quad (1)$$

The Markov chain is known as stationary if the above conditional probability is independent of t . The temporal evolution of a Markov chain is captured by an $n \times n$ matrix of transition probabilities, shown as follows:

$$a_{ij} = \Pr\{X_t = s_j | X_{t-1} = s_i\}, A = [a_{ij}]. \quad (2)$$

Thus a_{ij} is the probability that the Markov chain is in state s_j at the next time instant, given that it is in the state s_i at the current time instant. Obviously the matrix A is row-stochastic; that is,

$$a_{ij} \geq 0 \forall i, j, \text{ and } \sum_{j=1}^n a_{ij} = 1 \forall i \quad (3)$$

Hence the vector all one's is a column eigenvector of A . A standard result defines that each matrix also has at least one row eigenvector whose components are all nonnegative. Such a vector, i.e., a nonnegative vector π such that $\pi = A\pi$ and such that the components of π add up to 1, is called a stationary distribution. Under some additional conditions, like the irreducibility of the Markov chain, there is only single stationary distribution. Note that if the Markov chain is started off at time $t = 0$ with the initial state X_0 distributed according to π , then X_t is distributed according to π at all future times.

Hidden Markov Models

HMMs are among the foremost vital techniques for protein fold recognition. Earlier HMM approaches, like SAM and HMMer, designed an HMM for a query with its homologous sequences and so used this HMM to attain sequences with known structures in the PDB using the Viterbi algorithmic rule, an instance of dynamic programming ways. This could be viewed as a sort of profile-sequence alignment. More recently, profile-profile methods are shown to considerably improve the sensitivity of fold recognition over profile-sequence, or sequence-sequence, methods. In the HMM version of profile-profile ways, the HMM for the query is aligned with the prebuilt HMMs of the guide library. This manner of profile-profile alignment is additionally computed using standard dynamic programming methods.

In a Markovian sequence, the character showing at position t solely depends on the k preceding characters, k being the order of the Markov chain. Hence, a Markov chain is totally outlined by the set of probabilities of every character given the past of the sequence during a k -long window: the transition matrix. Within the hidden Markov model, the transition matrix can be modified along the sequence. The selection of the transition matrix is ruled by another Markovian process, usually called the hidden process. Hidden Markov models are thus particularly helpful to represent sequence heterogeneity. These models are utilized in predictive approaches: some algorithmic programs just like the Viterbi algorithm and also the forward-backward procedure permit to recover which transition matrix was used on the determined sequence[11].

One of the main motivations for studying Markov chains is that in some sense they have a finite explanation. The Markov property says that the latest measurement X_{t-1} have all the information contained in all the past

measurements. Thus, once the observer measures X_{t-1} , he can throw all past measurements without any loss of information. Now what happens if a stochastic process is not Markovian? Hidden Markov models (HMM's) are somewhat more general models for stochastic processes that still retain the "finite memory" feature of Markov chains. Specifically, suppose we now have two sets $X := \{s_1, \dots, s_n\}$ called the state space, and $Y := \{r_1, \dots, r_m\}$ called the output space. Suppose $\{X_t\}$ is a Markov chain. At each time t , the state X_t induces a probability distribution on the output space Y , as follows:

$$Pr\{Y_t = r_i | X_t = s_j\} = b_{ij}, \forall i, j. \quad (4)$$

The stochastic process $\{Y_t\}$ is said to obey a hidden Markov model (HMM). Thus a HMM is described by an $n \times n$ matrix A of state transition probabilities, and another $m \times n$ matrix B of readout probabilities [11].

3.2 Support Vector Machines

Support Vector Machine is supervised Machine Learning approach. Computer Vision is the vast area whereas Machine Learning is one of the application domains of Artificial Intelligence besides with pattern recognition, Natural Language Processing, Robotics. Supervised learning, Un-supervised learning, Semi-supervised learning and reinforcement learning are several types of Machine Learning.

Support Vector Machines (SVM) are learning systems that use a hypothesis space of linear function in a high dimensional feature space, trained with a learning algorithm from optimization theory that implements a learning bias derived from statistical learning theory. This learning strategy introduced by Vapnik and co-workers is a principled and very powerful technique that in the few years since its introduction has already outperformed various other systems in a large variety of applications. The learning machine is given a training set of examples (or inputs) with associated labels (or output values) in supervised learning. Mostly the examples are in the form of attributes vectors, so that the input space is a subset of R_n [6]. Once the attributes vectors are accessible, a number of sets of hypotheses could be selected for the problem. Among these, linear functions are best understood and simplest to apply. Traditional analytics and the classical neural networks literature have developed many methods for differentiating between two classes of instances using linear functions.

Vapnik, based on statistical learning theory, proposed a novel method known as support vectors machines to perform data classification and regression. Because of the high performance, SVM is receiving the attentions of more researchers in bioinformatics lately. Many learning algorithms including ANN implement the empirical risk minimization principle to learn the prediction model. The empirical risk $R_{emp}[f]$ in equation is given by the fitting error of the model f on the training data.

