# Project Report

M.S. student *Umberto Di Fabrizio*

*Abstract*— **The protein superfamily classification problem, which consists of determining the super-family membership of a given unknown protein sequence, is very important for a biologist. In this work, the classification task is tackled using several ANN (Backprop, CNN, LAMSTAR) which are compared with respect to their performances. The goal is to predict the family of novel protein sequences based on the sequence only. The results obtained are comparable to the state of the art and the methods can be further investigated to be adapted to a larger dataset of families.**

## I. INTRODUCTION

Bioinformatics has been growing in the last three decades[2] given the huge amount of biological data that is continuously gathered, mainly about DNA, RNA and proteins. The volume of data generated from project such as The Human Genome Project[3](1990-2003) has strengthen the collaboration between the computer scientist community and the biologist one.

One of the most challenging problem is to classify protein accordingly to their function or by their family or superfamily. Protein sequences are composed by an unique sequence of 20 amino acids which determines the protein function. They carry out fundamental roles for the cells functions ( basically represent the blueprint of the cell), infact they determine the shape and the structure of the cell.

Each protein encode a certain function which depends on its structure and amino acids sequence but can only be completely understood with experiments. Those experiments are costly and slow thus they cannot keep pace with the amount of information available and which needs to be annotated.

The challenge that has been tackled is to classify proteins into functional or structural existing superfamilies so that the annotation process can be automated.

A protein superfamily is a set of proteins for which common ancestry can be inferred so they possed sequence or structural homology. In a superfamily classification, an unlabeled protein sequence may belong to any of the superfamily from a set of known superfamilies. This classification is enormously useful because similar protein sequences exhibit almost the same biological structure and function, more importantly one of the main reason is treating and preventing genetic disease as well as drug discovery, prediction of molecular function and medical diagnosis.

## II. BACKGROUND

Several methods have been investigated in order to solve the superfamily classification problem: determining the superfamily membership of a given unknown protein sequence. BLAST (Basic Local Alignment Search Tool) [4] is a tool that uses direct modelling, performing a search of homologie between sequences. This software ex-

plores the local alignment in pairs to measure the similarity between sequences. The classification is done based on the alignment which had the greatest punctuation.

Another method that uses direct modelling is the HMM Hidden Markov Models that is widely used for probabilistic modelling of family of proteins. It uses probabilistic values to score how much an unkown protein belongs to a given family.

A Fuzzy ARTMAP model, a machine learning method has been proposed used to classify the protein sequence[8].

The use of Neural Networks to tackle this problem has been successfully presented in the work by C.H.Wu at al.[5][6][7] and an introductory survey of neural networks applied to genome problems can be found in the book by C.H.Wu and Jerry W. McLarty[9]. The book explains the basic idea of neural networks presenting the different kinds of network that can be used and gives meaningful example on how to use them.

The process of protein classification using ANN can be divided in two parts: (i) pre-processing: protein sequence encoding; (ii) processing the protein classification using the ANN.

The first step is very challenging, the reason is that protein sequences have different lengths and can even be very long (~1000) whilst the neural network has a fixed amount of inputs and can hardly handle missing inputs. The methods of encoding can be basically dived in two types[9]: direct or indirect. The direct encoding basically translate each amino acid into a binary vector[1] accordingly to the one-hot encoding, this means that given that we have 20 possible amino acids then each of them will be represented by a vector of 20 bits with only one position at '1' and all the rest with '0'. Of course this kind of encoding is not feasible in the common case because a protein sequence with 300 amino acid will be translated with 20*300=6000 binary inputs.

