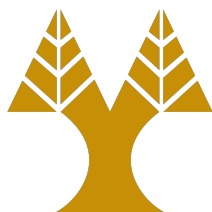


Thesis Dissertation

**PROTEIN SECONDARY STRUCTURE PREDICTION
USING CONVOLUTIONAL NEURAL NETWORKS
AND HESSIAN FREE OPTIMISATION**

Panayiotis Leontiou

UNIVERSITY OF CYPRUS



COMPUTER SCIENCE DEPARTMENT

May 2020

UNIVERSITY OF CYPRUS
COMPUTER SCIENCE DEPARTMENT

**Protein Secondary Structure Prediction Using Convolutional Neural
Networks And Hessian Free Optimisation**

Panayiotis Leontiou

Supervisor

Dr. Chris Christodoulou

Thesis submitted in partial fulfilment of the requirements for the award of
Bachelor degree in Computer Science at University of Cyprus

May 2020

Acknowledgments

First of all, I would like to express my opinion on academic research. Academic research is of major importance, especially thesis dissertation projects like this one, because students are given the opportunity to attempt something they have not tried before, discover something new. Even if their attempt is not successful it does not really matter, what matters is to learn how to research and analyse information. The research and analytical skills could be very valuable, not only for those who are aiming for an academic career, but also for those who are pursuing an industry career.

So far we had our teachers or professors to guide us through this world full of knowledge. Now, we have to guide ourselves through this maze of information and decide what to learn, in order to improve ourselves. As Albert Einstein once said, "Once you stop learning, you start dying", and one of the best ways to keep learning is through research. Research forces you to look for information from different sources and combine them with your way of thinking to reach some conclusions. Usually some experiments are required, and whether those are for a new evolutionary algorithm or a new pancake recipe is based on our own priorities. The first one could change the entire world, while the second could change the way your friends and family see your cooking skills. What both have in common is the continuous effort for personal improvement.

At this point, I would like to thank my advisor Dr. Chris Christodoulou, not only for his support on my related research, his kindness and motivation, but also for giving me the opportunity to work on a very interesting problem, using machine learning techniques. His guidance played a major role for the completion of this thesis dissertation project. If I could go back in time and choose a different topic or advisor I would choose the same, as they both helped me learn a lot of new things, that I would have not learned otherwise.

I would also like to thank Dr. Michalis Agathokleous and the master student Andreas Dionysiou for providing me with all the necessary data files and additional implementation advice for my project, based on their own experience on this machine learning problem.

Finally, I would like to thank my family for the continuous support, no matter what decisions I take. Even if they could not help me directly with this project, they helped me indirectly with their love and exceptional cooking skills.

Abstract

This dissertation attempts to solve the protein secondary structure prediction problem, a topic that has been concerning both Computer Science and Biology fields for decades.

Proteins are highly complex substances which are included in all living organisms. Proteins are not only of great nutritional value but are also involved in the chemical processes essential for life. The study of protein structures and functions can contribute to improved food supplements, drugs and antibiotics. In addition, the study of existing proteins could possibly help treat diseases and solve numerous biological problems, like covid-19 which, at the moment of writing, threatens human life on earth.

Even though there is a lot of information about the primary structure of millions of proteins, for most of them there is no information about their secondary or tertiary structure. The reason behind that is the extremely high cost, in both money and time, of the current state-of-the-art methods and instruments that are used for protein structure determination. As a result, computational algorithms and techniques, which are cheaper and faster, are essential for predicting the secondary and tertiary structures of proteins.

In the past, there were several attempts to solve the PSSP problem with Convolutional Neural Network (CNNs) and some of them managed to achieve very good results, 81% per residue Q3 accuracy [1]. Furthermore, an attempt with a simple Feed Forward Neural Network (FFNN), trained with the Hessian Free Optimisation (HFO) algorithm, managed to reach 80.4% Q3 accuracy [2]. These results are very close to the best results reported so far for the PSSP problem (84-85%), and the combination of these techniques was the motivation behind this dissertation project.

For the purpose of this dissertation, a CNN was trained with a variation of the HFO algorithm to predict the secondary structure of proteins (PSSP), which has never been attempted before. The original HFO algorithm could not be used, because of the complex structure of CNNs, instead a variation, known as the Subsampled Hessian Newton (SHN) method [3], was used. The results of this combination, for the CB513 dataset, were an overall per residue Q3 accuracy of 78.20% for a single fold and 81.80% for 10-fold cross-validation with ensembles, random forest and external rules filtering, while the SOV score was 75.67 and 78.98, respectively. Moreover, the SHN method did not require much tuning of the hyper parameters, which made the training process much faster compared to other state-of-the-art methods. As regards the PISCES dataset, the Q3 accuracy was 79.88% for a single fold and 83.02% for 5-fold cross-validation with ensembles, random forest and external rules filtering, while the SOV score was 76.67 and 82.64, respectively.

Contents

1	Introduction	1
1.1	Protein Secondary Structure Prediction problem	2
1.2	The Importance of PSSP	3
1.3	Previous Research on PSSP	4
2	Background	10
2.1	Biology Background	11
2.1.1	The Biological Role of Proteins	11
2.1.2	Amino Acids	11
2.1.3	Protein Structures	16
2.1.3.1	Primary Structure	16
2.1.3.2	Secondary Structure	19
2.1.3.3	Tertiary Structure	19
2.1.3.4	Quaternary Structure	20
2.2	Artificial Neural Networks Background	21
2.2.1	Origins of Artificial Neural Networks	21
2.2.2	Variations of Artificial Neural Networks and Optimizers	22
2.2.2.1	McCulloch and Pitts (McP)	22
2.2.2.2	Multi-Layer Perceptron (MLP)	26
2.2.2.3	Recurrent Neural Network (RNN)	30
2.2.2.4	Convolutional Neural Network (CNN)	30
2.2.2.5	Line Search	33
2.2.2.6	Conjugate Gradient (CG)	34
2.2.2.7	Newton's Method	35
2.3	Hessian Free Optimisation (HFO)	38
2.3.1	Intro to HFO	38
2.3.2	Analysis of HFO	40
2.3.3	Hessian-Vector Multiplication evaluation	40

3	Data Manipulation	43
3.1	PSSP Metrics	44
3.2	Protein Databases and DSSP	45
3.3	Dataset Format	46
3.4	Data Encoding and MSA profiles	47
3.5	CB513 and PISCES Datasets	48
3.6	Dataset preprocessing with MSA profiles	49
3.7	Significant neighboring amino acids	51
3.8	Training/ Testing Set and Cross Validation	52
3.9	Ensembles	53
3.10	Filtering	54
3.10.1	External Rules	54
3.10.2	Support Vector Machines	55
3.10.3	Decision Trees	57
3.10.4	Random Forests	58
4	Implementation	61
4.1	A new approach for the PSSP problem	62
4.2	CNN and HFO combination	63
4.3	Subsampled Hessian Newton (SHN) Method	64
4.4	Network Implementation	66
5	Experiments and Results	67
5.1	Experiments for Implementation Evaluation	68
5.2	Experiments with CB513 dataset	70
5.2.1	Fine Tuning of Hyper Parameters	70
5.2.2	10-fold Cross-Validation and Ensembles Results	72
5.2.3	CNN and SVM Combination	74
5.2.4	Filtering Results for CB513	75
5.2.5	Additional experiments with CB513	79
5.2.6	Final results for CB513	79
5.3	Experiments with PISCES dataset	80
5.3.1	5-fold Cross-Validation and Ensembles Results	80
5.3.2	Filtering Results for PISCES	81
5.3.3	Final Results for PISCES	83
5.4	Best Results for CB513 and PISCES	84

6 Conclusion and Future Work	87
6.1 Conclusions	88
6.2 Suggestions for Future Work on PSSP	90
References	93
Appendices	A-1
Appendix A Excluded proteins from CB513	A-1
Appendix B Excluded proteins from PISCES	B-1
Appendix C Convert datasets to Matlab files	C-1
Appendix D CB513 dataset pre-processing	D-1
Appendix E PISCES dataset pre-processing	E-1
Appendix F Python Implementation	F-1
Appendix G Ensembles Program	G-1
Appendix H External Rules Program	H-1
Appendix I SOV calculation	I-1
Appendix J Calculation of Q3 accuracy	J-1
Appendix K Data pre-processing for filtering	K-1
Appendix L Training Filtering Methods	L-1
Appendix M All filtering methods on CB513	M-1
Appendix N View filtering results of CB513	N-1
Appendix O All filtering methods on PISCES	O-1
Appendix P View filtering results of PISCES	P-1

List of Figures

1.1	Number of publications for PSSP per year [6].	5
2.1	The structure of amino acids [27].	12
2.2	The 20 standard amino acids [28].	13
2.3	An example for condensation reaction [29].	14
2.4	The Central Dogma of Molecular Biology: DNA makes RNA makes protein [30].	14
2.5	Example of the central dogma, which illustrates the first few amino acids for the alpha subunit of hemoglobin [23].	15
2.6	The amino acids specified by each codon [31].	15
2.7	All four protein structures.	17
2.8	The first amino acid of the A chain is glycine (Gly), whereas, the last is asparagine (Asn) [32].	18
2.9	The diagram shows the substitution in a small part of the hemoglobin β chain, where the amino acid at position seven, glutamate, is replaced by valine, in the sickle cell hemoglobin [32].	18
2.10	The diagram illustrates the shapes of the two main types of the secondary structure of proteins, the α -helix and the β -strand [32].	19
2.11	The diagram indicates some of the chemical interactions that determine the proteins' tertiary structure [32].	20
2.12	Structure of a Biological Neuron [33].	21
2.13	McCulloch and Pitts artificial neuron [35].	23
2.14	The step or heaviside function.	23
2.15	Decision lines for AND gate (a) and OR gate (b).	24
2.16	The OR gate is linearly separable while the XOR gate is not.	25
2.17	Multi-Layer Perceptron Neural Network with one hidden layer.	26
2.18	Decision regions based on the number of hidden layers.	27
2.19	RNN variations, Jordan network (left), Elman network (right).	30
2.20	A CNN example for digit image classification.	32
2.21	Example of max and average pooling.	32

2.22	Example of zero padding.	33
2.23	Gradient Descent (left) vs Conjugate Gradient (right) on a 2D problem. .	34
2.24	Newton's method in a first degree polynomial problem [23].	36
2.25	Local Quadratic approximations [23].	37
2.26	The Hessian matrix of the error function with respect to the weights. . . .	38
3.1	Protein representation example for protein 1bdsA_1-43.	46
3.2	Process of MSA profiling	47
3.3	Example of the encoded form of an MSA file [24].	49
3.4	The image shows the MSA file (before collapsing into a single file) [24]. .	50
3.5	The encoding of the new file, after combining the MSA files into a single file [24].	50
3.6	MSA record for a sequence of 6 amino acids.	51
3.7	Modified MSA record for a sequence of 6 amino acids.	51
3.8	An example of input data representation for a window size of 15 (or k = 7) amino acids [57].	52
3.9	10-fold cross validation	53
3.10	Results of different kernels for a 3-class classification problem.	55
3.11	SVM example of a linearly separable problem.	56
3.12	SVM projecting a problem in a higher dimension.	56
3.13	Example of simple decision tree [60].	57
3.14	Example of random forest prediction [60].	58
3.15	Distribution of the outcomes of 10000 simulations for each option [60]. .	59
3.16	Node splitting in a decision tree and a random forest model [60].	60
5.1	The test loss for each iteration compared to the test loss of the saved model.	69
5.2	The test Q3 accuracy after each iteration compared to the test Q3 accuracy of the saved model.	69
5.3	Hyper parameters and methods used that resulted in the best overall Q3 accuracy and best overall SOV score for CB513 dataset.	84
5.4	CM for CB513 fold 0 of single CNN model.	85
5.5	CM for CB513 fold 0 of ensembles model with random forest.	85
5.6	Hyper parameters and methods used that resulted in the best overall Q3 accuracy and best overall SOV score for PISCES dataset.	85
5.7	CM for PISCES fold 4 of single CNN model.	86
5.8	CM for PISCES fold 4 of ensembles model with random forest.	86

List of Tables

1.1	Methods used for PSSP in chronological order.	9
2.1	Types of proteins and their function [26].	12
2.2	Truth table for OR gate.	24
2.3	List of the most popular activation functions.	27
2.4	Derivatives and Hessians of typical loss function.	42
3.1	Table with the secondary structure abbreviations, grouped in 8 and 3 classes	45
4.1	Previous studies on Newton methods [3].	65
5.1	Hyper parameters for CNN for all experiments.	70
5.2	Q3 accuracy results for GNsize for fold 5 of CB513.	71
5.3	Tuning the C hyper parameter for fold 5 of CB513.	71
5.4	Tuning the batch size (bsize) hyper parameter for fold 5 of CB513.	72
5.5	Hyper parameters for trained models.	72
5.6	Q3 and SOV results for 10-fold cross validation for the CB513 dataset. . .	73
5.7	Q3 and SOV results for ensembles (with 5 experiments per fold) cross validation for the CB513 dataset.	73
5.8	Hyper parameters for SVM filtering.	74
5.9	Hyper parameters for Random Forest filtering.	75
5.10	Q3 accuracy and SOV score for ensembles (with 5 executions per fold) and external rules filtering for CB513 dataset.	75
5.11	Q3 accuracy and SOV score for ensembles (with 5 executions per fold), external rules and SVM filtering for CB513 dataset.	76
5.12	Q3 accuracy and SOV score for ensembles and SVM filtering for CB513 dataset.	76
5.13	Q3 accuracy and SOV score for ensembles, SVM and external rules fil- tering for CB513 dataset.	76
5.14	Q3 accuracy and SOV score for ensembles, external rules and decision tree filtering for CB513 dataset.	77

5.15	Q3 accuracy and SOV score for ensembles and decision tree filtering for CB513 dataset.	77
5.16	Q3 accuracy and SOV score for ensembles, decision tree and external rules filtering for CB513 dataset.	78
5.17	Q3 accuracy and SOV score for ensembles, external rules and random forest filtering for CB513 dataset.	78
5.18	Q3 accuracy and SOV score for ensembles and random forest filtering for CB513 dataset.	78
5.19	Q3 accuracy and SOV score for ensembles, random forest and external rules filtering for CB513 dataset.	78
5.20	Results for fold 0 of CB513 with the ensembles method applied before and after the filtering methods.	79
5.21	10-fold Cross validation, Q3 accuracy and SOV score for all methods for CB513 dataset.	79
5.22	Hyper parameters for SHN method, used for all PISCES experiments. . .	80
5.23	Q3 accuracy and SOV score for 5-fold cross validation for PISCES dataset.	80
5.24	Q3 accuracy and SOV score for ensembles method (with 5 trained models per fold) for PISCES dataset.	81
5.25	Q3 accuracy and SOV score for ensembles with external rules filtering for PISCES dataset.	81
5.26	Q3 accuracy and SOV score for ensembles with external rules and decision tree filtering for PISCES dataset.	82
5.27	Q3 accuracy and SOV score for ensembles with decision tree filtering for PISCES dataset.	82
5.28	Q3 accuracy and SOV score for ensembles with decision tree and external rules filtering for PISCES dataset.	82
5.29	Q3 accuracy and SOV score for ensembles with external rules and random forest filtering for PISCES dataset.	83
5.30	Q3 accuracy and SOV score for ensembles with random forest filtering for PISCES dataset.	83
5.31	Q3 accuracy and SOV score for ensembles with random forest and external rules filtering for PISCES dataset.	83
5.32	5-fold cross-validation, Q3 accuracy and SOV score for all methods for the PISCES dataset.	84
A.1	Excluded CB513 proteins due to zeroed MSA profiles.	A-1
B.1	Excluded PISCES proteins due to missing MSA profiles.	B-1
B.2	Excluded PISCES proteins due to corrupted or zeroed MSA profiles (1-120).	B-2

B.3	Excluded PISCES proteins due to corrupted or zeroed MSA profiles (121-240).	B-3
B.4	Excluded PISCES proteins due to corrupted or zeroed MSA profiles (241-341).	B-4

Chapter 1

Introduction

1.1	Protein Secondary Structure Prediction problem	2
1.2	The Importance of PSSP	3
1.3	Previous Research on PSSP	4

1.1 Protein Secondary Structure Prediction problem

Proteins are highly complex substances which are present in all living organisms. There are over 30,000 unique proteins in the human body, which are responsible for performing specific functions that are essential for life. The word protein is derived from the Greek word 'πρωτεϊος', which means 'of the first quality', 'in the lead' or 'holding first place', and their significance was recognised in the early 19th century by chemists.

These substances consist of smaller units, called amino-acids, which are organic compounds connected to each other, forming long chains. The differences between two proteins are based on their sequence of amino acids, which determines their structure and function. The interactions between the amino acids of a protein are responsible for the fold of the protein into a specific three-dimensional structure, which, under specific conditions, remains the same. This structure determines the function of each protein.

The study of protein structures and functions can contribute to improved food supplements, drugs and antibiotics. In addition, the study of existing proteins can help treat diseases and solve numerous biological problems with the help of modern technology, which is significantly cheaper and more efficient than a few years ago.

A hierarchical approach has been established for analysing the structure of proteins more effectively and observe their different forms. These forms are separated into four distinct categories, the primary, the secondary, the tertiary and the quaternary structure. The primary structure is just a linear sequence of amino acids, that are ordered based on where they appear in the unfolded protein. The secondary structure illustrates how the local parts of a protein are organised in a two-dimensional space. The tertiary structure, which determines the specific function of a protein, has a three-dimensional shape, formed when the amino acid chain is folded. Finally, the quaternary structure is formed when multiple tertiary structures are folded together and also has three-dimensional shape.

Even though the primary structure for millions of proteins is well documented, for most of them the secondary and tertiary structures are unknown, only for a small fraction of these proteins the secondary and tertiary structures are currently available. The research and the experimental determination of the secondary and tertiary structures of a protein are not only time consuming but also an extravagant process. More specifically, in order to determine the tertiary structure of proteins, expensive and tedious methods must be used, such as X-ray crystallography and nuclear magnetic resonance (NMR). The shape of a protein is completely determined by its primary structure, about 70% of the secondary structure is affected by the interactions of the nearby amino acids of the backbone, while the other

30% is affected by more distant interactions [4]. This made prediction techniques and implementations more appealing over the experimental methods, since they have high success rates on the prediction of secondary and tertiary structure of proteins, they cost significantly less and require considerably less time than the experimental methods.

One such prediction method is ab initio prediction, which tries to predict any of the three structures based only on the primary structure and without taking into consideration any patterns. This method is divided into two distinct cases. In the first case, the folding process is simulated or minimisation of the free energy of the polypeptide is attempted, and only the primary structure of the protein is used (no other known structures). On the other hand the second case attempts to predict the structure of a protein using already known and existing protein structures [4]. This thesis is concentrated entirely on the second prediction method, and more specifically on the use of Neural Networks (NN) to predict the secondary structure of proteins. These algorithms are designed based on computational statistics and mathematical optimization techniques. These optimisation techniques help computer systems learn hidden patterns and idiosyncrasies of data, which then gives them the ability to predict and classify new data.

To sum up, because of the extreme costs in both money and time of experimental methods, it is not possible to experimentally determine the structure of all proteins. In this thesis Convolutional Neural Networks (CNN) will be used in combination with the Hessian Free Optimisation (HFO) algorithm in order to predict the secondary structure of proteins.

1.2 The Importance of PSSP

The solution of the PSSP problem is very important because the secondary structure is essential in order to determine the tertiary structure, which gives information about the functions of a protein. The experimental methods used for determining the tertiary structure of proteins are extremely expensive in both time and money, which led to the study of just a small portion of known proteins. As a result, the scientific community has information about the functions of just a small subset (a few thousands) of proteins, compared to the millions of proteins that exist.

Furthermore, this means that the PSSP can help identify the tertiary structure of a protein with higher accuracy and less effort. It is very important to note that the functions of a protein are based on the 20 amino acids that compose a protein, which is the main reason why the research in this field is very important. Understanding how these molecules fold around space, assemble and function can help to understand why people are getting older,

why they suffer from dangerous diseases and viruses (such as cancer), how can a cure for a disease be found (like the cure for covid-19), and other ‘difficult to answer’ questions.

The proteins’ functions are related with their structure, which depends on both the physical and chemical parameters of these molecules. Bioinformatics is an interdisciplinary field that develops methods and software tools for understanding biological data. It combines knowledge from biology, computer science, information engineering, mathematics and statistics to analyse and interpret biological data.

1.3 Previous Research on PSSP

Researchers from different fields have been working on this problem for more than six decades. A wide variety of machine learning algorithms have been designed specifically for this problem and have achieved accuracy $>90\%$ [5], based on the Q3 score (Equation 1.1.). Additional structural templates from databases, which are called sequence-based structural similarity of proteins, were used in order to achieve accuracy higher than 88%. The additional information boosts significantly the learning process as well as the performance of these algorithms. The three-state accuracy for machine learning algorithms, that are not relying on the structural templates, is currently around 82-85%, which is good for such a complex problem. However, considering the theoretical limit of the three-state prediction which is around 88-90%, there is still room for improvement.

Figure 1.1 shows the number of publications per year for the PSSP problem as well as the cumulative number of publications, between 1973 and 2015. According to the graph the cumulative number of publications for the PSSP problem increased significantly between 1973 and 2015. More specifically, between 1973 and 1989 there were less than 5 publications for the PSSP problem per year. In 1990, the PSSP problem started to become more popular and the number of publications increased considerably to 8, while the cumulative number of publications was around 50. During the next two decades, the PSSP problem gained much popularity, probably because in that period there were some major breakthroughs, which helped to increase the three-state accuracy considerably. The popularity of PSSP dropped substantially in 2010 and for the following 5 years the interest for this problem was relatively moderate. A small selection of PSSP publications are mentioned below.

Feedforward Fully Connected Neural Network (FFNN) [7]: A fully connected Neural Network with local input window (usually of 13 amino acids with orthogonal encoding) and just one hidden layer. The output of the network was one of the three categories

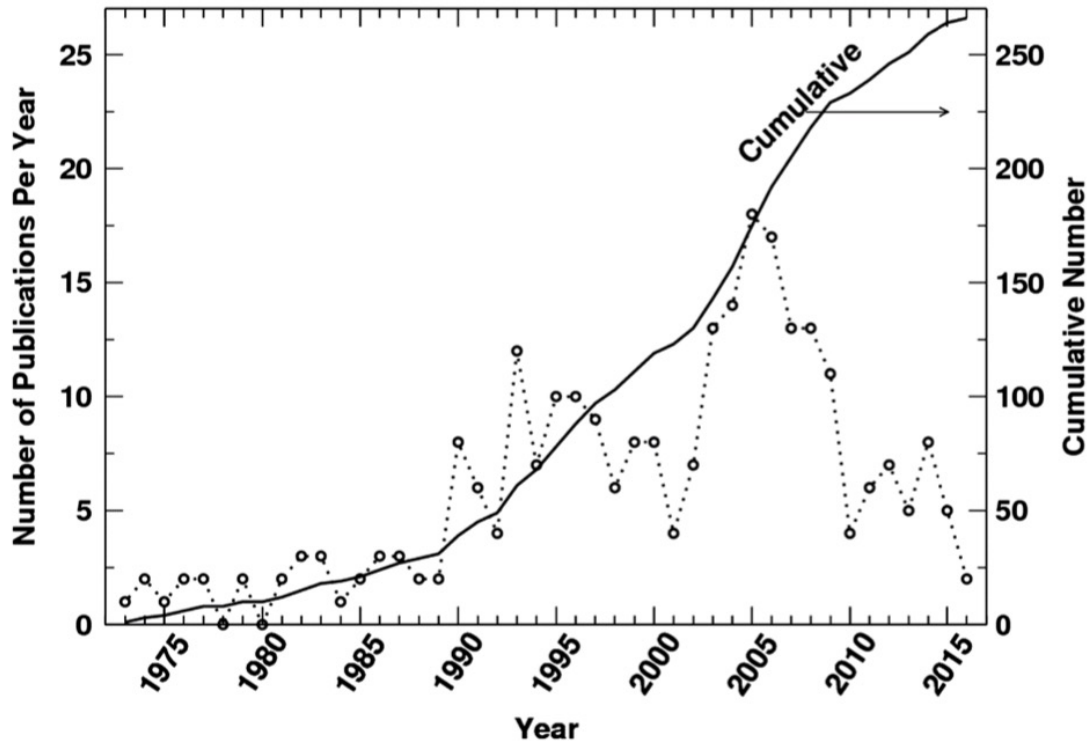


Figure 1.1: Number of publications for PSSP per year [6].

of the secondary structure of proteins (helix, pleated or other) based on the amino acid located in the centre of the input window. A secondary network was also used in this implementation to improve the output of the previous network. This method had issues with overfitting.

PHD: predicting 1D protein structure by profile based neural networks [8]: The structure of the network was the same with the Feed Forward Fully Connected Neural Network of Qian and Sejnowski [7], with the addition of techniques that deal with the overfitting problem. Two methods were used to counter overfitting, early stopping (terminating the training process before it starts to overfit) and ensemble average (training different networks at the same time with different data and learning methods). Furthermore, the multiple alignment technique was used in the input data, to take advantage of evolutionary information.

Gene-finding Programs (NNSSP) [9]: This Neural Network uses the ‘nearest neighbour’ method to group the sequences of amino acids based on their similarities and compare them with other sequences, that their secondary structure is known. Following that, the network tries to predict the secondary structure of other proteins that their secondary structure is not known.

Discrimination of Secondary structure Class (DSC) [10]: This algorithm groups the output data of the network and by using simple linear static methods attempts to predict the secondary structure of proteins.

PREDATOR [11]: It was implemented in a Neural Network which takes as input a sequence of amino acids and tries to predict the secondary structure based on possible hydrogen bonds that may exist in the output sequence.

Consensus [12]: In this method a Neural Network was used that took as input the multiple alignment with additional information about the protein (rather than just a simple sequence of amino acids). This Network attempts to locate similarities between the input sequence with other amino-acid sequences (similarities in genetic code, evolutionary history and common biological functions) in order to predict the secondary structure.

Bidirectional Recurrent Neural Network (BRNN) - Backpropagation ([13], [14]): This algorithm uses a Neural Network that takes as input a window with a sequence of amino acids and attempts to predict the secondary structure of the amino acid located in the centre of the input window, based on the amino-acids that precede and follow it in the input chain using bidirectional recursion. It is important to note that this algorithm had some of the best results in the PSSP problem at the time it was conceived, with 76% success rate.

Logical Analysis of Data (LAD) [15]: This method, which uses a machine learning algorithm, was implemented to identify properties of amino acids, and therefore, additional information about the homogeneity of proteins, which could help the prediction of the secondary structure of proteins. According to this method, the most important property that affects the helix class is molecular weights, for the pleated class is the mean ambient hydrophobicity, while for the other forms is the polarity.

Multiagent Secondary Structure Predictor with Postprocessing (MASSP3) [16]: This implementation attempts to solve this problem by using two distinct sections. The first section is based on a hybrid structure, which combines genetic and neural techniques, while the second section consists of a Multilayer Perceptron (MLP), which takes as input the output of the first section. The results of this method were fairly good.

Two-Stage method [17]: This approach uses two stages, the first identifies instabilities in how the protein folds into space and attempts to classify the different parts of the protein, while the second splits the proteins into sequences (3 to 7 residues) and tries to predict the secondary structure of these sequences.

Evolutionary method for learning HMM structure [18]: In this research genetic algorithms were used, which can dynamically change the parameters of a Hidden Markov Model (HMM) (since the construction of a HMM is very complicated) and build it dynamically, so that it can predict the secondary structure of the input sequences.

Cascade Bidirectional Recurrent Neural Network (BRNN) [19]: This implementation focused on the long range dependencies between the input data, which plays a major role in the folding of a protein and the correlation between adjacent secondary structures. In this article, the authors refer to the correlation of the secondary structure of an amino acid as regards to secondary structure of the adjacent amino acids. Two BRNNs are used, with the second taking its input from the output of the first BRNN. This method, although, it had relatively good results, could not outperform previous approaches.

Protein Secondary Structure Prediction Using Deep Convolutional Neural Fields [20]: This approach used a Deep Convolutional Neural Fields (DeepCNF), which is an extension of Deep Learning to Conditional Neural Field (CNF) (a combination of Conditional Random Fields (CRF) and shallow neural networks). The DeepCNF is much more powerful than the CNF, since it can model both the complex sequence-structure relationship (from a deep hierarchical architecture) and the interdependence between adjacent secondary structure tags. Based on the experimental results, the DeepCNF can reach prediction accuracy of about 84%, using the protein datasets CASP and CAMEO, surpassing existing methods of predicting the secondary structure of proteins. The DeepCNF networks can also be used to predict other properties of proteins, such as contact number, solvent accessibility and disorder regions.

Protein Secondary Structure Prediction with the use of Convolutional Neural Networks for Image Object Recognition [21]: The purpose of this research was to identify how Convolutional Neural Networks (CNN) can help in solving the PSSP problem. These type of networks take advantage of the spatial structure of the input data, which seems very promising. Furthermore, they manage input data of problems with sequences or problems that use the parameter of space, better, like image processing. This method could only reach an accuracy of about 40%, because there were problems in the representation of input data in the CNN, which prevented the network from learning effectively.

Capturing non-local interactions by long short-term memory bidirectional recurrent neural networks for improving prediction of protein secondary structure, backbone angles, contact numbers and solvent accessibility [22]: Unlike other methods that try to capture short to intermediate interactions between amino acid residues, this approach used Long Short-Term Memory (LSTM) Bidirectional Recurrent Neural Networks (BRNNs) to cap-

ture long range interactions. This method reported some of the best results so far with approximately 84% Q3 accuracy.

Protein Secondary Structure Prediction Using Bidirectional Recurrent Neural Networks (BRNN) and Hessian Free Optimisation (HFO) ([23], [2]): This dissertation was undertaken by a past Computer Science student of University of Cyprus in the context of his diplomatic research. This dissertation showed that simple Feed Forward Neural Networks (FFNNs) can be trained with the powerful second-order learning algorithm, Hessian Free Optimisation (HFO), to predict the secondary structure of proteins. This approach (FFNN with HFO) had very good results as regards the training time of the network and (Q3) accuracy, which was about 80.4% (using the PISCES dataset). The HFO does not require much tuning of the hyper parameters, which makes training much faster than other state of the art methods. The use of HFO seems very promising since it reduces the training time of the network and at the same time offers very good results.

Prediction of Secondary Structure of Proteins using Gabor filters and Support Vector Machines ([24], [1]): This dissertation was conducted by a past Computer Science student of University of Cyprus during his diplomatic research. This thesis project, was focused on the use of Convolutional Neural Network (CNN) with Gabor Filters and Support Vector Machines (SVMs) for filtering. The combination of a CNN with SVMs had very good results with about 81% (Q3) accuracy for the PSSP problem (using the PISCES dataset). A technique was also used to convert the primary structure of proteins from one dimension into two dimensions, since the CNN needs two dimensional input data to be trained.

Sixty-five years of the long march in protein secondary structure prediction: the final stretch? [6]: This paper focused on some of the state-of-the-art methods that are used to predict the secondary structure of proteins and compared them, using the same independent test sets. The reported results ranged from 77.1% to 82.3%. The best results (82.3% Q3 accuracy) were achieved by the DeepCNF [20]. In addition, this paper mentioned alternatives to discrete three-state secondary structure prediction (with eight-state prediction) and noted that the theoretical limit of secondary structure prediction is around 88%. This limit is very close to the best results reported so far (84%), which means that it is a matter of time for the PSSP problem to reach a plateau (where there will be no further improvements in Q3 accuracy).

MUFold-SS [25]: In this research a new deep learning architecture was suggested for the PSSP problem, the Deep inception-inside-inception (Deep3I) network. This network was implemented as a software tool, named MUFOLD-SS, which takes as input a specifically designed array of data, based on the primary structure of the proteins. This array includes

information for each amino acid and general information about the protein. The structure of MUFOLD-SS allows the extraction of information related to local and general interactions, between the amino acids, which made the predictions more accurate. This tool has outperformed other techniques used on the PSSP problem, with an accuracy of approximately 86.49%.

Table 1.1 shows the Q3 accuracy of the aforementioned methods, used on the PSSP problem, in chronological order.

NO.	METHOD	Q3 ACCURACY (%)
1	Feedforward Fully Connected NN (Qian και Sejnowski, 1988)	63.30
2	PHD (Rost, 2001; Rost και Sander, 1993)	71.40
3	NNSSP (Salamov και Soloveyev, 1997)	68.41
4	DSC (King και Sternberg, 1996)	71.95
5	PREDATOR (Frishman και Argos, 1997)	68.60
6	Consensus (Cuff και Barton, 1999)	72.70
7	BRNN – Backpropagation (Baldi et al., 1999)	76.00
8	LAD (Jacek et al., 2005)	70.60
9	MASSP3 (Giuliano et al., 2005)	76.10
10	Evolutionary method for learning HMM structure (Won et al., 2007)	65.00
11	Two-Stage method (Fadime et al., 2007)	74.10
12	Cascade BRNN (Jinmiao και Narendra, 2007)	74.38
13	Deep Convolutional Neural Fields (Wang et al., 2016)	83.00
14	Convolutional Neural Networks (Pavlidis, 2016)	40.00
15	LSTM-BRNN (Heffernan et al., 2017)	84.00
16	MUFold-SS (Fang et al., 2018)	86.49
17	Feed Forward NN with HFO (Charalambous et al., 2020)	80.40
18	Convolutional Neural Network with SVM filtering (Dionysiou et al., 2020)	81.00

Table 1.1: Methods used for PSSP in chronological order.

Chapter 2

Background

2.1	Biology Background	11
2.1.1	The Biological Role of Proteins	11
2.1.2	Amino Acids	11
2.1.3	Protein Structures	16
2.1.3.1	Primary Structure	16
2.1.3.2	Secondary Structure	19
2.1.3.3	Tertiary Structure	19
2.1.3.4	Quaternary Structure	20
2.2	Artificial Neural Networks Background	21
2.2.1	Origins of Artificial Neural Networks	21
2.2.2	Variations of Artificial Neural Networks and Optimizers	22
2.2.2.1	McCulloch and Pitts (McP)	22
2.2.2.2	Multi-Layer Perceptron (MLP)	26
2.2.2.3	Recurrent Neural Network (RNN)	30
2.2.2.4	Convolutional Neural Network (CNN)	30
2.2.2.5	Line Search	33
2.2.2.6	Conjugate Gradient (CG)	34
2.2.2.7	Newton's Method	35
2.3	Hessian Free Optimisation (HFO)	38
2.3.1	Intro to HFO	38
2.3.2	Analysis of HFO	40
2.3.3	Hessian-Vector Multiplication evaluation	40

2.1 Biology Background

2.1.1 The Biological Role of Proteins

Proteins are large macromolecules or biomolecules, that perform a variety of functions within organisms. Some of these functions are deoxyribonucleic acid (DNA) replication, responding to stimuli, providing structure to cells and organisms, catalysing metabolic reactions, and transporting molecules from one location to another. Proteins consist of hundreds or even thousands of smaller units, called amino acids, which are organic compounds that contain amine (NH₂) and carboxyl (COOH) functional groups.

The consumption of food, which contains proteins, is one of the main sources of proteins for the human body. The digestive system breaks down the consumed food into amino acids, which enter the blood stream. In order to perform a variety of functions, the cells of the human body gather amino acids from the blood stream to create all the essential proteins. If there is a shortage of amino acids in the blood stream, probably because of a poor diet with less proteins, the immune system will become weak, causing dizziness, exhaustion or even serious diseases. That happens because in order to create the necessary proteins for the human body, the cells need enough amino acids, otherwise they will not be able to support the needs of the entire human body.

In order to aid in the development of food supplements, drugs and antibiotics, it is mandatory to first understand the base structure and function of each protein. Research or studies on existing proteins could help solve numerous biological problems and treat diseases. This is considerably easier nowadays, with the help of the current technology, which is faster and computationally stronger than ten years ago.

The most important functions of proteins are displayed in table 2.1 and these reveal the significance of proteins, for all living organisms.

2.1.2 Amino Acids

Amino acids are organic compounds which contain amine (NH₂) and carboxyl (COOH) functional groups. Each amino acid has its specific side chain (R group), which is an atom or group of atoms that replace one or more hydrogen atoms on the parent chain of a hydrocarbon, which turns into a moiety of the resultant new molecule (Figure 2.1). The main elements of an amino acid are carbon (C), hydrogen (H), oxygen (O) and nitrogen (N), however, other elements can also be found in the side chains of some amino acids.

Type	Function Description	Example
Defense	Defense proteins help organisms fight infection, heal damaged tissue, and evade predators.	<i>Antibodies</i>
Enzyme	Enzymes build and break down molecules. They are critical for growth, digestion, and many other processes in the cell. Without enzymes, chemical reactions would happen too slowly to sustain life.	<i>Lactase</i>
Messenger	Messenger proteins transmit signals to coordinate biological processes between different cells, tissues, and organs.	<i>Growth Hormone</i>
Motor	Motor proteins keep cells moving and changing shape. They also transport components around, inside cells.	<i>Dynein, Kinesin</i>
Regulatory	Regulatory proteins bind DNA to turn genes on and off.	<i>Androgen, Estrogen</i>
Sensory	Sensory proteins help humans learn about their environment. They help them detect light, sound, touch, smell, taste, pain, and heat.	<i>Opsin</i>
Signaling	Signaling proteins allow cells to communicate with each other.	<i>Insulin</i>
Storage	Storage proteins store nutrients and energy-rich molecules for later use.	<i>Gluten</i>
Structural	Structural proteins strengthen cells, tissues, organs, and more.	<i>Collagen</i>
Transport	Transport proteins move molecules and nutrients around the body, in and out of cells.	<i>Hemoglobin</i>

Table 2.1: Types of proteins and their function [26].

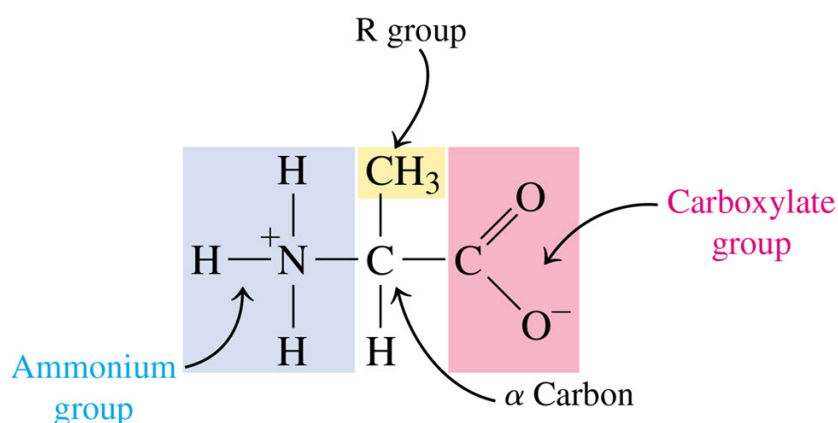


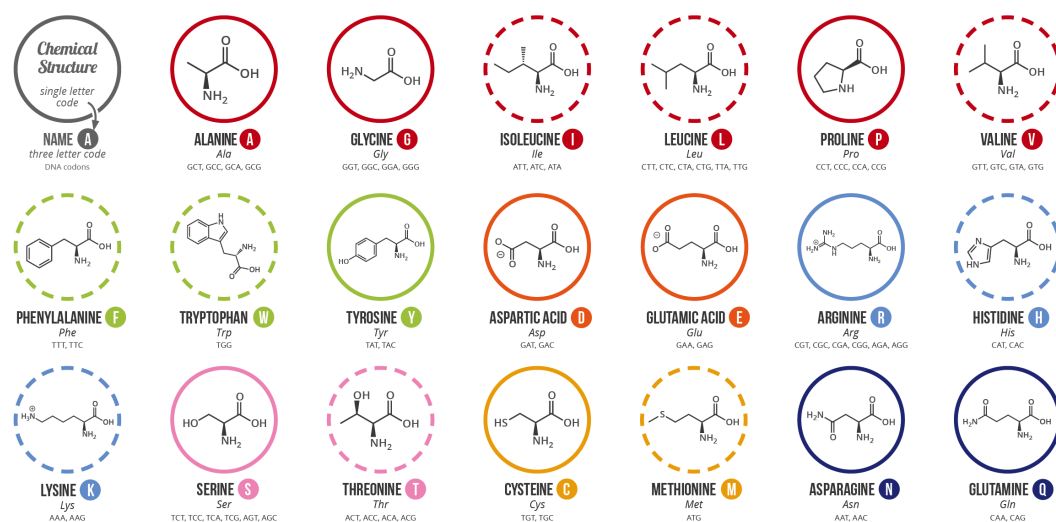
Figure 2.1: The structure of amino acids [27].