$$R_{emp}[f] = \frac{1}{l} \sum_{i=1}^l |f(x_i) - y_i| \quad (5)$$

Vapnik proved a type of error estimate below:

$$R[f] \leq R_{emp}[f] + capacity \quad (6)$$

Where $R[f] = \int |f(x) - y| dP(x, y)$ is the expected error of the model f on the test samples drawn from the distribution $P(x, y)$. The training samples and the test samples are assumed to have the same distribution function. SVM tries to control both the empirical error and the model complexity, which is controlled by the *capacity* term in equation (6), at the same time. Using the structure risk minimization principle, SVM finds a balance between the fitting power and the model complexity of the learning function. Frequently, SVM can avoid the overfitting problem commonly encountered in ANN and other learning algorithms. SVM was originally designed for binary classification. Since proteins have three different types of secondary structures according to our reduction method above, some modification of the SVM usage is necessary. Researchers have proposed different methods to solve the multi-class problem. One of them is to combine several binary classifiers to construct the tertiary classifier. This type of solution will be called the combination method. The other type of solution is to solve the multi-class problem directly by extending the original SVM theory. We will call the latter type the decomposition method. The open source software BSVM was used to train and predict data. A radial basis function (RBF) kernel was adopted for the BSVM, and we used a soft margin to handle the noises. This left us with two hyper-parameters C (the regularization parameter that controls the weight of the fitting error) and γ (the width of the Gaussian function) to determine. BSVM provides a tool to determine these values optimally. It was found that the best result was obtained when $C = 1.5$ and $\gamma = 0.15$ [6].

Protein secondary structure prediction has been satisfactorily performed by machine learning approaches such as support vector machines. Most secondary

structure prediction programs purpose only alpha helix and beta sheet structures and summarize all other structures in the random coil pseudo class. However, such an assignment often ignores current local ordering in so-called random coil regions. Signatures for this ordering are different dihedral angle pattern. Olav Zimmermann proposes a multi-step Support Vector Machine (SVM) procedure as a different technique to predict directly dihedral regions for each residue as this leads to a large amount of structural information that in dihedral regions without alpha helices or beta sheets is higher than those from secondary structure prediction programs [21].

3.3 Neural Network

Neural network is used in several fields of study and a high degree of attention, made some boost progress. The behaviour of the neural network depends mainly on two aspects: one is the topology of the network, a network learning rules. Neural network comprises of many nodes mesh structure, each node in the network structure has some assigned values. Neural network generally includes three levels i.e input layer, hidden layer and output layer. The organization form of internal nodes based on neural network, neural network can be divided into distinct types of structures. Multilayer neural network consists of input layer, hidden layer and output layer; this builds a multi-layer neural network. Analysis shows that by comparing with the single layer network, multilayer neural network has better capability to process information, especially for complicated information processing ability.

This technique is based on the operation of synaptic connections in neurons of the brain, where input is processed in various levels and generalized to a final output. The neural network is often “trained” to generalize definite input signals to a desired output. In the secondary structure prediction of neural networks, input is a sliding window with 13-17 residue sequence. Information from the central amino acid of each input window is adjusted by a weighting factor, calculates and sent to a next level, termed the hidden layer, where the signal is transferred into a number near to either 1 or zero, and then transported to three output units representing each of the possible secondary structures. The output units each weigh the signal again and sum them, and convert them into either a 1 or a 0. An output signal near to 1 for a secondary structure unit demonstrates that the secondary structure of that unit is predicted and an output signal near to 0 indicated that it is not predicted. Neural network are trained by adjusting the values of the weights that

modify that signals using a training set of sequences with known structure.

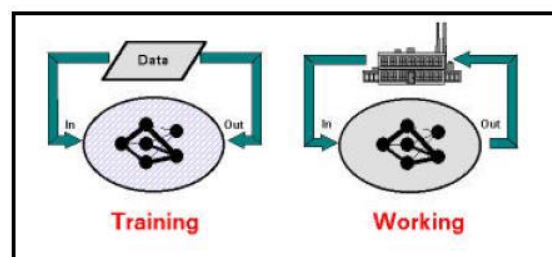


Fig.4- Training and Working Phase for Supervised Learning [8]

Neural Networks used in the Protein Secondary Structure Prediction

Neural networks are employed in protein prediction for regarding twenty years, and this previous work focussed on predicting the secondary structures of proteins. So as to evaluate the accuracy of these approaches, they all use the three state accuracy, or so- 31 referred to as Q3 scores. Here, the Q3 scores refer to the proportion of properly predicted residues within the three sub-structures, the helix, the strand and also the coil. The earliest use of neural networks for protein secondary structure prediction was in 1988, and it absolutely was introduced by Qian and Sejnowski. They employed a fully connected multi-layer feed-forward neural network. [28]

Artificial Neural Network - Supervised Training Algorithms.

There are many several neural network models and algorithms. Supervised learning infers the function from the supervised training data, and every example in the training data consists of an input and a target output. The supervised learning approach analyses the link between the inputs and target outputs within the training set, and produces an inferred function. The function is then employed by the network to predict the output for a given input. The prediction is predicated on the learned relationship between known amino acid sequences and 3D shapes, which is a method of supervised learning. A commonly employed algorithm for the multi-layer feed-forward neural network is the backpropagation algorithm. Its learning process can be divided into two sub-processes. First, the data feeds forward from the input layer to the output layer. Second, the error is back-propagated from the output layer to the input layer to raise the difference in between the actual output and the target output. Together

with the backpropagation algorithm, some other optimization algorithms can be employed as the learning method. Therefore, we employ five learning methods, namely, the Gradient descent BP algorithm, the Gradient descent BP algorithm with momentum, the Conjugate gradient descent algorithm with Fletcher-Reeves updates, the BFGS Quasi-Newton algorithm and the Levenberg-Marquadt (LM) algorithm. The structure of multilayer feed-forward neural networks and the feed-forward process is presented in the next section. [28]