The indirect encoding tries to extract useful features from the protein sequence in order to give meaningful inputs to the neural networks, one example is usually the n-gram hashing method[5]. This ngram method computes residue frequencies, the basic idead is that if we have a sequence *Seq*='ACACTGAC' then the possible bi-gram are 'AA', 'AC', 'AT', 'AG', 'CC', 'CA, 'CT', 'CG', 'TT', 'TA, 'TC', 'TG', 'GG', 'GA', 'GT', 'GC', and for each of them the occurrences are counted, usually the approach involves also the mapping of the original sequence to a smaller alphabet. A summary table on the encoding pros and cons is shown in Table I

---

[1]other techniques will be presented later

|  | Pros | Cons |
|---|---|---|
| Direct | Keeps the sequence order | Encoding is too large |
| Indirect | Independent of length of the biosequence | It is hard to find meaningful feature, this requires domain knowledge. The order of the sequence is not mantained |

## III. OUTLINE

The idea is to exploit the power of the Deep Learning Neural Network in order to classify the protein in its family.

In section DATA COLLECTION it is explained where the data was collected from and which tools have been used, in section DATA EN-CODING the encoding method to transform the sequence to numbers for the ANN is discussed. Sections BACKPROPAGATION CONVOLU-TIONAL NEURAL NETWORK, LAMSTAR discuss the design and structure of the ANNs employed to solve the problem, together with the choices that were made to optimize the execu-tion time and the accuracy of the nets. Finally section RESULTS presents the results and in CONCLUSION the contribution of the work are summarized and compared to the state of the art techniques.

## IV. DATA COLLECTION

The objective is to obtain for each superfamily or family the list of all the proteins that belong to that family and their complete protein sequence, an example of data is shown in table II.

This step is particularly tricky because there is not a database of protein sequences and fam-ilies, or better there are several but they all have to be queried manually through a web interface. For this reason the data is scraped from the web sites (http://prosite.expasy.org/, http://www.uniprot.org/uniprot/) with the use of the Kimono platform[10] and Rscripts. Automa-tizing this step is particularly useful in case any new superfamily data has to be collected.

With this process 4 families have been chosen: *ADH_SHORT*, *ABC_TRANSPORTER_2*, *G_TR_2*, *globin*, and for each of them 600 protein se-quences have been collected. The overall dataset is thus made by 2400 protein sequences and it is divided between training(90%) and testing(10%).

## V. DATA ENCODING

The sequences of proteins are encoded through the indirect method. The purpose is to extract as many information as possible from the protein without creating too many inputs for the neural network. The encoding adopted is a bigram hash-ing method as presented in [9] but the 20 amino acids are not mapped into 6 macro categories(as suggested), in order to keep as much information as possible. In this way we have 20x20 possible bigrams although there are also 3 special char-acters ('X','B','Z') that represent uncertainty of some amino acid base (e.g. X represent when it was not possible to detect the right amino acid), for this reason the total amount of inputs generated is 23x23=529. Notice that since we are using bigrams the number of inputs is the square of the number of possible amino acids, this is particular useful when we will reshape the input into a 23x23 matrix for the CNN.

A python script extracts the frequency of bigrams and normalize it on the length of the protein sequence, the information that we obtain represent a global feature of the sequence, an example of bigram frequency is shown in fig 1.

## VI. BACKPROPAGATION

The input of the neural network is the encoded protein sequence, the output is the family of the protein.

The network structure is shown in 2: it is composed by an input layer of 6 neurons, an hidden layer of 5 neurons and an output layer of 4 neurons. The layer are fully connected so the total amount of weights to train are $529 * 6 + 6 * 5 + 5 * 4 = 3224$. This particular configuration has been chosen for its simplicity, the network infact has been reduced to the min-imum number of neurons required to have an high accuracy so that the execution time is mini-mized. In all the implementation that are presented