Even though there are about 500 known amino acids, only 20 of them appear in genetic code and are considered as the standard amino acids (Figure 2.2). Amino acids can be classified in many different ways, according to the core structural functional groups' locations (alpha (α), beta (β), gamma (γ), delta (δ)), based on the polarity, pH level or on the side chain group type. Amino acids also participate in a number of other processes, such as neurotransmitter transport and biosynthesis. Short chains of amino acids (30 or less) linked by peptide bonds form peptides, and long, continuous, and unbranched peptide chains form polypeptides. Proteins consist of one or more polypeptides arranged in a biologically functional way.

A GUIDE TO THE TWENTY COMMON AMINO ACIDS

AMINO ACIDS ARE THE BUILDING BLOCKS OF PROTEINS IN LIVING ORGANISMS. THERE ARE OVER 500 AMINO ACIDS FOUND IN NATURE - HOWEVER, THE HUMAN GENETIC CODE ONLY DIRECTLY ENCODES 20. 'ESSENTIAL' AMINO ACIDS MUST BE OBTAINED FROM THE DIET, WHILST NON-ESSENTIAL AMINO ACIDS CAN BE SYNTHESISED IN THE BODY.

Chart Key: ● ALIPHATIC ● AROMATIC ● ACIDIC ● BASIC ● HYDROXYLIC ● SULFUR-CONTAINING ● AMIDIC ○ NON-ESSENTIAL ○ ESSENTIAL



Note: This chart only shows those amino acids for which the human genetic code directly codes for. Selenocysteine is often referred to as the 21st amino acid, but is encoded in a special manner. In some cases, distinguishing between asparagine/aspartic acid and glutamine/glutamic acid is difficult. In these cases, the codes asx (B) and glx (Z) are respectively used.

Figure 2.2: The 20 standard amino acids [28].

The process in which chains of amino acids are linked together is called condensation reaction (Figure 2.3). During this reaction, as the amino group of one amino acid joins the carboxyl group of a neighbouring amino acid, a water molecule is extracted, what is left of each amino acid is called amino acid residue.

Each amino acid can be represented by one or three characters from the English alphabet, so it is possible to represent a sequence of amino acids using a sequence of characters. Any change in this sequence, no matter how small it is, can lead to a completely different protein, which will have its own properties and functionalities.

The proteins in an organism are assembled based on its genes, also known as the DNA.

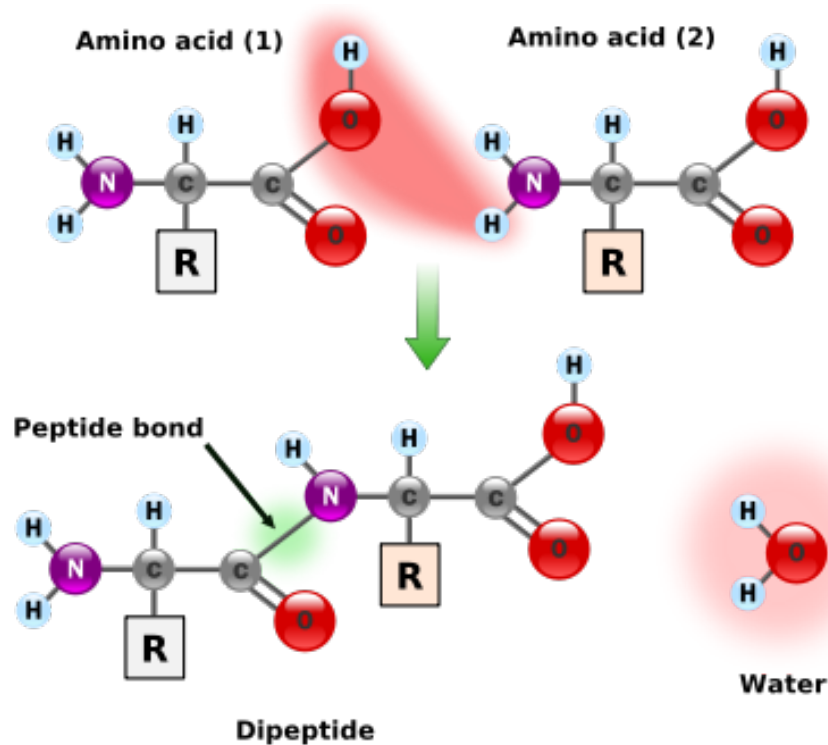


Figure 2.3: An example for condensation reaction [29].

In particular, the nucleotide sequence of a gene, which encodes a protein, specifies the unique amino acid sequence of that protein. For instance, there are around 30,000 genes in the human genome, and each one encodes one unique protein. According to The Central Dogma of Molecular Biology, the ‘DNA makes RNA and RNA makes protein’ (Figure 2.4). The first stage, ‘DNA makes RNA’, is called transcription, while the second stage, ‘RNA makes protein’, is called translation.

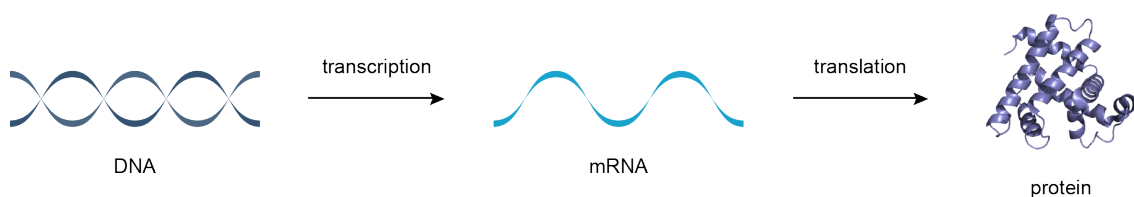


Figure 2.4: The Central Dogma of Molecular Biology: DNA makes RNA makes protein [30].

A sequence of three adjacent nucleotides composing the genetic code is called codon and designates an amino acid. There are four (4) unique nucleotides (adenine - A, uracil - U, guanine - G, and cytosine - C), which means that the maximum number of triplets that can be formed is sixty four ($4^3 = 64$). However, only twenty (20) amino acids can be encoded naturally. This means that some codons do not encode any amino acids or that

some amino acids can be described by multiple codons. Codons that do not encode any amino acids are called stop codons, because they are used as a termination signal in the translation process, signalling the release of the translated polypeptide or protein. Figure 2.5 shows an example of the translation stage, from DNA to protein, while figure 2.6 presents the table of codons, with the amino acid or the stop signal they encode.

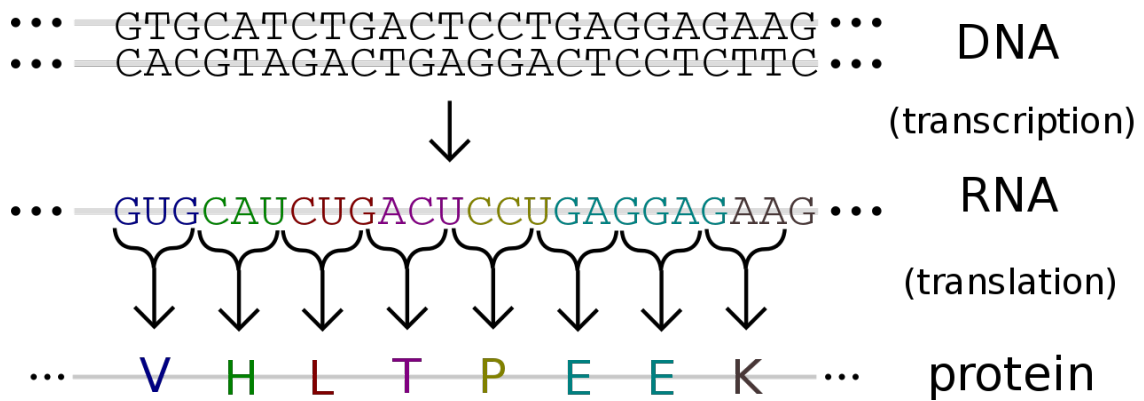


Figure 2.5: Example of the central dogma, which illustrates the first few amino acids for the alpha subunit of hemoglobin [23].

		Second nucleotide					
		U	C	A	G		
First nucleotide	U	UUU Phe UUC UUA Leu UUG	UCU UCC Ser UCA UCG	UAU Tyr UAC UAA STOP UAG STOP	UGU Cys UGC UGA STOP UGG Trp	U	Third nucleotide
	C	CUU CUC Leu CUA CUG	CCU CCC Pro CCA CCG	CAU His CAC CAA Gln CAG	CGU CGC Arg CGA CGG	C	
	A	AUU Ile AUC AUA AUG Met	ACU ACC Thr ACA ACG	AAU Asn AAC AAA Lys AAG	AGU Ser AGC AGA Arg AGG	A	
	G	GUU GUC Val GUA GUG	GCU GCC Ala GCA GCG	GAU Asp GAC GAA Glu GAG	GGU GGC Gly GGA GGG	G	

Figure 2.6: The amino acids specified by each codon [31].

2.1.3 Protein Structures

Protein structures range in size, from tens to several thousands of amino acids, and are categorised hierarchically into four distinct tiers, the primary, the secondary, the tertiary and the quaternary structure (Figure 2.7). This hierarchical approach was established to facilitate the observations of the various phases of protein formation. The number and type of amino acids of a protein are not enough, since the order and layout of their amino acids plays a major role because they determine the three-dimensional structure and hence the function of the protein.

2.1.3.1 Primary Structure

The primary structure of proteins is the sequence of amino acids in the polypeptide chain. This structure is determined by the gene, which is a sequence of nucleotides in deoxyribonucleic acid (DNA) or ribonucleic acid (RNA), corresponding to the protein. The sequence of a protein defines the structure and function of the protein and is unique to that protein. For example, the pancreatic hormone insulin is composed of 51 amino acids in 2 peptide chains, A chain has 21 amino acids while B chain has 30 amino acids, as shown in figure 2.8. The amino-acid sequences, in both chains, are unique to insulin and have a specific order. In each chain there are three-letter abbreviations, which represent the names of the amino acids. These are displayed in the order that are present and illustrate the primary structure of insulin.

The unique sequence for every protein is determined by the gene encoding of the protein. If the nucleotide sequence of the gene's coding region is changed, a different amino acid might be added to the growing polypeptide chain, which would change the protein structure and function. For instance, in sickle cell anemia (a hereditary disease that affects the red blood cells), a single amino acid substitution (valine in the β chain substitutes the amino acid glutamic) in the hemoglobin β chain, changes the protein structure and function (Figure 2.9). A hemoglobin molecule is comprised of two alpha and two beta chains, each consisting of about 150 amino acids. Therefore, the molecule has about 600 amino acids. The structural difference between the sickle cell molecule (which dramatically decreases life expectancy) and a normal hemoglobin molecule is just one of the 600 amino acids. As a result of this small change in the chain, hemoglobin molecules form long fibres that distort the biconcave, or disc-shaped, red blood cells and causes them to assume a crescent or 'sickle' shape, which clogs blood vessels and leads to myriad serious health issues such as breathlessness, dizziness, headaches, and abdominal pain, for those affected by this disease.

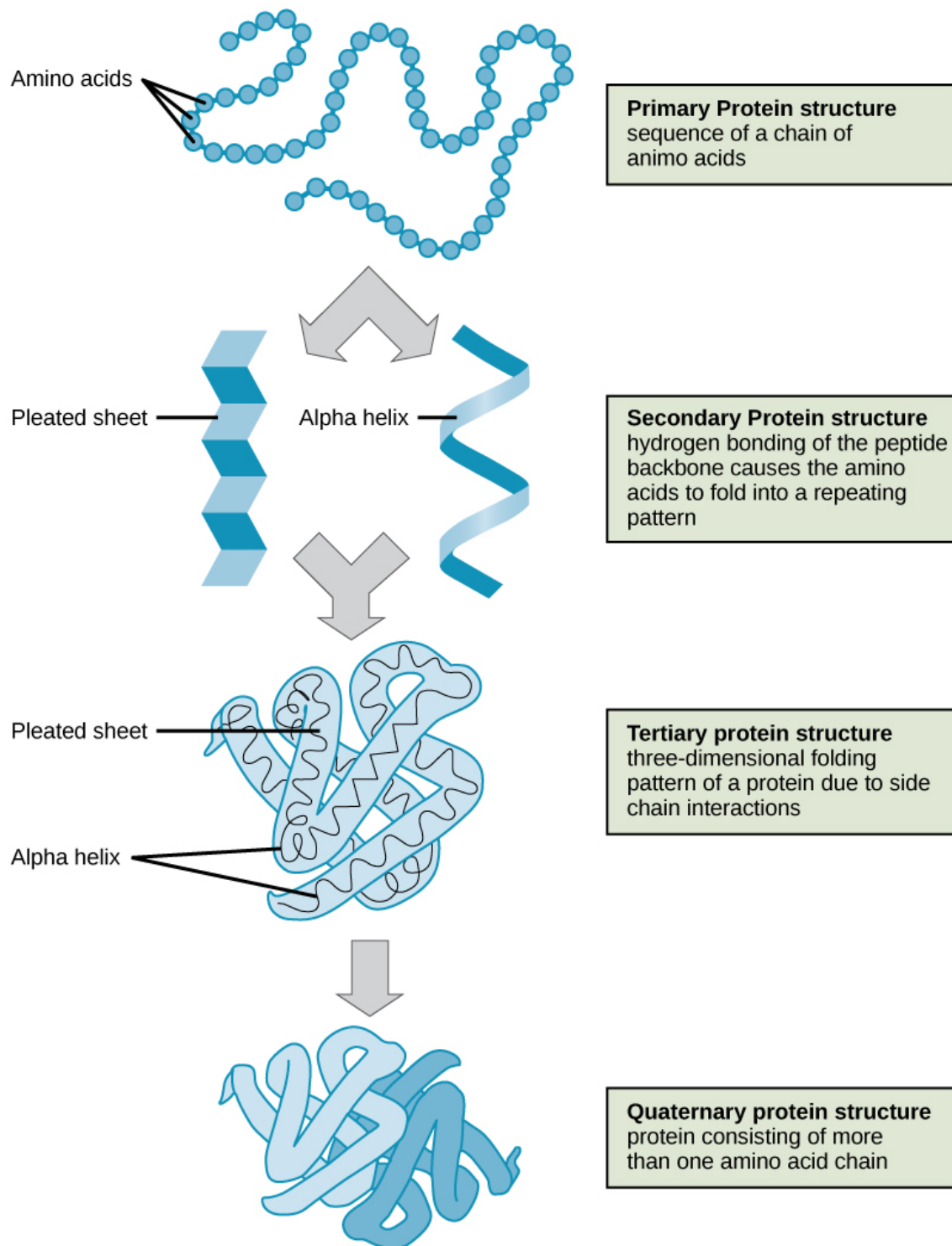


Figure 2.7: All four protein structures.

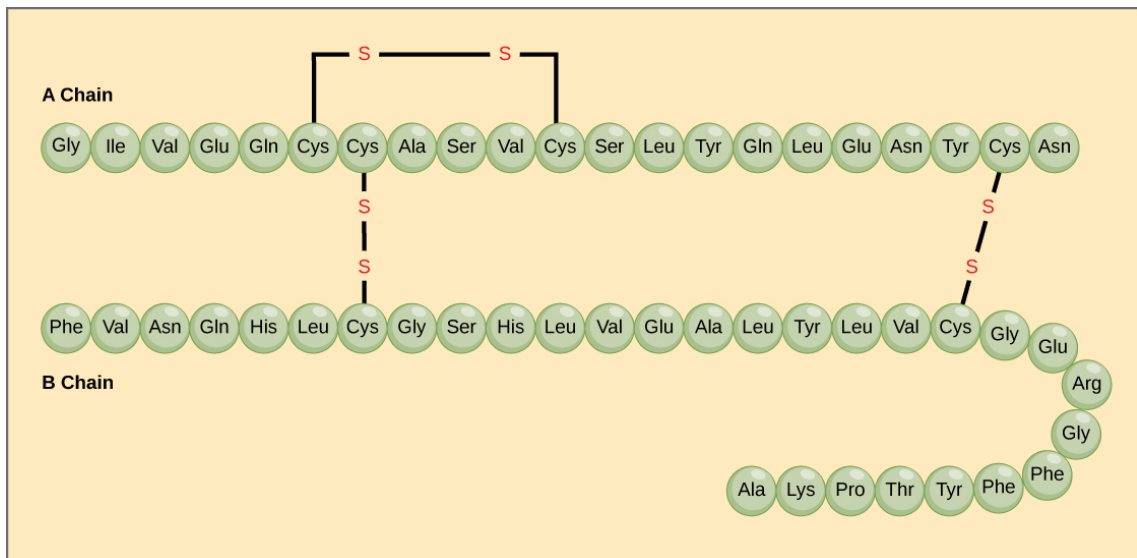


Figure 2.8: The first amino acid of the A chain is glycine (Gly), whereas, the last is asparagine (Asn) [32].

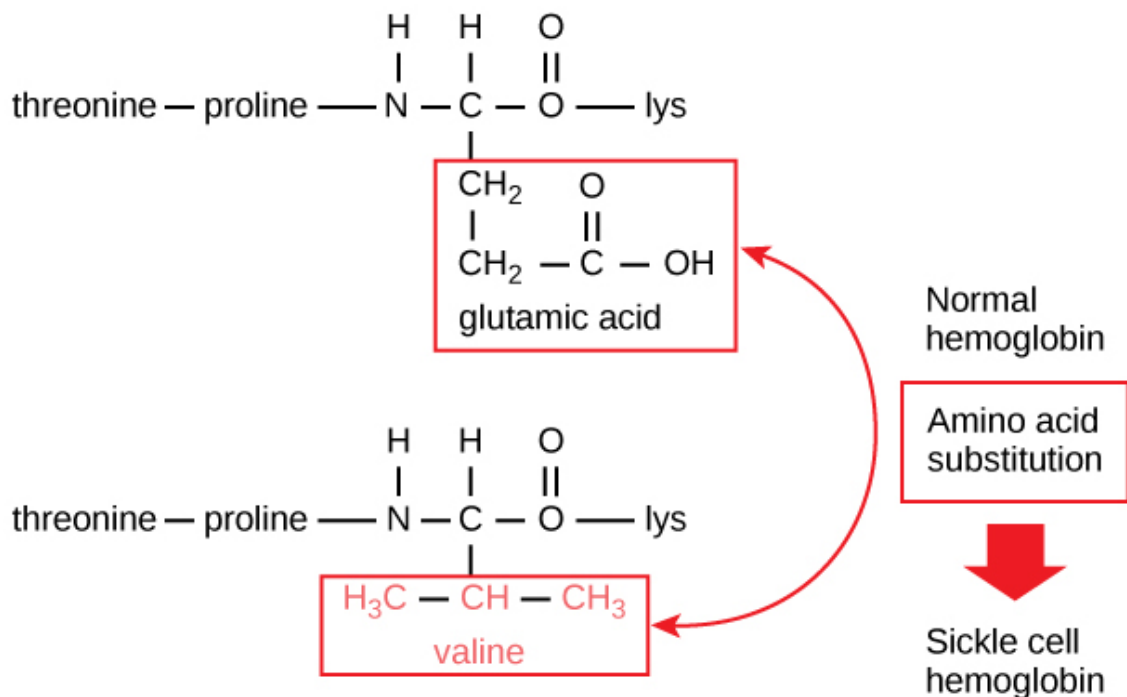


Figure 2.9: The diagram shows the substitution in a small part of the hemoglobin β chain, where the amino acid at position seven, glutamate, is replaced by valine, in the sickle cell hemoglobin [32].

2.1.3.2 Secondary Structure

The secondary structure of the protein refers to the local folding of the polypeptide in some regions and are defined by patterns of hydrogen bonds between the main-chain peptide groups. There are two main distinct categories of the secondary structure, the α -helix and the β -strand or β -sheets. Both of these are held in shape by hydrogen bonds, which form between carbonyl and amino groups in the peptide backbone. Certain amino acids have a propensity to form an α -helix, while others have a propensity to form a β -pleated sheet. The α -helix and β -pleated sheet structures are in most globular and fibrous proteins and play an important structural role.

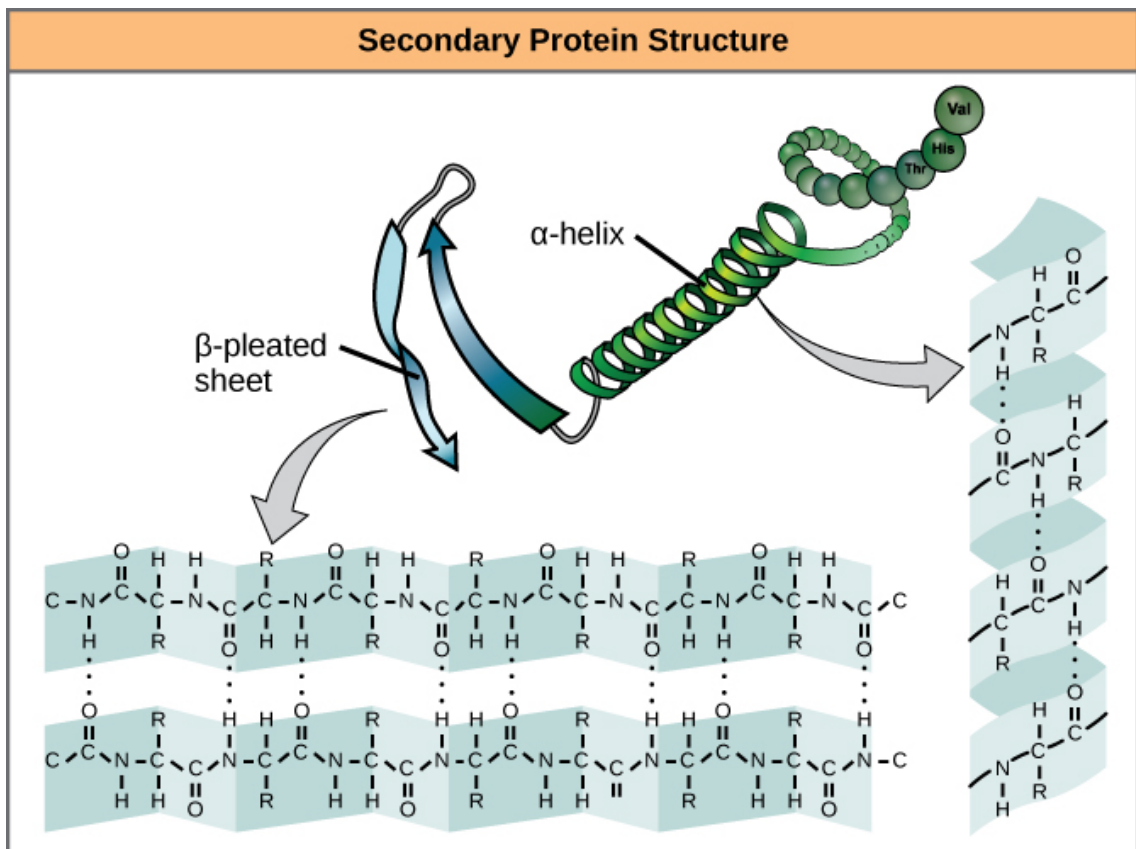


Figure 2.10: The diagram illustrates the shapes of the two main types of the secondary structure of proteins, the α -helix and the β -strand [32].

2.1.3.3 Tertiary Structure

The tertiary structure of proteins refers to a three-dimensional structure of monomeric and multimeric protein molecules. This structure is determined by a variety of chemical interactions on the polypeptide chain, such as ionic bonding, hydrophobic interactions, hydrogen bonding, and disulfide linkages (Figure 2.11). The protein's complex three-dimensional tertiary structure is created by the interactions among R groups. For instance,

R groups with like charges repel each other and those with unlike charges are attracted to each other (ionic bonds). The only covalent bond that forms during protein folding is the disulfide linkages, which are formed by interactions between cysteine side chains, in the presence of oxygen. As regards hydrophobic interactions, during the protein folding stage, the non-polar amino acids' hydrophobic R groups lie in the protein's interior, whereas, the hydrophilic R groups lie on the outside. Once a protein loses its three-dimensional shape, it may no longer be functional. The tertiary structure of a protein highly depends on the characteristics of its secondary structure, which is formed based on the order and layout of the amino acids (primary structure) of the protein.

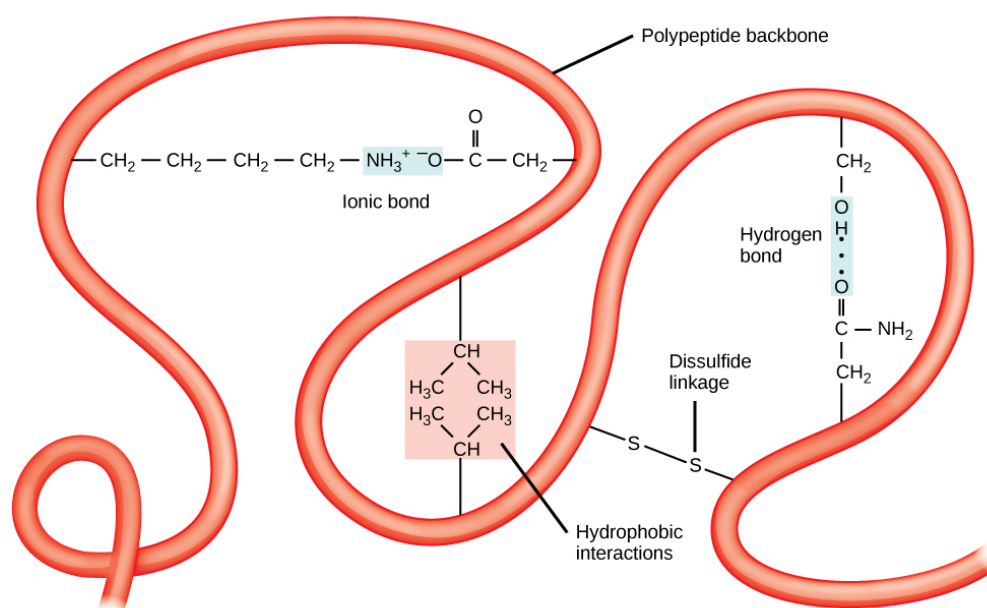


Figure 2.11: The diagram indicates some of the chemical interactions that determine the proteins' tertiary structure [32].

2.1.3.4 Quaternary Structure

The quaternary structure of a protein is the three-dimensional structure consisting of the aggregation of two or more individual polypeptide chains (subunits) that operate as a single functional unit (multimer). For example, insulin (which is a globular protein) has a combination of hydrogen and disulfide bonds, which cause it to mostly clump into a ball shape. Insulin starts out as a single polypeptide and after forming the disulfide linkages that hold the remaining chains together, it loses some internal sequences in the presence of post-translational modification. Silk (which is a fibrous protein), on the other hand, has a β -pleated sheet structure, which is the result of hydrogen bonding between different chains. A representation of the quaternary structure can be found in figure 2.7.

2.2 Artificial Neural Networks Background

2.2.1 Origins of Artificial Neural Networks

Artificial neural networks (ANNs) are computing systems which are inspired by the biological neural network that exists in the brains of humans and animals. The term ‘neural’ comes from the basic functional units of the human nervous system, called ‘neurons’ or ‘nerve cells’. These are located in various parts of the human body, like the brain which contains about 10^{11} neurons that are connected to 10^4 other neurons.

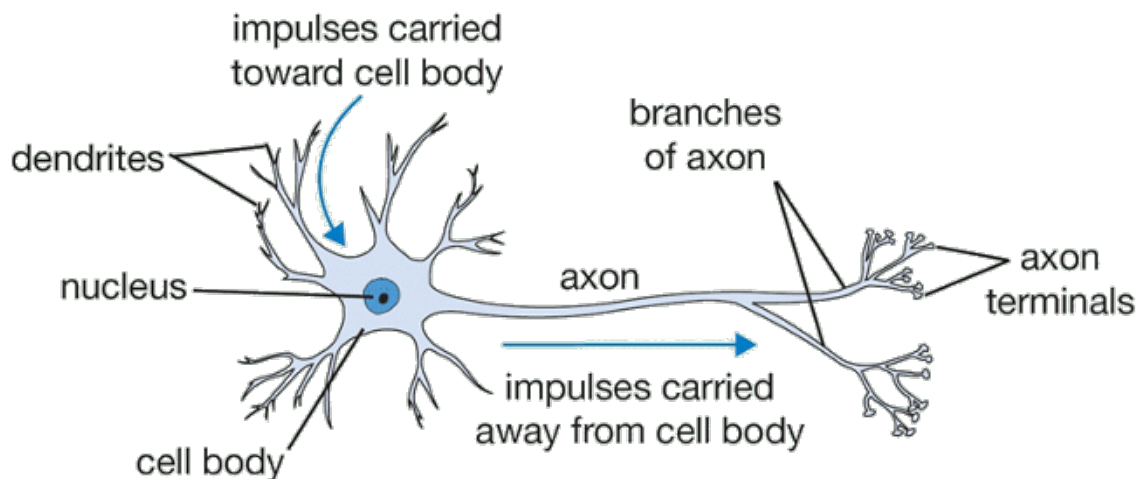


Figure 2.12: Structure of a Biological Neuron [33].

A biological neural network is a collection of neurons that receive, process and transmit information between each other, through electrical and chemical signals via specialized connections called synapses. It consists of three main components, the cell body, the axons and the dendrites. Figure 2.12 shows the direction of the impulses when a signal is carried towards or away from a neuron. The neuron receives signals from other neurons through dendrites. The body of the neuron adds all the incoming signals and calculates the input of the neuron. If the sum exceeds a certain threshold value the neuron triggers and the signal is transmitted through the axon to the other neighbouring neurons. Axon terminals are the connection point between brain neurons. The signal's strength, which is transmitted from one neuron to another, depends on the interconnection force of the neurons. The human nervous system is like an extremely high-connectivity network, which has trillions of neurons and billions of connections between them.

An Artificial Neural Network (ANN) has the same architecture with the biological neural network. An ANN has nodes that represent artificial neurons, a simplified version of biological neurons in terms of functionalities, and connections (edges) instead of synapses.

These connections are responsible for transmitting signals between the connected artificial neurons. ANNs have a similar behaviour with the biological neural network, but as they became more and more popular, the idea of replicating the human brain faded away.

The increasing demand for solving specific tasks, led to the development of various implementations of ANN, and some of them were based on the initial concept of biological neural networks. For instance, an ANN called Recurrent Neural Network (RNN) was based on the concept of short term memory and is used to recognise patterns, where the previous features can help predict the next ones. Another variation of ANNs is the Convolutional Neural Network (CNN) [34], which is used in this dissertation. The CNN is able to recognise patterns in two-dimensional (or three-dimensional) data, like images and videos, and feed the extracted features to a fully connected feed-forward Multi-Layer Perceptron (MLP) to classify the initial input data. There are many other variations of ANNs that were designed for specific tasks like speech translation or recognition, natural language processing, clustering or even playing video games. Some of these variations will be discussed in the following section.

2.2.2 Variations of Artificial Neural Networks and Optimizers

2.2.2.1 McCulloch and Pitts (McP)

The first ANN model was suggested by Warren McCulloch and Walter Pitts in 1943 [35] (Figure 2.13). The design of this artificial neuron was very simple and was based on a single biological neuron of the human brain. An input vector performs multiplications with the weight values and provides the signals to the artificial neuron. Then, the artificial neuron sums those signals and transfers the result to a threshold function, also known as step or heaviside function (Figure 2.14), which does not provide enough information about how close or how far the target output is. The output signal of the model was 1, if the value exceeded the a specific threshold value, otherwise the output signal was 0, which means this model can be used only for binary classification.

The inputs are classified based on the weights of the connections and the threshold value (Equation 2.1, where y is the output of the network, x the input vector, w the weight vector, $w \cdot x$ the dot product and s the threshold). For instance, for a two-dimensional input vector, in a simple two-dimensional scenario, the decision line can be calculated with equation 2.2.

$$y = \begin{cases} 1 & \text{if } w \cdot x > s \\ 0 & \text{otherwise} \end{cases} \quad (2.1)$$

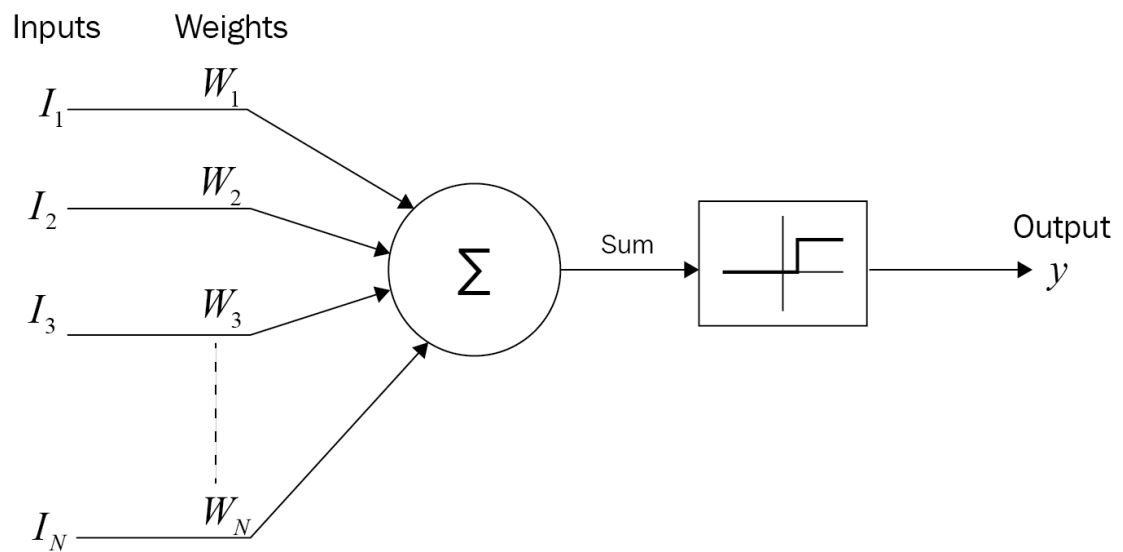


Figure 2.13: McCulloch and Pitts artificial neuron [35].

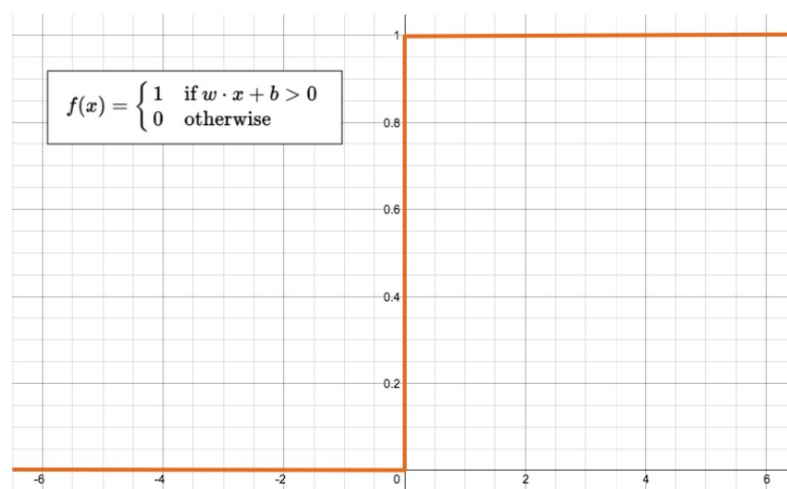


Figure 2.14: The step or heaviside function.

$$x_2 = -\left(\frac{w_1}{w_2}\right)x_1 + \left(\frac{s}{w_2}\right) \quad (2.2)$$

If the goal is to classify the OR gate (Table 2.2) with a McP model, the model could use infinite different ways to solve the problem. For instance, the model could have weights of $W = [2, 2]$ and a threshold value $S = 1$. In figure 2.15, (b) illustrates the decision line for the OR gate, where inputs above the line are classified as Class 1 while inputs below the line are classified as Class 0. The equation for this decision line is $x_2 = -x_1 + 0.5$.

x1 (input)	x2 (input)	y (output)
0	0	0
0	1	1
1	0	1
1	1	1

Table 2.2: Truth table for OR gate.

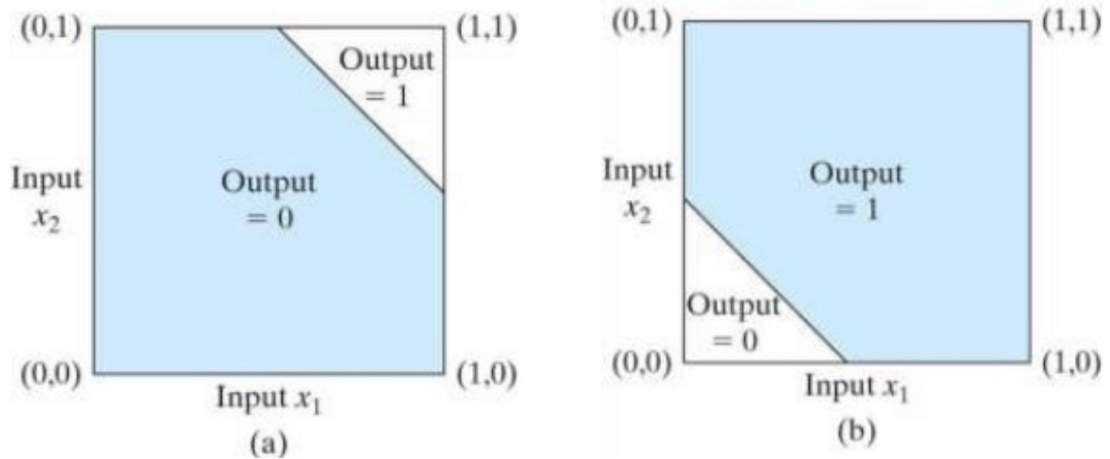


Figure 2.15: Decision lines for AND gate (a) and OR gate (b).

The training phase of McP neurons requires the input and target output to be presented to the network, which calculates the actual output for the given input and adjusts the weights accordingly. For example, if the output is 1 but the target output is 0 the weights are modified, while in the case where both the output and target output are 0 the weights remain the same. This process is also known as the Perceptron Learning algorithm [36] (Algorithm 1).

Perceptron Learning Algorithm

1. Initialize weights and threshold randomly.
2. Present input and desired output.
3. Calculate actual output (Equation 2.1).
4. Adapt weights:

if output 0, should be 1: $w_i(t+1) = w_i(t) + \eta \cdot x_i(t)$

if output 1, should be 0: $w_i(t+1) = w_i(t) - \eta \cdot x_i(t)$

if output is correct : $w_i(t+1) = w_i(t)$

where $0 \leq \eta \leq 1$ the learning rate, controlling the adaptation rate.

Algorithm 1: Perceptron Learning Algorithm.

This algorithm was thought to be very promising, but after a while it was proven that the perceptron algorithm could only solve problems with linearly separable patterns. In these problems, a straight line or hyperplane, which separates the patterns, can be found in space, like the OR gate problem which was mentioned earlier (Table 2.2). However, this algorithm cannot solve problems that require more than one straight lines or hyperplanes to separate the different classes, not even simple ones like the XOR gate problem (Figure 2.16). Except from that, there was no indication on how close to the target output was the predicted output because of the binary (either 1 or 0) output of the heaviside function. This problem was the main motivation for developing more sophisticated networks and algorithms, some of which will be discussed subsequently.