➤ **A Multilayer feed-forward Neural Network**

A Multilayer feed-forward neural network comprises of various simple, highly interconnected computing units, called neurons, which are coordinated in layers. Recall that every neuron acts a simple task of information processing, i.e. converting received inputs into processed outputs. The neurons are connected via linking edges. The knowledge learned between network inputs and outputs is achieved and stored as edges weights and biases, according to the strength of the links between different nodes. Although the information is shared at each node, overall the neural network acts well and efficiently. The architecture of a 3-layer feed-forward neural network is shown in Figure 5. The neurons in this network are coordinated into three layers (i.e. input layer, hidden layer and output layer), and the neurons in every layer are fully connected to neurons in the adjacent layer. In the feed-forward network, the data flows in one direction, from the input layer to the output layer, without feedback from the output layer. [28]

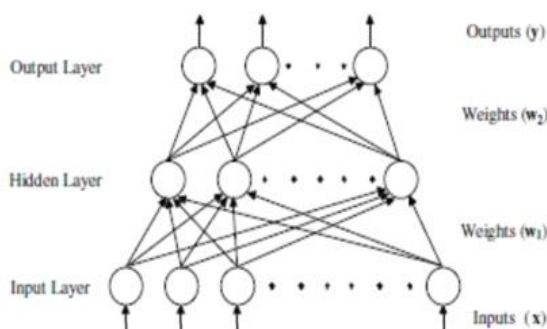


Fig.5- Multi-layer feed-forward neural networks [28]

Neurons in the input layer are static, since they simply get data and pass it on to the neurons in the hidden layer, without any data conversion. Neurons in the hidden layer are fully connected with neurons in both the input and output layers, and are critical for neural networks to learn the mapping relationships between inputs and outputs. After getting information from neurons in the

input layer, the hidden neurons process the information using the transfer function, and then propagate the processed information to the output layer for further processing to generate the outputs. Although neural networks may have more than one hidden layer, most applications only use one. Thus, the multi-layer feed-forward neural network architecture is represented by the number of layers, the number of nodes in every layer, and the transfer function employed in every layer, since the nodes in the similar layer use the similar transfer function. However, there is no widely accepted procedure to resolve the architecture of an MLP, like the number of hidden layers and the number of neurons in every layer. Therefore, it is a 22 complex process to construct a neural network. In general, the number of hidden layers and number of neurons in every layer depends on the categories of transfer functions and learning algorithms and the problems that required to be solved, and is usually resolved empirically. Due to the complication of establishing a neural network, the cost of overly huge neural networks may be more, particularly when the model structure is huge and the input has a large number of dimensions. In order to know how a neural network works, we first need to understand how the neurons in the hidden layers and output layers process data; this mechanism is shown in Figure 5. Information is prepared in two steps at each neuron in the hidden and output layers. In the first step, inputs are multiplied by the weights, related to every corresponding edge, and then gather to form a weighted sum. In the second step, the neuron uses the transfer function to transfer the sum into the output. In several cases there is an additional step between the formation of the weighted sum and the transformation, as few networks may add a bias onto the weighted sum before it is converted by the transfer function. The bias is a constant, which helps to increase the flexibility of the network.[28]

There are various options for the transfer function, but only a few are commonly used in practice. In general, the transfer function is bounded and non-decreasing. The logistic function is usually used transfer function, particularly in the hidden layers, because it is simple, non linear, bounded and monotonically increasing.

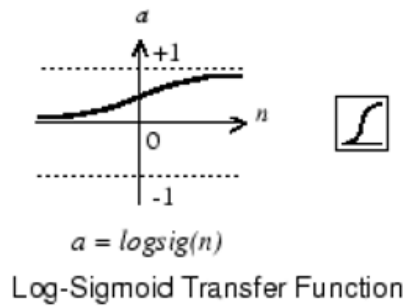


Fig.6- The Sigmoid (logistic) Function [28]

If the logistic function is employed in the hidden layer and the linear function is employed in the output layer, the neural network structure for a feed-forward neural network may be written as:

$$y = b_2 + \sum_{j=1}^q w_{1j} f \left(\sum_{i=1}^p w_{ij} x_i + b_i \right) \quad (7)$$

where y is the predicted output and $\{x_i, i = 1, 2, \dots, p\}$ are the input tuples to the network, p is the number of nodes in the input layer, q is the number of nodes in the hidden layer, $\{w_{ij}, i = 0, 1, 2, \dots, p; j = 1, 2, \dots, q\}$ show the weights located between the input and hidden layers, $\{w_{ij}, i = 0, 1, 2, \dots, p; j = 1, 2, \dots, q\}$ show the weights between the hidden and output layers, b_1 and b_2 are biases in the hidden and output layers, respectively, and f is the logistic function described above. After determining the neural network topology, we train the neural network with the training set, which comprises of the input tuples and the known output tuples. Then the input values are weighted and added at the hidden layer, and the weighted sum is converted through a specific transfer function to form the inputs to the output layer or the inputs to another hidden layer. The similar operation is conducted on the next layer until the network developed its outputs. The actual outputs give by the network can be compared with the desired (target) outputs to actuate how closely they match. This is measured by a score function, the mean squared errors (MSE), which is denoted by E . The target of training a neural network is to develop the score function, namely the mean squared error (MSE), to attain its global minimal. During the developing processes, the set of weights and biases keeps restoring. Thus, the network training process is actually a nonlinear optimization problem. [28]