| Family | Sequence |
|---|---|
| ABC_TRANSPORTER_2 | MAEAPAKKLTVSATEVAVEIVNMNKWYGDFHVLRDINLKVMRGERIVIAGPSGSGKSTMI RCINRLEEHQKGKIVVDGTELTNDLKKIDEVRREVGMVFQHFNLFPHLTILENCTLAPIW VRKMPKKQAEEVAMHFLKRVKIPEQANKYPGQLSGGQQQRVAIARSLCMNPKIMLFDEPT SALDPEMIKEVLDTMVGLAEEGMTMLCVTHEMGFARQVANRVIFMDQGQIVEQNEPAAFF DNPQHERTKLFLSQILH |
| ABC_TRANSPORTER_2 | MDKTRQTELVRWLKQHSTSAKRWLRISMLLGVVSGLLIIAQAWFLAVILQALIMEHTPRE QLLTPFILLLAVFVLRALLTVIRERVGFRCGQVVRQEVRNMVLNKLQALGPVWVKGKPAG SWATIVLEQIEDMQEYYSRYLPQMYLAGIIPIMILIAIFPFNWAAALILFATAPLIPIFM ALVGLGAADANRRNFVALGRLSGSFLDRLRGLDTLRLFFREKAEVQQIRESTEDFRSRTM EVLRMAFLSSGVLEFFASISIAIVAVYFGFSYLGELNFGSYGLPVTMFAGFLALILSPEF FQPLRDLGTYYHAKAQAVGAAESLVTLLESDGEQKTETGDKTPQDKPIQIEANKLEIYSH DGQRLVGPLDFTIEPQQRIAVFGQSGAGKSSLLNLLLGFLPYKGSIKINGDELKELCPDK WRALIGWVGQNPHLPEQTLIENICLGKPTASEAEIQQAIDDAYVSEFLPMLPDGLNTRLG DYAARLSVGQAQRVAVARTLLKPSRILLLDEPAASLDAHSEKRVMHTLNQLAQQQTTIMV THLLEETVNYDQIWVMANGQIIQRGHYAQLSQSEGPFARLLAHRSEEL |
| ADH_SHORT | MMDWNNKNVVYVGGFSGFGYQVCQMMMKKPMKHLIVCSRMENVEMLKKLQAINTSVKVMF VQMNIADYASIVKGVKQVIGHVGHVDVLINGVGGLADKDVETTVAVNLTGLINTTLMFMP YMDKTQSGHGGMVVSISSVYGLEPGPAFSVYSAAKHGGIGFTRSMADEHLYHKTGVAFMC ICPAMTSTELMMNKRDMNWMKWVPHSEEMWKMVMDAKMQTPEECAVNMMTAMEQAKNGAI YICSTSGMKEITPTVYMH |
| ADH_SHORT | MTNRLQGKVALVTGGASGVGLEVVKLLLGEGAKVAFSDINEAAGQQLAAELGERSMFVRH DVSSEADWTLVMAAVQRRLGTLNVLVNNAGILLPGDMETGRLEDFSRLLKINTESVFIGC QQGIAAMKETGGSIINMASVSSWLPIEQYAGYSASKAAVSALTRAAALSCRKQGYAIRVN SIHPDGIYTPMMQASLPKGVSKEMVLHDPKLNRAGRAYMPERIAQLVLFLASDESSVMSG SELHADNSILGMGL |

(Backprop,CNN,LAMSTAR) the objective is to have an accuracy higher than 95% minimizing the execution time that is notoriously the bottle neck of most ANN applications.

The 4 output neurons represent the 4 chosen families using the one-hot encoding, and a result is considered correct only is each neuron has an absolute error less than 0.1.

The network is initialized with random weights between -0.5 and 0.5 and the learning rate set at 1, this decision is the consequence of the tests presented in the RESULTS section. The backprop-agation algorithm runs on its stochastic (or on-line) version of the gradient descendent algorithm which means that weights are updated at each iteration, this solution has been used instead of batch or minibatch learning because it is well suited for the problem domain and leads to a faster convergence.

## VII. CONVOLUTIONAL NEURAL NETWORK

The convolutional neural network has the same output of the backprop and the same input, although the input has been reshaped to assume the matrix form needed by the CNN.