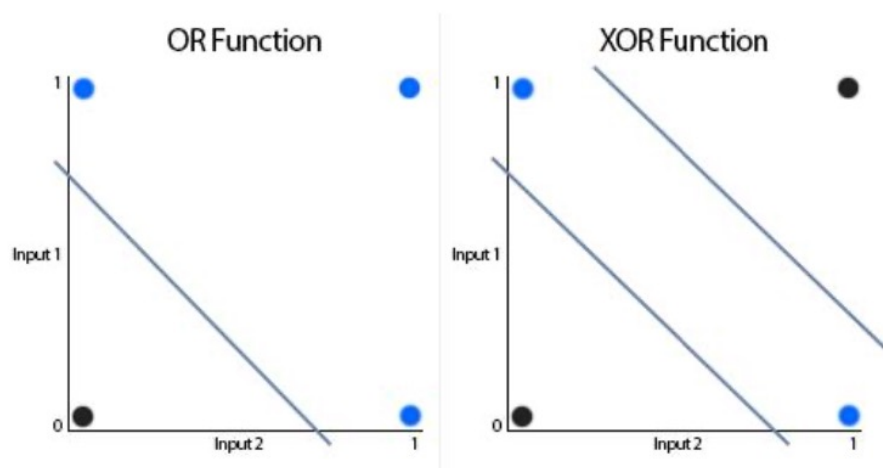


Figure 2.16: The OR gate is linearly separable while the XOR gate is not.

2.2.2.2 Multi-Layer Perceptron (MLP)

Multi-Layer Perceptron (MLP) neural networks are currently the most popular and well-known variants of ANNs. They consist of multiple, slightly modified, McCulloch and Pitts neurons, which form layered feed forward networks (Figure 2.17). McP neurons use a specific threshold activation function (step function) while MLP neurons can use any arbitrary activation function (Table 2.3). This is the reason why McP can only perform binary classification, while MLP can perform regression or classification, depending on the selected activation function. Furthermore, activation functions provide an indication to the network whether the outputs are closer or further of the expected outputs, which helps the network adjust the weights accordingly, to improve predictions.

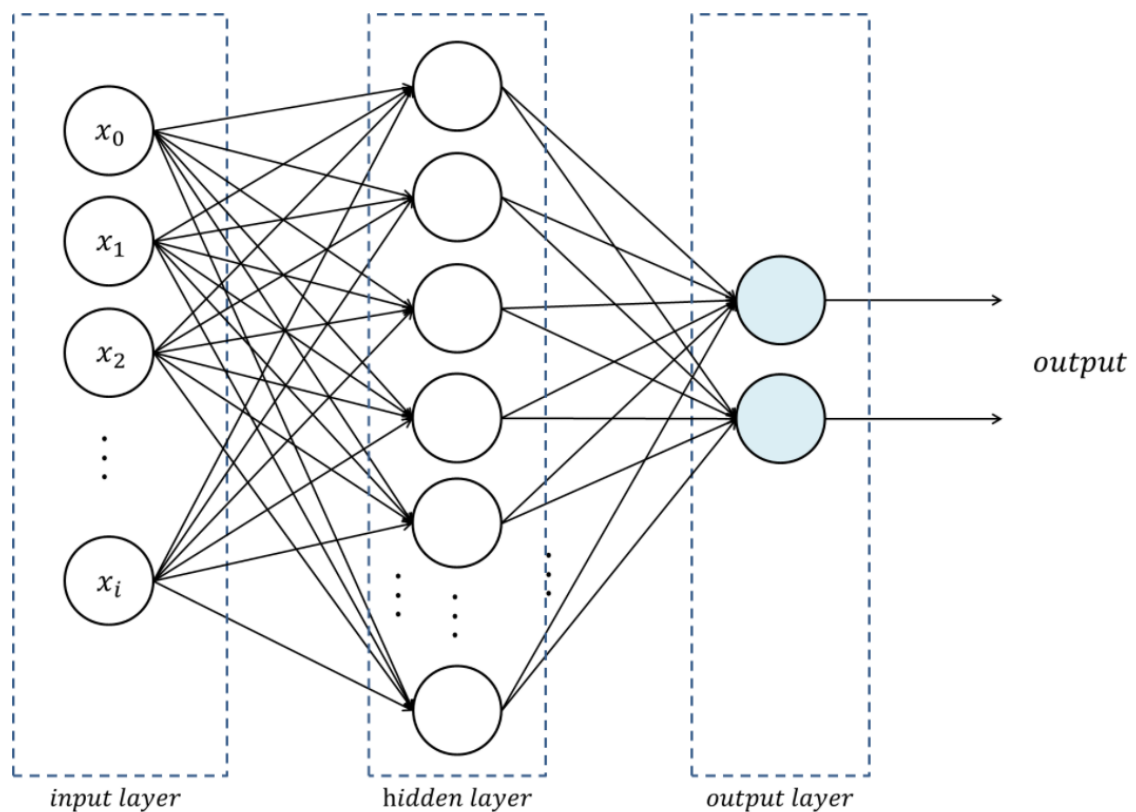


Figure 2.17: Multi-Layer Perceptron Neural Network with one hidden layer.

An MLP neural network consists of an input layer, one or more hidden layers and an output layer. The hidden and output layers are active, while the input layer is not active (only forwards the input data to the network). Each layer has one or more neurons and an independent neuron unit, also known as 'bias', which has a constant input value of one (1). The role of the bias unit is to help the network adapt more effectively to the provided data. The number of hidden layers is very important as it specifies the possibilities of the network and processes the biggest amount of information during the training (learning)



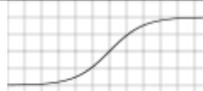

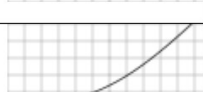
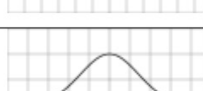
Name	Plot	Equation	Derivative	Range
Heaviside		$f(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$	$f'(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ ? & \text{if } x = 0 \end{cases}$	$\{0,1\}$
Logistic / Sigmoid		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	$(0,1)$
TanH		$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$f'(x) = \frac{1}{x^2 + 1}$	$(-1,1)$
Rectified linear unit (ReLU)		$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{otherwise} \end{cases}$	$f'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{otherwise} \end{cases}$	$[0,\infty)$
SoftPlus		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$	$(0,\infty)$
Gaussian		$f(x) = e^{-x^2}$	$f'(x) = -2xe^{-x^2}$	$(0,1)$

Table 2.3: List of the most popular activation functions.

phase (Figure 2.18). The neurons of the first hidden layer determine the number of decision lines that can separate the patterns into classes. The second hidden layer gives the ability to form arbitrary complex decision shapes, which are able to separate any classes, so there is no need for more than two hidden layers in a neural network [37].

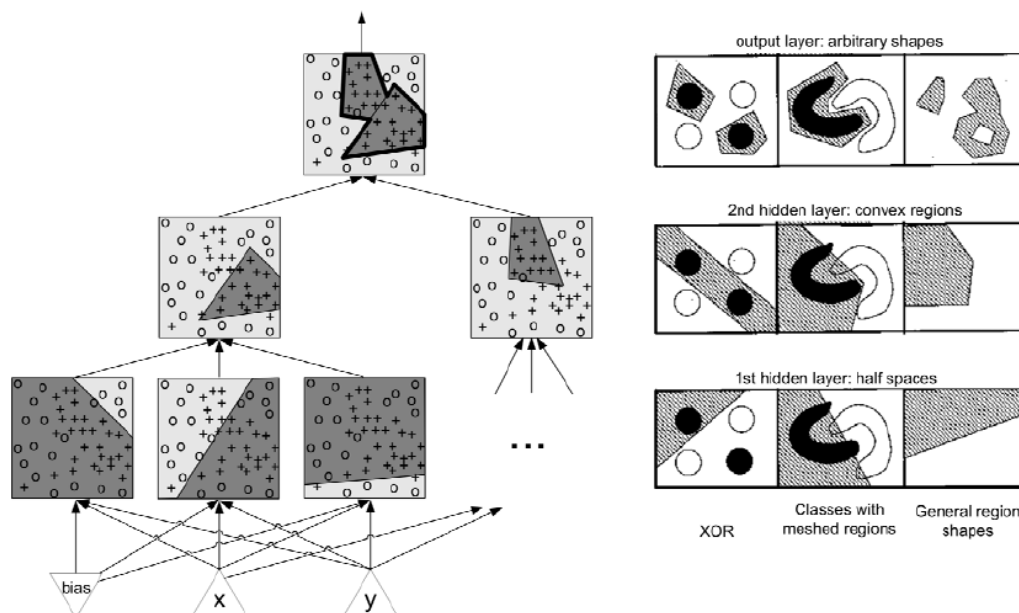


Figure 2.18: Decision regions based on the number of hidden layers.

The calculation process of the networks' output is very similar with the one used in McP. The input layer forwards the input values to the first hidden layer, which calculates the sum of the bias and the dot product of the weights and the input vector, and then passes that value to the activation function (Equation 2.3, where y is the output of a single neuron, x the input vector for that neuron, w the weight vector, $w^T \cdot x$ the dot product, b the threshold and f the arbitrary activation function). The output signals of the activation function are then fed as inputs of the next hidden or output layer, which then repeats this process until there are no more layers to pass the signals.

$$y = f(w^T x + b) \quad (2.3)$$

Gradient Descent

Gradient descent is one of the most popular optimisation algorithms for training ANNs. It is considered a mathematical optimization algorithm that is able to minimize a function by iteratively moving in the direction of steepest descent, which is defined as the negative of the gradient. An error function is used to calculate how successful the network predicting the classes of the input patterns was, like the mean squared error (MSE) function (Equation 2.4, where t is the target output, o the actual output, p denotes the pattern and j the neuron). The objective is to minimise the error value, which is the difference between the target and actual outputs. By adjusting the weight vectors according to the negative of the derivative of the error value, of the current pattern, with respect to each weight (Equation 2.5), where n is the learning rate), at some point the correct classifications will be maximized.

$$E = \frac{1}{2} \sum_j (t_{pj} - o_{pj})^2 \quad (2.4)$$

$$\Delta w_{ij} = -n \frac{\partial E_p}{\partial w_{ij}} \quad (2.5)$$

Backpropagation Algorithm (BP)

In order to calculate the error and use gradient descent to minimize it, both target and predicted outputs must be known. In the output layer this is fine as both values are available, however, in the hidden layers the target values are unknown, which means that only the weights between the last hidden layer and the output can be adjusted. To solve this issue, the backpropagation algorithm was suggested, which propagates the error from the output layer back to the last hidden layer, which then does the same until all the weights are updated (Algorithm 2, where δ_{ij} is the error signal, y_{ij} is the actual output and d_{ij} is

the target output of neuron i of layer j . The δ_{ik} is the same as δ_{ij} but for the previous iteration of the algorithm).

Backpropagation

```
Repeat:
  For each pattern :
    // Forward Pass
    Calculate the output
    // Backward Pass
    For each layer j, starting at the output:
      For each unit i:
        // Compute the error
        If output neuron:  $\delta_{ij} = y_{ij}(1 - y_{ij})(d_{ij} - y_{ij})$ 
        If hidden neuron:  $\delta_{ij} = y_{ij}(1 - y_{ij}) \sum \delta_{ik} \cdot W_{jk}$ 
        For each weight to this unit:
          Compute and apply  $\Delta w$ 
      Compute total error
      Increment epoch counter
  Until small enough error or epoch counter exceeded
```

Algorithm 2: The Backpropagation algorithm.

More specifically, to update all the weights two passes are required, a forward pass to calculate the error based on the given input pattern, and a backward pass, where the error is back propagated to the previous layers and all the weights are updated respectively. The entire process is repeated for every pattern, until all patterns have been passed into the network (one epoch), which is also known as the online update mode. There are two alternatives, the batch and mini-batch modes. The first feeds the network with all the patterns at once and gets cumulative updates for the weights, which usually helps the network learn more effectively. However, if the input datasets are too big and cannot fit into memory, this method cannot be used. The second method is a combination of the online and batch mode and can be used for big datasets. This method takes the input patterns and splits them into smaller chunks, called mini-batches, then it feeds the network with one mini-batch at a time. The size of the mini-batch can be adjusted to ensure that there are no ‘out of memory’ issues, which makes this more flexible compared to the other two methods. The goal is to feed the network all the input patterns several times until the error reaches a specified value or until a number of epochs (when all the patterns have been fed into the network) has passed.

2.2.2.3 Recurrent Neural Network (RNN)

The Recurrent Neural Network (RNN) is a variation of MLP, which instead of feeding the input forward to the next layers, it uses recurrent inputs. Recurrent inputs are the output signals from the hidden layer or the output layer, which are fed into a previous layer or even to the same layer. This technique creates a form of ‘memory’ for the network, since the output depends on both the current input and the input from the previous iterations. This makes RNNs great for dynamic problems, like timeseries or sequence predictions.

There are two main versions of RNNs, the Jordan RNN [38] and the Elman RNN [39] (Figure 2.19). The main difference between the two versions is that the first transfers its output to a context layer, also known as state units, which then feeds the network along with the new input patterns. The second variation, on the other hand, feeds the hidden layer output to a context layer, also known as context units, which is fed back to the hidden layer.

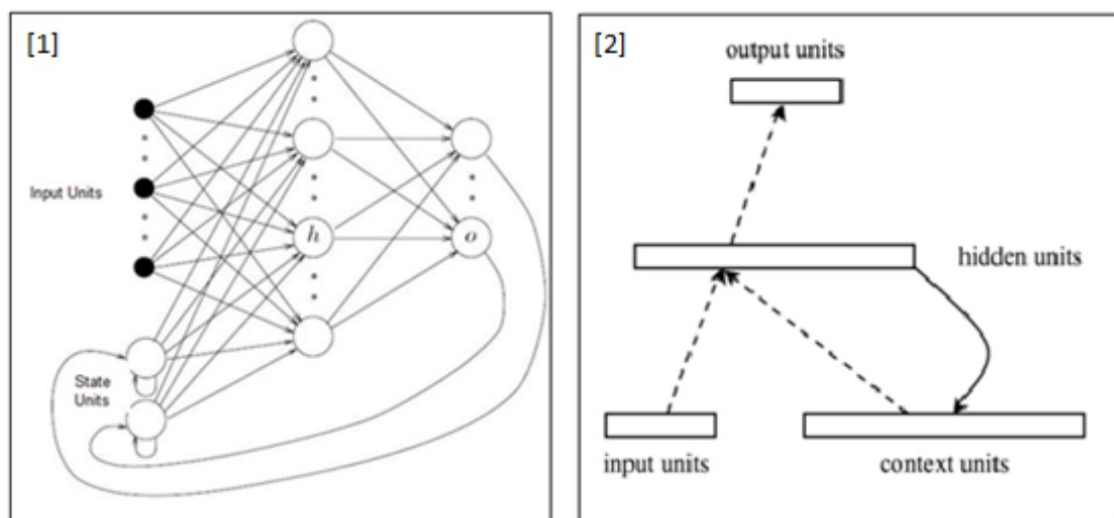


Figure 2.19: RNN variations, Jordan network (left), Elman network (right).

2.2.2.4 Convolutional Neural Network (CNN)

Description of Convolutional Neural Networks

A Convolutional Neural Network (CNN) is a class of deep artificial neural networks, which is most commonly applied to analysing visual imagery. Its application ranges from image and video recognition, recommender systems, image classification to medical image analysis, and natural language processing (NLP). CNNs are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers. Multilayer perceptrons (MLPs), which are usually fully connected networks (each

neuron in a layer is connected with all the neurons in the next layer), are prone to overfitting data. CNNs on the other hand take advantage of the hierarchical pattern in data and assemble more complex patterns using smaller and simpler patterns. Therefore, on the scale of connectedness and complexity, CNNs are on the lower extreme. The architecture of a CNN is designed so that it can take advantage of the two-dimensional (2D) structure of an input image (or any other 2D input, such as speech signals).

Architecture of Convolutional Neural Networks

A Convolutional Neural Network (CNN) consists of an input and an output layer, as well as multiple hidden layers. The hidden layers of a CNN typically consist of a series of convolutional layers. The activation function used is commonly a Rectifier Linear Unit (RELU) layer, and is subsequently followed by additional convolutions, pooling layers and a fully connected layer. The final fully connected layer, which is usually a multilayer perceptron (MLP) network, uses the backpropagation learning algorithm for training. The input of a convolutional layer is an image of size $d \times d \times c$, where d is the height and width of an image and c is the number of channels of the input image (e.g. an RGB image has $c = 3$). A convolutional layer has k filters (or kernels) of size $m \times m \times n$, where m is smaller than the dimensions of the image and n can be either the same as the number of channels c or smaller (may vary for each kernel). Convolutional layers convolve the input, which leads to the creation of k feature maps of size $d - m + 1$, and pass their output to the next layer. Subsequently, each feature map is sampled, typically averaging or maximizing above the same areas in feature maps of size $p \times p$ (where p is between 2, for small images, and usually does not exceed 5, for larger images). A bias and a sigmoid nonlinearity is applied to each feature map, prior or after the pooling layer.

Figure 2.20 illustrates an example of a CNN which is used to classify images of handwritten digits. The diagram shows the different layers of a CNN (convolution, max pooling, multilayer perceptron) and the feature maps that are extracted from each image (small squares). At the end of the CNN (right hand side), there is a fully connected network (MLP) which is used to classify the input image [40].

A pooling layer between the hidden convolutional layers is a common tactic for classic CNN architectures. The pooling layers are mainly used to reduce the dimensions of the output of each layer, the number of parameters and the complexity of the network, which consequently reduces the total computation time of the network. This practice also prevents the network from overfitting (adapting to the training data, making it less effective at predicting new data patterns), which can be determined by observing the training and test error values. The pooling layers are independent from the other layers and they pro-

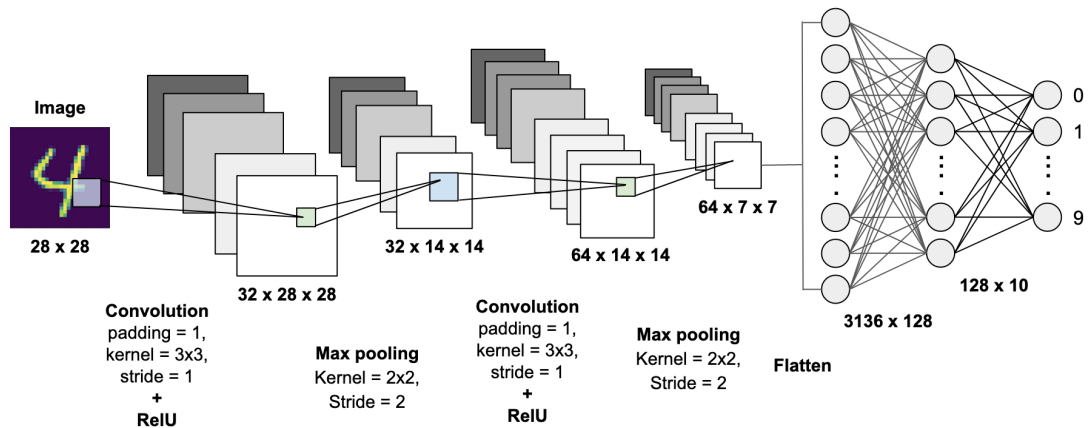


Figure 2.20: A CNN example for digit image classification.

process the output of each kernel (filter) separately. Even though, there are different types of pooling layers, max pooling, min pooling (which is the opposite of max pooling), average pooling and L2-normalization pooling, the max pooling technique seems to work better than the rest [41]. This technique, as its name suggests, takes the max value from each filtered result and returns it. Figure 2.21 illustrates an example of max and average pooling, where the kernel size (filter size) is 2×2 and the stride (how many slots to skip) is two (2). The applied filters can be distinguished by their colors.

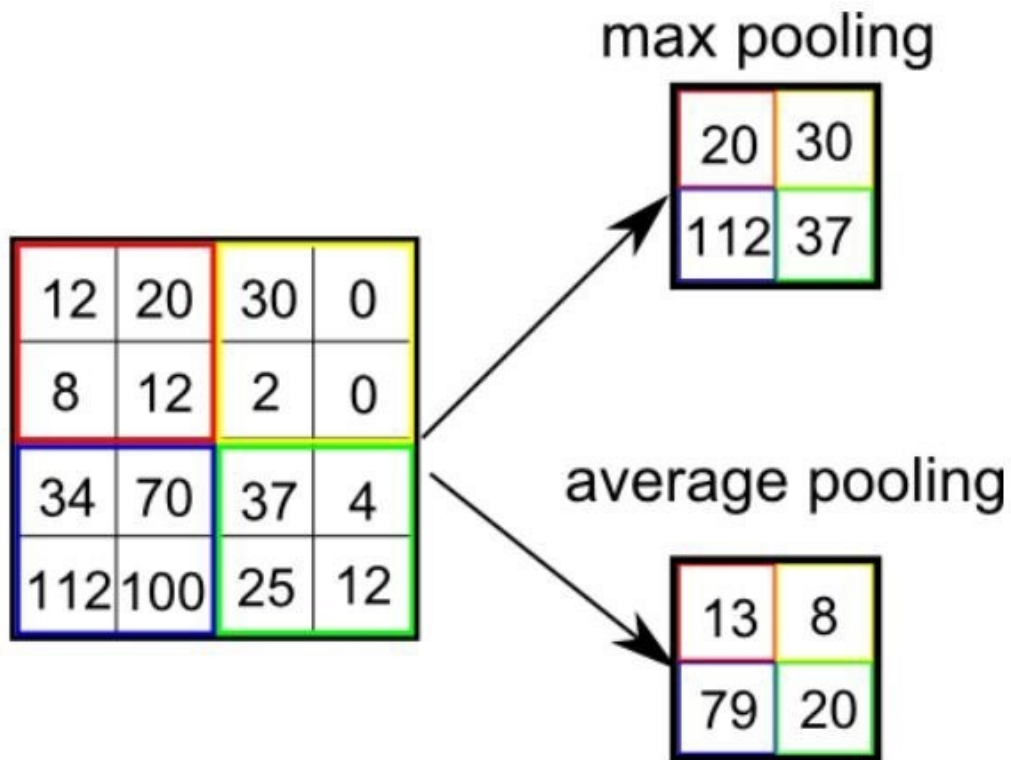


Figure 2.21: Example of max and average pooling.

The input of a CNN is a 3D array, also called 3D tensor. For instance, if the purpose of a CNN is to identify objects in a 50×50 pixels picture, the input would be a 3D tensor with shape $50 \times 50 \times 3$. That is because each pixel is represented by three values, one for red, one for green and one for blue (RGB). In the PSSP problem, a 2D tensor is enough, as a 2D tensor is the same with a 3D tensor where the third dimensions has size one. The shape of this tensor will be $L \times 20$, where L is the number of lines of the input file and 20 represents the 20 known amino acids. An example with visualizations and more details on how CNNs work can be found here [42].

One of the main advantages of CNNs, is the fact that they can extract features from complex sequences, due to the small number of synaptic weights. For example, if the input array size is $28 \times 28 \times 3$ (RGB) and the kernel size is 5×5 , then each neuron of the convolution layer will be connected with an area of the input array with shape $5 \times 5 \times 3$. This means that each neuron has 76 ($5 \times 5 \times 3 = 75 + 1 = 76$) synaptic weights, which can extract features and adapt to complex input data.

In some cases, it is necessary to add zero padding (append zero values) around the input data (like a frame for the input array). The number of rows and columns of zeroes is variable, which makes it possible to control the dimensions of the output of a hidden convolution layer, using zero padding. Figure 2.22 illustrates an example of zero padding, where two borders of zeros are placed around the $32 \times 32 \times 3$ input.

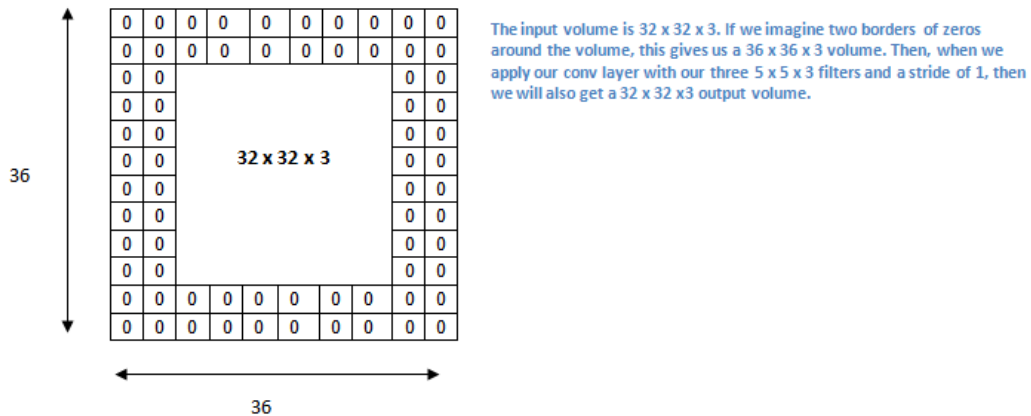


Figure 2.22: Example of zero padding.

2.2.2.5 Line Search

Line search is one of the basic iterative approaches, used to find a minimum x^* of an objective function. For an ANN, x represents the weights of the network, while the objective function represents the error function. Equation 2.6 illustrates the essential components

to calculate the next iteration of x . The step size determines the size of the step of x in that direction. Line search, in each iteration, attempts to find the best step size, which can minimise the objective function in a specific search direction. On the other hand, gradient descent requires a learning rate which determines how small or how big is each step. If the step is too small the learning process will take significantly more time and can lead the network to a local minimum (instead of the global minimum, which is the desired outcome). If the step is too big then it is very likely that the objective function will jump far away from the desired minimum.

$$x_{n+1} = x_n + a_n d_n \quad (2.6)$$

Therefore, applying the optimal step size is very important, as it can prevent the network from moving further away from the minimum. In order to find the step size, a naive approach was to move along a search direction in small steps and after each step calculate the error, if the error starts increasing then stop and change direction [43]. However, this approach is not very efficient, robust or accurate compared to other variations of line search [44].

2.2.2.6 Conjugate Gradient (CG)

The Conjugate gradient algorithm (Algorithm 3), unlike gradient descent, in each iteration changes the direction to prevent the network from becoming counterproductive (reversing the progress). In addition to that, in an N-dimensional problem, the CG algorithm is guaranteed to find a solution in N steps, since in every CG step the network obtains the minimum of that direction. Figure 2.23 compares the CG and the gradient descent algorithm on the same two-dimensional problem. Conjugate gradient managed to converge in just two steps, while gradient descent required several steps.

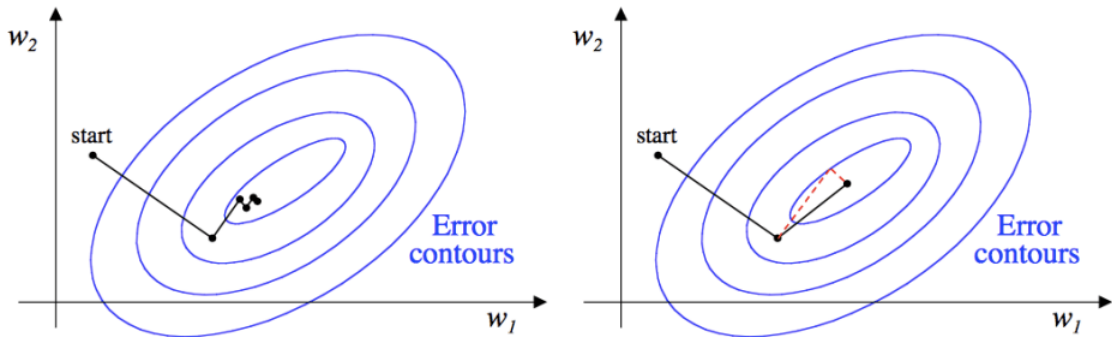


Figure 2.23: Gradient Descent (left) vs Conjugate Gradient (right) on a 2D problem.

Conjugate Gradient

1. Initialize weight vector w_0 randomly, set $i=0$
2. Evaluate the gradient vector g_i , and set the initial search direction $d_i = -g_i$
3. Use Line Search to find best step size a , which minimizes the function $f(w_i + ad_i)$
4. Update weights $w_{i+1} = w_i + ad_i$
5. Test stopping conditions
6. Evaluate new gradient vector g_{i+1}
7. Evaluate new search direction $d_{i+1} = -g_{i+1} + \beta_i d_i$, where β_i is given by one of:

$$\beta_i = \frac{g_{i+1}^T (g_{i+1} - g_i)}{g_i^T g_i} \quad (\text{Polar and Ribiere})$$

$$\beta_i = \frac{g_{i+1}^T g_{i+1}}{g_i^T g_i} \quad (\text{Fletcher and Reeves})$$

8. Set $i=i+1$ and go to step 3

Algorithm 3: Conjugate Gradient Algorithm [45].

2.2.2.7 Newton's Method

An iterative method, originally used to find approximations of the roots of real-valued functions, is currently used in optimisation problems to find the maximum or minimum of a function and is known as the Newton's Method. The derivative of a function at a maximum or a minimum point is zero, which makes it possible to find local maxima and minima by using the Newton's Method on the derivative of the optimisation function.

Newton's Method is considered a second-order optimisation algorithm, since it requires information about the second derivative of the optimisation function. Compared to first-order optimisation algorithms (like gradient descent), second-order optimisation methods can achieve faster and more accurate convergence to the minimum of a function.

In a simple first-degree polynomial (Figure 2.24), 1D problem, of a function $f(x)$ and a sub-optimal initial solution x_0 , Newton's method suggests the following:

1. Set $x_i = x_0$
2. Find the equation of the tangent at x_i
3. Find the point x_{i+1} at which the tangent line intersects with the x-axis
4. Find the projection of x_{i+1} on $f(x)$
5. Set $x_i = x_{i+1}$ and repeat from 2 until $f(x_i) < \text{threshold}$

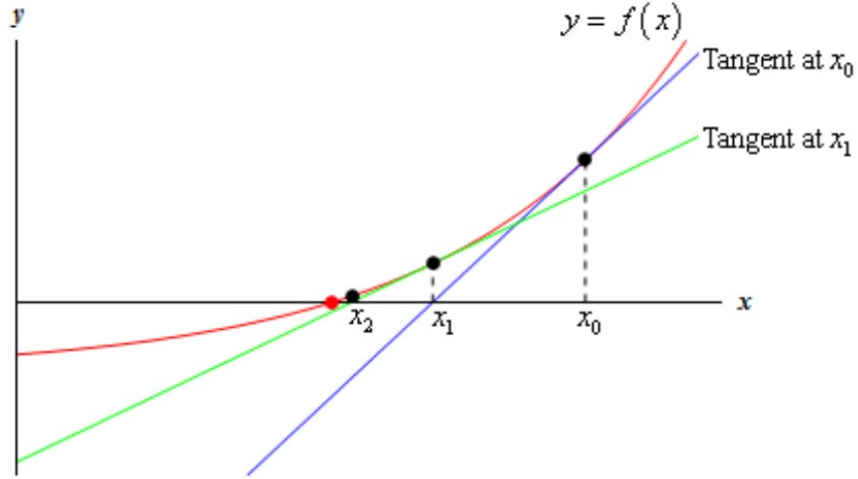


Figure 2.24: Newton's method in a first degree polynomial problem [23].

The equation of a point-slope line is:

$$y - y_1 = m(x - x_1) \quad (2.7)$$

In 2.7 the derivative can be used instead of the slope m and this can be rewritten as:

$$f(x) - f(x_1) = f'(x)(x - x_1) \quad (2.8)$$

Since x_1 is the point of interaction on x-axis, $f(x_1) = 0$ which gives the update rule for x for optimizing the function as:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \quad (2.9)$$

The previous simple example was used just to provide the intuition behind the method of finding the roots of a function. In optimisation theory, this method actually approximates the function $f(x)$ with a local quadratic function around x and moves towards the minimum of that approximated function with iterative steps. This process is repeated until a specified error threshold is reached or after a certain number of iterations has passed. The quadratic approximations around the weights at each iteration are shown in figure 2.25.

For the approximation of the function $f(x)$, the second-order Taylor expansion (second series Taylor approximation) is being utilized.

$$f(x_0 + x) \approx f(x_0) + f'(x_0)x + \frac{f''(x_0)}{2}x^2 \quad (2.10)$$

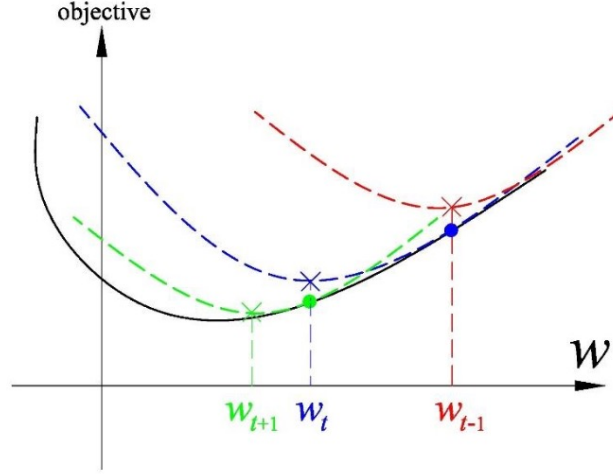


Figure 2.25: Local Quadratic approximations [23].

In order for $f(x_0 + x)$ to be a minimum, an optimal x value must be specified. Newton's method takes the derivative of the Taylor series and sets it equal to zero (Equation 2.11).

$$\frac{d \left(f(x_0) + f'(x_0)x + f''(x_0)\frac{x^2}{2} \right)}{dx} = f'(x_0) + f''(x_0)x = 0 \Rightarrow x = -\frac{f'(x_0)}{f''(x_0)} \quad (2.11)$$

This x is just the absolute minimum of the local approximation of $f(x)$ around the initial solution of x_0 and not the absolute minimum of $f(x)$. For the minimum of the objective function this process must be repeated multiple times, until it eventually converges to a minimum. The final update rule for optimizing the function $f(x)$ for a 1D problem is given by the equation 2.12 .

$$x_{n+1} = -\frac{f'(x_n)}{f''(x_n)} \quad (2.12)$$

This algorithm, however, can work only for objective functions with a single dimension ($f: \mathbb{R} \rightarrow \mathbb{R}$).

If the objective function, has multiple dimensions ($f: \mathbb{R}^n \rightarrow \mathbb{R}$), the algorithm must be modified by replacing derivatives with gradients and second derivatives with Hessians (the matrix of second partial derivatives, figure 2.26)

$$x_{n+1} = -\frac{\nabla f(x_n)}{H(f)(x_n)} \quad (2.13)$$

Equation 2.13 is the final update rule, which is the one cited as the Newton's method.

$$H = \nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \end{pmatrix}$$

Figure 2.26: The Hessian matrix of the error function with respect to the weights.

The Newton's Method seems very efficient computationally because it calculates the quadratic approximation around the solution and immediately finds the minimum of that curvature, instead of fitting a plane to the solution, like the Gradient Descent algorithm. The problem is that it can become computationally impossible to calculate and store the entire hessian matrix of the function, as the parameters increase. Because of that, the standard Newton's method cannot be applied and used in Artificial Neural Networks, which have thousands or even millions of parameters. There are some variations of this algorithm, however, which can be used with ANNs. One such variation is the Hessian Free Optimization algorithm [46] which, instead of calculating and storing the entire Hessian matrix, calculates an approximation that requires less computational resources and does not have to be stored. This algorithm will be discussed in more detail in section 2.3 below.

2.3 Hessian Free Optimisation (HFO)

2.3.1 Intro to HFO

As mentioned in the previous section, the Newton's method, as a second order optimization algorithm, can achieve faster and more accurate convergence to the minimum of a function, compared to first order algorithms, like the gradient descent. In high dimensional problems, first order optimization algorithms can be extremely slow or ineffective due to a problematic phenomenon, called Vanishing Gradient. This phenomenon can be described as a state where the updates for the first layers of a network are very close to zero, because of the backpropagation of the error and the decreasing gradient. As a result, the front layers have almost no information to adjust their weights, which means that the training process becomes slower or even ineffective.

On the other hand, second order optimization algorithms, like Newton's method, calculate the curvature of the error surface (Hessian Matrix) which significantly improves each step

of the optimisation process. What makes these algorithms so efficient is the fact that they attempt to find a quadratic curve that tightly fits at each point, which helps them find the minimum of that curvature immediately, unlike first order algorithms which select a fitting plane and then calculate the next step. However, these second order algorithms have some limits. For instance, in case of a big ANN (with thousands to millions of parameters) sometimes it may not be possible to calculate the Hessian Matrix, due to the extremely high memory requirements. Because of that, several variations of Newton’s method were suggested, like Newton-CG, CG-Steihaug, Newton-Lanczos [47], and Truncated Newton [48], but their applications on machine learning and neural networks have been either extremely limited or not effective at all [49].

The Hessian Free Optimization (HFO) algorithm [46] is a variation of Newton’s method, which uses the local quadratic approximations to generate the suggested updates. Unlike other Newton’s variations, HFO managed to lift the memory constraints, which made it an effective optimisation algorithm for ANNs. This algorithm, instead of calculating and storing the entire Hessian Matrix (H), calculates the dot product of H with an arbitrary vector u (Hu). It takes advantage of mathematical techniques, like finite differences, which computationally costs the same as a single gradient calculation. This means that HFO can calculate the dot products of the Hessian with arbitrary vectors, instead of using the Hessian, and it can optimize the local quadratic objective approximations by using the conjugate gradient (CG) algorithm, to compensate for not having the Hessian Matrix. As mentioned in section 2.2.2.6, the CG method requires N iterations to converge (where N is the number of the network’s parameters), but there are various stopping criteria that allow early termination (after significant progress is made), which reduce the total training time.

Even though, the Hessian Matrix is not calculated in HFO, there are no approximations, as the Hu product is computed accurately. In the standard Newton’s method the approximated quadratic is fully optimized, while the HFO does not perform complete optimization with the un-converged CG algorithm [46] and this is the only difference between the two approaches. The difference between the accuracy of Newton’s method and the HFO with the not fully converged CG is that small that makes it insignificant, where the benefits in terms of efficiency of the HFO (by not calculating the full Hessian Matrix) are obvious.

It is important to note that instead of the Hu product, the Gu product is used, where G is the Gauss-Newton Matrix (an approximation of the Hessian Matrix) [50]. It might look pointless to use an approximation instead of the actual matrix, however, the Gauss-Newton matrix bypasses possible problems that can occur with the use of the Hessian,

which could make it completely ineffective during the training process. Even for the cases where these problems do not appear, the G matrix provides better results, in terms of search directions, which lead to lower memory consumption (about half) and higher running speeds compared to the H matrix.

2.3.2 Analysis of HFO

A detailed analysis of how HFO works was described by Charalambous [23] for anyone interested into diving deeper into this variation of Newton's method. As it was mentioned before, when the H matrix is used, some issues can occur. One of the most important problems is the lack of the utilization of the CG algorithm, on a quadratic model with a non-positive definite curvature matrix, since the Hessian matrix in some cases is non-positive definite. To deal with this issue, the Gauss-Newton matrix is used, which is guaranteed to always be positive semi-definite and is an approximation of the Hessian matrix. Except from that, the Gauss-Newton matrix usually outperforms the Hessian matrix in terms of efficiency.