➤ Optimization Algorithms

Recall that back-propagation (BP) is an approach that propagates the errors backward from network outputs to inputs. During this operation, weights and biases involved in the networks are updated to decrease the distance between the actual predicted values and the desired values. Actually, the goal is to develop the score function. There are five optimization algorithms. These are the gradient descent BP, the gradient descent with momentum BP, the conjugate gradient algorithm with Fletcher-Reeves updates, the BFGS Quasi-Newton algorithm, and the Levenberg Marquadt (LM) algorithm. The first three methods are depends upon the gradient of the score function, and the other two rely on the Hessian Matrix of the score function. Recall that the optimization operation target is to develop the score function, and we value the MSE function as the score function in our work. In the MSE function, the predicted value achieved from neural networks comprises by a function of the parameter θ , that is within the model and required to be updated to lowers the error function. In a neural network, the parameters will be the weights and biases linked to the network. We typically values iterative improvement search methods for this optimization process, and local information to guide the local search. The iterative local optimization algorithm may be broken down into three steps:

1. Initialize: Choose the initial values for the vector of parameter θ . In general these are chosen randomly.

2. Iterate: Starting from $i = 0$, let $\theta^{i+1} = \theta^i + \lambda^{i-1} \bar{v}$, where \bar{v} is the direction for the next step, and λ represents the step size. In our research, \bar{v} is a direction in which the score function is minimized.

3. Convergence: Repeat step 2 until the score function attained a local minimum. Actually, all the following optimization algorithms provide the required information to determine \bar{v} , the direction for the next step. To find the local minimum, we count the gradient and the Hessian of the targeted function when the function is twice continuously differentiable. [28]

i. Gradient Descent BP

The gradient descent BP, the most basic BP algorithm, is referred to as the steepest descent. Among all directions we may move, the steepest descent direction $-\nabla_{\theta} E_{(\theta)}$ is the direction on which the target function (the score function MSE) decreases most fastly. The advantage of this

approach is that it needs only the calculation of the gradient $-\nabla_{\theta} E(\theta)$, but not the second derivatives (Hessian). If the selected learning rate is small enough, the gradient descent is guaranteed to be confluent.

ii. Gradient Descent with momentum BP

This method depends on the gradient descent BP, but it accelerates its concurrency by adding a momentum. This is because there is always a trade-off when selecting the learning rate: if it is too short the convergence will take too much time, while if it is too high we will step too far and miss the minimum. When the momentum term equal to 0, this optimization algorithm is equals to the gradient descent BP, and when the momentum term is greater than 0, the algorithm will accelerate the convergence in the present direction.

3.4 RBF Neural Networks Model

Radial basis functions were introduced in 1985 by Powell and Broomhead. Lowe was the first to exploit the use of Radial basis functions in neural networks. Radial basis function neural networks comprises of 3 layers, as shown in figure

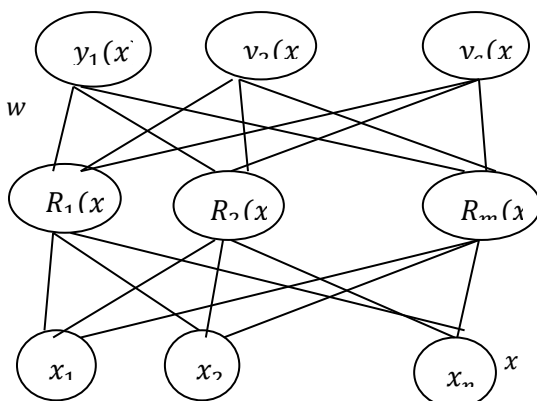


Fig. 7- Radial basis function neural networks

Radial basis function neural networks are feed-forward networks. RBF networks stand for a class of neural network models, in which the hidden units are activated according to the distance between the input units. RBF networks join two different types of learning: supervised and unsupervised. At this time we used supervised learning. A special class of functions is used in RBF networks to perform a nonlinear transformation on a given network input. These functions give the RBF network its name. Radial basis functions are defined by the fact that their response varies (decreases or increases) monotonically with distance from a mid point. A typical radial function is the Gaussian function that is given by

$$R_t(x) = \exp\left(-\frac{\|x-c_i\|^2}{2\sigma_i^2}\right), i = 1,2,3, \dots, m \quad (8)$$

The weight correction can be calculated by back-propagation. The hidden neurons serve as computing units that functions a nonlinear transformation on the input vector by means of radial basis functions that is responsible as a basis for this transformation. Referring to every neuron comprises two parameters: The neuron specific center and the radius of the RBF, which is similar to all hidden neurons. The nonlinear transformation of the input vector is carried out in every hidden neuron by counting the response of the radial function. The Euclidean distance between the input vector and the center vector should be counted [27].

The Prediction Method Based on RBFNN

Basically there are twenty kinds of amino acids. In the first step of the analysis each amino acid might be converted into binary code of 21 units (each unit representing 1 of the amino acids) before it is delivered to neural networks. In the second step twenty units are valued to enter the percentage of every amino acid in the multiple sequence alignment at that position.