The network structure is shown in 3: the input is a 23x23 matrix of bigram frequencies, then the convolutional layer produces 2 feature maps by convolving the input with 2 kernels (filters) of 7x7 weights, this leads to 2 maps of 17x17 neurons. The convolution is implemented with the formula:

$$(f * g)[n] = 1 \sum_{m=-M}^{M} f[n-m]g[m]$$

The result of the convolution is used as input for the pooling layer which is implemented as a SpatialMaxPooling accordingly to the torch documentation[11], this produces two matrix of 8x8 neurons, then the 2x8x8=128 neurons are reshaped to a vector form and fully connected to the output layer which is composed by 4 neurons. **Convolution operation:** The neurons of a convolutional layer are linked to the previous ones through an operation of convolution. In contrast with the neurons of a fully connected layers, that can be disposed in vectors, neurons of convolutional layers are disposed in matrices. In order to compute the convolution, a matrix kernel (the weights of the convolution) is put over an area of

```
(A, A)    0.781250
(A, B)    0.000000
(A, C)    0.390625
(A, D)    0.000000
(A, E)    1.562500
(A, F)    0.000000
(A, G)    0.781250
(A, H)    0.390625
(A, I)    2.734375
(A, K)    0.390625
(A, L)    0.000000
(A, M)    0.000000
(A, N)    0.000000
(A, P)    0.390625
(A, Q)    0.390625
(A, R)    0.000000
(A, S)    0.000000
(A, T)    0.000000
(A, V)    0.781250
(A, W)    0.000000
(A, X)    0.000000
(A, Y)    0.390625
(A, Z)    0.000000
(B, A)    0.000000
(B, B)    0.000000
(B, C)    0.000000
(B, D)    0.000000
(B, E)    0.000000
(B, F)    0.000000
(B, G)    0.000000
          ...
(Y, S)    0.390625
(Y, T)    0.781250
(Y, V)    0.000000
(Y, W)    0.000000
(Y, X)    0.000000
(Y, Y)    0.000000
(Y, Z)    0.000000
(Z, A)    0.000000
(Z, B)    0.000000
(Z, C)    0.000000
(Z, D)    0.000000
(Z, E)    0.000000
(Z, F)    0.000000
(Z, G)    0.000000
(Z, H)    0.000000
(Z, I)    0.000000
(Z, K)    0.000000
(Z, L)    0.000000
(Z, M)    0.000000
(Z, N)    0.000000
(Z, P)    0.000000
(Z, Q)    0.000000
(Z, R)    0.000000
(Z, S)    0.000000
(Z, T)    0.000000
(Z, V)    0.000000
(Z, W)    0.000000
(Z, X)    0.000000
(Z, Y)    0.000000
(Z, Z)    0.000000
```
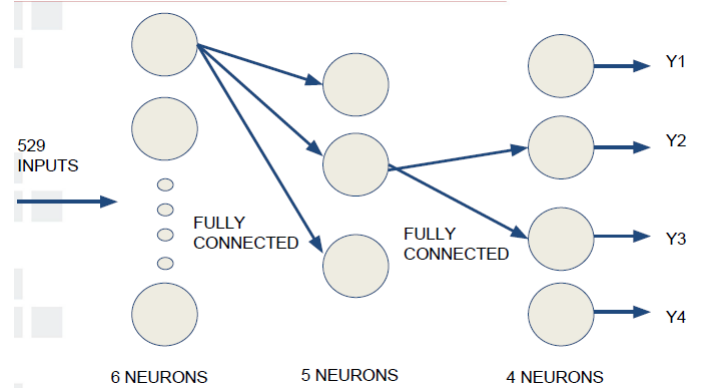
Fig. 1. Frequency bigrams for protein P00334.



Fig. 2. BackPropagation network structure.

neurons of the previous layer and each input is multiplied with the corresponding kernel weight (the filter) to generate the output, this operations is shown in figure 4.

**Pooling operation:** The pooling used is the max-pooling which means that for each submatrix (2x2) of the input layer the pooling layer extracts only one value that is the maximum of the 4, in this way the matrix dimension is reduced.