2.3.3 Hessian-Vector Multiplication evaluation

As mentioned in previous sections, instead of an explicit evaluation of the Hessian matrix, dot products with the Hessian and arbitrary vectors are performed in HFO, which computationally cost the same with a gradient calculation. If the Hessian is considered as the Jacobian matrix of the gradient, based on the definition of directional derivatives, the $H(w)v$ product is the directional derivative of the gradient $\nabla f(w)$ in the direction v (Equation 2.14).

$$H(w)v = \lim_{\varepsilon \rightarrow 0} \frac{\nabla f(w + \varepsilon v) - \nabla f(w)}{\varepsilon} \quad (2.14)$$

In practice, finite-differences suffer from numerical errors, which are troublesome for training ANNs. To counter this issue, a method called 'Forward Differentiation' was proposed [51] and was adapted for ANN training [52]. The main idea was to repeat the chain rule for the value of each node of the gradient, and in order to do that an $R_v(x)$ operator was defined to denote the directional derivative of x in the direction v .

$$R_v X = \lim_{\varepsilon \rightarrow 0} \frac{X(w + \varepsilon v) - X(w)}{\varepsilon} = \frac{\partial X}{\partial w} v \quad (2.15)$$

The R operator is a derivative operator, so it obeys the usual rules of differentiation (2.16):

$$\begin{aligned}
R_v(X + Y) &= R_vX + R_vY && \text{linearity} \\
R_v(XY) &= (R_vX)Y + XR_vY && \text{product rule} \\
R_v(h(X)) &= (R_vX)h'(X) && \text{chain rule}
\end{aligned} \tag{2.16}$$

If these rules are repeated recursively in the gradient calculation algorithm the Hv product will be computed. The algorithm for a simple gradient evaluation is illustrated in algorithm 4 (where $L(y_l; t_l)$ is one of the loss functions of table 2.4), while the algorithm 5 (where $L(y_l; t_l)$ is one of the loss functions of table 2.4) shows the modified version, where the differentiation rules are used to calculate the $H(w)v$ product. The algorithm 6, illustrates how the $G(w)v$ product is calculated and obviously it is simpler than algorithm 5 [49].

```

input:  $y_0; \theta$  mapped to  $(W_1, \dots, W_{\ell-1}, b_1, \dots, b_{\ell-1})$ .
for all  $i$  from 0 to  $\ell - 1$  do
     $x_{i+1} \leftarrow W_i y_i + b_i$ 
     $y_{i+1} \leftarrow s_{i+1}(x_{i+1})$ 
end for
 $dy_\ell \leftarrow \partial L(y_\ell; t_\ell) / \partial y_\ell$  ( $t_\ell$  is the target)
for all  $i$  from  $\ell - 1$  downto 0 do
     $dx_{i+1} \leftarrow dy_{i+1} s'_{i+1}(x_{i+1})$ 
     $dW_i \leftarrow dx_{i+1} y_i^\top$ 
     $db_i \leftarrow dx_{i+1}$ 
     $dy_i \leftarrow W_i^\top dx_{i+1}$ 
end for
output:  $\nabla f(\theta)$  as mapped from  $(dW_1, \dots, dW_{\ell-1}, db_1, \dots, db_{\ell-1})$ .

```

Algorithm 4: Algorithm for computing the gradient of a FFNN [49].

```

input:  $v$  mapped to  $(RW_1, \dots, RW_{\ell-1}, Rb_1, \dots, Rb_{\ell-1})$ 
 $Ry_0 \leftarrow 0$  (since  $y_0$  is not a function of the parameters)
for all  $i$  from 0 to  $\ell - 1$  do
     $Rx_{i+1} \leftarrow RW_i y_i + W_i Ry_i + Rb_i$  (product rule)
     $Ry_{i+1} \leftarrow Rx_{i+1} s'_{i+1}(x_{i+1})$  (chain rule)
end for
 $Rdy_\ell \leftarrow R \left( \frac{\partial L(y_\ell; t_\ell)}{\partial y_\ell} \right) = \frac{\partial \{ \partial L(y_\ell; t_\ell) / \partial y_\ell \}}{\partial y_\ell} Ry_\ell = \frac{\partial^2 L(y_\ell; t_\ell)}{\partial y_\ell^2} Ry_\ell$ 
for all  $i$  from  $\ell - 1$  downto 0 do
     $Rdx_{i+1} \leftarrow Rdy_{i+1} s'_{i+1}(x_{i+1}) + dy_{i+1} R \{ s'_{i+1}(x_{i+1}) \}$  (product rule)
     $\phantom{Rdx_{i+1} \leftarrow} = dy_{i+1} s''_{i+1}(x_{i+1}) Rx_{i+1}$  (chain rule)
     $RdW_i \leftarrow Rdx_{i+1} y_i^\top + dx_{i+1} Ry_i^\top$  (product rule)
     $Rdb_i \leftarrow Rdy_i$ 
     $Rdy_i \leftarrow RW_i^\top dx_{i+1} + W_i^\top Rdx_{i+1}$  (product rule)
end for
output:  $H(w)v$  as mapped from  $(RdW_1, \dots, RdW_{\ell-1}, Rdb_1, \dots, Rdb_{\ell-1})$ .

```

Algorithm 5: Algorithm for computing the $H(w)v$ product in a FFNN [49].

input: $RW_1, \dots, RW_{\ell-1}, Rb_1, \dots, Rb_{\ell-1}$.
 $Ry_0 \leftarrow 0$ (y_0 is not a function of the parameters)
for all i **from** 1 **to** $\ell - 1$ **do**
 $Rx_{i+1} \leftarrow RW_i y_i + W_i R y_i + Rb_i$ (product rule)
 $Ry_{i+1} \leftarrow Rx_{i+1} s'_{i+1}(x_{i+1})$
end for
 $Rdy_\ell \leftarrow \frac{\partial^2 L(y_\ell; t_\ell)}{\partial y_\ell^2} Ry_\ell$
for all i **from** $\ell - 1$ **downto** 1 **do**
 $Rdx_{i+1} \leftarrow Rdy_{i+1} s'_{i+1}(x_{i+1})$
 $RdW_i \leftarrow Rdx_{i+1} y_i^\top$
 $Rdb_i \leftarrow Rdx_{i+1}$
 $Rdy_i \leftarrow RW_i^\top dx_{i+1}$
end for
output: $(RdW_1, \dots, RdW_{\ell-1}, Rb_1, \dots, Rb_{\ell-1})$.

Algorithm 6: Algorithm for computing the G(w)v product in a FFNN [49].

Name	$L(z; t)$	$\nabla L(z; t)$	$L''(z; t)$
Squared error	$\frac{1}{2} \ p - t\ ^2$	$-(p - t)$	I
Cross-entropy error	$-t \log p - (1 - t) \log(1 - p)$	$-(p - t)$	$\text{diag}(p(1 - p))$
Cross-entropy error (multi-dim)	$-\sum_i [t]_i \log[p]_i$	$-(p - t)$	$\text{diag}(p) - pp^\top$

Table 2.4: Derivatives and Hessians of typical loss function.

Chapter 3

Data Manipulation

3.1	PSSP Metrics	44
3.2	Protein Databases and DSSP	45
3.3	Dataset Format	46
3.4	Data Encoding and MSA profiles	47
3.5	CB513 and PISCES Datasets	48
3.6	Dataset preprocessing with MSA profiles	49
3.7	Significant neighboring amino acids	51
3.8	Training/ Testing Set and Cross Validation	52
3.9	Ensembles	53
3.10	Filtering	54
3.10.1	External Rules	54
3.10.2	Support Vector Machines	55
3.10.3	Decision Trees	57
3.10.4	Random Forests	58

3.1 PSSP Metrics

The Protein Secondary Structure Prediction (PSSP) problem concentrates on predicting, as accurately as possible, the secondary structure of proteins based on their primary structure. In this thesis project supervised learning methods are utilised, which require both the input data (primary structure) and output data (secondary structure) to train an Artificial Neural Network (ANN) to make predictions. Supervised learning is like learning with a teacher. The input is presented to the ANN which attempts to predict the output and then receives feedback on whether its predictions were correct or not. This way the network can adjust the weights accordingly to improve the prediction results. Both primary and secondary structure data must be encoded in a way that can be fed into the network.

To get an indication of how good are the predictions of the trained models two metrics were used, the per residue Q3 accuracy and the Segment Overlap (SOV), which are commonly used for the PSSP problem. The Q3 accuracy measures the number of correctly classified amino acids divided by the number of total amino acids (Equation 3.1, where n is the number of amino acid residues and m_i takes the value of 1 if the predicted value of the i^{th} amino acid residue is correct and 0 otherwise). The Segment Overlap (SOV) [53] score is used to measure how good are the predicted results for each class and the general structure of the entire protein. More specifically, unlike Q3, SOV considers the size of continuous overlapping segments and assigns extra allowance to longer continuous overlapping segments (instead of just checking the individual positions, like Q3).

$$Q = 100 \frac{1}{n} \sum_{i=1}^n m_i \quad (3.1)$$

For instance, if the target secondary structure of a protein consists of four (4) helices followed by two (2) coils and then another four (4) helices and the prediction has only ten (10) helices, Q3 and SOV will produce different accuracy values. The Q3 accuracy will be 80%, as eight of the ten amino acids were predicted correctly, while the SOV score would be just 48. Even though the original SOV score was not a percentage, a modified definition of SOV [54] was suggested, which fixed this issue with normalization techniques.

$$SOV_{\alpha} = \frac{1}{N_{\alpha}} \sum_{s_{\alpha}} \frac{\min OV(s_1, s_2) + \delta(s_1, s_2)}{\max OV(s_1, s_2)} \quad (3.2)$$

The SOV score for the α -helix can be calculated with equation 3.2, where s_1 and s_2 are the actual and predicted segments of the secondary structure of the α -helices, respectively.

The s_a is the number of segment pairs (s_1, s_2) , where s_1 and s_2 have at least one common residue α -helix. The $\min OV(s_1, s_2)$ is the length of the overlap between s_1 and s_2 , and the $\max OV(s_1, s_2)$ is the length of the total area for which one of the s_1 and s_2 has one residue of type α -helix. The N_α is the total number of residues of type α -helix. The calculation of $\delta(s_1, s_2)$ is based on equation 3.3.

$$\delta(s_1, s_2) = \min \left\{ \begin{array}{l} \max OV(s_1, s_2) - \min OV(s_1, s_2) \\ \min OV(s_1, s_2) \\ \text{int}(0.5 \times \text{len}(s_1)) \\ \text{int}(0.5 \times \text{len}(s_2)) \end{array} \right\} \quad (3.3)$$

3.2 Protein Databases and DSSP

There are several protein databases, like the iProClass (Protein Information Resource), PDBe (Protein Data Bank in Europe), PDBj (Protein Data Bank in Japan) and RCSB (RCSB Protein Data Bank), which include various information about millions of proteins. This information includes protein names, length, structures (primary, secondary, tertiary and quaternary) and other biological information related to proteins. The protein information included in the datasets of the PSSP problem was extracted from these databases.

Secondary Structure	8 class code	3 class code
α-helix	H	H
3-helix	G	
π-helix	I	
β-strand	E	E
β-bridge	B	
β-turn	T	C
bend	S	
Random coil	C	

Table 3.1: Table with the secondary structure abbreviations, grouped in 8 and 3 classes

The Dictionary of Secondary Structure of Proteins (DSSP) [55] defined a standardized format for categorising the secondary structures of proteins. According to this format, there are eight (8) distinct classes of secondary structures, based their shape, which are represented by a capital letter of the English alphabet. These are the α – *helix* (H), 3-helix

(G), π -helix (I), β -strand (E), β -bridge (B), β -turn (T), bend (S), and random coil (C) for residues which are not in any of the other conformations (Table 3.1). Usually these eight (8) categories are grouped into three (3) more general categories, which describe the shape of a specific local segment of the protein. For the purpose of this dissertation, the 3-class classification is used, which includes the helix (H) conformations, containing the first three categories (H, G, I), the sheet (E) conformations, containing the next two categories (E, B), and the Coil (C) conformations, containing the rest categories (T, S, C).

3.3 Dataset Format

The protein datasets, that were used for training, had records of a 3-line format per protein. An example is shown in figure 3.1, where the first line contains the protein name, the second line the primary structure and the third line the secondary structure of the protein. The protein name can be used to combine the primary and secondary structures with the Multiple Sequence Alignment (MSA) [8] profiles. The primary structure corresponds to the sequence of amino acids of each protein and each letter represents one amino acid. The secondary structure, located in the third line, is the target output which must be predicted by the network, and each letter represents the class of each amino acid (based on table 3.1).

	Protein Name	Primary Structure	Secondary Structure
13	1bdsA_1-43		
14	AAPCFCSGKPGRGDLWILRGTCPPGGYGYTSNCYKWPNICYPH		
15	CCCCCCCCCCCCCEEECCCCCCCCCCCCCEEECCCCCCCC		

Figure 3.1: Protein representation example for protein 1bdsA_1-43.

To prepare training and validation datasets, a python program was developed, which creates files with comma-separated values (CSV files) based on the input datasets (in the form of figure 3.1) and MSA profiles, which will be discussed in the following sections. This program gives the ability to process multiple datasets at once by including the names of the datasets in the 'datasets' variable, located at the top section of the program. Moreover, if the MSA profiles for some proteins are missing, the program will ignore these proteins and print their names on the screen. This program for CB513 can be found in appendix D and for PISCES can be found in appendix E.

3.4 Data Encoding and MSA profiles

It is suggested the input and output data, that is used to train ANNs and most machine learning algorithms, to be normalised before they are used in training. The new encoded data should consist of real values between zero and one (0,1), or between minus one and one (-1,1), according to the selected activation function. The reason behind this suggestion is to speed up the learning process and help the network reach convergence faster.

For the PSSP problem a suggested encoding method is to use the Multiple Sequence Alignment (MSA) profiles along with the protein datasets (that include protein names, primary and secondary structures). These MSA profiles, in general, include information about about DNA and RNA protein sequences, and are very popular in the field of Bioinformatics. In many cases, the proteins that are selected to create the MSA profile have an evolutionary relationship with each other and are descended from a common ancestor. Because of that, these proteins are presumed to have the same secondary and tertiary structure [8]. The amino acids of these proteins are aligned together and are encoded in a way such that each position of their sequence represents the probability of the appearance of each amino acid, to form an MSA profile. An example of this alignment process is illustrated in figure 3.2.



Figure 3.2: Process of MSA profiling

However, the alignment of three or more biological sequences is extremely hard and time consuming when done manually. Because of that, computational algorithms have been developed to analyse and align these amino acid sequences. These algorithms use heuristics to find an approximation of the alignment, as the optimal alignment is computationally expensive. Furthermore, the MSA files that are created for each protein contain N rows (where N is the number of amino acids of the protein) and 20 columns, where each column represents the probability of each amino acid (from the 20 known types) appearing

in that specific position in the protein sequence.

In the example of figure 3.2, the highlighted column would have 89% (8/9) for the V amino acid and 11% (1/9) for the E amino acid, while all other amino acids would have zero (0) values. The values of amino acids for each line must add up to 100 and before they are used for training an ANN they must be normalised in the range (0,1), which can be easily done by dividing them with 100. Therefore, the ANN would be able to adapt to the structure of the data more easily.

3.5 CB513 and PISCES Datasets

In general, in order to successfully train a prediction model the datasets, which will be used for the training phase, must be preprocessed. In this phase data selections and data cleaning techniques are performed. There are various datasets for the PSSP problem, that have been created and preprocessed over the years. For the purpose of this dissertation, two widely used datasets were selected, the CB513 dataset [12] and the PISCES dataset [56]. These two datasets were chosen because they have been used for the PSSP problem by many researchers, which makes the comparison of the results possible and gives an indication of how well trained is the neural network. To prevent the network from memorising the order of the input patterns, a good practice is to shuffle the input patterns on each epoch, and therefore get better prediction results.

Initially, the smaller dataset was used, CB513, which has 513 unique proteins, from which eight (8) were excluded (these can be found in Appendix A), due to the fact that their MSA files included only zeros. This dataset required less time to train and helped to identify whether the neural network was able to learn how to predict the secondary structure of proteins or not. In the next phase, the bigger dataset was used, PISCES, which consists of about 8500 sequences, from which 341 were excluded because their MSA files were either corrupted or zeroed and another 16 were excluded due to missing MSA profiles (all of them are shown in Appendix B). The bigger dataset (PISCES) was utilized because in many machine learning problems, by feeding the model with more data, the predictions become more accurate.

Finally, the PISCES dataset was not in the expected form, that was mentioned before (protein name, primary structure, secondary structure) but fortunately a team of University of Cyprus implemented a Java program, which was able to convert the PISCES data into the expected format. These new modified PISCES datasets were provided by Dionysiou ([1], [24]), who also worked on the PSSP problem in the past.

3.6 Dataset preprocessing with MSA profiles

Convolutional Neural Networks (CNNs) expect their input as a two-dimension (2D) or three-dimension (3D) array, so in order to train a CNN to predict the secondary structure of proteins, the training and test data must be presented in the form of 2D or 3D arrays. The input representation method, which will be used, is the same with the one used by [24].

Over recent years, the Multiple Sequence Alignment (MSA) profiles approach was used by many researchers. In the Bioinformatics sector, the sequence alignment is a well known approach, that refers to sequences of DNA, RNA or proteins. In general, this method attempts to find similarities between these types of sequences, which can usually define some biological association, leading to a better understanding of the biological mechanisms. An example of an MSA file is illustrated in figure 3.3.

```
0 0 0 100 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
29 3 15 0 9 0 0 0 25 3 2 1 0 0 1 7 4 0 1 0
0 0 0 0 0 0 0 0 1 4 1 10 0 2 3 28 3 43 3 2
37 5 58 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
12 4 4 0 8 0 45 0 0 0 0 14 0 0 3 2 0 7 0 0
3 3 31 11 0 0 0 30 0 0 0 15 0 0 7 0 0 0 0 0
2 0 11 7 10 13 47 0 0 0 0 1 0 0 6 2 1 0 0 0
10 1 2 2 0 0 0 1 1 17 23 17 1 0 2 3 5 3 0 11
1 1 2 4 0 0 0 10 1 6 13 34 0 0 1 15 0 9 0 2
5 4 3 1 1 0 1 4 8 6 10 16 0 1 4 3 3 12 6 13
```

Figure 3.3: Example of the encoded form of an MSA file [24].

Since the input data of a CNN must consist of 2D (or 3D) arrays, these MSA files must be visualised in a way so that they can be used to successfully train a CNN. In order to achieve that, all the MSA files were combined in a single file and the desired output label was added at the end of each record.

For instance, the information included in two MSA files, named ‘1bdoa_77-156’ and ‘1bfga_19-144’, before the two files were combined together, is shown in figure 3.4. Figure 3.5 shows the new encoded file (where all spaces were replaced with commas ‘,’), after combining the two MSA files, where the red line separates the data of the first file from the data of the second file. Each line now has twenty one (21) numbers (columns) instead of twenty (20), as the predicted class was added (C:0, E:1, H:2) at the end of each

line, according to the dataset files (with the protein name, primary and secondary structures), described in section 3.3. This means that both PSSP datasets and MSA profiles were utilized to create the new dataset files.

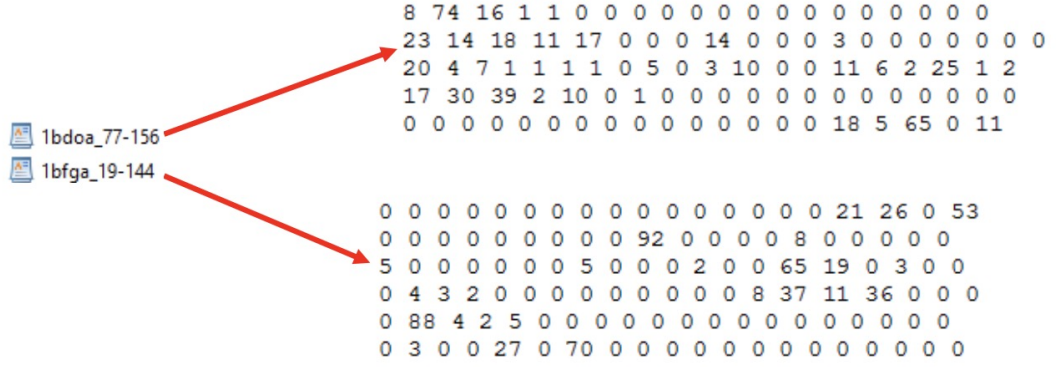


Figure 3.4: The image shows the MSA file (before collapsing into a single file) [24].

By using the above technique, a new file was created, which included all the data from the MSA profiles and the desired labels of the secondary structure for each record (Figure 3.5). The same method was used for both training and test datasets. This new representation with the MSA profiles can be used to successfully train a CNN, since these files can be presented in a 2D table format. The CNN will receive as input one record (one line) of the new dataset at a time and will attempt to predict the secondary structure, using the output class representation mentioned earlier (C:0, E:1, H:2).

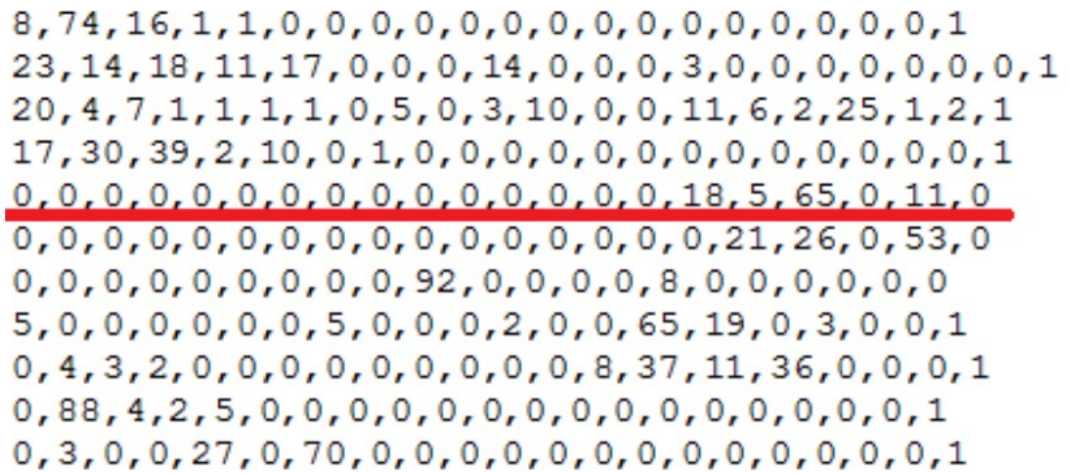


Figure 3.5: The encoding of the new file, after combining the MSA files into a single file [24].

3.7 Significant neighboring amino acids

The sequence of amino acids plays a major role, as it determines the interactions that take place and the folds that are formed in the secondary structure of a protein. The secondary structure of an amino acid is significantly affected by the adjacent amino acids (previous and next amino acids), according to the distance between them (short distance usually means bigger impact, while longer distance means less impact) [20].

The training and test datasets were modified to take advantage of this interaction between the neighboring amino acids. For each record, except from the information about each amino acid (from the MSA profiles) and the expected output class, described in the previous section, the information of k -neighboring amino acids was added (where k is an integer variable). For example, if k is equal to one, each record will consist of the MSA records of the left amino acid, the MSA records of the current amino acid, the MSA records of the right amino acid and the target label (class) of the current amino acid (which is located in the middle). If an amino acid is not preceded (first amino acid in a sequence) or followed (last amino acid in a sequence) by another amino acid, zero values are added instead (zero padding) to ensure that all records have the same length. An example where k is one (1) is illustrated in figure 3.6, for a sequence of six (6) amino acids.

[illegible]

Figure 3.6: MSA record for a sequence of 6 amino acids.

The new modified dataset has 61 ($20 \times 3 + 1$) numbers for each record, as shown in figure 3.7.

[illegible]

Figure 3.7: Modified MSA record for a sequence of 6 amino acids.

As soon as this new representation is fed into the CNN it is rearranged into a 2D array, where each line includes the MSA profile vector for each amino acid and the last value will be the target label (class). Figure 3.8 illustrates an example for the data representation method after it was rearranged into a 2D array, for a window size of 15 amino acids (or 'plus7'), where each row represents the vector of the MSA profile for the specific amino acid and the SS label represents the class (H, E, C) of the middle amino acid.

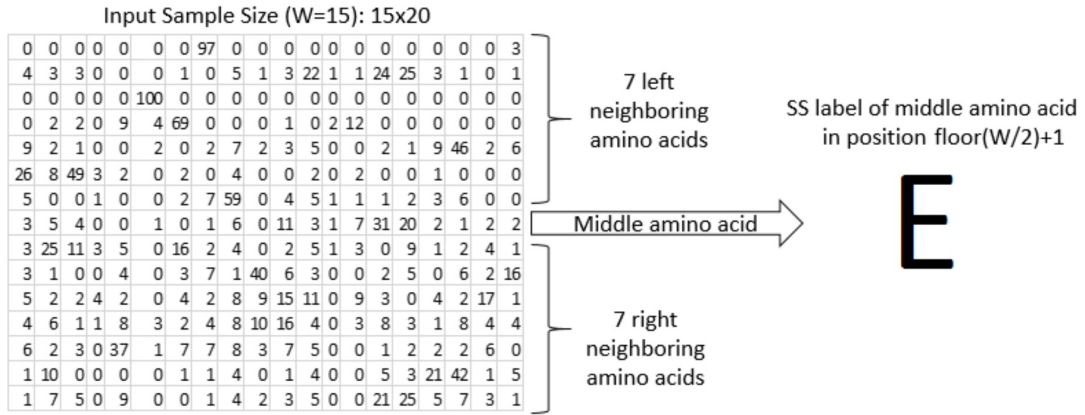


Figure 3.8: An example of input data representation for a window size of 15 (or $k = 7$) amino acids [57].

In order to create these modified datasets two Python programs were developed, that use the CB513 (Appendix D) and PISCES (Appendix E) datasets along with the MSA profiles to prepare the modified datasets, according to the provided 'plus' value, which can take any positive integer value (>0). In the previous example the plus value was one (1), as the neighboring amino acids were one for each side (one left and one right) of each amino acid. These programs prepare multiple datasets (or folds) at once and print on the screen the names of proteins for which the MSA files were not found.

This technique can improve the accuracy of a CNN [24], as the CNN can identify the neighboring amino acids, which can affect the secondary structure of an amino acid.

3.8 Training/ Testing Set and Cross Validation

To train an ANN a specific set of data is required, called training dataset, which is used for training the model so that it can extract features from the input patterns and classify these patterns into a number of classes. However, it is very important to ensure that the model was able to generalize the extracted knowledge so that it can predict patterns that has not 'seen' before. For this reason, another set of data is used, called test dataset,

which is completely different from the training dataset, and its purpose is to measure the effectiveness of the network to classify new data, that has never seen before. In general, a good rule of thumb is to split the entire dataset into 80% for the training dataset and 20% for the test dataset (80-20 rule). However, in different problems, other splitting criteria can be used, which may lead to better results.

Sometimes this method is not enough to test the ability of a network to predict new data, since the accuracy depends on a specific test dataset. A method that can be used to address this issue is to evenly split the data into N folds and train N different models. Each model will have a unique fold selected as the test dataset and the rest N-1 folds will be used as the training dataset. This method is called N-fold cross validation (Figure 3.9) and the cross validation accuracy is equal to the average test accuracy of all models.

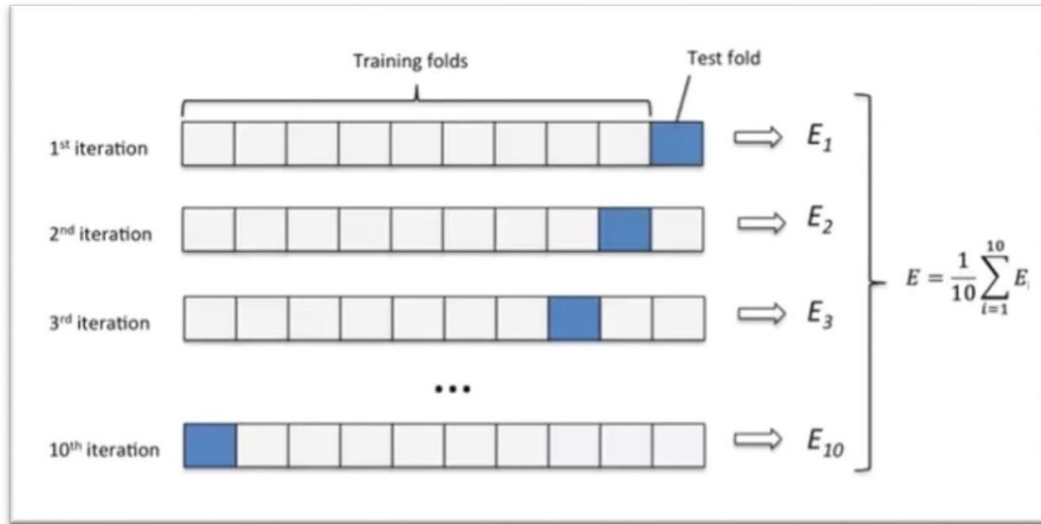


Figure 3.9: 10-fold cross validation

3.9 Ensembles

Ensemble learning is a method which can be used to improve the performance of a machine learning model. According to this method, instead of training just a single model, multiple models can be trained and then their results can be combined somehow, to improve the final results.

Even though, there are various ensemble methods, ranging from simple to advanced and more sophisticated methods, for the purpose of this dissertation a relatively basic approach was used, also known as averaging ensemble method. This method basically calculates the average of the outputs of its models. In the PSSP problem, for instance, if there are

five different trained models the following steps are applied. First, for each of the input-output pairs, the output of each of the five models is calculated and classified in one of the three classes (H, E, C). Then, the results of each of the five models are compared with the method ‘winner takes all’ and the class that had the most appearances is chosen as the final class for that input. If there is a tie between some of the classes, an arbitrary class (from those) is selected as the final class.

This ensemble method, even if it is very simple, it can remove random errors from the models, which can lead to improved results. More advanced ensemble methods might have a bigger impact on the predicted results, but at the cost of computation resources, as these are usually more complex.

3.10 Filtering

3.10.1 External Rules

Post-processing filtering is an additional method that is used to improve the accuracy of a model. The applied filtering method can be problem specific or more generic, with the use of different learning algorithms. Both methods were used in this dissertation, which affected, by a small amount, the final accuracy (Q3 Score) and the quality of the predictions (SOV score).

The first filtering method was based on a set of external rules, that are specific for the PSSP problem. These rules are based on empirical observations and were used to ‘fix’ the quality of the results (SOV), rather than improving the overall accuracy (Q3).

The external rules applied are (where H, E and C are the three possible classes):

1. Single ‘H’ or ‘E’ are replaced with ‘C’
2. Sequence ‘HEEH’ is replaced with ‘HHHH’
3. Sequence ‘HEH’ is replaced with ‘HHH’
4. Sequence ‘!HH!’ is replaced with ‘!CC!’

These simple rules can be applied extremely fast and can increase the SOV score, while sometimes can slightly drop the Q3 accuracy.

3.10.2 Support Vector Machines

In 1995, Cortes and Vapnik suggested the use of Support Vector Machines (SVMs), in Machine Learning. Initially SVMs were used for binary classification problems and their purpose was to find hyperplanes that best divide a dataset into classes [58]. If the data cannot be separate linearly, SVMs attempt to map the data into a higher dimension using a non-linear kernel function. These kernel functions are very effective and efficient as they just compute inner products. This transformation to a higher dimension is more likely to make the data linearly separable. Figure 3.10 illustrates four popular SVMs that are currently used, along with their kernels.

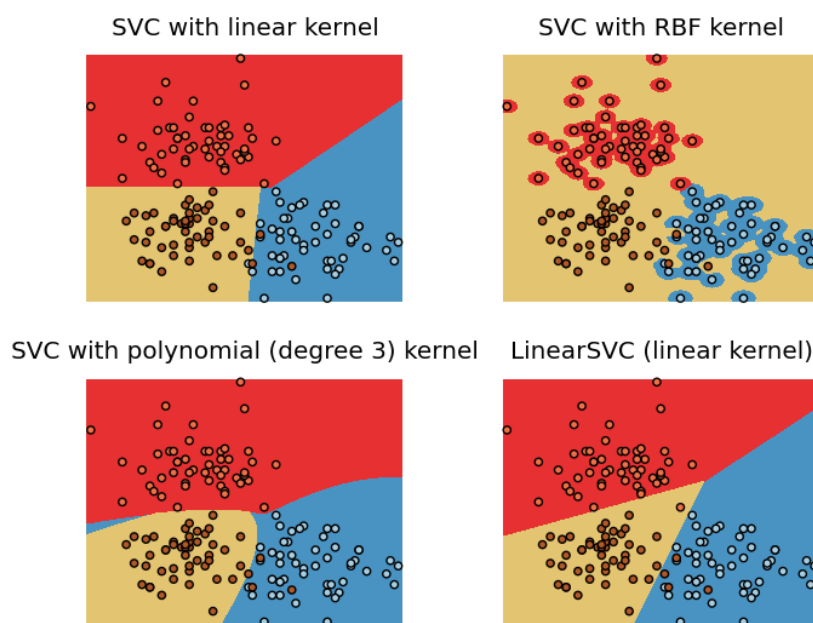


Figure 3.10: Results of different kernels for a 3-class classification problem.

In order to best separate two linearly separable classes by finding the optimal hyperplane, SVMs attempt to maximize the distance between the points, that are closer to the hyperplane, for each different class. The points that are located near the limits of this separation are called support vectors and the points that are located in the area of multiple classes (overlapping classes) are not taken into consideration, in order to create a more generalized model. Figure 3.11 shows three possible separation lines (A, B and C), where the blue star and red circle that are connected with them are considered the support vectors and line C is considered the optimal hyperplane. In figure 3.12, even though the separation of the initial data (left plot) is hard, if they are projected in a higher dimension they

can be easily separated by a hyperplane (right plot).

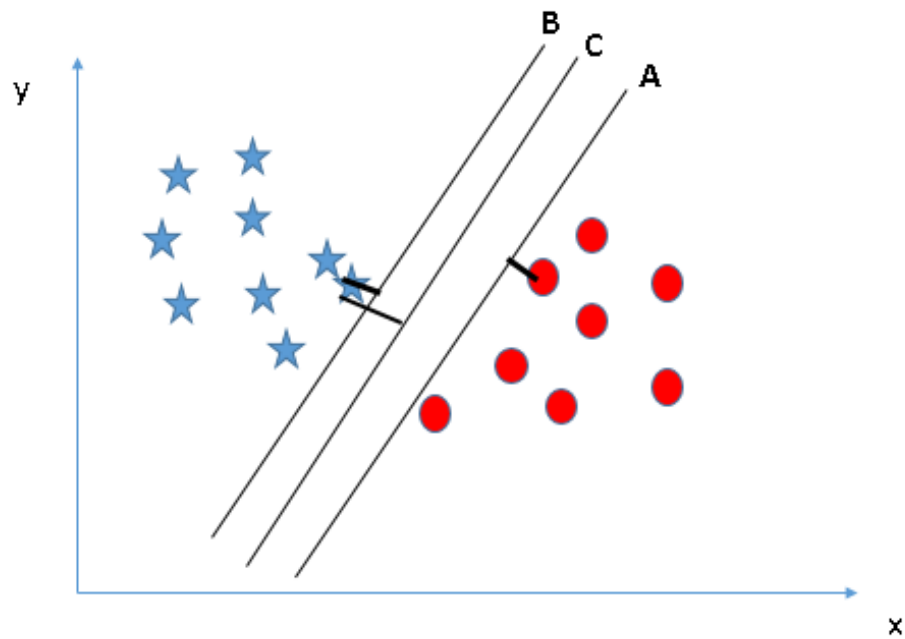


Figure 3.11: SVM example of a linearly separable problem.

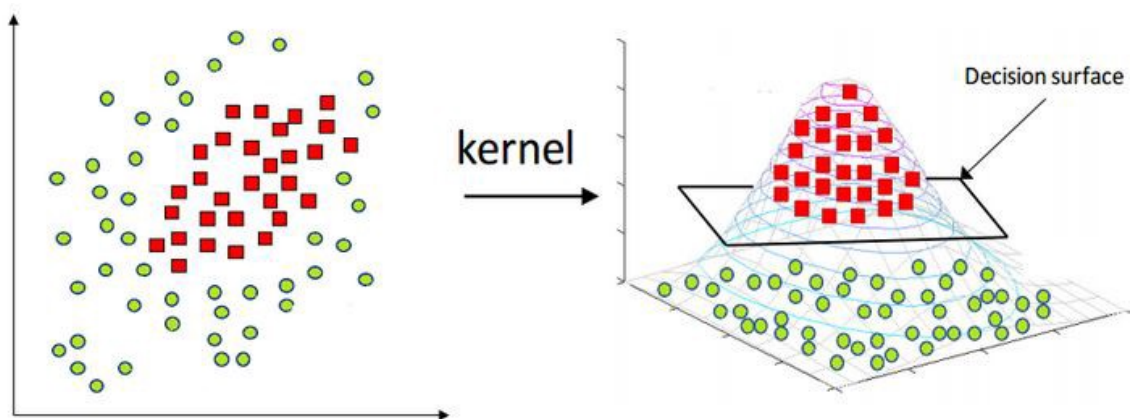


Figure 3.12: SVM projecting a problem in a higher dimension.

In particular, SVMs had very good filtering results for the PSSP problem [59]. More specifically, they were used by Dionysiou [24] and Dionysiou et al. [57] and their good final results makes them very promising.

3.10.3 Decision Trees

Decision trees are most commonly known for their use in operations research, and more specifically in decision analysis, but are also a popular tool for machine learning. They can be used to identify a strategy that is most likely to reach a target goal. A decision tree is defined as a support tool, with a tree-like shape, which models decisions and possible consequences, including resource costs, chance events outcomes and utility.

The best way to explain how a decision tree operates is through a simple example. A scenario, where a dataset contains numbers with different features, is illustrated in figure 3.13. There are two 1s and five 0s, which represent the two classes. The goal is to separate the data using their features, which are color (red or blue) and whether the number is underlined or not.

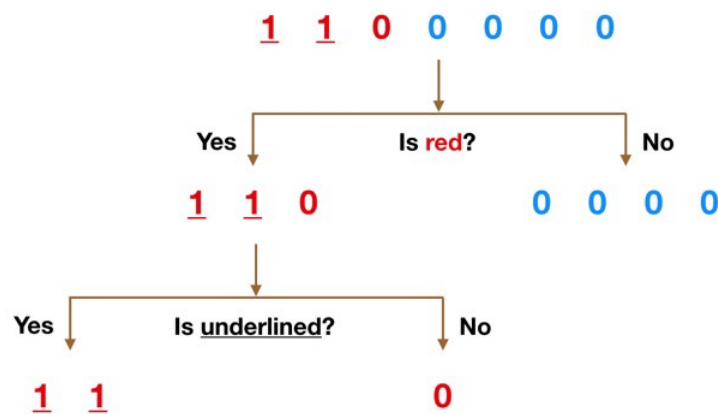


Figure 3.13: Example of simple decision tree [60].

Obviously, the color feature can be used to split the data, as only one of the 0s is red, while the rest are blue. The question ‘Is it red?’ can be used to split the first node. A node in a tree is like a point where the path splits into two branches, where the data that meet the criteria go under the ‘Yes’ branch and ones that do not go under the ‘No’ branch, as shown in figure 3.13. The ‘No’ branch contains only blue 0s that are not underlined, which means no further splits should be made. On the other hand, the ‘Yes’ branch contains data that have different features, so the question ‘Is it underlined?’ can be used to split the red data. The two underlined 1s go under the ‘Yes’ subbranch, while the not-underlined 0 goes under the ‘No’ subbranch. At this point no further splits of the data are required.

Even though, in real life examples the data will not be as clean as the one used in this example, the applied logic of a decision tree remains the same. A decision tree will decide at each node which feature can split the observations into two groups in a way that

the differences are maximised, while maximising the similarities between the members of each subgroup.

3.10.4 Random Forests

A random forest is a classification algorithm that consists of a large number of individual decision trees that function as an ensemble. Each individual tree, outputs a class prediction and the class with the most votes is selected as the prediction of the random forest (Figure 3.14).