The training of the cluster layer executes first. Once finished, the output layer is subjected to supervised learning as employed in the feed-forward network with back-propagation. To train the network, various parameters required to be specified. The maximum learning rate parameter (L_{max}) ranging between 0 and 1, the minimum learning rate parameter (L_{min}) ranging between 0 and 1 and the maximum number of epochs used for training (N). They are used to construct the learning gradient. This gradient is required to update the weights of a cluster node whose input weight vector has the least squared Euclidean distance to the input vector pattern. The maximal learning rate also serves as the initial learning rate. In updating the weight of every dimension, the following formula is applied:

$$W_{new} = W_{old} + [L \times (I_{value} - W_{old})] \quad (9)$$

W_{old} means the weight before updating. W_{new} means the updated weight. L is learning rate. I_{value} is input pattern value. After all input patterns have been run via the network, the learning rate itself is updated via:

$$L_{new} = G \times L_{old} \quad (10)$$

L_{old} means the learning rate before updating. L_{new} means the updated learning rate. G is learning gradient. The gradient is calculated by the following formula:

$$G = L_{max} - E_{complete} \times (L_{max} - L_{min}) \quad (11)$$

$E_{complete}$ means epochs completed. The weight correction w_{ij} of neuron i and j is calculated by back-propagation:

$$w_{ij} = L_{constant} \times G_{local} \times I_j \quad (12)$$

The learning rate constant $L_{constant}$ and the non-normalized minimum average squared error of BP algorithm must be specified. G_{local} means local Gradient. I_j means input signal of neuron j [27].

3.5 Backpropagation network

The artificial neural network, which is made by multi-neuron via a certain rule and is a structure of temporal hierarchical network, can be used to adaptive non-program information processing depends on working mode of brain. Its basic principle and characteristics are mutual supervision and mutual cooperation between neurons, continuous-time nonlinear, the overall role of network, paralleled distributive treatment with high scale, high robust and learning and memory function respectively[24].

Generation of BP neural network

BP neural network is the most representative and usually applied learning algorithm in artificial neural network which is the acyclic multi-level network training algorithm. 3-layered network input layer, hidden layer, and output layer can be selected considered network topology because correctness of network and expression ability can't always be enhanced when hidden layer and its neurons can be enhanced in BP neural network. Generally, for an amino acid, its amino acid residues have statistical correlation which could affect the secondary structure of amino acid. Thus, the input window of neural network can be designed. If 9 amino acids VRKKRWACD can be inputted, the number of neurons is 9x21 in input layer which is coding bit rate of amino acid. The output layer consisted of 3 neurons corresponding to three secondary structures of protein alpha helix, beta sheet, and gamma crimp. Comparing 3 results of output layer, alpha helix coding is 100, beta sheet coding is 010, and gamma crimp is 001. Adaptive adjustment strategy can be employed in this algorithm[24]. BP network can be defined with Mat lab as follows

```
Net=newff(TEMP,[30,3],
{'tansing','purelin'},'traingd');
```

Protein Structure Prediction using Back Propagation Neural Network

The algorithm which is employed to train the ANN having 3-layer is as follows:

- Initialize the random weights in the network (often randomly)
- Do
 - For each example e in the training set
 - O = neural-net-output(network, i) ; forward pass
 - T = teacher output for i
 - Calculate error (T - O) at the output units
 - Compute delta_wh for all weights from hidden layer to output layer ; backward pass
 - Compute delta_wi for all weights from input layer to hidden layer ; backward pass continued
 - Update the weights in the network until all examples classified correctly or stopping criterion satisfied
- Return the network [14]

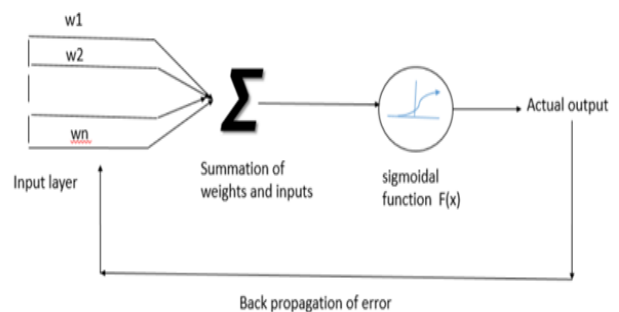


Fig. 8- artificial neural network model with back propagation [1]

Training Parameters

The performance parameters are learning algorithm, transfer function and the number of cycles the network consumes, to be converged.

- Epochs: This decides for how many cycles we would like to train the network. For example, if the number of epochs is 100 then the whole training data will be presented 100 times.
- Learning Algorithm: the learning algorithm is LEARN_GDM (Learning by Gradient Descent Method).
- Min_grad: This decides the acceptable error that we would choose. If this MSE (mean squared error) is reached the network has converged and training will stop. Generally this value is 1e-5.
- Transfer Function: The transfer function selected is Log Sigmoidal function (LOGSIG)

• Gradient descent (GDM) was used to lower the mean squared error between network output and the actual error rate. After initializing all the parameters, the network is ready to be trained. Repeated experiments were performed to get the neural network converged. Weights were initialized to random values and networks were run until at least one of the following termination conditions was satisfied

1. Maximum Epoch
2. Minimum Gradient
3. Performance Goal

If the MSE reaches the set value the network has converged and training stops. If the network does not converge and the number of epochs has reached the set value, then the network has not converged. We will have to retrain the network, by changing some of the parameters or the training algorithm or the network architecture. Once, the network has been trained will be simulated with a set of inputs that the network has not seen. We have classified the data sets into 2 parts. i.e training set and testing set which is not used in the training process, and is used to test and then we have simulated our results with these datasets. Almost 2/3rd of the total dataset can be taken as training set and 1/3rd of the rest can be taken as test set. This is done through the analysis of the correctness attained via testing against these set. [14]

3.6 Approaches for Protein Tertiary Structure Prediction

Presently 3 approaches are followed for the estimating of the tertiary structure of proteins. These are 1) Homology modelling, 2) Threading and 3) Ab initio structure prediction.