The network has been trained with a learning rate = 0.01 during 1 epoch those values have been selected after computing the performances for the learning rate in the range [0.01,0.05] and epochs between [1,5]. The number of filters (convolution kernels) has been set to 2 because it was the minimum number that led to high accuracy.

The weights and bias of the network are initialized ramdomly by the framework used : Torch[5]. The weights and biases are then trained using the back-propagation algorithm (extended for convolutoinal and subsampling layers) and accordingly to the Minimum Squared Error (MSE) error function. Given the layer described in the previous section we can count the total number of weights and biases, we will se how the fact that the convolution shares weights will lead to a small number of weights to be trained: $2*7*7+2*8*8 = 226$ which is much less than the 3224 needed for backpropagation.

## VIII. LAMSTAR

The LAMSTAR[12]

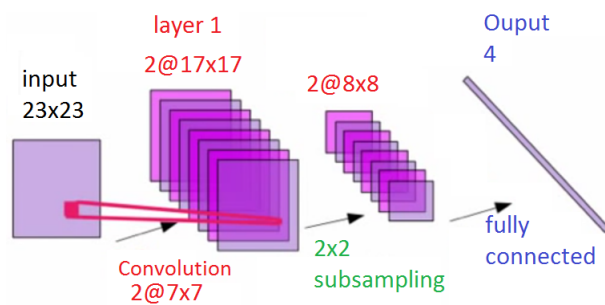Fig. 3. CNN network structure.



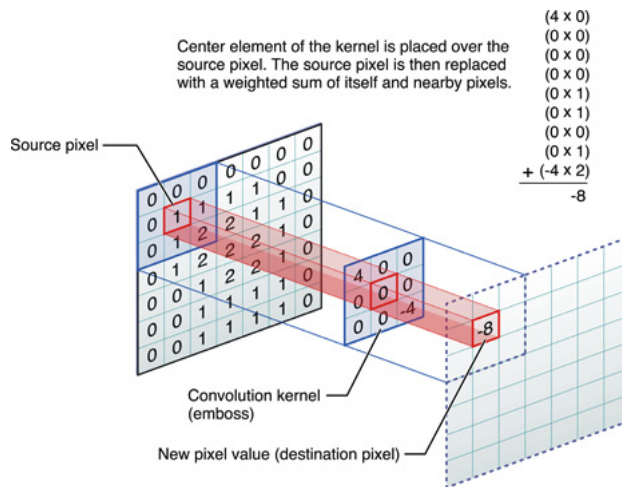Fig. 4. Convolution operation example.

## IX. RESULTS

## X. CONCLUSION

## XI. CODE

The code: R scripts and python scripts for the data collection and preprocessing, and Java for the ANN coding is available at https://github.com/umbertoDifa/protein-family-classification-ann

## REFERENCES

[1] D.Graupe, Principles of Artificial Neural Networks, 3nd ed., Advanced Series in Circuits and System-Vol.7

[2] Efficient Feature Selection and Classification of Protein Sequence Data in Bioinformatics,Muhammad Javed Iqbal et al., 2014

[3] https://www.genome.gov/12011238

[4] Basic local alignment tool,S. Altschul et al., 1990

[5] Protein classification artificial neural system, C. H. Wu et al, 1992

[6] Neural Networks for Full-Scale Protein Sequence, C. H. Wu et al, 1995

[7] Motif identification neural design for rapid and sensitive protein family search, C. H. Wu et al,1996

[8] Multi-class Protein Sequence Classification Using Fuzzy ARTMAP,Shakir Mohamed et al.,2006

[9] Neural Networks and Genome Informatics

[10] https://www.kimonolabs.com/

[11] https://github.com/torch/nn/blob/master/doc/convolution.md#nn.SpatialMaxPooling

[12] Large memory storage and retrieval (LAMSTAR) network, D.Graupe,1999