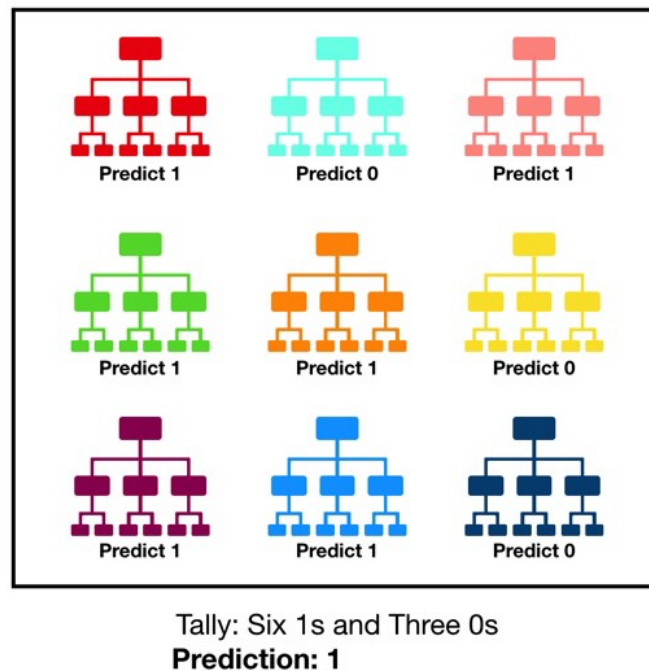


Figure 3.14: Example of random forest prediction [60].

The reason, a random forest model works so well, in data science, is because a large number of unrelated models, that operate as a group, can outperform any of the individual models. One of the most important things in a random forest is the low correlation between the individual models (trees), since the trees ‘correct’ each others’ errors, as long as they do not make the same mistakes in the same direction. In order for a random forest to have good predictive results, there must be an actual signal that helps the models adapt to the features of the data and the correlation, between the predictions of the individual trees, must be as low as possible.

In order to illustrate why uncorrelated predictions are so important, a simple example will be used. In a gambling game a uniformly distributed random generator is used to produce a number between 1 and 100. If the number is above 40 the player wins and earns money

based on the bet amount, which means the player has 60% chance to win. The player has three options, play 100 games betting \$1 per game (choice1), play 10 games betting \$10 per game (choice2) or play 1 game and bet \$100 (choice3). Below are the expected values for all three options:

$$ExpectedValue(choice1) = (0.60 * 1 + 0.40 * (-1)) * 100 = 20$$

$$ExpectedValue(choice2) = (0.60 * 10 + 0.40 * (-10)) * 10 = 20$$

$$ExpectedValue(choice3) = 0.60 * 100 + 0.40 * (-100) = 20$$

It is obvious that all options have the same expected value, which makes it difficult to choose. A visualization of a Monte Carlo simulation could reveal the distributions of the available options. Figure 3.15 illustrates the distribution of the outcome of 10000 simulations for each of the three options. The three options, even though they share the same expected value, they have completely different outcome distributions. With the first option (choice1) there is 97% chance to make money, while for the other two options (choice2 and choice3) the chance to make money is 63% and 60%, respectively. It seems that the more the \$100 bet is split up, the higher the chance for the player to make money, as each game does not dependent on the other games.

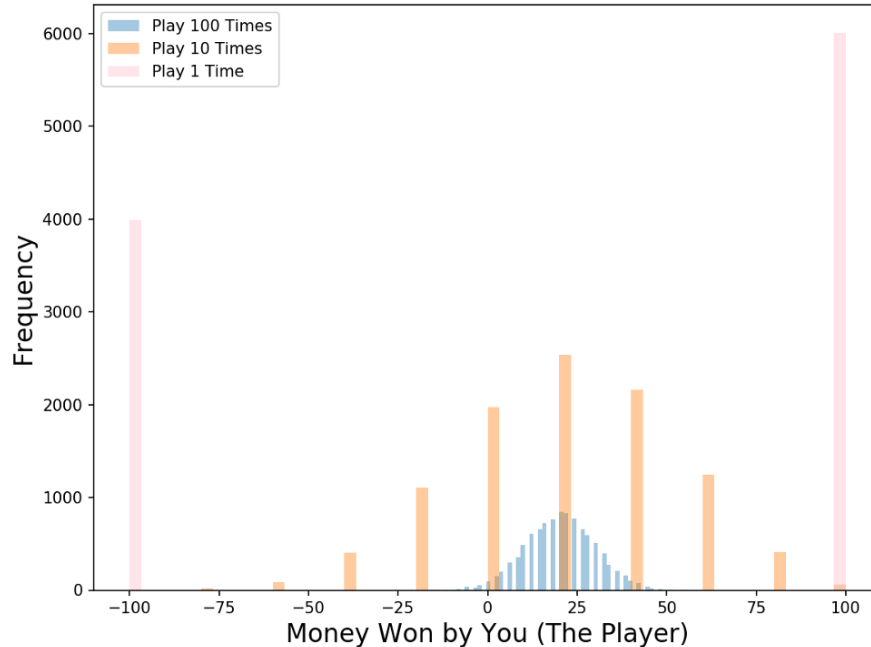


Figure 3.15: Distribution of the outcomes of 10000 simulations for each option [60].

A random forest works in the same way, with the game mentioned above. The higher the number of uncorrelated trees, the higher the chance of making correct predictions.

To ensure that each individual tree is uncorrelated with the other trees, a random forest uses two methods, bagging and feature randomness. The first one (bagging), randomly selects a sample from the dataset for each individual decision tree, instead of using the entire dataset. The second method (feature randomness), restricts the number of features that can be used to split a node in each decision tree, by selecting a random subset of the available features. This increases the variation between the individual trees of the model, which results in lower correlation.

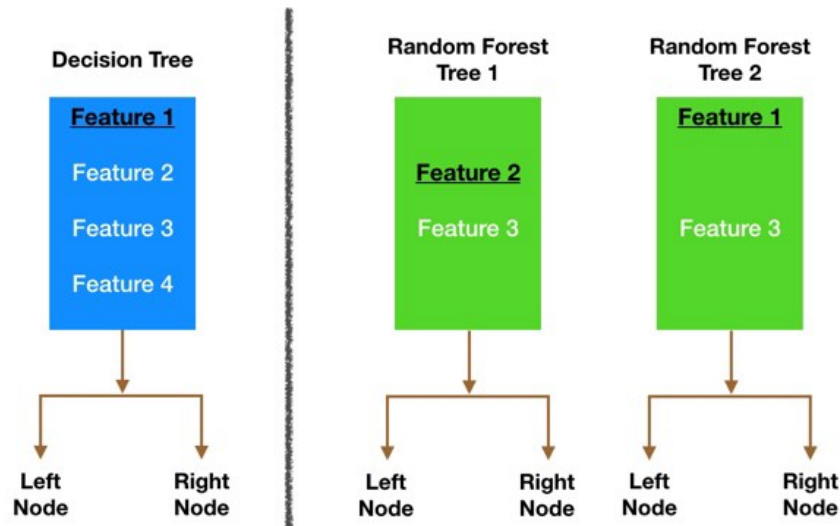


Figure 3.16: Node splitting in a decision tree and a random forest model [60].

In order to make things easier to understand, an example will be illustrated. Figure 3.16 shows a decision tree (blue) and two trees from a random forest (green), where both models can separate the data based on four features. The decision tree chose the Feature 1 to best separate the data into groups. The first tree of the random forest (Tree 1) could only choose between Features 2 and 3, which were selected randomly, to split the data, while the second tree (Tree 2) could only choose between Features 1 and 3. Even though Feature 1 was the best splitting option, only Tree 2 could use it, since it was not included in the available features of Tree 1.

To sum up, bagging helps to create trees that are trained on different sets of data, while feature randomness forces them to use different features to make decisions.

Chapter 4

Implementation

4.1	A new approach for the PSSP problem	62
4.2	CNN and HFO combination	63
4.3	Subsampled Hessian Newton (SHN) Method	64
4.4	Network Implementation	66

4.1 A new approach for the PSSP problem

The PSSP problem can be considered a classification problem, which means that an ANN can be used to predict the secondary structure of proteins. For this dissertation, a Convolutional Neural Network (CNN) was selected because, according to previous attempts [1], it managed to produce very good results (>80%), which makes it very promising.

A thesis dissertation should explore new methods or combinations, in order to provide some value to the world of research. Because of that, the use of just a CNN for the PSSP problem would be a poor choice as it has already been used before and it would not help the research community or the researchers that are involved with the PSSP problem. The main idea was to combine a CNN with the Hessian Free Optimisation (HFO) algorithm (a second order optimiser) to predict the secondary structure of proteins, which has never been attempted before, mainly because of the complexity of second order optimisation algorithms. This optimiser, combined with a simple Feed Forward Neural Network (FFNN), managed to achieve great results for the PSSP problem [2], with more than 80% Q3 accuracy.

Usually CNNs, because of their complexity, contain thousands of parameters and that makes the training process very time consuming on a Central Processing Unit (CPU). For small datasets, like CB513, a few hours (around 4-12 hours depending on the selected settings) would be enough. However, for bigger datasets, like PISCES, the training process could take days to complete. To speed up the training process a Graphics Processing Unit (GPU) was utilized, from the Google's Colab cloud service. Colab is a free Jupyter notebook environment that runs entirely in the cloud, does not require any setup and supports many popular machine learning libraries (paid services are also available). A Jupyter Notebook, also known as the IPython Notebook, is an interactive computational environment based on the web (usually ending with the extension '.ipynb') [61]. This allows users to combine code, comments, graphical visualizations and multimedia, in an interactive document, which can be run via a web browser, hosted on a local machine or even a remote server.

The classification model of this dissertation was implemented on a notebook to ensure portability, remove machine constraints (requirements), as this can be run entirely in Colab even with an 'old' machine using just a web browser. In addition to that, a notebook makes it easier to interact with the program and it comes with some of the most popular machine learning libraries and frameworks, like TensorFlow [62], PyTorch [63] and Scikit-learn [64], which are pre-installed and ready to use.

4.2 CNN and HFO combination

The purpose of this dissertation was to combine a CNN with HFO and train it to predict the secondary structure of proteins. The PyTorch machine learning framework [63] was initially used, along with the fastai library, to implement a CNN and train it for the PSSP problem. This part was successful and the results were around 70% Q3 accuracy, without tuning the hyper parameters. The next step was to implement the HFO algorithm in PyTorch, since the available implementation was written in pure Python and could not interact with the CNN. An alternative option was to implement the CNN from scratch (in pure Python), which would significantly drop the efficiency as the PyTorch framework (and most machine learning frameworks) has its functions written in C++, which is much more efficient.

The implementation of HFO was probably the most difficult task of this dissertation. The HFO implementation which was already available was very complex, which made this task even harder. After many failed attempts, the HFO algorithm managed to train a Feed Forward Neural Network (FFNN) to predict the XOR gate (a toy problem used check if a network is learning effectively). The next step was to try the new implementation on the PSSP, and fortunately the FFNN with HFO managed to extract some patterns from the proteins which resulted in around 72% Q3 accuracy (without any tuning of the hyper parameters). The final step of the implementation was to combine the two sections, the CNN and the FFNN with the HFO optimiser.

Unfortunately, this was not as simple as it seemed at the beginning of this dissertation. The HFO algorithm was specifically designed for a FFNN and not a CNN, which made the updates totally ineffective. An alternative approach could be to use a different optimisation algorithm, like gradient descent, to train the CNN layers and then train the FFNN at the end with the HFO. However, this practice seems pointless, since the purpose of HFO is to replace a different optimizer, not depend on it. Another approach could be to ignore the CNN layers and train only the FFNN with HFO, which would be a waste of resources, since the effectiveness of the CNN layers would not be utilized.

At this point, an ordinary dissertation project would possibly come to an end, as the main purpose was to attempt to combine the HFO algorithm with a CNN. The conclusion was that this was not possible because the HFO was designed explicitly for a FFNN and not a CNN. However, this is not an ordinary dissertation, so despite the tight margins of available time, an alternative approach was pursued with the help of additional research. This seemed to be a great decision as a recent article [3], published in January 2020,

explained why the HFO algorithm is not compatible with CNNs and suggested a variation of HFO, specifically designed for CNNs. This new method, called Subsampled Hessian Newton (SHN) method [3], is discussed in the following section.

4.3 Subsampled Hessian Newton (SHN) Method

There are several studies on Newton methods for training deep ANNs ([65], [66], [67], [46], [68], [69], [70]), but almost all of them used fully connected FFNNs. The Newton methods are very complicated and that is possibly the main reason why CNNs have not been utilized in those studies. Apart from this article [3], there is no evidence, or published documents, that describe how the Newton methods can be applied in deep learning (CNNs) effectively. This made Gradient Descent, and its variations, the most popular optimisation algorithms for CNNs, although the Newton's methods are more robust, more efficient and require less tuning of the hyper parameters (for FFNNs at least).

A new variation of HFO was suggested by Wang et al. (2020) for CNNs which is explained in high detail [3] (mathematical proof included). This new method was labeled as the Subsampled Hessian Newton (SHN) method (Algorithm 7, where (35) is Equation 4.1, (36) is Equation 4.2 and (37) is Equation 4.3). Due to the high complexity of the proof and extensive explanation of this algorithm, it is better to refer to the original paper for a better understanding of the transition from the HFO to the SHN algorithm.

```

Given initial  $\theta$ . Calculate  $f(\theta)$ ;
while  $\nabla f(\theta) \neq \mathbf{0}$  do
    Choose a set  $S \subset \{1, \dots, l\}$ ;
    Compute  $\nabla f(\theta)$  and the needed information for Gauss Newton matrix-vector products;
    Approximately solve the linear system in (36) by CG to obtain a direction  $\mathbf{d}$ ;
     $\alpha = 1$ ;
    while true do
        Compute  $f(\theta + \alpha \mathbf{d})$ ;
        if (35) is satisfied then
            break;
        end
         $\alpha \leftarrow \alpha/2$ ;
    end
    Update  $\lambda$  based on (37);
     $\theta \leftarrow \theta + \alpha \mathbf{d}$ ;
end

```

Algorithm 7: A subsampled Hessian Newton method for CNNs [3].

$$f(\theta + \alpha \mathbf{d}) \leq f(\theta) + \eta \alpha \nabla f(\theta)^T \mathbf{d} \quad (4.1)$$

$$(G + \lambda I)d = -\nabla f(\theta) \quad (4.2)$$

$$\lambda_{\text{next}} = \begin{cases} \lambda \times \text{drop} & \rho > \rho_{\text{upper}} \\ \lambda & \rho_{\text{lower}} \leq \rho \leq \rho_{\text{upper}} \\ \lambda \times \text{boost} & \text{otherwise} \end{cases} \quad (4.3)$$

The memory consumption of the Newton method depends on the size of data, which makes it difficult to handle large datasets. To counter this issue, the SHN method uses a subset S of the training data to obtain the subsampled Gauss-Newton matrix, which is used to approximate the Hessian matrix. This technique not only reduces the execution time per iteration (with a slightly less accurate direction) but also decreases the memory usage considerably.

For instance, at the m th convolutional layer for the Gauss-Newton matrix-vector products only the following matrices must be stored:

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(S^{m,i})^T}, \forall i \in S \quad (4.4)$$

For the gradient evaluations and the activation function the whole training data is required, so the independent results over all instances for each mini-batch must be summed. If the index set $\{1, \dots, N\}$ of data is split to R equal-sized subsets S_1, \dots, S_R and the result for each subset is calculated, then to find the final output all the subset results must be accumulated. The utilization of subsets can effectively decrease the memory consumption (Wang et al., 2020, section 3.5 and section 5) [3].

Types of Neural Networks	
LeCun et al. [15]	Fully-connected
Martens [20]	Autoencoder
Martens and Sutskever [21]	Fully-connected, Recurrent
Kiros [9]	Fully-connected, Autoencoder
Wang et al. [29, 30]	Fully-connected
Botev et al. [1]	Autoencoder

Table 4.1: Previous studies on Newton methods [3].

Table 4.1 illustrates some of the previous studies on Newton methods, performed on different types of ANNs. Other studies investigated the use of second-order optimization

methods for training CNNs, however, those are different from the Newton method considered in this dissertation.

4.4 Network Implementation

For the purpose of this dissertation, the implementation of a Convolutional Neural Network (CNN) with the Subsampled Hessian Newton (SHN) method was used, which was implemented in Python by Wang et al. [3] and can be found here [<https://github.com/cjlin1/simpleNN>]. Many optimisation tricks were applied to reduce memory consumption and to improve efficiency, which are discussed in that paper. The Python implementation used the Tensorflow [62] machine learning framework and is slightly different from the one used in [3], which was implemented in Matlab.

The initial implementation, which the paper [3] used for the experiments, used Matlab. Consequently, the input datasets used a matlab format (.mat), which was transferred to the Python version. The input files must contain a ‘y’ variable (of size $N \times 1$), which includes all the labels of the target class, and a ‘Z’ variable (of size $N \times M$), which includes all the features. Since the datasets were already preprocessed with a specific format which could be easily adapted to the matlab format, a script was implemented which was responsible to convert the text files (.txt) to matlab files (.mat) (Appendix C). This script can be found at [https://gitlab.com/perf.ai/pssp_project/-/blob/master/datasets2mat.sh].

In addition to the above, the implementation was modified so that it could be executed in a Jupyter notebook [61] and the datasets were uploaded to a public Gitlab repository, to be easily accessible. All the necessary scripts, programs, data files and instructions were uploaded in that repository, which can be found here [https://gitlab.com/perf.ai/pssp_project/-/tree/master].

Further modifications were made to the Python implementation to adapt it to the PSSP problem and improve the results, as the initial version was not very effective for this particular problem. For all experiments of this dissertation a free Colab machine was utilized (to use one visit [<https://colab.research.google.com/notebooks/welcome.ipynb>]). For information according TensorFlow visit [<https://www.tensorflow.org/>] and for PyTorch visit [<https://pytorch.org/>]. As regards the fastai library, which can be found here [<https://www.fast.ai/>], a very informative course is available at [<https://course.fast.ai/part2>], that describes how to create more advanced neural networks.

Chapter 5

Experiments and Results

5.1	Experiments for Implementation Evaluation	68
5.2	Experiments with CB513 dataset	70
5.2.1	Fine Tuning of Hyper Parameters	70
5.2.2	10-fold Cross-Validation and Ensembles Results	72
5.2.3	CNN and SVM Combination	74
5.2.4	Filtering Results for CB513	75
5.2.5	Additional experiments with CB513	79
5.2.6	Final results for CB513	79
5.3	Experiments with PISCES dataset	80
5.3.1	5-fold Cross-Validation and Ensembles Results	80
5.3.2	Filtering Results for PISCES	81
5.3.3	Final Results for PISCES	83
5.4	Best Results for CB513 and PISCES	84

5.1 Experiments for Implementation Evaluation

Many experiments have been performed, in order to find the optimal hyper parameters for the Convolutional Neural Network (CNN) and the Subsampled Hessian Newton (SHM) optimiser. The initial implementation [3] was already tested on some well-known benchmarks, like MNIST and CIFAR10 problems, which proved that the network was able to learn effectively. Because of that, the experiments of this dissertation were focused more on the PSSP problem. Initially, the model was trained with the CB513 dataset, which is relatively small, to identify the best hyper parameters and then additional experiments were performed on the bigger dataset, PISCES.

To ensure that each trained model has the best possible accuracy, during the training process after each iteration the model (all the weights) with the highest test accuracy was saved to an output file. This file can be then loaded to predict the test data and display the Q3 accuracy. This practice ensures that the model does not overfit to the training data and is able to predict new, never seen before data samples.

Figure 5.1 shows the test loss after each iteration for a CNN model trained with fold 5 of CB513. The red line illustrates the test loss after each iteration, while the green line illustrates the test loss of the saved model. The test loss after iteration 5, fluctuated within a narrow margin of about 0.05, while the test loss of the saved model followed a downward trend until it reached a plateau. Figure 5.2 displays the test accuracy for the same model for the first 35 iterations. The red line represents the test accuracy after each iteration, while the green line represents the test accuracy of the saved model. According to the line graph (Figure 5.2), the test accuracy dropped slightly in iteration 20, while the test accuracy of the saved model remained the same. This proves that at any iteration the saved model has the best possible test accuracy, which does not drop throughout the entire training process. In addition to that, the two line graphs (Figure 5.1 and 5.2) confirm that the model is able to train effectively and manages to converge in about twenty iterations.

For each experiment the following steps were performed. First the global parameters for the datasets were set (plus_var: the number of neighboring amino acids added, ds_num: the fold number of the dataset, dataset: 'CB513' or 'PISCES' to choose between the two PSSP datasets). Then, the appropriate dataset was retrieved from the Gitlab repository and all necessary functions were loaded. In the next phase, the hyper parameters were selected to prepare the model for the training phase. As soon as the training process was finished, a new file was created, which had all weights of the model with the best test accuracy. This file was used in the final step, in which the saved model was loaded and

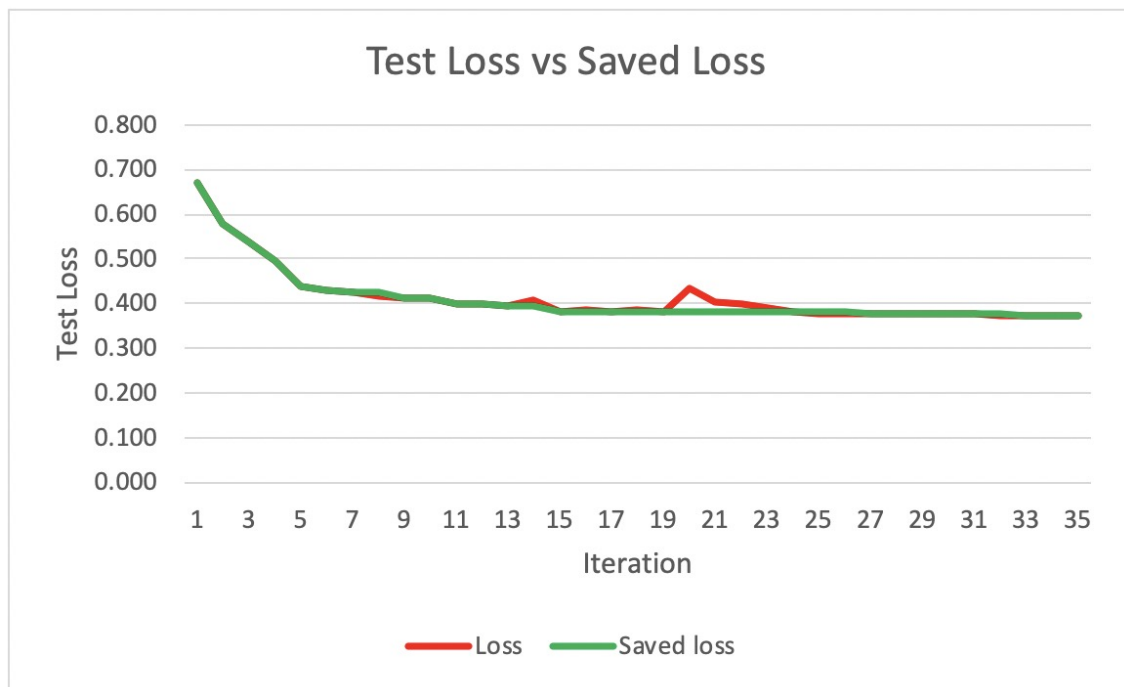


Figure 5.1: The test loss for each iteration compared to the test loss of the saved model.

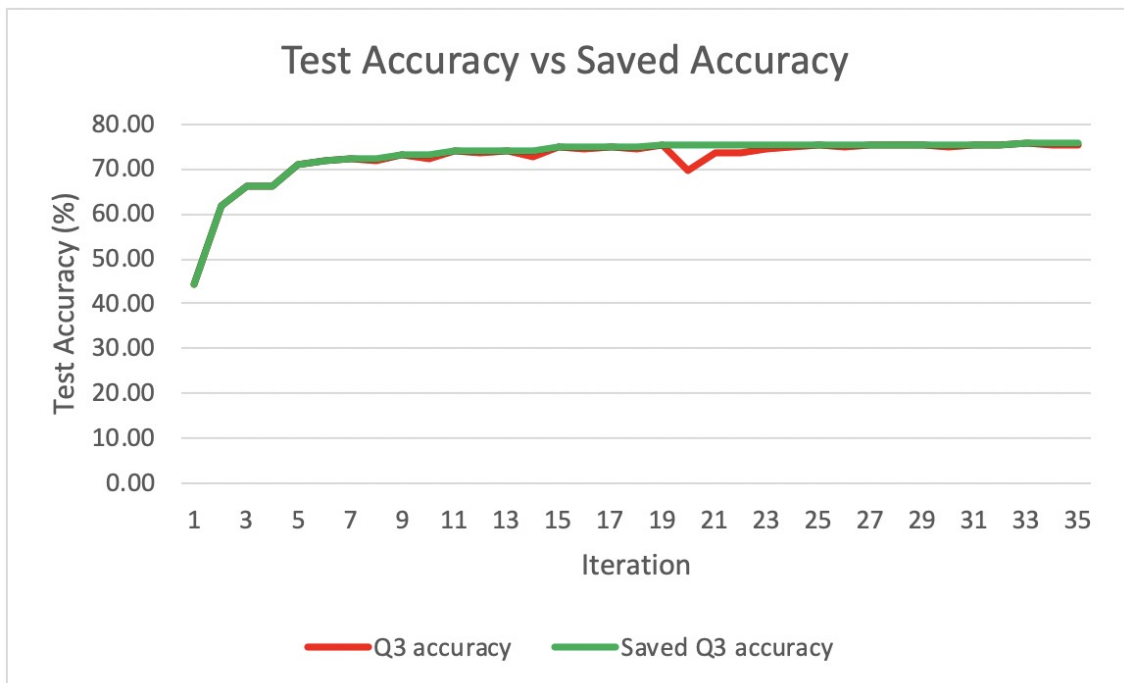


Figure 5.2: The test Q3 accuracy after each iteration compared to the test Q3 accuracy of the saved model.

was used to create two files with the predicted secondary structures of the proteins for the test and train datasets, respectively.

In order to check the efficiency of the Subsampled Hessian Newton (SHN) method, the Gradient Descent algorithm could be used to train CNN models, with the same structure as the one used for the SHN method. This would make it possible to compare the two optimisation algorithms, both in terms of accuracy (for the PSSP problem) and fine tuning of their hyper parameters. The Stochastic Gradient Descent (SGD) algorithm is already implemented and can be selected, as an alternative optimisation method.

5.2 Experiments with CB513 dataset

5.2.1 Fine Tuning of Hyper Parameters

In order to find the best hyper parameters for the network, experiments must be performed where each time only one hyper parameter is altered (the rest remain the same). Table 5.1 illustrates the hyper parameters used for the CNN, where the last layer has only three (3) neurons (one for each possible class). For each combination of hyper parameters five different models were trained and the average Q3 accuracy was saved in an excel file. For hyper parameter tuning, fold 5 was selected (for all the experiments), because it was observed that models trained with this fold performed very poorly, compared to the other folds. The motivation behind this was to maximize the performance of the hardest-to-learn fold with the hope that this would increase the overall Q3 accuracy and SOV score of the cross-validation.

	Type	Kernel size	Number of Filters	Activation Function
CNN Layer 0	Convolutional Hidden Layer	3 x 3	64	ReLU
CNN Layer 1	Convolutional Hidden Layer	3 x 3	64	ReLU
CNN Layer 2	Convolutional Hidden Layer	3 x 3	64	ReLU
CNN Layer 3	Fully Connected MLP	-	-	SOFTMAX

Table 5.1: Hyper parameters for CNN for all experiments.

First of all, the CB513 datasets were prepared with ‘plus7’ amino acids (seven left and seven right neighbouring amino acids were added, for each amino acid). This selection was based on [24], as this number (7) of neighbouring amino acids seemed very promising. The next step was to choose the number of samples used in the subsampled Gauss-Newton matrix (GNsize). Six different values were tested, as shown in table 5.2, while all

other parameters were selected randomly or based on the default values of the implementation. According to table 5.2, the best value for the GNsize was 2048 with approximately 75.54% Q3 accuracy.

GNsize	C	CNN layers	bsize	Max Iterations	Dimensions	Q3 Accuracy
50	0.01	4	8192	50	15 20 1	73.89%
100	0.01	4	8192	50	15 20 1	75.20%
200	0.01	4	8192	50	15 20 1	75.13%
512	0.01	4	8192	50	15 20 1	75.25%
1024	0.01	4	8192	50	15 20 1	75.40%
2048	0.01	4	8192	50	15 20 1	75.54%

Table 5.2: Q3 accuracy results for GNsize for fold 5 of CB513.

For the following experiments GNsize was selected to be equal to 2048 (GNsize = 2048). After that, the C value had to be determined so the same process was repeated but this time the C values were examined. Table 5.3 illustrates the Q3 accuracy results of the models in relation to the C value. It is obvious that the best option was 0.01 with 75.54% accuracy, so C was set to this value for the following experiments (C = 0.01).

GNsize	C	CNN layers	bsize	Max Iterations	Dimensions	Q3 Accuracy
2048	0.01	4	8192	50	15 20 1	75.54%
2048	0.05	4	8192	50	15 20 1	75.36%
2048	0.10	4	8192	50	15 20 1	75.02%
2048	0.50	4	8192	50	15 20 1	75.32%
2048	1.00	4	8192	50	15 20 1	75.29%

Table 5.3: Tuning the C hyper parameter for fold 5 of CB513.

This process was replicated for the batch size (bsize), which usually is set based on the memory constraints. It is very important to note that if the model cannot begin the training process, it is probably because this value was set too high. In this case, lowering the bsize value can fix the issue. Table 5.4 shows the results for bsize, however, it is not clear which one is the best, as most of them are very close to each other. For the purpose of this dissertation, the value 12288 was selected as the batch size (bsize = 12288) to reduce the training time of the model.

GNsize	C	CNN layers	bsize	Max Iterations	Dimensions	Q3 Accuracy
2048	0.01	4	1024	50	15 20 1	75.56%
2048	0.01	4	2048	50	15 20 1	75.39%
2048	0.01	4	4096	50	15 20 1	75.40%
2048	0.01	4	8192	50	15 20 1	75.54%
2048	0.01	4	10240	50	15 20 1	75.57%
2048	0.01	4	12288	50	15 20 1	75.64%

Table 5.4: Tuning the batch size (bsize) hyper parameter for fold 5 of CB513.

5.2.2 10-fold Cross-Validation and Ensembles Results

In order to check whether the results of a model are good just for a specific test dataset or whether the trained network is a good prediction model, additional techniques must be utilized. One such technique is cross-validation, which was described in section 3.8. More specifically, a 10-fold cross-validation was used for the CB513 dataset to validate the model’s ability to generalize.

Table 5.5 shows the hyper parameters for all the trained models, which were used for the cross-validation of CB513.

GNsize	C	CNN layers	bsize	Max Iterations	Dimensions
2048	0.01	4	12288	50	15 20 1

Table 5.5: Hyper parameters for trained models.

The cross-validation results for the CB513 dataset are shown in table 5.6. This table displays the overall Q3 accuracy and overall SOV score for the best trained model for each fold. In addition, the Q3 accuracy and SOV scores for each of the three classes (H, E, C) are shown separately, as well as the average results for all folds (cross-validation values).

According to table 5.6, the best trained model achieved 78.20% overall Q3 accuracy and 75.67 overall SOV score, while the cross-validation results were 77.25% and 72.91, respectively. Even though the optimisation for the hyper parameters was based on fold 5, which had the lowest Q3 accuracy as expected, the results for all the other folds were considerably better. It is obvious that most of the models had trouble identifying the class ‘E’ and that is why the QE accuracy for all folds are substantially lower than the QH and the QC accuracy. Most models were able to predict, to some extent, the class ‘C’, as the QC accuracy for the 10-fold cross-validation was approximately 82.13%.

	Q3	QH	QE	QC	SOV	SOVH	SOVE	SOVC
Fold0	78.20%	80.56%	68.94%	81.71%	75.67	81.18	72.21	72.32
Fold1	76.25%	78.82%	64.59%	81.38%	73.02	70.56	72.02	74.00
Fold2	77.85%	81.19%	65.77%	81.15%	73.26	75.80	69.85	70.43
Fold3	77.85%	80.75%	66.33%	81.70%	74.33	73.70	67.38	72.28
Fold4	77.97%	80.44%	65.36%	82.36%	73.38	74.11	68.25	70.52
Fold5	75.77%	79.45%	61.01%	81.05%	71.60	70.57	65.69	71.05
Fold6	77.91%	75.81%	65.82%	85.59%	74.42	74.04	70.42	74.58
Fold7	76.74%	75.78%	67.67%	82.37%	68.37	69.53	72.39	68.62
Fold8	76.82%	77.10%	69.57%	80.82%	72.61	67.49	73.48	72.57
Fold9	77.13%	79.63%	62.46%	83.15%	72.43	80.30	66.98	72.42
Average	77.25%	78.95%	65.75%	82.13%	72.91	73.73	69.87	71.88

Table 5.6: Q3 and SOV results for 10-fold cross validation for the CB513 dataset.

Table 5.7 shows the results for the cross validation of the ensembles method, where five (5) CNNs were trained with the SHN method, using the CB513 dataset. For the experiments of this dissertation, multiple models were trained for each fold (about 7-10) and the five (5) models, which formed the best ensembles model, were selected for the final ensembles model of each fold.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	78.46	80.86	68.99	82.07	75.19	79.61	73.17	72.87
Fold1	76.49	79.05	64.36	81.90	72.76	71.39	72.17	74.38
Fold2	78.19	81.23	66.77	81.41	75.03	76.87	70.88	71.48
Fold3	78.15	80.87	66.63	82.16	74.70	74.46	67.75	72.54
Fold4	78.15	80.92	65.60	82.26	74.95	77.38	66.96	71.64
Fold5	75.97	77.96	62.59	81.75	71.24	69.26	65.35	70.96
Fold6	77.91	75.86	66.77	85.08	74.25	74.23	71.83	73.94
Fold7	76.94	76.11	68.00	82.40	69.53	69.98	73.40	68.12
Fold8	77.03	77.33	68.55	81.71	73.72	70.17	72.32	72.77
Fold9	77.29	79.63	65.24	81.98	74.12	82.34	68.70	72.90
Average	77.46	78.98	66.35	82.27	73.55	74.57	70.25	72.16

Table 5.7: Q3 and SOV results for ensembles (with 5 experiments per fold) cross validation for the CB513 dataset.

Usually, the ensembles method is more effective when there is high variance between the trained models, because each trained model explores a different space of the dataset and it learns how to predict based on different features. That is why, it is better to combine results from different machine learning models. The combination of these models can create a new model with more accurate predictions than any of the separate models, at the cost of time and processing power, but this is not guaranteed.

In this case (Table 5.7), the new ensemble model managed to outperform every single CNN model for all folds (the best CNN model for each fold is shown in table 5.6). The increase in Q3 accuracy was relatively small, probably because all models were trained with the same hyper parameters, which resulted in less variance. The boost in accuracy

could be better if models with different hyper parameters or different types of models were used for creating the ensembles model. For instance, a Convolutional Neural Network, a Recurrent Neural Network, a Bidirectional Recurrent Neural Network, a Feed Forward Neural Network and a Long Short-Term Memory model could be trained and then combined with the ensembles method. This combination could have a greater impact on the accuracy of the new ensembles model.

5.2.3 CNN and SVM Combination

As mentioned in section 3.10.2, the final attempt to improve the results was to use Support Vector Machines (SVMs), which managed to improve the results of past PSSP studies [59]. More specifically, an SVM will be used to filter the output data from the CNN, with the ambition that the Q3 accuracy and SOV score could be improved.

In order to train the SVM, a window (of odd size) will be used to extract information from the prediction file created by the CNN, which will be used as the input features of SVM, while the expected output will be the secondary structure of the middle amino acid. A python program (prepare_SVM_files.py) was used to prepare the datasets for the SVM based on the given window size (Appendix K). The SVM was also implemented as a python program (train_SVM.py) (using Scikit-learn machine learning library [64]), which exploits the output datasets from ‘prepare_SVM_files.py’ to train an SVM and create a new output file with the new filtered predictions (Appendix L). Both programs can be found in this Gitlab repository [https://gitlab.com/perf.ai/pssp_project/-/tree/master].

For the experiments of this dissertation both networks will be utilized and the final results will be compared with the results of a standalone CNN. The CNN will use as input the data described in section 3.6, while the SVM will take as input the output of the CNN, in order to filter the results. The table 5.8 shows the hyper parameters used for all the experiments.

Kernel	C	Decision Function Shape	Degree	Shrinking	Tol	Gamma
RBF (Radial Basis Function)	10	ovr	3	TRUE	0.001	0.1

Table 5.8: Hyper parameters for SVM filtering.

The same technique could be used with almost any other classification model. In this dissertation, except from SVM filtering, Decision Trees and Random Forests [71] (which are basically ensembles of decision trees) were used, as alternative filtering methods, in combination with external rules, which were explained in section 3.10.1. The parameters

used for the random forest filter are shown in table 5.9. For the decision tree filter, the only non-default parameter used was the maximum depth (max_depth) parameter which was set to twenty (20).

n_estimators	max_depth	random_state	min_samples_split	min_samples_leaf
100	25	42	2	1

Table 5.9: Hyper parameters for Random Forest filtering.

5.2.4 Filtering Results for CB513

The chosen window size for the filtering methods, for the CB513 experiments, was thirteen (13), because it produced relatively good filtering results without a major impact on the total filtering time.

Table 5.10 shows the results for the Q3 accuracy and SOV score after applying the external rules to the ensembles model. It seems that the Q3 accuracy increased only by a tiny amount, while the SOV score rose by 1.33, which is relatively good.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	78.83	79.47	68.27	84.41	77.35	77.96	72.07	75.74
Fold1	76.22	77.30	62.99	83.41	72.57	70.78	71.79	72.44
Fold2	78.17	79.83	65.77	83.18	76.68	76.98	71.88	72.92
Fold3	78.26	80.06	65.79	83.57	75.34	73.92	67.42	72.66
Fold4	78.07	80.04	63.89	83.73	74.78	77.09	65.58	70.69
Fold5	75.89	76.64	61.06	83.38	72.42	68.98	65.97	70.56
Fold6	78.07	74.93	65.52	86.74	76.35	76.02	71.33	75.46
Fold7	76.97	75.10	66.99	83.75	70.87	70.16	72.37	68.07
Fold8	77.46	76.46	67.76	83.98	76.72	73.16	72.52	73.06
Fold9	77.30	78.88	63.79	83.47	75.69	84.14	68.71	73.19
Average	77.52	77.87	65.18	83.96	74.88	74.92	69.96	72.48

Table 5.10: Q3 accuracy and SOV score for ensembles (with 5 executions per fold) and external rules filtering for CB513 dataset.

Table 5.11 illustrates the Q3 accuracy and SOV score for the ensembles model after applying external rules and SVM filtering. The external rules filtering usually offers a significant boost in the overall SOV score and sometimes a slight drop in the Q3 accuracy.

The order in which the filters are applied can produce different results, so the same filters could be applied in various ways (different orders). For this purpose, a bash script was created, which applies the filtering methods in various orders and creates an output file with all the results (Appendix M). The results for the ensembles model with SVM filtering are shown in table 5.12.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	79.24	80.52	72.51	82.13	76.80	77.09	74.23	73.38
Fold1	76.75	79.00	73.33	77.14	72.87	72.04	75.10	69.87
Fold2	78.43	81.84	77.56	75.79	77.24	77.12	72.87	71.79
Fold3	78.56	82.40	72.19	78.78	74.17	72.84	70.41	69.30
Fold4	78.42	82.39	76.89	75.78	74.91	77.74	69.21	68.31
Fold5	76.38	77.75	71.02	78.26	72.31	70.23	68.51	68.56
Fold6	78.05	77.63	76.11	79.34	77.53	76.56	74.23	74.05
Fold7	77.44	77.07	78.18	77.31	69.65	69.72	73.08	67.02
Fold8	77.16	77.60	76.62	77.09	75.49	70.53	75.06	69.80
Fold9	77.42	80.71	73.23	76.91	76.07	85.19	73.08	70.92
Average	77.79	79.69	74.76	77.85	74.70	74.91	72.58	70.30

Table 5.11: Q3 accuracy and SOV score for ensembles (with 5 executions per fold), external rules and SVM filtering for CB513 dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	80.44	84.00	74.47	81.14	75.95	82.45	76.98	71.19
Fold1	79.00	83.48	75.44	77.79	75.55	77.70	78.09	72.51
Fold2	80.03	84.21	80.86	75.79	78.29	77.47	74.29	73.83
Fold3	79.44	83.10	74.22	79.17	75.61	74.67	73.21	71.32
Fold4	79.19	84.23	79.34	74.76	75.98	78.05	73.74	68.63
Fold5	77.84	80.86	74.82	77.25	73.15	69.51	70.32	70.94
Fold6	79.57	79.37	77.88	80.57	77.43	75.99	76.77	73.39
Fold7	78.72	79.17	80.67	77.35	71.46	72.70	75.24	68.51
Fold8	79.42	79.11	76.23	81.54	77.19	71.45	77.70	73.76
Fold9	79.01	82.50	76.59	77.33	77.31	85.96	75.42	72.35
Average	79.27	82.00	77.05	78.27	75.79	76.60	75.18	71.64

Table 5.12: Q3 accuracy and SOV score for ensembles and SVM filtering for CB513 dataset.