3.6.1 Homology Modelling

This is the simplest and very reliable approach. The observation that proteins with same sequences tend to fold into same structures forms the basis for this method. It can be noticed that even proteins with 25% sequence identity fold into same structures. This approach does not work for remote homologs (< 25% pair wise identity). The method for homology modelling may be briefly defined as: given a query sequence Q, and a sequence database of known protein structures, find a protein P such that P has same sequence as to Q and return P's structure as an approximation to Q' structure. The following are the main steps in homology modelling:

- 1) Finding known structures linked to the query sequence whose structure has to be modeled
- 2) Aligning the query sequence to the templates
- 3) Constructing variable side-chains and main-chains and
- 4) Model refinement, assessing the model built and choosing the most native conformations.

3.6.2 Threading

Threading is an approach for fold recognition. This is employed for sequences having sequence identity $\leq 30\%$. In this approach, given a sequence and the set of folds available in the Protein Data Bank (PDB) the target is to see if the sequence can accept one of the folds of known structure. This approach takes merit of the knowledge of existing structures and the principles by which they are stabilized. Fold assignment and alignment are attained by threading the sequence via every structure in a library of all known folds.

3.6.3 Ab initio (de novo) structure prediction

While homology modelling and threading needs knowledge of known structures, ab initio structure prediction has no limitations like above approaches. It starts with the estimation that the real structure of a protein is at the global free energy minimal. In previous years a particularly successful approach called Rosetta has been developed by Baker and colleagues (Simons et al, 1997). This approach has assimilated information got from known structures and is depends upon a picture of protein folding in which small segments of the protein chain flicker between distinct local structures consistent with their local sequence, and folding to the real state occurs when these local segments are oriented so that low free energy interactions are made throughout the protein. [22]

Applications and limitations of methods for structure prediction

Each of the above methods described provide structural description to different extent. While homology modelling can give atomic level details of the target protein, threading can help only to know the fold of the protein. Baker and Sali (2001) have describes the accuracy and application of protein structure models with examples. Large and medium level homology models with sequence identity > 30% are convenient in refining functional prediction like ligand binding. Low accuracy models of several of the ribosomal proteins were

favorable in building the molecular model for whole yeast ribosome. The correctness and applicability of models produced by ab initio methods are in general of lesser accuracy compared to models obtained from either homology modelling or threading. These are convenient in predicting functional relationships from structural similarity and for identification of patches of conserved surface residues.

3.7 PSP using Bio-Inspired algorithms

Bio-inspired methods are simpler than standard GAs, producing results in less time with less cost. Swarm Intelligence (SI) that is one of the research trends, inspires from animal colonies, their way of life and their behaviour, to solve the optimization problems. We will discuss these algorithms in this section.

3.7.1 Honey Bee Algorithms

One of the algorithms under SI group, is marriage in honey bees and inspires from reproduction process of honey bees. Honey bee colony comprises of a queen, broods, drones, and workers. It considers workers as local search methods that enhance the results and queen as the answer. At the end, the enhanced and mutated brood will be taken as the queen if it has a good fitness value. Children would be produced from crossover of the drones and queens. Some modification can be done for replacing a child with the queen to enhance the result and rescue from trapping in local minimums. ECEPP/3 is employed for fitness function and torsion angles are employed for representation of proteins. The approach has 3 main parts. In the first part, a random structure is comprised as the queen then few random structures are created as drones in the second part. In the last part one of the drones can be chosen to take part in the crossover with the queen and produce a brood. This brood can be mutated and enhanced and will replace the queen if it has good fitness value. This modified approach brings good result but it is tested only on a one small structure. Besides, it is slow and computationally complicated.

Another method in this group is honey bee that is employed for PSP problem. Food sources are possible conformations and structures for every sequence. These sources are presented as sequences of torsion angles that are first produced randomly. Then scout bees search for the sources and compute their energy values by ECEPP/2 function. Values will be sorted and conformations with small energies will be selected. These good conformations can be the next generation in addition to new randomly produced conformations. These processes continue until

the best conformation is produced. The main difficulty is that it is also only tested on a one small protein.

Artificial bee colony technique is used in for protein conformational space to find the real conformation that has the smallest energy value. Dihedral angles are employed to represent the protein and ECEPP/2 as the fitness function. At first, N food sources are produced randomly and their energy is computed using the fitness function. Then every structure is enhanced with the SMMP package, which uses torsion angles to simulate peptides and proteins linearly. After enhancing, structures are sorted and the best half is chosen to make the half of next generation. In addition, new sources near the best half are produced and enhanced with SMMP. If the energy value of these new sources is better than the second half of parents, they will be selected for the next generation. This process goes on until the best conformation is attained [9].

3.7.2 Ant Colony Algorithms

Another group of methods used for PSP is the ant colony algorithm that is inspired from the social life of ants and food foraging. It was first proposed by Dorigo et al. to resolve the complicated problems. Ants communicate with each other and find food conformations via pheromone trails. Some research has been done to employ this method for PSP problem. In the ant colony algorithm is used based on 2D HP model that first produces a random pheromone matrix, then few ants are comprised as conformations. Each ant chose a random position on the sequence as the initial point and each moment adds amino acids in both directions, according to the position of two previous amino acids and pheromone matrix. Movement direction is one of the directions of forward, right or left. When the first step is done, next step is the local search only for sources with good fitness values. These conformations are changed, initialize from a random first position on the sequence until the energy values are minimized. At last the pheromone matrix of these selected sources can be updated to be valued in the next iteration. Finally, the source with the lesser energy is called the native structure. The accuracy of this algorithm is not good because it only comprised three degrees of folding while minimizing the computational complication. Another different work is done in that only adds amino acids in construction phase from one direction in contrast to. Also the local search only comprises the best sources not all of them as in. Both these two works estimate the longer sequences in more time and produced results are not needed. In the ant colony, method is used based on 3D

HP model and five distinct directions are considered; forward, right, left, up and down[9].

3.7.3 Particle Swarm Optimization Algorithms

In comparison to the methods that are valued to find the near optimal solutions, PSO produces acceptable result. Each particle of the PSO is one solution of the search space. The target of this algorithm is that the whole particle reaches the optimal global point. In this way the speed of every particle changes dynamically, comprises the last move and the neighbor particles [3]. So position and velocity vectors for the *i*th particle in the *D* dimensional search space are $X_i = (X_{i1}, \dots, X_{id})$ and $V_i = (V_{i1}, \dots, V_{id})$ respectively. Speed changes give each particle the capability to search around the best local or global point. The overall equations to update the position and velocity recursively are as follows:

$$V_{id}^{k+1} = V_{id}^k + c_1 \times rand(.) \times (P_{id} - x_{id}^k) + c_2 \times rand(.) \times (P_{gd} - x_{id}^k).$$

$$X_{id}^{k+1} = X_{id}^k + V_{id}^{k+1}$$

Each particle in PSP is one possible structure for the protein and position vectors are the vectors of angles. Velocity Vectors tells the rate of getting nearer to the answer and changing these angles. It's also used depends on various models, of which are successful but none of them is used to compute the real tertiary structures. Most of these implementations are serial except for parallel implementation in that depends upon distributed systems [9].

3.7.4 Clonal Selection Algorithms

Burnet introduced Clonal Selection Algorithm (CSA) which inspires from nature that immune system chooses specific cells for reproduction and clonal selection. CSA can be employed for PSP problem; randomly chosen angles are the probable sources of the first generation. Then every individual is cloned and mutation is applied to all of them to produce children and the structure with good fitness, among every parent and its children is chosen. This process is continued until the native structure is attained. Also to accelerate this process, Zhu et al. employed Graphical Processing Unit (GPU). They converted fitness function calculations to the GPU to be counted for every individual in parallel. The speed up for sequences of length 13 and 21, were 10.94 and 16.4 respectively. But the sequences being studied are not real and are Fibonacci sequence [9].

3.8 Framework of Machine Learning-Based Methods

This section discussed the mechanism of protein fold recognition by machine learning-based algorithms. The overall process in protein fold recognition by machine learning-based algorithms comprises 2 phases (Figure 9): (1) model training; and (2) prediction. In the first phase (model building), query protein sequences are first submitted into a pipeline of feature representation, in which sequences of distinct lengths are encoded with fixed-length feature vectors by feature descriptors. The mostly employed feature descriptors include Amino Acid Composition (AAC), Pseudo AAC, Functional Domain (FunD), Position Specific Scoring Matrix (PSSM)-based descriptors, Secondary Structure-based descriptors, and Autocross-covariance (ACC) transformation. When the resulting feature representations show few irrelevant features or redundant features, an alternative step is mostly performed to choose the optimal feature subsets, which can yield the best performance, from the resulting feature representations. Subsequently, the feature vectors are put into a pre-selected classification method to train a estimation model. Typical classification methods often employed in model building such as SVM, Random Forest (RF), Naïve Bayes (NB), and Logistic Regression (LR). The first phase is completed in this step.

In the second phase (prediction), uncharacterized query proteins are first submitted into the similar pipeline of feature representation as in the first phase. Notice that if feature optimization of the generated feature representation is performed in the first phase, feature optimization might also be performed in the second phase; otherwise, the resulting feature vectors are put into the trained prediction model, wherein the protein fold class to which the query proteins belong is predicted feature vectors by feature descriptors.

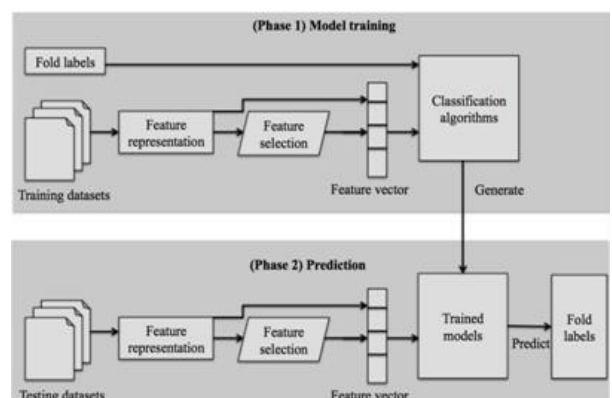


Fig. 9: Framework of machine learning-based methods for protein fold recognition. [26]

4. ACKNOWLEDGEMENT

I would like to thank my guide Dr. BALJIT SINGH KHEHRA, Professor and Head of Computer Science and Engineering department, for providing necessary infrastructure for the research work and helping me to write this paper.