Table 5.13 illustrates the Q3 accuracy and SOV score for the ensembles model with SVM and external rules filtering. According to tables 5.10 and 5.11, the SVM filtering improved the overall Q3 accuracy by a small amount, but decreased the overall SOV score slightly. The impact of SVM filtering was significant, for both overall Q3 accuracy and overall SOV score (Tables 5.7 - before, and 5.12 - after), with 79.27% and 75.79, respectively. It seems that if SVM filtering is applied before the external rules, the cross validation results are substantially better, with approximately 1.75% increase in overall Q3 accuracy and about 1.65 growth in overall SOV score.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	80.55	83.59	73.18	82.45	78.55	82.86	76.21	74.03
Fold1	78.89	82.56	74.36	78.90	75.88	76.96	76.14	71.84
Fold2	80.06	83.46	80.17	76.90	78.63	77.80	73.64	73.68
Fold3	79.36	82.37	72.61	80.49	76.15	74.00	70.71	72.27
Fold4	79.51	83.93	78.66	76.13	77.79	78.19	73.33	70.94
Fold5	77.75	80.24	73.46	78.23	74.11	71.05	69.18	70.31
Fold6	79.95	78.86	76.70	82.39	78.83	78.00	76.35	73.88
Fold7	78.60	78.33	80.10	77.99	72.13	72.23	74.79	68.43
Fold8	79.22	78.39	75.10	82.33	77.93	70.94	75.55	73.30
Fold9	79.09	82.23	75.84	78.17	78.01	86.06	74.59	73.46
Average	79.30	81.40	76.02	79.40	76.80	76.81	74.05	72.21

Table 5.13: Q3 accuracy and SOV score for ensembles, SVM and external rules filtering for CB513 dataset.

The results for the ensembles model with the external rules and decision tree filtering are shown in table 5.14. Table 5.15 illustrates the Q3 accuracy and SOV score for the ensembles model with decision tree filtering, while table 5.16 displays the results for the ensembles model with decision tree and external rules filtering. The decision tree filtering improved the results significantly (Tables 5.7 - before, and 5.15 - after), reaching 81.69% overall Q3 accuracy and 75.93 overall SOV score. According to tables 5.14 and 5.16, when the decision tree filtering is applied before the external rules, the prediction results of the model are considerably better. More specifically, there is an increase of about 2.27% in the overall Q3 accuracy and approximately 5.11 in the overall SOV score.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	80.18	80.77	78.16	80.88	76.08	75.57	74.92	72.28
Fold1	78.50	79.88	74.59	79.83	72.27	72.92	76.51	69.75
Fold2	78.90	81.09	82.54	75.00	75.95	74.95	75.64	71.06
Fold3	79.92	82.02	75.30	80.65	73.81	70.73	73.27	70.12
Fold4	79.46	81.65	78.36	78.13	73.09	75.81	72.17	67.85
Fold5	78.26	78.96	75.10	79.45	73.23	69.87	72.26	69.85
Fold6	79.45	77.89	78.00	81.30	75.40	72.63	75.64	73.74
Fold7	79.14	78.00	78.69	80.23	68.38	71.51	73.84	66.06
Fold8	79.08	77.10	78.32	81.25	74.59	68.77	74.50	71.09
Fold9	78.58	80.82	76.77	77.66	74.43	80.86	74.58	70.19
Average	79.15	79.82	77.58	79.44	73.72	73.36	74.33	70.20

Table 5.14: Q3 accuracy and SOV score for ensembles, external rules and decision tree filtering for CB513 dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	82.47	84.72	81.56	81.26	76.30	80.35	79.15	72.91
Fold1	82.25	84.86	81.50	80.76	76.29	80.57	83.03	72.36
Fold2	82.01	84.68	84.54	78.24	77.79	79.46	79.07	72.41
Fold3	80.98	83.79	79.43	79.44	73.83	74.74	77.94	69.90
Fold4	81.33	84.34	83.69	77.50	76.02	76.71	78.33	70.07
Fold5	81.34	83.51	80.59	80.13	74.44	73.88	75.34	70.47
Fold6	81.39	81.61	82.91	80.45	77.52	77.07	78.86	74.10
Fold7	81.60	81.31	83.27	80.90	73.31	78.21	79.10	70.13
Fold8	81.94	81.24	82.72	82.10	77.45	70.24	81.65	73.14
Fold9	81.55	82.79	82.50	79.95	76.33	81.65	80.33	71.13
Average	81.69	83.29	82.27	80.07	75.93	77.29	79.28	71.66

Table 5.15: Q3 accuracy and SOV score for ensembles and decision tree filtering for CB513 dataset.

Table 5.17 shows the results for the ensembles model with external rules and random forest filtering. The Q3 accuracy and SOV score for the ensembles model with random forest filtering are displayed in table 5.18, while the results for the ensembles model with random forest and external rules filtering are presented in table 5.19. The boost of random forest filtering was great (Tables 5.7 - before, and 5.18 - after), since it increased the overall Q3 accuracy to 81.75% and the overall SOV score to 76.33. It is obvious that the results are better when the random forest filtering is applied before the external rules filtering (Tables 5.17 and 5.19).

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	82.67	84.00	79.33	83.57	80.17	81.83	78.42	76.82
Fold1	82.43	83.62	79.61	83.24	78.58	81.72	77.98	74.73
Fold2	82.18	83.42	82.29	80.99	80.58	79.29	75.54	75.87
Fold3	81.12	82.48	77.03	82.19	77.39	76.58	74.07	73.36
Fold4	81.69	83.79	81.48	79.97	78.24	78.87	75.89	71.63
Fold5	81.44	82.35	78.78	82.21	76.76	75.13	73.45	72.30
Fold6	81.84	80.93	81.19	82.81	80.92	79.78	79.62	76.50
Fold7	81.57	80.18	81.40	82.68	75.84	78.85	78.15	72.00
Fold8	81.92	79.68	81.31	84.21	80.85	73.83	79.49	76.04
Fold9	81.71	81.53	81.05	82.24	79.00	82.81	79.48	74.14
Average	81.86	82.20	80.35	82.41	78.83	78.87	77.21	74.34

Table 5.16: Q3 accuracy and SOV score for ensembles, decision tree and external rules filtering for CB513 dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	80.00	81.91	79.66	78.73	76.32	76.51	76.54	70.72
Fold1	78.29	81.31	75.39	77.79	72.81	76.14	75.82	69.33
Fold2	78.55	83.14	83.79	71.63	74.11	75.06	74.38	67.22
Fold3	79.89	83.63	76.79	78.38	73.58	71.49	73.57	68.60
Fold4	79.35	83.60	80.32	75.17	72.80	77.21	72.32	65.82
Fold5	78.34	81.15	75.61	77.74	72.73	69.35	71.73	69.45
Fold6	79.11	79.15	78.47	79.40	76.42	76.08	74.53	72.51
Fold7	78.91	80.89	79.48	77.16	68.24	72.67	73.65	64.55
Fold8	79.07	79.42	78.83	78.91	74.94	71.56	75.57	69.61
Fold9	78.78	82.68	77.11	76.33	75.58	85.09	74.33	69.10
Average	79.03	81.69	78.55	77.12	73.75	75.12	74.24	68.69

Table 5.17: Q3 accuracy and SOV score for ensembles, external rules and random forest filtering for CB513 dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	82.27	85.81	83.74	78.73	76.39	82.40	79.65	71.18
Fold1	82.47	86.29	83.55	78.97	77.02	81.79	83.20	72.90
Fold2	82.18	87.41	86.47	75.16	77.99	80.07	78.40	71.31
Fold3	81.36	85.17	80.74	78.45	74.72	77.20	76.71	69.18
Fold4	81.61	86.51	84.86	75.68	74.66	77.52	77.64	68.22
Fold5	81.07	84.67	81.55	78.14	74.04	73.18	74.84	69.70
Fold6	81.12	82.20	83.68	79.03	77.96	78.38	78.53	73.18
Fold7	81.66	83.99	84.91	78.20	74.18	79.35	80.48	69.76
Fold8	82.17	82.68	83.85	80.75	79.30	74.91	82.88	74.00
Fold9	81.55	84.95	84.94	76.72	77.06	85.18	80.14	70.15
Average	81.75	84.97	83.83	77.98	76.33	79.00	79.25	70.96

Table 5.18: Q3 accuracy and SOV score for ensembles and random forest filtering for CB513 dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold0	82.29	85.01	81.73	80.53	79.98	84.08	78.58	73.99
Fold1	82.62	85.42	81.67	81.10	78.76	82.96	78.73	74.28
Fold2	82.38	86.26	84.41	77.78	81.07	81.60	76.53	76.04
Fold3	81.42	84.02	78.29	80.91	77.83	77.43	73.99	73.54
Fold4	81.74	85.85	82.53	77.78	78.32	80.52	75.29	71.37
Fold5	81.07	84.09	79.63	79.61	76.72	76.00	72.81	71.34
Fold6	81.35	81.90	81.79	80.73	80.63	80.30	79.57	75.56
Fold7	81.42	82.94	83.04	79.43	75.96	79.43	78.51	71.51
Fold8	81.93	81.12	82.33	82.40	80.97	75.23	79.86	75.37
Fold9	81.73	83.95	83.14	79.01	79.56	85.64	78.75	74.14
Average	81.80	84.06	81.86	79.93	78.98	80.32	77.26	73.71

Table 5.19: Q3 accuracy and SOV score for ensembles, random forest and external rules filtering for CB513 dataset.

5.2.5 Additional experiments with CB513

Since the order of the filtering methods matters, the order of the ensembles method could also lead to different results. A few experiments were performed with fold 0 of CB513, where the ensembles method was applied after the various filtering methods. According to table 5.20, applying the ensembles after the filtering methods, leads to better results. Apart from that, it is possible to combine all the filtering methods into one ensembles model, which was not possible in the scenario where the ensembles method was applied first. The new ensembles model, which combined 15 models (5 models with external rules and SVM filtering, 5 models with external rules and decision tree filtering, and 5 models with external rules and random forest filtering), had the highest SOV score. Further experiments have not been performed because of the shortage of time.

Method Used (for fold 0 of CB513)	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Ensembles + External Rules + SVM	79.24	80.52	72.51	82.13	76.80	77.09	74.23	73.38
Ensembles + External Rules + Decision Tree	80.18	80.77	78.16	80.88	76.08	75.57	74.92	72.28
Ensembles + External Rules + Random Forest	80.00	81.91	79.66	78.73	76.32	76.51	76.54	70.72
External Rules + SVM + Ensembles	79.49	79.93	72.57	83.13	76.56	77.30	74.36	73.41
External Rules + Decision Tree + Ensembles	81.55	80.48	78.49	84.12	77.05	75.55	77.35	75.07
External Rules + Random Forest + Ensembles	81.52	81.70	81.01	81.68	76.59	76.91	76.67	74.01
Ensembles (SVM + Decision Tree + Random Forest)	81.21	80.98	77.26	83.64	77.16	76.53	76.81	75.17

Table 5.20: Results for fold 0 of CB513 with the ensembles method applied before and after the filtering methods.

5.2.6 Final results for CB513

After collecting all the results for all filtering methods, the 10-fold cross validation method (average) was used to combine the results for all folds. These results are presented in table 5.21, which makes it easier to compare the different filtering methods. According to table 5.21, the best results for CB513, in terms of overall Q3 accuracy and overall SOV score, came from the ensembles model with the random forest and external rules filtering. This model managed to reach 81.80% Q3 (per residue) accuracy and 78.98 SOV score, which is very close with the current state-of-the-art results (84-85% Q3 accuracy).

METHOD	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
CROSS-VALIDATION	77.25	78.95	65.75	82.13	72.91	73.73	69.87	71.88
ENSEMBLES (5 EXPERIMENTS / FOLD)	77.46	78.98	66.35	82.27	73.55	74.57	70.25	72.16
ENSEMBLES + EXTERNAL RULES	77.52	77.87	65.18	83.96	74.88	74.92	69.96	72.48
ENSEMBLES + EXTERNAL RULES + SVM	77.79	79.69	74.76	77.85	74.70	74.91	72.58	70.30
ENSEMBLES + SVM	79.27	82.00	77.05	78.27	75.79	76.60	75.18	71.64
ENSEMBLES + SVM + EXTERNAL RULES	79.30	81.40	76.02	79.40	76.80	76.81	74.05	72.21
ENSEMBLES + EXTERNAL RULES + DECISION TREE	79.15	79.82	77.58	79.44	73.72	73.36	74.33	70.20
ENSEMBLES + DECISION TREE	81.69	83.29	82.27	80.07	75.93	77.29	79.28	71.66
ENSEMBLES + DECISION TREE + EXTERNAL RULES	81.86	82.20	80.35	82.41	78.83	78.87	77.21	74.34
ENSEMBLES + EXTERNAL RULES + RANDOM FOREST	79.03	81.69	78.55	77.12	73.75	75.12	74.24	68.69
ENSEMBLES + RANDOM FOREST	81.75	84.97	83.83	77.98	76.33	79.00	79.25	70.96
ENSEMBLES + RANDOM FOREST + EXTERNAL RULES	81.80	84.06	81.86	79.93	78.98	80.32	77.26	73.71

Table 5.21: 10-fold Cross validation, Q3 accuracy and SOV score for all methods for CB513 dataset.

5.3 Experiments with PISCES dataset

The PISCES dataset is much bigger than the CB513 dataset and the experiments for this dataset required a lot more time. Because of that, most hyper parameters used in PISCES experiments were derived from the CB513 experiments. This will probably have an impact on the prediction capabilities of the final model, but further experiments could not be made due to the shortage of time. The hyper parameters used for the CNN in the PISCES experiments are shown in table 5.1. The max epochs (max iterations) were increased from 50 to 100 because the model required more epochs to learn the bigger dataset. Table 5.22 shows the hyper parameters that were used to train all PISCES models.

GNsize	C	CNN layers	bsize	Max Iterations	Dimensions
2048	0.01	4	12288	100	15 20 1

Table 5.22: Hyper parameters for SHN method, used for all PISCES experiments.

5.3.1 5-fold Cross-Validation and Ensembles Results

For the PISCES dataset a 5-fold cross validation was used, where seven (7) models were trained per fold and the best one was selected. The main reason a 5-fold cross validation was chosen, instead of a 10-fold, was to make the results comparable with past studies on PSSP, that used the PISCES dataset. Table 5.23 displays the Q3 accuracy and SOV score of the best model for each fold.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	79.22	84.45	69.90	79.39	76.00	78.84	77.47	71.60
Fold2	79.41	84.19	71.02	79.59	76.47	77.80	76.85	72.44
Fold3	79.60	84.33	70.47	79.87	76.48	78.23	76.09	72.16
Fold4	79.88	84.62	71.06	80.09	76.67	79.04	77.30	72.55
Fold5	79.41	84.38	70.04	79.96	76.59	80.37	76.65	72.26
Average	79.50	84.39	70.50	79.78	76.44	78.86	76.87	72.20

Table 5.23: Q3 accuracy and SOV score for 5-fold cross validation for PISCES dataset.

To create the ensembles model five (5) from the seven (7) trained models were selected, so that the Q3 accuracy of the new model was maximized. Table 5.24 presents the results for Q3 accuracy and SOV score of the new ensembles model, for each fold of PISCES.

The comparison between table 5.23 and 5.24 reveals that there is a similar issue with the CB513 dataset. There is not enough variance between the trained models, which results in only a small improvement in overall Q3 accuracy (0.30%) and SOV score (0.63).

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	79.55	84.62	70.19	79.88	76.80	79.44	77.84	72.12
Fold2	79.74	84.24	71.05	80.36	77.01	78.08	77.09	73.05
Fold3	79.83	84.47	70.43	80.36	76.92	78.49	76.42	72.50
Fold4	80.15	84.65	71.30	80.60	77.34	79.48	77.61	73.18
Fold5	79.72	84.46	70.28	80.52	77.29	80.75	77.16	72.95
Average	79.80	84.49	70.65	80.34	77.07	79.25	77.22	72.76

Table 5.24: Q3 accuracy and SOV score for ensembles method (with 5 trained models per fold) for PISCES dataset.

5.3.2 Filtering Results for PISCES

A bigger dataset, like PISCES, can help the model to learn more effectively the patterns of the data, and that is why the cross validation results (Table 5.23) are better compared to the CB513 results (Table 5.6), but at the same makes it very difficult to use SVM filtering. SVMs are usually very effective for small datasets, however, on big datasets the memory scales quadratically with the number of data points, which makes them very difficult to train and impractical. Several attempts were made to train an SVM with samples from the PISCES dataset, but the results were worse than the results without the SVM filtering. Because of that, the SVM filtering was not applied in any of the PISCES experiments. If a good sample is extracted from the PISCES dataset, it might be possible to train an SVM for filtering the results of models trained with the PISCES dataset. However, for this dissertation such technique could not be found.

Table 5.25 shows the results for each of the five folds after applying the external rules to the ensembles model. It seems there was a slight increase (0.06) to the overall Q3 accuracy and a considerable increase in the overall SOV score (1.18), which was expected, since external rules are used mainly to improve the overall SOV score.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	79.64	83.79	69.08	81.61	78.06	79.90	77.61	72.95
Fold2	79.79	83.40	70.07	81.93	78.04	78.36	76.89	73.57
Fold3	79.91	83.66	69.40	81.98	78.15	78.91	76.11	73.17
Fold4	80.22	83.81	70.24	82.25	78.53	79.85	77.40	73.88
Fold5	79.75	83.54	69.23	82.14	78.47	80.75	76.93	73.63
Average	79.86	83.64	69.60	81.98	78.25	79.55	76.99	73.44

Table 5.25: Q3 accuracy and SOV score for ensembles with external rules filtering for PISCES dataset.

The results for the ensembles model with external rules and decision tree filtering are displayed in table 5.26. The decision tree filtering improved the Q3 accuracy by 0.80% and affected slightly the SOV score.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	80.37	82.94	76.74	79.90	77.76	77.91	80.07	72.75
Fold2	80.78	83.01	77.19	80.68	78.39	77.28	79.61	73.97
Fold3	80.62	82.85	77.28	80.20	78.01	77.17	78.95	72.86
Fold4	80.92	83.16	77.71	80.46	78.46	78.52	80.16	73.35
Fold5	80.62	82.93	76.95	80.47	78.47	79.25	79.25	73.97
Average	80.66	82.98	77.17	80.34	78.22	78.03	79.61	73.38

Table 5.26: Q3 accuracy and SOV score for ensembles with external rules and decision tree filtering for PISCES dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	81.69	84.45	81.54	78.98	78.68	78.99	83.19	73.64
Fold2	82.10	84.55	81.06	80.24	79.59	78.85	82.70	75.32
Fold3	81.79	84.39	81.15	79.46	78.95	78.93	82.11	73.97
Fold4	82.13	84.57	81.93	79.75	79.52	79.58	82.93	74.97
Fold5	81.93	84.31	80.65	80.31	79.52	80.80	82.03	75.50
Average	81.93	84.45	81.27	79.75	79.25	79.43	82.59	74.68

Table 5.27: Q3 accuracy and SOV score for ensembles with decision tree filtering for PISCES dataset.

Table 5.27 shows the results, for each fold of PISCES, for the ensembles model after applying the decision tree filtering, while table 5.28 illustrates the results for the ensembles model with decision tree and external rules filtering. It is obvious that applying the decision tree filtering before the external rules produces significantly better results. Moreover, the decision tree filtering improved the results of the ensembles model considerably (81.93% overall Q3 accuracy and 79.25 overall SOV score).

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	81.75	83.12	80.08	81.32	81.28	80.67	82.87	75.91
Fold2	82.10	83.32	79.67	82.33	81.57	79.71	82.33	76.84
Fold3	81.86	83.26	79.52	81.75	81.11	80.18	81.32	75.86
Fold4	82.16	83.34	80.47	81.91	81.48	80.77	82.41	76.54
Fold5	81.91	83.06	79.20	82.37	81.55	82.02	81.58	76.98
Average	81.96	83.22	79.79	81.94	81.40	80.67	82.10	76.43

Table 5.28: Q3 accuracy and SOV score for ensembles with decision tree and external rules filtering for PISCES dataset.

The Q3 accuracy and SOV score for the ensembles model with external rules and random forest filtering are presented in table 5.29. A comparison between table 5.26 and table 5.29 shows that random forest filtering is clearly more effective than the decision tree filtering, for the PISCES dataset.

Table 5.30 illustrates the results for the ensembles model with the random forest filtering, while table 5.31 illustrates the results for the ensembles model with random forest and external rules filtering. According to tables 5.29 and 5.31, when random forest filtering

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	80.62	84.11	77.92	78.66	77.99	79.28	80.40	72.05
Fold2	80.98	84.01	78.49	79.40	78.43	78.26	80.05	72.99
Fold3	80.85	83.89	78.52	79.04	78.37	78.60	79.33	72.39
Fold4	81.15	84.19	78.99	79.26	78.83	79.74	80.48	72.93
Fold5	80.88	83.97	77.99	79.49	78.86	80.65	79.88	73.42
Average	80.90	84.03	78.38	79.17	78.50	79.31	80.03	72.76

Table 5.29: Q3 accuracy and SOV score for ensembles with external rules and random forest filtering for PISCES dataset.

is applied before external rules, the final results are much better (2.12% for Q3 accuracy and 4.14 for SOV score). Furthermore, if the external rules are applied after the random forest filtering, the overall Q3 accuracy drops slightly, while the overall SOV score grows by a small amount (Tables 5.30 and 5.31).

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	82.86	85.98	83.74	79.19	80.97	82.15	84.83	75.21
Fold2	83.23	86.12	83.81	79.99	81.71	81.80	84.59	76.44
Fold3	83.00	85.78	83.87	79.63	81.08	81.62	84.29	75.63
Fold4	83.31	86.09	84.36	79.86	81.68	82.54	84.66	76.32
Fold5	83.19	86.02	83.51	80.19	81.89	83.91	84.39	76.83
Average	83.12	86.00	83.86	79.77	81.47	82.40	84.55	76.09

Table 5.30: Q3 accuracy and SOV score for ensembles with random forest filtering for PISCES dataset.

	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
Fold1	82.79	85.21	82.27	80.64	82.37	82.79	83.71	76.45
Fold2	83.13	85.36	82.42	81.31	82.74	82.09	83.57	77.41
Fold3	82.89	85.01	82.24	81.05	82.25	82.08	82.66	76.59
Fold4	83.18	85.27	82.90	81.21	82.82	82.77	83.58	77.28
Fold5	83.09	85.21	81.96	81.64	83.00	84.35	83.08	77.79
Average	83.02	85.21	82.36	81.17	82.64	82.82	83.32	77.10

Table 5.31: Q3 accuracy and SOV score for ensembles with random forest and external rules filtering for PISCES dataset.

5.3.3 Final Results for PISCES

The 5-fold cross validation (average) was applied for all filtering methods and the results are shown in table 5.32. The best results according to this table, came from the ensembles model after applying the random forest and external rules filtering. This method managed to reach 83.02% overall Q3 (per residue) accuracy and 82.64 overall SOV score, which is very good, considering that the state-of-the-art results were around 84-85% Q3 accuracy.

METHOD	Q3 (%)	QH (%)	QE (%)	QC (%)	SOV	SOVH	SOVE	SOVC
CROSS-VALIDATION	79.50	84.39	70.50	79.78	76.44	78.86	76.87	72.20
ENSEMBLES (5 EXPERIMENTS / FOLD)	79.80	84.49	70.65	80.34	77.07	79.25	77.22	72.76
ENSEMBLES + EXTERNAL RULES	79.86	83.64	69.60	81.98	78.25	79.55	76.99	73.44
ENSEMBLES + EXTERNAL RULES + DECISION TREE	80.66	82.98	77.17	80.34	78.22	78.03	79.61	73.38
ENSEMBLES + DECISION TREE	81.93	84.45	81.27	79.75	79.25	79.43	82.59	74.68
ENSEMBLES + DECISION TREE + EXTERNAL RULES	81.96	83.22	79.79	81.94	81.40	80.67	82.10	76.43
ENSEMBLES + EXTERNAL RULES + RANDOM FOREST	80.90	84.03	78.38	79.17	78.50	79.31	80.03	72.76
ENSEMBLES + RANDOM FOREST	83.12	86.00	83.86	79.77	81.47	82.40	84.55	76.09
ENSEMBLES + RANDOM FOREST + EXTERNAL RULES	83.02	85.21	82.36	81.17	82.64	82.82	83.32	77.10

Table 5.32: 5-fold cross-validation, Q3 accuracy and SOV score for all methods for the PISCES dataset.

5.4 Best Results for CB513 and PISCES

In this section the hyper parameters for the best models will be displayed, along with the filtering methods used. For both CB513 and PISCES datasets the hyper parameters used to train the CNN models are shown in table 5.1. Figure 5.3 shows the hyper parameters for the trained models (CNN with SHN), the hyper parameters for the random forest filtering and the order of the applied techniques that produced the best final results for the CB513 dataset. These techniques resulted in an overall 81.80% Q3 (per residue) accuracy and an overall 78.98 SOV score. The confusion matrix for fold 0 of CB513 of a single CNN model trained with SHN method is shown in figure 5.4. The confusion matrix for the same fold for the ensembles model (combination of 5 trained models) with random forest filtering is shown in figure 5.5. It seems that the miss-predictions for classes ‘E’ and ‘H’ are less in the ensembles model, while the miss-predictions for class ‘C’ are slightly more, compared to the single CNN model.

MODEL HYPER PARAMETERS

GNsize	C	CNN layers	bsize	Max Iterations	Plus
2048	0.01	4	12288	50	7

FILTERING HYPER PARAMETERS

n_estimators	max_depth	random_state	min_samples_split	min_samples_leaf	Window
100	25	42	2	1	13

METHODS USED (IN ORDER)



Figure 5.3: Hyper parameters and methods used that resulted in the best overall Q3 accuracy and best overall SOV score for CB513 dataset.

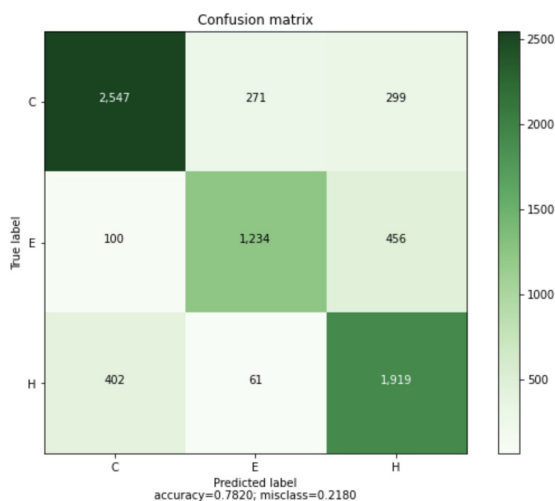


Figure 5.4: CM for CB513 fold 0 of single CNN model.

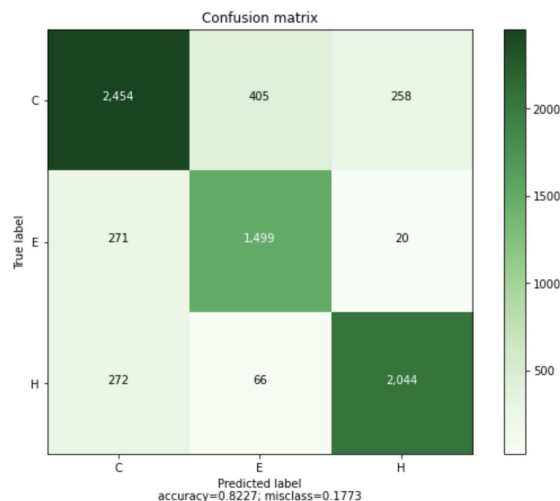


Figure 5.5: CM for CB513 fold 0 of ensembles model with random forest.

Figure 5.6 shows the hyper parameters for the trained models (CNN with SHN), the hyper parameters for the random forest filtering and the order of the applied methods, which produced the best final results for the PISCES dataset. These methods resulted in an overall 83.02% Q3 (per residue) accuracy and an overall 82.64 SOV score, which is very close with the current state-of-the-art results (84-85%).

MODEL HYPER PARAMETERS

GNsize	C	CNN layers	bsize	Max iterations	Plus
2048	0.01	4	12288	100	7

FILTERING HYPER PARAMETERS

n_estimators	max_depth	random_state	min_samples_split	min_samples_leaf	Window
100	25	42	2	1	19

METHODS USED (IN ORDER)



Figure 5.6: Hyper parameters and methods used that resulted in the best overall Q3 accuracy and best overall SOV score for PISCES dataset.

The confusion matrix for a single CNN model trained with SHN method, using fold 4 of PISCES, is displayed in figure 5.7. Figure 5.8 illustrates the confusion matrix for the ensembles model (combination of 5 CNN models) with random forest filtering, for

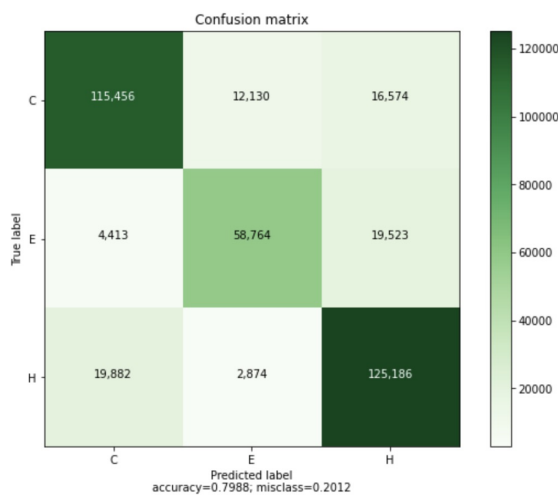


Figure 5.7: CM for PISCES fold 4 of single CNN model.

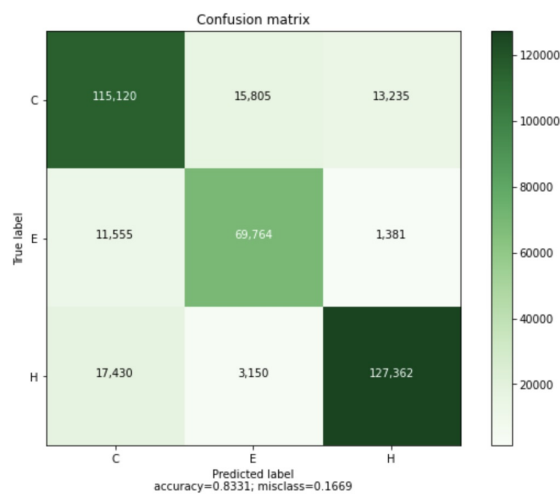


Figure 5.8: CM for PISCES fold 4 of ensembles model with random forest.

the same fold (fold 4). It is obvious that the correct predictions for classes ‘E’ and ‘H’ were increased, after applying the ensembles and filtering methods, while the correct predictions for class ‘C’ decreased by a small amount, compared to the single CNN model.

There are some hyper parameters that could still be modified and potentially improve the accuracy of the single CNN model with SHN. However, due to the shortage of available time, additional experiments could not be performed. In the next chapter, there will be some suggestions for further research regarding Convolutional Neural Networks with the Subsampled Hessian Newton method and the Protein Secondary Structure Prediction (PSSP) problem.

Chapter 6

Conclusion and Future Work

6.1	Conclusions	88
6.2	Suggestions for Future Work on PSSP	90

6.1 Conclusions

The initial purpose of this dissertation was to combine a Convolutional Neural Network (CNN) with the Hessian Free Optimisation (HFO) algorithm in order to train a model that predicts the Secondary Structure of Proteins (PSSP), given its primary structure, by exploiting the MSA profiles. The attempt to solve the PSSP problem was very important, since the experimental methods that are currently available are extremely expensive in both money and time. The ability to infer (predict) the secondary structure of proteins based on the primary structure could be very beneficial for the manufacture of pharmaceutical drugs, food complements and antibiotics. Except from that, the secondary structure of proteins could be used to determine the tertiary and quaternary structures, which could help scientists define the exact functionality of the studied proteins, and possibly provide an indication on how dangerous diseases, like cancer or covid-19, can be cured.

The attempt to combine the CNN with HFO was not fruitful because of the nature of HFO, which was specifically designed for Feed Forward Neural Networks (FFNNs). However, a new variation of HFO, called Subsampled Hessian Newton (SHN) method, was utilized to train a CNN for the PSSP problem. The results for the CB513 dataset were extremely promising with about 78.20% Q3 accuracy for a single fold (fold 0), 77.25% using 10-fold cross validation and approximately 77.46% using the ensembles method with 5 trained CNN models. The highest overall Q3 accuracy was 81.80% (Table 5.19) and was achieved by combining the ensembles model with random forest filtering and external rules. For the CB513 dataset, the overall SOV score for the 10-fold cross validation was 72.91, while the same figure increased to 73.55, after applying the ensembles method (with 5 trained models). The highest overall SOV score was achieved by combining the ensembles model with random forest and external rules filtering (Table 5.19), which was approximately 78.98.

Even though some studies reported results of 84-85% Q3 accuracy, the datasets used for training, were much larger than CB513, which means that they cannot be compared directly with the results of this dissertation. The results for CB513, however, can be compared with the results of [24] and [23], who also used the CB513 dataset (with the same 10-fold cross validation). The comparison between their results and the best results of this dissertation, for CB513 dataset, shows that the single CNN with SHN (78.20% Q3 accuracy) outperformed both the CNN with SGD (76.47% Q3 accuracy) [24], and the BRNN with HFO (77.01% Q3 accuracy) [23]. The same applies for the 10-fold cross validation, where the CNN with SHN achieved 77.25% (Table 5.6), the CNN with SGD achieved 75.16% and the BRNN with HFO achieved 75.80% overall Q3 accuracy. Moreover, the

ensembles model with random forest and external rules filtering (81.80% Q3 accuracy) managed to outperform the best results reported by [24] (80.40% overall Q3 accuracy) and [23] (78.15% overall Q3 accuracy).

For the PISCES dataset, the best overall Q3 accuracy of a single CNN with SHN was 79.88% (Table 5.23), while the overall Q3 accuracy for the 5-fold cross validation was 79.50%. The highest Q3 accuracy for the 5-fold cross validation was 83.12% and was achieved with the ensembles model with the random forest filtering (Table 5.30). The best overall SOV score for a single model was 76.67 and for the 5-fold cross validation was 76.44. The highest overall SOV score achieved was 82.64, with the combination of the ensembles model (with 5 trained CNNs), random forest and external rules filtering (Table 5.31).

The PISCES results can be compared with [1] and [2], where the same 5-fold cross validation was used. The overall Q3 accuracy for the 5-fold cross validation was 79.50% in this dissertation (for PISCES), which was slightly lower than [2] (79.57%) and considerably lower than [1] (80.65%). This means that there is still room for improvement for the models trained with the PISCES dataset. However, the final results for PISCES dataset were better (83.02% overall Q3 accuracy and 82.64 overall SOV score) compared to [1] (80.98% overall Q3 accuracy and 77.26 overall SOV score) and [2] (80.37% overall Q3 accuracy and 76.71 overall SOV score).

Based on the results from the various filtering methods, it seems that the order in which the filtering is applied plays a major role on the final outcome. That was expected, as the filtering of each method applies its own ‘corrections’ to the results, which could ‘reveal’ or ‘hide’ the errors for the next filtering method. For instance, in a sequence of seven (7) amino acids, where the predicted secondary structure is ‘HHHHEHE’, two (hypothetical) filtering methods could be applied. The first filtering method replaces sequences of ‘HHEH’ with ‘HHHH’, while the second method replaces ‘EHE’ with ‘EEE’. If the first method is applied, the sequence would become ‘HHHHHHE’ and after the second method it would remain the same, as there are no corrections that can be made. On a different scenario, the second filtering method could be applied first, which would result in a new sequence ‘HHHHEEEE’, that would remain the same after applying the first method. If the desired output was ‘HHHHHHHH’ the two approaches would result in completely different results. In the first scenario (1st method + 2nd method), the correct results would be 6 out of 7 (85.7%), while in the second scenario (2nd method + 1st method), the correct results would be 4 out of 7 (57.1%). This simple example shows that the order of the filtering methods applied can produce different final results, so various combinations

should be tested.

In addition to the above, the best filtering technique depends on the dataset and the output of the initial machine learning model. Even if two models are trained with the same machine learning architecture (e.g. CNN, FFNN, MLP, etc.) but with different optimisation algorithms (e.g. SGD, HFO, SHN, etc.), there is no guarantee that the boost from any of filtering methods will be the same for both models. The filtering results can vary between datasets and that can be observed from the filtering results of CB513 and PISCES datasets. It seems that there is no clear approach or a ‘best filtering’ method that will guarantee better results for all the machine learning models. Consequently, different filtering methods should be applied and the ones that produce the best results should be selected.

As mentioned earlier, some of the hyper parameters have not been exploited in this dissertation, not to a higher extend at least. This means there is still room for improvement for the prediction results. Some of these hyper parameters are the ‘plus’ parameter (window for CNN), the number of convolution layers, the number of filters and the kernel size of the CNN, as well as the number of trained models used in ensembles method. In addition to that, for the PISCES dataset, the hyper parameters used, to train the models, could be tuned in the same way that were tuned for the CB513 dataset.

6.2 Suggestions for Future Work on PSSP

Over the past years, multiple machine learning algorithms were utilized to predict the secondary structure of proteins, given their primary structure. However, none of these techniques managed to reach the maximum theoretical limit for the 3-class prediction of the PSSP problem, which is around 88-90% Q3 accuracy. That leaves the question whether a single machine learning model can even manage to reach such high accuracy, for such a complex problem (PSSP).

Maybe it is time to look for other alternative methods, like training multiple models and then combining their predicted results. One such method is called stacking ensemble method. This technique is similar to the ensembles method used in this dissertation, with the only difference that, instead of training the same type of model multiple times, it suggests to train different types of machine learning models, like K-means, Decision Trees, Support Vector Machines, Naive Bayes, Logistic Regression, and variations of Neural Networks. After the selected models are trained, they must be used to predict the test (or validation) dataset, and the predictions must be stored in an output file. Then a logistic regression model could be utilized, to learn how to best combine the predictions from

each of the separate models. This method does not guarantee that the stacked ensemble results will be better than the results of all separate models, however, even in that case, the model with the highest accuracy could be used instead of the stacked ensemble.

As many data scientists claim, ‘the answers are in the data’. Given this it is possible that the reason behind the accuracy limitations of the prediction models are related with the input data. Therefore, another suggestion would be to use different datasets to train the models or perform some modifications to the datasets in order to help the network extract the most important features. In addition to that, a separate dataset could be used only for tuning the hyper parameters and another one just for testing the model. According to the final results, it was obvious that the models were able to predict the classes ‘H’ and ‘C’ with higher success rate than the ‘E’ class. This means that the network could not extract all necessary features to be able to predict the ‘E’ class accurately. This phenomenon was observed because the datasets were not balanced.

It was observed that the accuracy of the predictions for various proteins with different lengths was not the same. A statistical analysis, on the final results, could help identify for which proteins the model had higher or lower accuracy. This could give an indication on which types of proteins the accuracy should be improved, in order to increase the overall accuracy of the model.