REFERENCES

- [1] Shivani, B. Arushi, M. Deepali, "Design and Implementation of an Algorithm to Predict Secondary Structure of Proteins using Artificial Neural Network", International Journal of Emerging Research in Management & Technology ISSN: 2278-9359, Vol. 2, No. 12, pp. 34-40, 2013.
- [2] B. Hemashree, S.K. Kandarpa, "Protein Structure Prediction using Artificial Neural Network", Special Issue of International Journal of Computer Applications (0975 - 8887) on Electronics, Information and Communication Engineering - ICEICE, No.3, pp. 22-24, 2011.
- [3] B. Wenzheng, C. Yiming, C. Yuehui, "Multiple Feature Fusion Protein Tertiary Structure Prediction", International Conference on Information Science and Cloud Computing Companion, pp. 751-756, 2013.
- [4] C. Jianlin, N. T. Allison, and B. Pierre, "Machine Learning Methods for Protein Structure Prediction", IEEE reviews in biomedical engineering, VOL. 1, pp. 41-49, 2008.
- [5] C. Nandini, "A Survey on Protein Sequence Classification with Data Mining Techniques", International Journal Of Scientific & Engineering Research, Vol. 7, No. 7, pp. 1442-1449, July 2016.
- [6] D. H. Shing, Y. Y. Chi, "Secondary Structure Prediction Using SVM and Clustering", Proceedings of the Fourth International Conference on Hybrid Intelligent Systems, 2004
- [7] F. Mhamdi, "Text mining, feature selection and data mining for proteins classification", International conference on information and communication technologies, pp. 457-458, April 2004.
- [8] I. S. Mohammad, A. Hakimeh, "RMSD Protein Tertiary Structure Prediction with Soft Computing", International Journal of Mathematical Sciences and Computing, Vol. 2, pp. 24-33, 2016.
- [9] K. Yasaman, F. Mahmood, K. Hamed, S. Hossein, "Protein Structure Prediction using Bio-Inspired Algorithm: a review", The 16th CSI International Symposium on Artificial Intelligence and Signal Processing, pp. 201-206, 2012.
- [10] M. Sonal, P. Yadunath and A. Anamika, "Classification of Protein Structure (RMSD \leq 6Å) using Physicochemical Properties", International Journal of Bio-Science and Bio-Technology, Vol.7, No.6, pp.141-150, 2015.
- [11] M. Vidyasagar, "Some Challenges in Computational Biology", European control conference, pp. 3364-3369, 2013.
- [12] N. Mathuriya, "Comparison of K-means and Backpropagation Data Mining Algorithms", International Journal of Computer Technology and Electronics Engineering, Vol. 2, No. 2, pp. 151-155, April 2012.
- [13] P. Mayuri, S. Hitesh, "Protein Secondary Structure Prediction Using Support Vector Machines (SVMs)", International Conference on Machine Intelligence Research and Advancement, pp. 594-598, 2013
- [14] P. Rojalina, D. Nilamadhab, R. Samita, "A Novel Approach for Protein Structure Prediction using Back Propagation Neural Network", International Journal of Computer Science And Technology, Vol. 3, Issue 2, pp. 600-603, 2012.
- [15] P. Sun and J. Zhang, "Improved Prediction Method of Protein Contact Based on RBF Neural Network", 3rd International conference on bioinformatics and biomedical engineering, pp. 1-4, 2009.
- [16] R. N. Chandrayani, K. Manali, "Bioinformatics: Protein Structure Prediction", 4th IEEEICNT, 2013
- [17] R. S. Prashant, S. Harish, B. Mahua and S. Anupam, "Quality assessment of modelled protein structure using physicochemical properties", Journal of Bioinformatics and Computational Biology, pp. 1-16, 2015.
- [18] S. Saha, "Application of Data Mining In Protein Sequence Classification", International Journal of Database Management Systems, Vol. 4, No. 5, pp. 103-118, October 2012.
- [19] V. Swati, S. K. Bithin, R. K. Santanu, "An Efficient Technique for Protein Classification Using Feature Extraction by Artificial Neural Networks", Annual IEEE India conference (INDICON), 2010.
- [20] W. Dianhui, L. K. Nung, D. S. Tharam, "Data Mining for Building Neural Protein Sequence Classification Systems with Improved Performance", IEEE Publications, pp. 1746-1751, 2003.
- [21] W. Bo, L. Yongkui, Y. Jian, L. Shuang, "Application Research of Protein Structure Prediction Based Support Vector Machine", International Symposium on Knowledge Acquisition and Modeling, pp. 581-584, 2008.
- [22] https://www.researchgate.net/publication/228894892_Approaches_to_Protein_Structure_Prediction_and_Their_Applications
- [23] W. J. Barry, "Backpropagation Neural Network: A tutorial", *Chemometrics and Intelligent Laboratory System* Elsevier Science Publishers B.V., Amsterdam, Vol. 18, pp. 115-155, 1993.
- [24] W. Jian, L. Jian-Ping, "Protein Secondary Structure Prediction Based on BP Neural Network and Quasi-Newton Algorithm" International conference on apperceiving computing and intelligent analysis (ICACIA) IEEE, pp. 128-131, 2008.
- [25] W. L. George, P. Marius, L. N. Michael, *Neural Network Analysis Of protein tertiary structure*, Tetrahedron Computer Methodology, Vol. 3, No. 3/4, pp. 191-211, 1990.
- [26] W. Leyi and Z. Quan, "Recent Progress in Machine Learning-Based Methods for Protein Fold Recognition", International Journal of molecular sciences, pp. 1-13, 2016.
- [27] Z. Zhen, J. Nan, "RADIAL BASIS FUNCTION METHOD FOR PREDICTION OF PROTEIN SECONDARY STRUCTURE", Proceedings of the Seventh International Conference on Machine Learning and Cybernetics, Kunming, Vol. 3, pp. 1379-1383, 2008.
- [28] https://www.ruor.uottawa.ca/bitstream/10393/23636/3/Zhao_Jin_g_2013_thesis.pdf