Different filtering methods can produce different results for different algorithms used for the PSSP and for different protein datasets, which means that experimentation, with various filtering methods, is essential in order to find the optimal filtering method. In this dissertation, only a few filtering methods were used, which leaves the door open for further experimentation with other filtering methods. The order, in which the filtering methods were applied, can affect the final results. Because of that, it is highly suggested to apply the filtering methods in different orders and choose the one that produces the best results. Moreover, the filtering methods could be applied before the ensembles method, which could lead to better results. Another approach could be to apply the external rules or other filtering methods multiple times (e.g. apply external rules, SVM filtering, ensembles method and then apply external rules for a second time).

In this dissertation, one of the goals was to train a CNN with the stochastic gradient descent (SGD) algorithm and compare the results with the SHN method. However, due to the lack of time these experiments have not been performed. These experiments could illustrate whether the SHN method could outperform the SGD algorithm for the PSSP problem. Theoretically, the CNN with SHN should require less time and iterations to train, compared to the CNN with SGD. The comparison between the CNN with SHN

and the CNN with SGD from [1] could not be made, because the implementations were different, the machines used for the experiments had different specifications and in [1] a CPU was used, instead of a GPU. According to some benchmark results from other problems [3], even though SHN performed almost the same with SGD in terms of accuracy, it was more robust than SGD in terms of hyper parameter tuning. This can still make SHN a better optimisation option, since the trained models that are required, in order to find the best hyper parameters, are significantly less than SGD. In addition, the total training iterations of SHN, for each model, are considerably less which means less training time.

References

- [1] A. Dionysiou, M. Agathocleous, C. Christodoulou, and V. Promponas, *Input representation of sequence to structure prediction problems for deep learning, (in preparation)*, 2020.
- [2] K. Charalambous, M. Agathocleous, C. Christodoulou, and V. Promponas, “Solving the protein secondary structure prediction problem with the hessian free optimization algorithm”, *IEEE Access*, 2020, Under review.
- [3] C.-C. Wang, K. Tan, and C.-J. Lin, “Newton methods for convolutional neural networks”, *ACM Transactions on Intelligent Systems and Technology*, vol. 11, pp. 1–30, 2020.
- [4] J. M. Berg, J. L. Tymoczko, and L. Stryer, *Biochemistry*, 5th. New York: NY: WH. Freeman, 2002.
- [5] C. Magnan and P. Baldi, “Sspro/accpro 5: Almost perfect prediction of protein secondary structure and relative solvent accessibility using profiles, machine learning and structural similarity.”, *Bioinformatics*, vol. 30, no. 18, pp. 2592–2597, 2014.
- [6] Y. Yang, J. Gao, J. Wang, R. Heffernan, J. Hanson, K. Paliwal, and Y. Zhou, “Sixty-five years of the long march in protein secondary structure prediction: The final stretch?”, *Briefings in Bioinformatics*, vol. 19, no. 3, pp. 482–494, 2018.
- [7] N. Qian and T. J. Sejnowski, “Predicting the secondary structure of globular proteins using neural network models”, *Journal of Molecular Biology*, vol. 202, no. 4, pp. 865–884, 1988.
- [8] B. Rost and C. Sander, “Improved prediction of protein secondary structure by use of sequence profiles and neural networks”, *Proceedings of the National Academy of Sciences, USA*, vol. 90, pp. 7558–7562, 1993.
- [9] A. Salamov and V. Soloveyev, “Ab initio gene finding in drosophila genomic dna”, *Genome Research*, vol. 10, no. 4, pp. 516–522, 2000.
- [10] R. D. King and M. J. Sternberg, “Identification and application of the concepts important for accurate and reliable protein secondary structure prediction”, *Protein Science*, vol. 5, no. 11, pp. 2298–2310, 1996.

- [11] D. Frishman and P. Argos, “Knowledge-based protein secondary structure assignment”, *Proteins: Structure, Function, and Bioinformatics*, vol. 23, no. 4, pp. 566–579, 1995.
- [12] J. A. Cuff and G. J. Barton, “Evaluation and improvement of multiple sequence methods for protein secondary structure prediction”, *Proteins: Structure, Function, and Bioinformatics*, vol. 34, no. 4, pp. 508–519, 1999.
- [13] P. Baldi, S. Brunak, P. Frasconi, G. Soda, and G. Pollastri, “Exploiting the past and the future in protein secondary structure prediction”, *Bioinformatics*, vol. 15, no. 11, pp. 937–946, 1999.
- [14] P. Baldi, S. Brunak, P. Frasconi, G. Pollastri, and G. Soda, *Bidirectional Dynamics for Protein Secondary Structure Prediction*. In: Sun R., Giles C.L. (eds) *Sequence Learning*, ser. Lecture Notes in Computer Science. Springer, Berlin, Heidelberg, 2000, vol. 1828, pp. 80–104.
- [15] J. Blazewicz, P. Hammer, and P. Lukasiak, “Predicting secondary structures of proteins”, *IEEE engineering in medicine and biology magazine*, vol. 24, no. 3, pp. 88–94, 2005.
- [16] G. Armano, A. Orro, and E. Vargiu, “Massp3: A system for predicting protein secondary structure”, *EURASIP Journal on Advances in Signal Processing*, vol. 2006, no. 17195, pp. 1–9, 2006.
- [17] F. U. Yüksektepe, O. Yılmaz, and M. Türkay, “Prediction of secondary structures of proteins using a two-stage method”, *Computers & Chemical Engineering*, vol. 32, no. 1–2, pp. 78–88, 2008.
- [18] K. J. Won, T. Hamelryck, A. Prugel-Bennett, and A. Krogh, “An evolutionary method for learning hmm structure: Prediction of protein secondary structure”, *BMC bioinformatics*, vol. 8, p. 357, Feb. 2007.
- [19] J. Chen and N. S. Chaudhari, “Cascaded bidirectional recurrent neural networks for protein secondary structure prediction”, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 4, no. 4, pp. 572–582, 2007.
- [20] S. Wang, J. Peng, J. Ma, and J. Xu, “Protein secondary structure prediction using deep convolutional neural fields”, *Scientific reports*, vol. 6, 2016.
- [21] P. Pavlidis, Πρόβλεψη δευτεροταγούς δομής των πρωτεϊνών με τη χρήση των *convolutional neural networks* για οπτική αναγνώριση αντικειμένων, University of Cyprus, Computer Science Department, Thesis Project, 2016.

- [22] R. Heffernan, Y. Yang, K. Paliwal, and Y. Zhou, “Capturing non-local interactions by long short-term memory bidirectional recurrent neural networks for improving prediction of protein secondary structure, backbone angles, contact numbers and solvent accessibility”, *Bioinformatics*, vol. 33, pp. 2842–2849, 2017.
- [23] K. Charalambous, *Protein secondary structure prediction using bidirectional recurrent neural networks and hessian free optimisation*, BSc Thesis, Department of Computer Science, University of Cyprus, 2018.
- [24] A. Dionysiou, Πρόβλεψη δευτεροταγούς δομής πρωτεϊνών με χρήση συνελικτικών νευρωνικών δικτύων σε συνδιασμό με φίλτρα *gabor* και *support vector machines*, BSc Thesis, Department of Computer Science, University of Cyprus, 2018.
- [25] C. Fang, Y. Shang, and D. Xu, “Mufold-ss: New deep inception–inside–inception networks for protein secondary structure prediction”, *Proteins: Structure, Function, and Bioinformatics*, 2018.
- [26] Learn.Genetics, *Genetic science learning center*, <https://learn.genetics.utah.edu/>, [Online; accessed April 21, 2020], 2018.
- [27] C. Simons, *Food science toolbox*, <https://cwsimons.com/wp-content/uploads/2017/12/AA.jpg>, [Online; accessed April 23, 2020].
- [28] A. Brunning, *A guide to the 20 common amino acids*, <http://www.compoundchem.com/wp-content/uploads/2014/09/20-Common-Amino-Acids-v3.png>, [Online; accessed April 23, 2020].
- [29] BioTopics, *Amino acid condensation*, <http://www.biotopics.co.uk/as/aminoon.html>, [Online; accessed April 27, 2020].
- [30] T. Brown and T. Brown Jr, *Nucleic acids book*, <https://www.atdbio.com/>, [Online; accessed April 25, 2020].
- [31] S. Clancy and W. Brown, “Translation: Dna to mrna to protein”, *Nature Education*, vol. 1, no. 1, p. 101, 2008.
- [32] M. A. Clark, M. Douglas, and J. Choi, *Biology 2e*, <https://openstax.org/books/biology-2e/pages/3-4-proteins>, [Online; accessed December 15, 2019], 2018.
- [33] A. Byun, *Convolutional neural networks for visual recognition*, <https://cs231n.github.io/neural-networks-1/>, [Online; accessed April 23, 2020].
- [34] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning”, *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [35] W. S. McCulloch and W. Pitts, “A logical calculus of the ideas immanent in nervous activity”, *The bulletin of mathematical biophysics*, vol. 5, no. 4, pp. 115–133, 1943.

- [36] F. Rosenblatt, “The perceptron: A probabilistic model for information storage and organization in the brain”, *Psychological Review*, pp. 65–386, 1958.
- [37] V. Kurková, “Kolmogorov’s theorem and multilayer neural networks”, *Neural Networks*, vol. 5, no. 3, pp. 501–506, 1991.
- [38] M. I. Jordan, *Serial order: A parallel distributed processing approach*, San Diego: University of California, Institute for Cognitive Science, Technical Report number AD-A-173989/5/XAB; ICS-8604, 1986.
- [39] J. L. Elman, “Finding structure in time”, *Cognitive Science*, vol. 14, no. 2, pp. 179–211, 1990.
- [40] K. Patel, *Mnist handwritten digits classification using a convolutional neural network (cnn)*, [Web; accessed December 16, 2019], 2019.
- [41] Y. L. Boureau, J. Ponce, and Y. LeCun, “A theoretical analysis of feature pooling in visual recognition”, in *In Proceedings of the 27th international conference on machine learning*, 2010, pp. 111–118.
- [42] T. Ganegedara, *Intuitive guide to cnns*, <https://towardsdatascience.com/light-on-math-machine-learning-intuitive-guide-to-convolution-neural-networks-e3f054dd5daa>, [Online; accessed April 27, 2020], 2018.
- [43] D. R. Hush and J. M. Salas, “Improving the learning rate of back propagation with the gradient reuse algorithm”, in *Proceedings of the IEEE 1988 International Conference on Neural Networks*, vol. 1, San Diego, CA, USA, 1988, pp. 441–447.
- [44] C. Charalambous, “Conjugate gradient algorithm for efficient training of artificial neural networks”, *Circuits, Devices and Systems, IEE Proceedings G*, vol. 139, no. 3, pp. 301–310, 1992.
- [45] C. M. Bishop, *Neural networks for pattern recognition*, Oxford University Press, Oxford, UK., 1995.
- [46] J. Martens, “Deep learning via hessian-free optimization”, in *Proceedings of the 27th International Conference on Machine Learning (ICML’10)*, Bottou, L. and Littman, M., (eds.), 2010, pp. 735–742.
- [47] S. G. Nash, “Newton-type minimization via the lanczos method”, *SIAM Journal on Numerical Analysis*, vol. 21, no. 4, pp. 770–788, 1984.
- [48] ———, “A survey of truncated-newton methods”, *Journal of Computational and Applied Mathematics*, vol. 124, no. 1, pp. 45–59, 2000.

- [49] J. Martens and I. Sutskever, *Training Deep and Recurrent Networks with Hessian-Free Optimization*. In: Montavon G., Orr G.B., Müller K.R. (eds) *Neural Networks: Tricks of the Trade*, ser. Lecture Notes in Computer Science. Springer, Berlin, Heidelberg, 2012, vol. 7700.
- [50] N. Schraudolph, “Fast curvature matrix-vector products for second-order gradient descent”, *Neural Computation*, vol. 14, no. 7, pp. 1723–1738, 2002.
- [51] R. E. Wengert, “A simple automatic derivative evaluation program”, *Communications of the ACM*, vol. 7, no. 8, pp. 463–464, 1964.
- [52] B. A. Pearlmutter, “Fast exact multiplication by the hessian”, *Neural Computation*, vol. 6, no. 1, pp. 147–160, 1994.
- [53] B. Rost, C. Sander, and R. Schneider, “Phd—an automatic mail server for protein secondary structure prediction”, *Bioinformatics*, vol. 10, no. 1, pp. 53–60, 1994.
- [54] A. Zemla, Č. Venclovas, K. Fidelis, and B. Rost, “A modified definition of sov, a segment-based measure for protein secondary structure prediction assessment”, *Proteins: Structure, Function, and Bioinformatics*, vol. 34, no. 2, pp. 220–223, 1999.
- [55] W. Kabsch and C. Sander, “Dictionary of protein secondary structure: Pattern recognition of hydrogen-bonded and geometrical features”, *Biopolymers*, vol. 22, no. 12, pp. 2577–2637, 1983.
- [56] G. Wang and R. L. Dunbrack Jr, “Pisces: A protein sequence culling server”, *Bioinformatics*, vol. 19, no. 12, pp. 1589–1591, 2003.
- [57] A. Dionysiou, M. Agathocleous, C. Christodoulou, and V. Promponas, “Convolutional neural networks in combination with support vector machines for complex sequential data classification”, *Artificial Neural Networks and Machine Learning - ICANN 2018, Lecture Notes in Computer Science*, ed. by V. Kurkova, Y. Manolopoulos, B. Hammer, L. Iliadis, I. Maglogiannis, Cham: Springer, vol. 11140, pp. 444–455, 2018.
- [58] C. Cortes and V. Vapnik, “Support-vector network”, *Machine Learning*, vol. 20, pp. 1–25, 1995.
- [59] P. Kountouris, M. Agathocleous, V. J. Promponas, G. Christodoulou, S. Hadjicostas, V. Vassiliades, and C. Christodoulou, “A comparative study on filtering protein secondary structure prediction”, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 9, no. 3, pp. 731–739, 2012.
- [60] T. Yiu, *Understanding random forest*, Torward Data Science, 2019. [Online]. Available: <https://towardsdatascience.com/understanding-random-forest-58381e0602d2>.

- [61] T. Kluyver, B. Ragan-Kelley, F. Pérez, B. Granger, M. Bussonnier, J. Frederic, K. Kelley, J. Hamrick, J. Grout, S. Corlay, P. Ivanov, D. Avila, S. Abdalla, and C. Willing, “Jupyter notebooks - a publishing format for reproducible computational workflows”, in *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, F. Loizides and B. Schmidt, Eds., IOS Press, 2016, pp. 87–90.
- [62] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, *TensorFlow: Large-scale machine learning on heterogeneous systems*, Software available from tensorflow.org, 2015. [Online]. Available: <https://www.tensorflow.org/>.
- [63] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library”, in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [64] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python”, *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [65] A. Botev, H. Ritter, and D. Barber, “Practical gauss-newton optimisation for deep learning”, in *In Proceedings of the 34th International Conference on Machine Learning*, 2017, pp. 557–565.
- [66] X. He, D. Mudigere, M. Smelyanskiy, and M. Taká, *Large scale distributed hessian free optimization for deep neural network*, 2017.
- [67] R. Kiros, *Training neural networks with stochastic hessian-free optimization*, arXiv preprint arXiv:1301.3641, 2013.
- [68] O. Vinyals and D. Povey, “Krylov subspace descent for deep learning”, in *In Proceedings of Artificial Intelligence and Statistics*, 2012, pp. 1261–1268.
- [69] C.-C. Wang, C.-H. Huang, and C.-J. Lin, “Subsampled hessian newton methods for supervised learning”, *Neural Computation*, vol. 27, no. 8, pp. 1766–1795, 2015.

- [70] C.-C. Wang, K.-L. Tan, C.-T. Chen, Y.-H. Lin, S. S. Keerthi, D. Mahajan, S. Sundararajan, and C.-J. Lin, “Distributed newton methods for deep learning”, *Neural Computation*, vol. 30, no. 6, pp. 1673–1724, 2018.
- [71] L. Breiman, “Random forests”, *Machine Learning*, vol. 45, pp. 5–32, 2001.

Appendix A

Excluded proteins from CB513

Table A.1 shows the names of the eight (8) proteins that were excluded from the CB513 dataset, because of zeroed MSA profiles.

No.	Protein
1	1coiA_1-29
2	1mctl_1-28
3	1tiiC_195-230
4	2erlA_1-40
5	1ceoA_202-254
6	1mrtA_31-61
7	1wfbB_1-37
8	6rlxC_-2-20

Table A.1: Excluded CB513 proteins due to zeroed MSA profiles.

Appendix B

Excluded proteins from PISCES

Table B.1 displays the identities of the sixteen (16) proteins that were excluded from the PISCES dataset, because their MSA profiles were missing.

No.	Protein	No.	Protein
1	1VPPX	9	4P6LA
2	3MLSP	10	3S0RA
3	4P2OP	11	1WFBA
4	1G0YI	12	1RPQW
5	4JQIV	13	4JO6Y
6	3UKWC	14	4H8LA
7	4H25C	15	1YODA
8	3SGRA	16	4KE2A

Table B.1: Excluded PISCES proteins due to missing MSA profiles.

Tables B.2, B.3 and B.4 illustrate the names of the PISCES proteins that were excluded from the PISCES dataset because their MSA profiles were corrupted or zeroed (according to [24]). In total those proteins were three hundred forty one (341).

No. Protein	No. Protein	No. Protein	No. Protein
1 3P51A	31 2OU5A	61 3U5WA	91 2R19A
2 1V96A	32 3L60A	62 3O6QA	92 3OP6A
3 3L7HA	33 3H0DA	63 3UV0A	93 2PW9A
4 4DHXA	34 1Q2HA	64 3NS4A	94 4GOFA
5 4F2LA	35 4I86A	65 2HQLA	95 3N7XA
6 2D7EA	36 2G7SA	66 3TE8A	96 2P9XA
7 4MTUA	37 2P63A	67 2O38A	97 2VS0A
8 4BSXA	38 3Q18A	68 4PF3A	98 2WG8A
9 4F87A	39 2FI1A	69 4H41A	99 4P49A
10 4MYVA	40 2Y5PA	70 3H16A	100 3ZC0A
11 1QV9A	41 4Q53A	71 2D59A	101 3B0FA
12 3FF5A	42 2Q3TA	72 1VR4A	102 2OX7A
13 4P2VA	43 1VPRA	73 2O1QA	103 4I16A
14 1WWPA	44 3DFUA	74 2HX5A	104 4PSFA
15 2D68A	45 3DNXA	75 2NPTA	105 4AP5A
16 2R85A	46 4GUCA	76 3CRYA	106 3K8RA
17 4MO0A	47 3Q6CA	77 2ERWA	107 3D3OA
18 4L3UA	48 4LTBA	78 3C4RA	108 3HLSA
19 4J5RA	49 3MD9A	79 2IP6A	109 1WPNA
20 2OL5A	50 3ESMA	80 3GO9A	110 2O99A
21 4KTWA	51 3H0NA	81 1YPYA	111 3I76A
22 3D33A	52 4HHVA	82 3E0RA	112 4LQBA
23 3PD7A	53 3M5QA	83 4F27A	113 4JX0A
24 3KVPA	54 4N74A	84 3PL0A	114 4R7RA
25 3QH6A	55 4KQIA	85 3IBWA	115 4GT9A
26 3VS8A	56 4IHQA	86 3TS9A	116 4F4WA
27 3CPTA	57 4K92A	87 3D55A	117 2UVPA
28 3GP6A	58 4Q7OA	88 3NOQA	118 1Z4RA
29 3O65A	59 4J1VA	89 4E6WA	119 2OU1A
30 1TU9A	60 4X33A	90 4AQOA	120 3F67A

Table B.2: Excluded PISCES proteins due to corrupted or zeroed MSA profiles (1-120).

No. Protein	No. Protein	No. Protein	No. Protein
121 1Q9UA	151 2R5SA	181 3ESLA	211 1U6HB
122 1WV3A	152 2FB0A	182 1U9LA	212 1UVQC
123 3ONQA	153 3L49A	183 3H35A	213 1ZVZB
124 2YF2A	154 2O4AA	184 1BB1A	214 1ZW2B
125 3Q0HA	155 2NS0A	185 1BB1B	215 2BPA3
126 3TUOA	156 2ZEXA	186 1BB1C	216 2BPTB
127 2E1FA	157 3UV1A	187 1C94A	217 2C5IP
128 3C0DA	158 3I4UA	188 1DPJB	218 2C5KP
129 3IV4A	159 3RK6A	189 1DTDB	219 2DS2A
130 2CVIA	160 4P78A	190 1F8VD	220 2ERLA
131 1SQWA	161 1UV7A	191 1GWMA	221 2GWWB
132 4BOQA	162 2HL7A	192 1HX6A	222 2HZSI
133 3W6SA	163 4U9OA	193 1KD8A	223 2MLTA
134 3TVQA	164 4BSVA	194 1KP6A	224 2P06A
135 3LQ9A	165 2C0NA	195 1KVEA	225 2PBDV
136 3BPQA	166 3I00A	196 1L2WI	226 2PLXB
137 2BDRA	167 3FH3A	197 1M3WA	227 2QUOA
138 3F43A	168 3NR5A	198 1M45B	228 2W4YA
139 3G21A	169 4E6SA	199 1M46B	229 2WBYC
140 4J91A	170 4LKUA	200 1MCTI	230 2WFUA
141 4K12A	171 3KTOA	201 1MQSB	231 2WFVA
142 2Q3SA	172 4QRNA	202 1MW5A	232 2WWXB
143 4QSGA	173 2V9PA	203 1O06A	233 2X5CA
144 2V73A	174 3F95A	204 1P9IA	234 2X5RA
145 2WVQA	175 1M1QA	205 1PJMA	235 2XF7A
146 4NUAA	176 3CSXA	206 1PJNA	236 2XZEQ
147 4G97A	177 4FX7A	207 1SVFB	237 3AJBB
148 3BRVA	178 3JQOA	208 1T01B	238 3C5TB
149 3O12A	179 4L2WA	209 1TQEX	239 3DT5A
150 3KUPA	180 3U97A	210 1TTWB	240 3E8YX

Table B.3: Excluded PISCES proteins due to corrupted or zeroed MSA profiles (121-240).

No. Protein	No. Protein	No. Protein	No. Protein
241 3FBLA	271 3UKXC	301 4H62V	331 4PNBA
242 3GP2B	272 3UL1A	302 4H7RA	332 4PNDA
243 3H0TC	273 3V62C	303 4H8MA	333 4PW1A
244 3HE4B	274 3V86A	304 4H8OA	334 4QMFA
245 3HE5A	275 3VU5B	305 4HB1A	335 4R0RA
246 3HE5B	276 3VU6B	306 4HBEA	336 4R80A
247 3L9AX	277 3VVIA	307 4HLBA	337 4R8TA
248 3LCNC	278 3W8VA	308 4HR1A	338 4RIQC
249 3LJMA	279 3W92A	309 4I7ZE	339 4TTLA
250 3M6ZA	280 3WKNE	310 4IICA	340 4UEBB
251 3NK4C	281 3WOEB	311 4J4AA	341 4W6YA
252 3OWTC	282 3WX4A	312 4JGLA	
253 3P06A	283 3WY9C	313 4JHKC	
254 3PLVC	284 3ZTAA	314 4KVTA	
255 3QKSC	285 4A94C	315 4KYTB	
256 3R46A	286 4BFHA	316 4LOOB	
257 3R4AA	287 4BLQA	317 4M1XA	
258 3RA3B	288 4BPLB	318 4M6BC	
259 3RF3C	289 4C1AA	319 4MGPA	
260 3RKLA	290 4CAYC	320 4MI8C	
261 3S1BA	291 4CU4B	321 4N3BB	
262 3S6PE	292 4CXFB	322 4N3CB	
263 3SHPA	293 4DACA	323 4OGQE	
264 3SJHB	294 4EHQG	324 4OQ9A	
265 3TQ2A	295 4F87A	325 4OYDB	
266 3TWEA	296 4FBWC	326 4OZKA	
267 3TZ1B	297 4FTBD	327 4PC0C	
268 3U4VA	298 4FZ0M	328 4PN8A	
269 3U4ZA	299 4G1AA	329 4PN9A	
270 3UC7A	300 4GVBB	330 4PNAA	

Table B.4: Excluded PISCES proteins due to corrupted or zeroed MSA profiles (241-341).

Appendix C

Convert datasets to Matlab files

The following bash script was created and used to convert the ‘.txt’ datasets (text files) to ‘.mat’ datasets (Matlab files).

```
1  #!/bin/bash
2  # This script finds all the testSet and trainSet files in the current directory,
3  # converts them to matlab datasets and saves them in the folder mat_datasets.
4  folderName="mat_datasets"
5  mkdir "$folderName"
6
7  runAll=""
8  datasets=$(ls | grep -e 'testSet' -e 'trainSet' | grep -v '\.mat')
9
10 for ds in $datasets
11 do
12     # echo "$ds"
13     loaded=$(echo "$ds" | sed "s:.txt:.")
14     # echo "$loaded"
15     outFile=$(echo "$folderName/$loaded.mat")
16     # echo "$outFile"
17     runMat="load _$ds; _y_=_loaded(1:end,_end); _Z_=_loaded(1:end,_1:end-1); _save _outFile _y _Z_ -v7.3; _clear;"
18     # echo "$runMat"
19     runAll="$runAll$runMat_"
20 done
21 runAll="$runAll_exit;"
22 # echo "$runAll"
23 /Applications/MATLAB_R2019a.app/bin/matlab -nodisplay -r "$runAll" > "$folderName/log.txt"
```

Appendix D

CB513 dataset pre-processing

This Python program prepares the CB513 datasets for training the Convolutional Neural Network. It was implemented for the purposes of this dissertation.

```
1  """
2  Uses the DATASETS files to create new datasets for CB513 based
3  on the specified number of amino acids (ADD_AMINO_ACIDS).
4  """
5
6  import os, sys
7
8  ADD_AMINO_ACIDS = 7 # 7 + 1 + 7 = 15 amino acids per row
9  DATASETS = ['trainSet0.txt', 'testSet0.txt',
10             'trainSet1.txt', 'testSet1.txt',
11             'trainSet2.txt', 'testSet2.txt',
12             'trainSet3.txt', 'testSet3.txt',
13             'trainSet4.txt', 'testSet4.txt',
14             'trainSet5.txt', 'testSet5.txt',
15             'trainSet6.txt', 'testSet6.txt',
16             'trainSet7.txt', 'testSet7.txt',
17             'trainSet8.txt', 'testSet8.txt',
18             'trainSet9.txt', 'testSet9.txt']
19
20  FOLDER_NAME = 'plus{0}_CB513'.format(ADD_AMINO_ACIDS)
21
22  if not os.path.exists(FOLDER_NAME):
23      os.makedirs(FOLDER_NAME)
24
25  protein_name = None
26  hssp_file = None
27  CATEGORIES = ['C', 'E', 'H']
28
29  def enumerate_cat(labels):
30      for i, cat in enumerate(CATEGORIES):
31          labels = labels.replace(cat, str(i))
32      return labels
33
34  def get_zero_lines(num_of_lines):
35      if (num_of_lines < 1):
36          return ""
37      zeros = (("0," * 20) + '\n') * num_of_lines
```

```

38     return zeros
39
40 for dataset_name in DATASETS:
41     # dataset_name = DATASETS[0]
42     print('Preparing_{0}..._Missing_hssp_files:'.format(dataset_name))
43     output_file = '{0}/plus{1}_{2}'.format(FOLDER_NAME, ADD_AMINO_ACIDS, dataset_name)
44
45     with open(dataset_name, 'r') as ds_f:
46         with open(output_file, 'w') as out_f:
47             line_num = 0
48             for line in ds_f:
49                 if line_num == 0:
50                     protein_name = line.split()[0]
51                     hssp_file = './msaFiles/{0}.hssp'.format(protein_name)
52                     # print(hssp_file)
53                     line_num += 1
54                 elif (line_num == 1):
55                     # print(line)
56                     line_num += 1
57                     continue
58                 else:
59                     labels = (line[:-1]).replace('!', '')
60                     label_nums = enumerate_cat(labels)
61                     label_index = 0
62                     # print(labels)
63                     try:
64                         with open(hssp_file, 'r') as hssp_f:
65                             buf = get_zero_lines(ADD_AMINO_ACIDS)
66                             buf_len = ADD_AMINO_ACIDS
67                             amino_count = 0
68                             for msa_line in hssp_f:
69                                 if (buf_len > 2 * ADD_AMINO_ACIDS):
70                                     temp = buf.replace('\n', '') + label_nums[amino_count]
71                                     out_f.write(temp)
72                                     out_f.write('\n')
73                                     buf = buf.split("\n", 1)[-1]
74                                     buf_len -= 1
75                                     amino_count += 1
76
77                                     modif_line = (msa_line[:-1]).replace('_', ' ')
78                                     buf = '{0}{1}\n'.format(buf, modif_line)
79                                     buf_len += 1
80
81
82                             for i in range(0, ADD_AMINO_ACIDS+1):
83                                 temp = buf.replace("\n", '') + label_nums[amino_count]
84                                 out_f.write(temp)
85                                 out_f.write('\n')
86                                 buf = buf.split("\n", 1)[-1]
87                                 buf = buf + get_zero_lines(1)
88                                 amino_count += 1
89                                 assert amino_count == len(label_nums)
90                             except Exception:
91                                 print(protein_name)
92                                 line_num = 0
93     print('Done_with_{0}_file!'.format(dataset_name))

```

Appendix E

PISCES dataset pre-processing

This Python program prepares the PISCES datasets for training the Convolutional Neural Network. It was implemented for the purposes of this dissertation.

```
1  """
2  Uses the DATASETS files to create new datasets for PISCES with
3  the specified number of neighboring amino acids (ADD_AMINO_ACIDS).
4  """
5
6  import os, sys
7
8  ADD_AMINO_ACIDS = 7 # 7 + 1 + 7 = 15 amino acids per row
9  DATASETS = ['trainSet1.txt', 'testSet1.txt',
10             'trainSet2.txt', 'testSet2.txt',
11             'trainSet3.txt', 'testSet3.txt',
12             'trainSet4.txt', 'testSet4.txt',
13             'trainSet5.txt', 'testSet5.txt']
14
15  FOLDER_NAME = 'plus{0}_PISCES'.format(ADD_AMINO_ACIDS)
16  MSA_FOLDER = 'msaFiles'
17
18  if not os.path.exists(FOLDER_NAME):
19      os.makedirs(FOLDER_NAME)
20
21  protein_name = None
22  hssp_file = None
23  CATEGORIES = ['C', 'E', 'H']
24
25  def enumerate_cat(labels):
26      for i, cat in enumerate(CATEGORIES):
27          labels = labels.replace(cat, str(i))
28      return labels
29
30  def get_zero_lines(num_of_lines):
31      zeros = ("0," * 20) + '\n' * num_of_lines
32      return zeros
33
34  for dataset_name in DATASETS:
35      # dataset_name = DATASETS[0]
36      print('Preparing_{0}..._Missing_hssp_files:'.format(dataset_name))
37      output_file = './{0}/plus{1}_{2}'.format(FOLDER_NAME, ADD_AMINO_ACIDS, dataset_name)
```

```

38
39 with open(dataset_name, 'r') as ds_f:
40     with open(output_file, 'w') as out_f:
41         line_num = 0
42         for line in ds_f:
43             if line_num == 0:
44                 protein_name = line.split()[0]
45                 hssp_file = '{0}/{1}.hssp'.format(MSA_FOLDER, protein_name)
46                 # print(hssp_file)
47                 line_num += 1
48             elif (line_num == 1):
49                 # print(line)
50                 line_num += 1
51             continue
52         else:
53             labels = line[:-1]
54             label_nums = enumerate_cat(labels)
55             label_index = 0
56             # print(labels)
57             try:
58                 with open(hssp_file, 'r') as hssp_f:
59                     buf = get_zero_lines(ADD_AMINO_ACIDS)
60                     buf_len = ADD_AMINO_ACIDS
61                     amino_count = 0
62                     for msa_line in hssp_f:
63                         if (buf_len > 2 * ADD_AMINO_ACIDS):
64                             temp = buf.replace('\n', '') + label_nums[amino_count]
65                             out_f.write(temp)
66                             out_f.write('\n')
67                             buf = buf.split("\n", 1)[-1]
68                             buf_len -= 1
69                             amino_count += 1
70
71                             modif_line = (msa_line[:-1]).replace('_', ',')
72                             buf = '{0}{1}\n'.format(buf, modif_line)
73                             buf_len += 1
74
75
76                     for i in range(0, ADD_AMINO_ACIDS+1):
77                         temp = buf.replace('\n', '') + label_nums[amino_count]
78                         out_f.write(temp)
79                         out_f.write('\n')
80                         buf = buf.split("\n", 1)[-1]
81                         buf = buf + get_zero_lines(1)
82                         amino_count += 1
83                     assert amino_count == len(label_nums)
84             except Exception:
85                 print(protein_name)
86                 line_num = 0
87     print('Done_with_{0}_file!'.format(dataset_name))

```

Appendix F

Python Implementation

The following code includes the implementation of the Convolutional Neural Network with the Subsampled Hessian Newton method. This program was used to perform all the experiments of this dissertation. Note that commands that begin with ‘!’ should be executed as bash commands. It is highly advised to use the notebook version of the implementation which can be found at [https://gitlab.com/perf.ai/pssp_project/-/blob/master/Notebooks/CNN_HFO.ipynb]. This implementation was based on the Python implementation from [3], however, several modifications have been made to improve the results of the CNN for the PSSP problem.

```
1  # -*- coding: utf-8 -*-
2  """shn_cnn_May22.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1KZtk3v3joX5pAUQJIbpGV9I-kmnddRjV
8  """
9
10 # plus_var=7
11 # ds_num=1
12 # dataset="PISCES"
13 plus_var=7
14 ds_num=5
15 dataset="CB513"
16
17 """## Imports """
18
19 # Commented out IPython magic to ensure Python compatibility.
20 # %load_ext autoreload
21 # %autoreload 2
22
23 # %matplotlib inline
24
25 !pip install hdf5storage
26
27 import pdb
28 import numpy as np
```

```

29 import tensorflow as tf
30 tf.compat.v1.disable_eager_execution()
31 import time
32 import math
33 import argparse
34 import os
35 import scipy.io as sio
36 import tensorflow.compat.v1 as tf
37 tf.disable_v2_behavior()
38 from tensorflow.python.client import device_lib
39 import pandas as pd
40 import hdf5storage
41
42 """## Get data ##"""
43
44 test_url="https://gitlab.com/perf.ai/pssp_project/-/raw/master/plus{0}_{1}/mat_datasets/plus{2}_testSet{3}.mat".format(str(
    ↪ plus_var), dataset, str(plus_var), str(ds_num))
45 train_url="https://gitlab.com/perf.ai/pssp_project/-/raw/master/plus{0}_{1}/mat_datasets/plus{2}_trainSet{3}.mat".format(str(
    ↪ plus_var), dataset, str(plus_var), str(ds_num))
46 TEST_FILE="plus" + str(plus_var) + "_testSet" + str(ds_num) + ".mat"
47 TRAIN_FILE="plus" + str(plus_var) + "_trainSet" + str(ds_num) + ".mat"
48
49 !echo "$test_url"
50 !echo "$train_url"
51
52 ![ -f "$TEST_FILE" ] && echo "$TEST_FILE_exist" || wget "$test_url"
53 ![ -f "$TRAIN_FILE" ] && echo "$TRAIN_FILE_exist" || wget "$train_url"
54
55 !ls
56
57 NEIGHBOURS = plus_var # number of amino-acids to add left and right
58 AMINO_ACID_LEN = 20
59 WINDOW = 2 * NEIGHBOURS + 1
60 TOTAL_AMINO_ACIDS = WINDOW * AMINO_ACID_LEN
61 TOTAL_COLS = TOTAL_AMINO_ACIDS + 1 # plus the secondary structure category
62 CATEGORIES = 3 # number of different classification categories
63 TOTAL_COLS
64
65 """## VGG ##"""
66
67 """
68 Codes are modifeid from PyTorch and Tensorflow Versions of VGG:
69 https://github.com/pytorch/vision/blob/master/torchvision/models/vgg.py, and
70 https://github.com/keras-team/keras-applications/blob/master/keras\_applications/vgg16.py
71 """
72
73 # import tensorflow.compat.v1 as tf
74 # tf.disable_v2_behavior()
75 # import numpy as np
76 # import pdb
77 from tensorflow.keras.applications.vgg16 import VGG16 as vgg16
78 from tensorflow.keras.applications.vgg19 import VGG19 as vgg19
79
80 __all__ = ['VGG11', 'VGG13', 'VGG16', 'VGG19']
81
82 def VGG(feature, num_cls):
83

```

```

84     with tf.variable_scope('fully_connected') as scope:
85         dim = np.prod(feature.shape[1:])
86         x = tf.reshape(feature, [-1, dim])
87
88         x = tf.keras.layers.Dense(units=4096, activation='relu', name=scope.name)(x)
89         x = tf.keras.layers.Dense(units=4096, activation='relu', name=scope.name)(x)
90         x = tf.keras.layers.Dense(units=num_cls, name=scope.name)(x)
91
92     return x
93
94 def make_layers(x, cfg):
95     for v in cfg:
96         if v == 'M':
97             x = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(x)
98         else:
99             x = tf.keras.layers.Conv2D(
100                 filters=v,
101                 kernel_size=[3, 3],
102                 padding='SAME',
103                 activation=tf.nn.relu
104             )(x)
105     return x
106
107 cfg = {
108     'A': [64, 'M', 128, 'M', 256, 256, 'M', 512, 512, 'M', 512, 512, 'M'],
109     'B': [64, 64, 'M', 128, 128, 'M', 256, 256, 'M', 512, 512, 'M', 512, 512, 'M'],
110     'D': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'M', 512, 512, 512, 'M', 512, 512, 512, 'M'],
111     'E': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 256, 'M', 512, 512, 512, 512, 'M',
112           512, 512, 512, 512, 'M'],
113 }
114
115 def VGG11(x_images, num_cls):
116     feature = make_layers(x_images, cfg['A'])
117     return VGG(feature, num_cls)
118
119 def VGG13(x_images, num_cls):
120     feature = make_layers(x_images, cfg['B'])
121     return VGG(feature, num_cls)
122
123 def VGG16(x_images, num_cls):
124     feature = make_layers(x_images, cfg['D'])
125     return VGG(feature, num_cls)
126
127 def VGG19(x_images, num_cls):
128     feature = make_layers(x_images, cfg['E'])
129     return VGG(feature, num_cls)
130
131 """## Net ##"""
132
133 # import tensorflow.compat.v1 as tf
134 # tf.disable_v2_behavior()
135 # import math
136 # import pdb
137 # from tensorflow.python.client import device_lib
138 # import numpy as np
139 # from net.vgg import *
140

```



```

141 def CNN_4layers(x_image, num_cls, reuse=False):
142     _NUM_CLASSES = num_cls
143     with tf.variable_scope('conv1', reuse=reuse) as scope:
144         conv = tf.keras.layers.Conv2D(
145             filters=64,
146             kernel_size=[3, 3],
147             padding='SAME',
148             activation=tf.nn.relu
149         )(x_image)
150
151     with tf.variable_scope('conv2', reuse=reuse) as scope:
152         conv = tf.keras.layers.Conv2D(
153             filters=64,
154             kernel_size=[3, 3],
155             padding='SAME',
156             activation=tf.nn.relu
157         )(conv)
158
159     with tf.variable_scope('conv3', reuse=reuse) as scope:
160         conv = tf.keras.layers.Conv2D(
161             filters=64,
162             kernel_size=[3, 3],
163             padding='SAME',
164             activation=tf.nn.relu
165         )(conv)
166
167     with tf.variable_scope('fully_connected', reuse=reuse) as scope:
168         dim = np.prod(conv.shape[1:])
169         flat = tf.reshape(conv, [-1, dim])
170         outputs = tf.keras.layers.Dense(units=_NUM_CLASSES, name=scope.name)(flat)
171
172     return outputs
173
174
175     # with tf.variable_scope('conv1', reuse=reuse) as scope:
176     # conv = tf.keras.layers.Conv2D(
177     #     filters=32,
178     #     kernel_size=[5, 5],
179     #     padding='SAME',
180     #     activation=tf.nn.relu
181     # )(x_image)
182     # pool = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(conv)
183     # # N x 16 x 16 x 32
184
185     # with tf.variable_scope('conv2', reuse=reuse) as scope:
186     # conv = tf.keras.layers.Conv2D(
187     #     filters=64,
188     #     kernel_size=[3, 3],
189     #     padding='SAME',
190     #     activation=tf.nn.relu
191     # )(pool)
192     # pool = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(conv)
193     # # N x 8 x 8 x 64
194
195     # with tf.variable_scope('conv3', reuse=reuse) as scope:
196     # conv = tf.keras.layers.Conv2D(
197     #     filters=64,

```

```

198 # kernel_size=[3, 3],
199 # padding='SAME',
200 # activation=tf.nn.relu
201 # )(pool)
202 # pool = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(conv)
203 # # N x 4 x 4 x 64
204
205 # with tf.variable_scope('fully_connected', reuse=reuse) as scope:
206 # dim = np.prod(pool.shape[1:])
207 # flat = tf.reshape(pool, [-1, dim])
208 # outputs = tf.keras.layers.Dense(units=_NUM_CLASSES, name=scope.name)(flat)
209
210 # return outputs
211
212 def CNN_7layers(x_image, num_cls, reuse=False):
213     _NUM_CLASSES = num_cls
214     with tf.variable_scope('conv1', reuse=reuse) as scope:
215         conv = tf.keras.layers.Conv2D(
216             filters=64,
217             kernel_size=[3, 3],
218             padding='SAME',
219             activation=tf.nn.relu
220         )(x_image)
221         conv = tf.keras.layers.Conv2D(
222             filters=64,
223             kernel_size=[3, 3],
224             padding='SAME',
225             activation=tf.nn.relu
226         )(conv)
227
228     with tf.variable_scope('conv2', reuse=reuse) as scope:
229         conv = tf.keras.layers.Conv2D(
230             filters=64,
231             kernel_size=[3, 3],
232             padding='SAME',
233             activation=tf.nn.relu
234         )(conv)
235         conv = tf.keras.layers.Conv2D(
236             filters=64,
237             kernel_size=[3, 3],
238             padding='SAME',
239             activation=tf.nn.relu
240         )(conv)
241
242     with tf.variable_scope('conv3', reuse=reuse) as scope:
243         conv = tf.keras.layers.Conv2D(
244             filters=64,
245             kernel_size=[3, 3],
246             padding='SAME',
247             activation=tf.nn.relu
248         )(conv)
249         conv = tf.keras.layers.Conv2D(
250             filters=64,
251             kernel_size=[3, 3],
252             padding='SAME',
253             activation=tf.nn.relu
254         )(conv)

```

```

255
256 with tf.variable_scope('fully_connected', reuse=reuse) as scope:
257     dim = np.prod(conv.shape[1:])
258     flat = tf.reshape(conv, [-1, dim])
259     outputs = tf.keras.layers.Dense(units=_NUM_CLASSES, name=scope.name)(flat)
260
261 return outputs
262
263 # with tf.variable_scope('conv1', reuse=reuse) as scope:
264 # conv = tf.keras.layers.Conv2D(
265 #     filters=32,
266 #     kernel_size=[5, 5],
267 #     padding='SAME',
268 #     activation=tf.nn.relu
269 # )(x_image)
270 # conv = tf.keras.layers.Conv2D(
271 #     filters=32,
272 #     kernel_size=[3, 3],
273 #     padding='SAME',
274 #     activation=tf.nn.relu
275 # )(conv)
276 # pool = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(conv)
277 # # N x 16 x 16 x 32
278
279 # with tf.variable_scope('conv2', reuse=reuse) as scope:
280 # conv = tf.keras.layers.Conv2D(
281 #     filters=64,
282 #     kernel_size=[3, 3],
283 #     padding='SAME',
284 #     activation=tf.nn.relu
285 # )(pool)
286 # conv = tf.keras.layers.Conv2D(
287 #     filters=64,
288 #     kernel_size=[3, 3],
289 #     padding='SAME',
290 #     activation=tf.nn.relu
291 # )(conv)
292 # pool = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(conv)
293 # # N x 8 x 8 x 64
294
295 # with tf.variable_scope('conv3', reuse=reuse) as scope:
296 # conv = tf.keras.layers.Conv2D(
297 #     filters=64,
298 #     kernel_size=[3, 3],
299 #     padding='SAME',
300 #     activation=tf.nn.relu
301 # )(pool)
302 # conv = tf.keras.layers.Conv2D(
303 #     filters=128,
304 #     kernel_size=[3, 3],
305 #     padding='SAME',
306 #     activation=tf.nn.relu
307 # )(conv)
308 # pool = tf.keras.layers.MaxPool2D(pool_size=[2, 2], strides=2, padding='valid')(conv)
309 # # pool = tf.layers.dropout(pool, rate=0.25, name=scope.name)
310 # # N x 4 x 4 x 128
311

```

```

312 # with tf.variable_scope('fully_connected', reuse=reuse) as scope:
313 #     dim = np.prod(pool.shape[1:])
314 #     flat = tf.reshape(pool, [-1, dim])
315 #     outputs = tf.keras.layers.Dense(units=_NUM_CLASSES, name=scope.name)(flat)
316
317 # return outputs
318
319 def CNN(net, num_cls, dim):
320
321     _NUM_CLASSES = num_cls
322     _IMAGE_HEIGHT, _IMAGE_WIDTH, _IMAGE_CHANNELS = dim
323
324     with tf.name_scope('main_params'):
325         x = tf.placeholder(tf.float32, shape=[None, _IMAGE_HEIGHT, _IMAGE_WIDTH, _IMAGE_CHANNELS], name='
            ↳ input_of_net')
326         y = tf.placeholder(tf.float32, shape=[None, _NUM_CLASSES], name='labels')
327
328     # call CNN structure according to string net
329     outputs = globals()[net](x, _NUM_CLASSES)
330     outputs = tf.identity(outputs, name='output_of_net')
331
332     return (x, y, outputs)
333
334 """## Utilities ##"""
335
336 # import numpy as np
337 # import math
338 # import scipy.io as sio
339 # import os
340 # import math
341 # import pdb
342
343 class ConfigClass(object):
344     def __init__(self, args, num_data, num_cls):
345         super(ConfigClass, self).__init__()
346         self.args = args
347         self.iter_max = args.iter_max
348
349         # Different notations of regularization term:
350         # In SGD, weight decay:
351         # weight_decay <- lr/(C*num_of_training_samples)
352         # In Newton method:
353         # C <- C * num_of_training_samples
354
355         self.seed = args.seed
356
357         if self.seed is None:
358             print('You_choose_not_to_specify_a_random_seed.'+'\
359                 'A_different_result_is_produced_after_each_run.')
360             elif isinstance(self.seed, int) and self.seed >= 0:
361                 print('You_specify_random_seed_{ }.'.format(self.seed))
362             else:
363                 raise ValueError('Only_accept_None_type_or_nonnegative_integers_for'+\
364                     '_random_seed_argument!')
365
366         self.train_set = args.train_set
367         self.val_set = args.val_set

```

```

368     self.num_cls = num_cls
369     self.dim = args.dim
370
371     self.num_data = num_data
372     self.GNsize = min(args.GNsize, self.num_data)
373     self.C = args.C * self.num_data
374     self.net = args.net
375
376     self.xi = 0.1
377     self.CGmax = args.CGmax
378     self._lambda = args._lambda
379     self.drop = args.drop
380     self.boost = args.boost
381     self.eta = args.eta
382     self.lr = args.lr
383     self.lr_decay = args.lr_decay
384
385     self.bsize = args.bsize
386     if args.momentum < 0:
387         raise ValueError('Momentum_needs_to_be_larger_than_0!')
388     self.momentum = args.momentum
389
390     self.loss = args.loss
391     if self.loss not in ('MSELoss', 'CrossEntropy'):
392         raise ValueError('Unrecognized_loss_type!')
393     self.optim = args.optim
394     if self.optim not in ('SGD', 'NewtonCG', 'Adam'):
395         raise ValueError('Only_support_SGD,_Adam_&_NewtonCG_optimizer!')
396
397     self.log_file = args.log_file
398     self.model_file = args.model_file
399     self.screen_log_only = args.screen_log_only
400
401     if self.screen_log_only:
402         print('You_choose_not_to_store_running_log._Only_store_model_to_{}'.format(self.log_file))
403     else:
404         print('Saving_log_to:{}'.format(self.log_file))
405         dir_name, _ = os.path.split(self.log_file)
406         if not os.path.isdir(dir_name):
407             os.makedirs(dir_name, exist_ok=True)
408
409         dir_name, _ = os.path.split(self.model_file)
410         if not os.path.isdir(dir_name):
411             os.makedirs(dir_name, exist_ok=True)
412
413     self.elapsed_time = 0.0
414
415 def read_data(filename, dim, label_enum=None):
416     """
417     args:
418         filename: the path where .mat files are stored
419         label_enum (default None): the list that stores the original labels.
420             If label_enum is None, the function will generate a new list which stores the
421             original labels in a sequence, and map original labels to [0, 1, ... number_of_classes-1].
422             If label_enum is a list, the function will use it to convert
423             original labels to [0, 1,..., number_of_classes-1].
424     """

```

```

425
426 # mat_contents = sio.loadmat(filename)
427 mat_contents = hdf5storage.loadmat(filename)
428 images, labels = mat_contents['Z'], mat_contents['y']
429
430 labels = labels.reshape(-1)
431 images = images.reshape(images.shape[0], -1)
432
433 _IMAGE_HEIGHT, _IMAGE_WIDTH, _IMAGE_CHANNELS = dim
434 zero_to_append = np.zeros((images.shape[0],
435     _IMAGE_CHANNELS*_IMAGE_HEIGHT*_IMAGE_WIDTH-np.prod(images.shape[1:]))
436 images = np.append(images, zero_to_append, axis=1)
437
438 # check data validity
439 if label_enum is None:
440     label_enum, labels = np.unique(labels, return_inverse=True)
441     num_cls = labels.max() + 1
442
443     if len(label_enum) != num_cls:
444         raise ValueError('The_number_of_classes_is_not_equal_to_the_number_of\
445     labels_in_dataset._Please_verify_them.')
446     else:
447         num_cls = len(label_enum)
448         forward_map = dict(zip(label_enum, np.arange(num_cls)))
449         labels = np.expand_dims(labels, axis=1)
450         labels = np.apply_along_axis(lambda x:forward_map[x[0]], axis=1, arr=labels)
451
452
453 # convert groundtruth to one-hot encoding
454 labels = np.eye(num_cls)[labels]
455 labels = labels.astype('float32')
456
457 return [images, labels], num_cls, label_enum
458
459 def normalize_and_reshape(images, dim, mean_tr=None):
460     _IMAGE_HEIGHT, _IMAGE_WIDTH, _IMAGE_CHANNELS = dim
461     images_shape = [images.shape[0], _IMAGE_CHANNELS, _IMAGE_HEIGHT, _IMAGE_WIDTH]
462
463     # images normalization and zero centering
464     images = images.reshape(images_shape[0], -1)
465
466     images = images/255.0
467
468     if mean_tr is None:
469         print('No_mean_of_data_provided!_Normalize_images_by_their_own_mean.')
470         # if no mean_tr is provided, we calculate it according to the current data
471         mean_tr = images.mean(axis=0)
472     else:
473         print('Normalize_images_according_to_the_provided_mean.')
474         if np.prod(mean_tr.shape) != np.prod(dim):
475             raise ValueError('Dimension_of_provided_mean_does_not_agree_with_the_data!_Please_verify_them!')
476
477     images = images - mean_tr
478
479     images = images.reshape(images_shape)
480     # Tensorflow accepts data shape: B x H x W x C
481     images = np.transpose(images, (0, 2, 3, 1))

```

```

482     return images, mean_tr
483
484
485 def predict(sess, network, test_batch, bsize):
486     x, y, loss, outputs = network
487
488     test_inputs, test_labels = test_batch
489     batch_size = bsize
490
491     num_data = test_labels.shape[0]
492     num_batches = math.ceil(num_data/batch_size)
493
494     results = np.zeros(shape=num_data, dtype=np.int)
495     infer_loss = 0.0
496
497     for i in range(num_batches):
498         batch_idx = np.arange(i*batch_size, min((i+1)*batch_size, num_data))
499
500         batch_input = test_inputs[batch_idx]
501         batch_labels = test_labels[batch_idx]
502
503         net_outputs, _loss = sess.run(
504             [outputs, loss], feed_dict={x: batch_input, y: batch_labels}
505         )
506
507         results[batch_idx] = np.argmax(net_outputs, axis=1)
508         # note that _loss was summed over batches
509         infer_loss = infer_loss + _loss
510
511     avg_acc = (np.argmax(test_labels, axis=1) == results).mean()
512     avg_loss = infer_loss/num_data
513
514     return avg_loss, avg_acc, results
515
516 """## Newton - CG """
517
518 # import pdb
519 # import tensorflow as tf
520 # import time
521 # import numpy as np
522 # import os
523 # import math
524 # from utilities import predict
525
526 def Rop(f, weights, v):
527     """Implementation of R operator
528     Args:
529         f: any function of weights
530         weights: list of tensors.
531         v: vector for right multiplication
532     Returns:
533         Jv: Jacobian vector product, length same as
534             the number of output of f
535     """
536     if type(f) == list:
537         u = [tf.zeros_like(ff) for ff in f]
538     else:

```

```

539     u = tf.zeros_like(f) # dummy variable
540     g = tf.gradients(ys=f, xs=weights, grad_ys=u)
541     return tf.gradients(ys=g, xs=u, grad_ys=v)
542
543 def Gauss_Newton_vec(outputs, loss, weights, v):
544     """Implements Gauss-Newton vector product.
545     Args:
546         loss: Loss function.
547         outputs: outputs of the last layer (pre-softmax).
548         weights: Weights, list of tensors.
549         v: vector to be multiplied with Gauss Newton matrix
550     Returns:
551         J'BJv: Gauss-Newton vector product.
552     """
553     # Validate the input
554     if type(weights) == list:
555         if len(v) != len(weights):
556             raise ValueError("weights_and_v_must_have_the_same_length.")
557
558     grads_outputs = tf.gradients(ys=loss, xs=outputs)
559     BJv = Rop(grads_outputs, weights, v)
560     JBJv = tf.gradients(ys=outputs, xs=weights, grad_ys=BJv)
561     return JBJv
562
563
564 class newton_cg(object):
565     def __init__(self, config, sess, outputs, loss):
566         """
567         initialize operations and vairables that will be used in newton
568         args:
569             sess: tensorflow session
570             outputs: output of the neural network (pre-softmax layer)
571             loss: function to calculate loss
572         """
573         super(newton_cg, self).__init__()
574         self.sess = sess
575         self.config = config
576         self.outputs = outputs
577         self.loss = loss
578         self.param = tf.compat.v1.trainable_variables()
579
580         self.CGiter = 0
581         FLOAT = tf.float32
582         model_weight = self.vectorize(self.param)
583
584         # initial variable used in CG
585         zeros = tf.zeros(model_weight.get_shape(), dtype=FLOAT)
586         self.r = tf.Variable(zeros, dtype=FLOAT, trainable=False)
587         self.v = tf.Variable(zeros, dtype=FLOAT, trainable=False)
588         self.s = tf.Variable(zeros, dtype=FLOAT, trainable=False)
589         self.g = tf.Variable(zeros, dtype=FLOAT, trainable=False)
590         # initial Gv, f for method minibatch
591         self.Gv = tf.Variable(zeros, dtype=FLOAT, trainable=False)
592         self.f = tf.Variable(0., dtype=FLOAT, trainable=False)
593
594         # rTr, cgtol and beta to be used in CG
595         self.rTr = tf.Variable(0., dtype=FLOAT, trainable=False)

```



```

596     self.cgtol = tf.Variable(0., dtype=FLOAT, trainable=False)
597     self.beta = tf.Variable(0., dtype=FLOAT, trainable=False)
598
599     # placeholder alpha, old_alpha and lambda
600     self.alpha = tf.compat.v1.placeholder(FLOAT, shape=[])
601     self.old_alpha = tf.compat.v1.placeholder(FLOAT, shape=[])
602     self._lambda = tf.compat.v1.placeholder(FLOAT, shape=[])
603
604     self.num_grad_segment = math.ceil(self.config.num_data/self.config.bsize)
605     self.num_Gv_segment = math.ceil(self.config.GNsize/self.config.bsize)
606
607     cal_loss, cal_lossgrad, cal_lossGv, \
608     add_reg_avg_loss, add_reg_avg_grad, add_reg_avg_Gv, \
609     zero_loss, zero_grad, zero_Gv = self._ops_in_minibatch()
610
611     # initial operations that will be used in minibatch and newton
612     self.cal_loss = cal_loss
613     self.cal_lossgrad = cal_lossgrad
614     self.cal_lossGv = cal_lossGv
615     self.add_reg_avg_loss = add_reg_avg_loss
616     self.add_reg_avg_grad = add_reg_avg_grad
617     self.add_reg_avg_Gv = add_reg_avg_Gv
618     self.zero_loss = zero_loss
619     self.zero_grad = zero_grad
620     self.zero_Gv = zero_Gv
621
622     self.CG, self.update_v = self._CG()
623     self.init_cg_vars = self._init_cg_vars()
624     self.update_gs = tf.tensordot(self.s, self.g, axes=1)
625     self.update_sGs = 0.5*tf.tensordot(self.s, -self.g-self.r-self._lambda*self.s, axes=1)
626     self.update_model = self._update_model()
627     self.gnorm = self.calc_norm(self.g)
628
629
630     def vectorize(self, tensors):
631         if isinstance(tensors, list) or isinstance(tensors, tuple):
632             vector = [tf.reshape(tensor, [-1]) for tensor in tensors]
633             return tf.concat(vector, 0)
634         else:
635             return tensors
636
637     def inverse_vectorize(self, vector, param):
638         if isinstance(vector, list):
639             return vector
640         else:
641             tensors = []
642             offset = 0
643             num_total_param = np.sum([np.prod(p.shape.as_list()) for p in param])
644             for p in param:
645                 numel = np.prod(p.shape.as_list())
646                 tensors.append(tf.reshape(vector[offset: offset+numel], p.shape))
647                 offset += numel
648
649             assert offset == num_total_param
650             return tensors
651
652     def calc_norm(self, v):

```

```

653     # default: frobenius norm
654     if isinstance(v, list):
655         norm = 0.
656         for p in v:
657             norm = norm + tf.norm(tensor=p)**2
658         return norm**0.5
659     else:
660         return tf.norm(tensor=v)
661
662     def _ops_in_minibatch(self):
663         """
664         Define operations that will be used in method minibatch
665         Vectorization is already a deep copy operation.
666         Before using newton method, loss needs to be summed over training samples
667         to make results consistent.
668         """
669
670     def cal_loss():
671         return tf.compat.v1.assign(self.f, self.f + self.loss)
672
673     def cal_lossgrad():
674         update_f = tf.compat.v1.assign(self.f, self.f + self.loss)
675
676         grad = tf.gradients(ys=self.loss, xs=self.param)
677         grad = self.vectorize(grad)
678         update_grad = tf.compat.v1.assign(self.g, self.g + grad)
679
680         return tf.group(*[update_f, update_grad])
681
682     def cal_lossGv():
683         v = self.inverse_vectorize(self.v, self.param)
684         Gv = Gauss_Newton_vec(self.outputs, self.loss, self.param, v)
685         Gv = self.vectorize(Gv)
686         return tf.compat.v1.assign(self.Gv, self.Gv + Gv)
687
688     # add regularization term to loss, gradient and Gv and further average over batches
689     def add_reg_avg_loss():
690         model_weight = self.vectorize(self.param)
691         reg = (self.calc_norm(model_weight))**2
692         reg = 1.0/(2*self.config.C) * reg
693         return tf.compat.v1.assign(self.f, reg + self.f/self.config.num_data)
694
695     def add_reg_avg_lossgrad():
696         model_weight = self.vectorize(self.param)
697         reg_grad = model_weight/self.config.C
698         return tf.compat.v1.assign(self.g, reg_grad + self.g/self.config.num_data)
699
700     def add_reg_avg_lossGv():
701         return tf.compat.v1.assign(self.Gv, (self._lambda + 1/self.config.C)*self.v
702             + self.Gv/self.config.GNsize)
703
704     # zero out loss, grad and Gv
705     def zero_loss():
706         return tf.compat.v1.assign(self.f, tf.zeros_like(self.f))
707     def zero_grad():
708         return tf.compat.v1.assign(self.g, tf.zeros_like(self.g))
709     def zero_Gv():

```

```

710         return tf.compat.v1.assign(self.Gv, tf.zeros_like(self.Gv))
711
712     return (cal_loss(), cal_lossgrad(), cal_lossGv(),
713            add_reg_avg_loss(), add_reg_avg_lossgrad(), add_reg_avg_lossGv(),
714            zero_loss(), zero_grad(), zero_Gv())
715
716     def minibatch(self, data_batch, place_holder_x, place_holder_y, mode):
717         """
718         A function to evaluate either function value, global gradient or sub-sampled Gv
719         """
720         if mode not in ('funonly', 'fungrad', 'Gv'):
721             raise ValueError('Unknown_mode_other_than_funonly_&_fungrad_&_Gv!')
722
723         inputs, labels = data_batch
724         num_data = labels.shape[0]
725         num_segment = math.ceil(num_data/self.config.bsize)
726         x, y = place_holder_x, place_holder_y
727
728         # before estimation starts, need to zero out f, grad and Gv according to the mode
729
730         if mode == 'funonly':
731             assert num_data == self.config.num_data
732             assert num_segment == self.num_grad_segment
733             self.sess.run(self.zero_loss)
734         elif mode == 'fungrad':
735             assert num_data == self.config.num_data
736             assert num_segment == self.num_grad_segment
737             self.sess.run([self.zero_loss, self.zero_grad])
738         else:
739             assert num_data == self.config.GNsize
740             assert num_segment == self.num_Gv_segment
741             self.sess.run(self.zero_Gv)
742
743         for i in range(num_segment):
744
745             load_time = time.time()
746             idx = np.arange(i * self.config.bsize, min((i+1) * self.config.bsize, num_data))
747             batch_input = inputs[idx]
748             batch_labels = labels[idx]
749             batch_input = np.ascontiguousarray(batch_input)
750             batch_labels = np.ascontiguousarray(batch_labels)
751             self.config.elapsed_time += time.time() - load_time
752
753             if mode == 'funonly':
754
755                 self.sess.run(self.cal_loss, feed_dict={
756                     x: batch_input,
757                     y: batch_labels,})
758
759             elif mode == 'fungrad':
760
761                 self.sess.run(self.cal_lossgrad, feed_dict={
762                     x: batch_input,
763                     y: batch_labels,})
764
765             else:
766

```

```

767         self.sess.run(self.cal_lossGv, feed_dict={
768             x: batch_input,
769             y: batch_labels})
770
771     # average over batches
772     if mode == 'funonly':
773         self.sess.run(self.add_reg_avg_loss)
774     elif mode == 'fungrad':
775         self.sess.run([self.add_reg_avg_loss, self.add_reg_avg_grad])
776     else:
777         self.sess.run(self.add_reg_avg_Gv,
778             feed_dict={self._lambda: self.config._lambda})
779
780
781     def _update_model(self):
782         update_model_ops = []
783         x = self.inverse_vectorize(self.s, self.param)
784         for i, p in enumerate(self.param):
785             op = tf.compat.v1.assign(p, p + (self.alpha - self.old_alpha) * x[i])
786             update_model_ops.append(op)
787         return tf.group(*update_model_ops)
788
789     def _init_cg_vars(self):
790         init_ops = []
791
792         init_r = tf.compat.v1.assign(self.r, -self.g)
793         init_v = tf.compat.v1.assign(self.v, -self.g)
794         init_s = tf.compat.v1.assign(self.s, tf.zeros_like(self.g))
795         gnorm = self.calc_norm(self.g)
796         init_rTr = tf.compat.v1.assign(self.rTr, gnorm**2)
797         init_cgtol = tf.compat.v1.assign(self.cgtol, self.config.xi*gnorm)
798
799         init_ops = [init_r, init_v, init_s, init_rTr, init_cgtol]
800
801         return tf.group(*init_ops)
802
803     def _CG(self):
804         """
805         CG:
806             define operations that will be used in method newton
807             Same as the previous loss calculation,
808             Gv has been summed over batches when samples were fed into Neural Network.
809         """
810
811         def CG_ops():
812
813             vGv = tf.tensordot(self.v, self.Gv, axes=1)
814
815             alpha = self.rTr / vGv
816             with tf.control_dependencies([alpha]):
817                 update_s = tf.compat.v1.assign(self.s, self.s + alpha * self.v, name='update_s_ops')
818                 update_r = tf.compat.v1.assign(self.r, self.r - alpha * self.Gv, name='update_r_ops')
819
820             with tf.control_dependencies([update_s, update_r]):
821                 rnewTrnew = self.calc_norm(update_r)**2
822                 update_beta = tf.compat.v1.assign(self.beta, rnewTrnew / self.rTr)
823                 with tf.control_dependencies([update_beta]):

```



```

881     idx = np.random.choice(np.arange(0, full_labels.shape[0]),
882                             size=self.config.GNsize, replace=False)
883
884     mini_inputs = full_inputs[idx]
885     mini_labels = full_labels[idx]
886
887     start = time.time()
888
889     self.sess.run(self.init_cg_vars)
890     cgtol = self.sess.run(self.cgtol)
891
892     avg_cg_time = 0.0
893     for CGiter in range(1, self.config.CGmax+1):
894
895         cg_time = time.time()
896         self.minibatch((mini_inputs, mini_labels), x, y, mode='Gv')
897         avg_cg_time += time.time() - cg_time
898
899         self.sess.run(self.CG)
900
901         mnewTrnew = self.sess.run(self.rTr)
902
903         if mnewTrnew**0.5 <= cgtol or CGiter == self.config.CGmax:
904             break
905
906         self.sess.run(self.update_v)
907
908     print('Avg_time_per_Gv_iteration:_{:.5f}_s\r\n'.format(avg_cg_time/CGiter))
909
910     gs, sGs = self.sess.run([self.update_gs, self.update_sGs], feed_dict={
911         self._lambda: self.config._lambda
912     })
913
914     # line_search
915     f_old = f
916     alpha = 1
917     while True:
918
919         old_alpha = 0 if alpha == 1 else alpha/0.5
920
921         self.sess.run(self.update_model, feed_dict={
922             self.alpha:alpha, self.old_alpha:old_alpha
923         })
924
925         prered = alpha*gs + (alpha**2)*sGs
926
927         self.minibatch(full_batch, x, y, mode='funonly')
928         f = self.sess.run(self.f)
929
930         actred = f - f_old
931
932         if actred <= self.config.eta*alpha*gs:
933             break
934
935         alpha *= 0.5
936
937     # update lambda

```

```

938     ratio = actred / prered
939     if ratio < 0.25:
940         self.config._lambda *= self.config.boost
941     elif ratio >= 0.75:
942         self.config._lambda *= self.config.drop
943
944     self.minibatch(full_batch, x, y, mode='fungrad')
945     f = self.sess.run(self.f)
946
947     gnorm = self.sess.run(self.gnorm)
948
949     summary = self.sess.run(merged)
950     train_writer.add_summary(summary, k)
951
952     # exclude data loading time for fair comparison
953     end = time.time()
954
955     end = end - self.config.elapsed_time
956     total_running_time += end - start
957
958     self.config.elapsed_time = 0.0
959
960     total_CG += CGiter
961
962     output_str = '{ } --iter_f:_{:.3f}_gl:_{:.5f}_alpha:_{:.3e}_ratio:_{:.3f}_lambda:_{:.5f}_#CG:_{ }_actred:_{:.5f}_prered:
        ↳ _{:.5f}_time:_{:.3f}'.\
963         format(k, f, gnorm, alpha, actred/prered, self.config._lambda, CGiter, actred, prered, end - start)
964     print(output_str)
965     if not self.config.screen_log_only:
966         print(output_str, file=log_file)
967
968     if val_batch is not None:
969         # Evaluate the performance after every Newton Step
970         if test_network == None:
971             val_loss, val_acc, _ = predict(
972                 self.sess,
973                 network=(x, y, self.loss, outputs),
974                 test_batch=val_batch,
975                 bsize=self.config.bsize,
976             )
977         else:
978             # A separat test network part has not been done...
979             val_loss, val_acc, _ = predict(
980                 self.sess,
981                 network=test_network,
982                 test_batch=val_batch,
983                 bsize=self.config.bsize
984             )
985
986     output_str = '\r\n{ } --iter_val_acc:_{:.3f}%_val_loss_{:.3f}\r\n'.\
987         format(k, val_acc*100, val_loss)
988     print(output_str)
989     if not self.config.screen_log_only:
990         print(output_str, file=log_file)
991
992     if val_acc > best_acc:
993         best_acc = val_acc

```

```

994         checkpoint_path = self.config.model_file
995         save_path = saver.save(self.sess, checkpoint_path)
996         print('Best_model_saved_in_{ }\r\n'.format(save_path))
997
998     if val_batch is None:
999         checkpoint_path = self.config.model_file
1000         save_path = saver.save(self.sess, checkpoint_path)
1001         print('Model_at_the_last_iteration_saved_in_{ }\r\n'.format(save_path))
1002         output_str = 'total_#CG_{ }_total_running_time_{:.3f}s'.format(total_CG, total_running_time)
1003     else:
1004         output_str = 'Final_acc:{:.3f}%_best_acc:{:.3f}%_total_#CG_{ }_total_running_time_{:.3f}s'\
1005             format(val_acc*100, best_acc*100, total_CG, total_running_time)
1006     print(output_str)
1007     if not self.config.screen_log_only:
1008         print(output_str, file=log_file)
1009     log_file.close()
1010
1011     """##Set Train Arguments##"""
1012
1013     # Arguments for HFO - PSSP dataset
1014     train_args = ("--optim_NewtonCG_--GNsize_2048_--C_0.01_--net_CNN_4layers_--bsize_12288_--iter_max_50_
1015         ↪ " +
1016         "--train_set_/" + TRAIN_FILE + "_--val_set_/" + TEST_FILE + "_--dim_" +
1017         str(WINDOW) + "_" + str(AMINO_ACID_LEN) + "_1").split()
1018
1019     # Arguments for SGD - PSSP dataset
1020     # train_args = ("--optim SGD --lr 0.01 --C 0.01 --net CNN_4layers --bsize 256 " +
1021     # "--train_set _/" + TRAIN_FILE + " --val_set _/" + TEST_FILE + " --dim " +
1022     # str(WINDOW) + " " + str(AMINO_ACID_LEN) + " 1").split()
1023
1024     """##Declare Train Function##"""
1025
1026     # import pdb
1027     # import numpy as np
1028     # import tensorflow as tf
1029     # tf.compat.v1.disable_eager_execution()
1030     # import time
1031     # import math
1032     # import argparse
1033
1034     # from net.net import CNN
1035     # from newton_cg import newton_cg
1036     # from utilities import read_data, predict, ConfigClass, normalize_and_reshape
1037
1038     def parse_args():
1039         parser = argparse.ArgumentParser(description='Newton_method_on_DNN')
1040         parser.add_argument('--C', dest='C',
1041             help='regularization_term,_or_so-called_weight_decay_where'+\
1042             'weight_decay=_lr/(C*num_of_samples)_in_this_implementation',
1043             default=0.01, type=float)
1044
1045         # Newton method arguments
1046         parser.add_argument('--GNsize', dest='GNsize',
1047             help='number_of_samples_for_estimating_Gauss-Newton_matrix',
1048             default=4096, type=int)
1049         parser.add_argument('--iter_max', dest='iter_max',
1050             help='the_maximal_number_of_Newton_iterations',

```



```

1050         default=100, type=int)
1051     parser.add_argument('--xi', dest='xi',
1052         help='the_tolerance_in_the_relative_stopping_condition_for_CG',
1053         default=0.1, type=float)
1054     parser.add_argument('--drop', dest='drop',
1055         help='the_drop_constants_for_the_LM_method',
1056         default=2/3, type=float)
1057     parser.add_argument('--boost', dest='boost',
1058         help='the_boost_constants_for_the_LM_method',
1059         default=3/2, type=float)
1060     parser.add_argument('--eta', dest='eta',
1061         help='the_parameter_for_the_line_search_stopping_condition',
1062         default=0.0001, type=float)
1063     parser.add_argument('--CGmax', dest='CGmax',
1064         help='the_maximal_number_of_CG_iterations',
1065         default=250, type=int)
1066     parser.add_argument('--lambda', dest='_lambda',
1067         help='the_initial_lambda_for_the_LM_method',
1068         default=1, type=float)
1069
1070     # SGD arguments
1071     parser.add_argument('--epoch_max', dest='epoch',
1072         help='number_of_training_epoch',
1073         default=500, type=int)
1074     parser.add_argument('--lr', dest='lr',
1075         help='learning_rate',
1076         default=0.01, type=float)
1077     parser.add_argument('--decay', dest='lr_decay',
1078         help='learning_rate_decay_over_each_mini_batch_update',
1079         default=0, type=float)
1080     parser.add_argument('--momentum', dest='momentum',
1081         help='momentum_of_learning',
1082         default=0, type=float)
1083
1084     # Model training arguments
1085     parser.add_argument('--bsize', dest='bsize',
1086         help='batch_size_to_evaluate_stochastic_gradient_Gv_etc.Since_the_sampled_data_\\
1087         \\for_computing_Gauss-Newton_matrix_and_etc._might_not_fit_into_memory_\\
1088         \\for_one_time,_we_will_split_the_data_into_several_segements_and_average\\
1089         \\over_them.',
1090         default=1024, type=int)
1091     parser.add_argument('--net', dest='net',
1092         help='classifier_type',
1093         default='CNN_4layers', type=str)
1094     parser.add_argument('--train_set', dest='train_set',
1095         help='provide_the_directory_of_.mat_file_for_training',
1096         default=None, type=str)
1097     parser.add_argument('--val_set', dest='val_set',
1098         help='provide_the_directory_of_.mat_file_for_validation',
1099         default=None, type=str)
1100     parser.add_argument('--model', dest='model_file',
1101         help='model_saving_address',
1102         default='./saved_model/model.ckpt', type=str)
1103     parser.add_argument('--log', dest='log_file',
1104         help='log_saving_directory',
1105         default='./running_log/logger.log', type=str)
1106     parser.add_argument('--screen_log_only', dest='screen_log_only',

```

```

1107         help='screen_printing_running_log_instead_of_storing_it',
1108         action='store_true')
1109     parser.add_argument('--optim', '--optim',
1110                         help='which_optimizer_to_use:_SGD,_Adam_or_NewtonCG',
1111                         default='NewtonCG', type=str)
1112     parser.add_argument('--loss', dest='loss',
1113                         help='which_loss_function_to_use:_MSELoss_or_CrossEntropy',
1114                         default='MSELoss', type=str)
1115     parser.add_argument('--dim', dest='dim', nargs='+', help='input_dimension_of_data,'+
1116                         'shape_must_be:_height_width_num_channels',
1117                         default=[32, 32, 3], type=int)
1118     parser.add_argument('--seed', dest='seed', help='a_nonnegative_integer_for_'
1119     'reproducibility', type=int)
1120
1121     args = parser.parse_args(args=train_args)
1122     return args
1123
1124
1125     args = parse_args()
1126
1127     def init_model(param):
1128         init_ops = []
1129         for p in param:
1130             if 'kernel' in p.name:
1131                 weight = np.random.standard_normal(p.shape)* np.sqrt(2.0 / ((np.prod(p.get_shape().as_list()[:-1]))))
1132                 opt = tf.compat.v1.assign(p, weight)
1133             elif 'bias' in p.name:
1134                 zeros = np.zeros(p.shape)
1135                 opt = tf.compat.v1.assign(p, zeros)
1136             init_ops.append(opt)
1137         return tf.group(*init_ops)
1138
1139     def gradient_trainer(config, sess, network, full_batch, val_batch, saver, test_network):
1140         x, y, loss, outputs, = network
1141
1142         global_step = tf.Variable(initial_value=0, trainable=False, name='global_step')
1143         learning_rate = tf.compat.v1.placeholder(tf.float32, shape=[], name='learning_rate')
1144
1145         # Probably not a good way to add regularization.
1146         # Just to confirm the implementation is the same as MATLAB.
1147         reg = 0.0
1148         param = tf.compat.v1.trainable_variables()
1149         for p in param:
1150             reg = reg + tf.reduce_sum(input_tensor=tf.pow(p,2))
1151         reg_const = 1/(2*config.C)
1152         batch_size = tf.compat.v1.cast(tf.shape(x)[0], tf.float32)
1153         loss_with_reg = reg_const*reg + loss/batch_size
1154
1155         if config.optim == 'SGD':
1156             optimizer = tf.compat.v1.train.MomentumOptimizer(
1157                 learning_rate=learning_rate,
1158                 momentum=config.momentum).minimize(
1159                     loss_with_reg,
1160                     global_step=global_step)
1161         elif config.optim == 'Adam':
1162             optimizer = tf.compat.v1.train.AdamOptimizer(learning_rate=learning_rate,
1163                 beta1=0.9,

```

```

1164         beta2=0.999,
1165         epsilon=1e-08).minimize(
1166         loss_with_reg,
1167         global_step=global_step)
1168
1169     train_inputs, train_labels = full_batch
1170     num_data = train_labels.shape[0]
1171     num_iters = math.ceil(num_data/config.bsize)
1172
1173     print(config.args)
1174     if not config.screen_log_only:
1175         log_file = open(config.log_file, 'w')
1176         print(config.args, file=log_file)
1177     sess.run(tf.compat.v1.global_variables_initializer())
1178
1179
1180     print('-----_initializing_network_by_methods_in_He_et_al._(2015)-----')
1181     param = tf.compat.v1.trainable_variables()
1182     sess.run(init_model(param))
1183
1184     total_running_time = 0.0
1185     best_acc = 0.0
1186     lr = config.lr
1187
1188     for epoch in range(0, args.epoch):
1189
1190         loss_avg = 0.0
1191         start = time.time()
1192
1193         for i in range(num_iters):
1194
1195             load_time = time.time()
1196             # randomly select the batch
1197             idx = np.random.choice(np.arange(0, num_data),
1198                                   size=config.bsize, replace=False)
1199
1200             batch_input = train_inputs[idx]
1201             batch_labels = train_labels[idx]
1202             batch_input = np.ascontiguousarray(batch_input)
1203             batch_labels = np.ascontiguousarray(batch_labels)
1204             config.elapsed_time += time.time() - load_time
1205
1206             step, _, batch_loss= sess.run(
1207                 [global_step, optimizer, loss_with_reg],
1208                 feed_dict = {x: batch_input, y: batch_labels, learning_rate: lr}
1209             )
1210
1211             # print initial loss
1212             if epoch == 0 and i == 0:
1213                 output_str = 'initial_f_(reg+_avg_loss_of_1st_batch):_{:.3f}'.format(batch_loss)
1214                 print(output_str)
1215                 if not config.screen_log_only:
1216                     print(output_str, file=log_file)
1217
1218             loss_avg = loss_avg + batch_loss
1219             # print log every 10% of the iterations
1220             if i % math.ceil(num_iters/10) == 0:

```

```

1221     end = time.time()
1222     output_str = 'Epoch_{:d}/{:d} \t loss_{:.4f} lr_{:.6f} elapsed_time_{:.3f}'\
1223         .format(epoch, i, num_iters, batch_loss, lr, end-start)
1224     print(output_str)
1225     if not config.screen_log_only:
1226         print(output_str, file=log_file)
1227
1228     # adjust learning rate for SGD by inverse time decay
1229     if args.optim != 'Adam':
1230         lr = config.lr/(1 + args.lr_decay*step)
1231
1232     # exclude data loading time for fair comparison
1233     epoch_end = time.time() - config.elapsed_time
1234     total_running_time += epoch_end - start
1235     config.elapsed_time = 0.0
1236
1237     if val_batch is None:
1238         output_str = 'In_epoch_{:d} \t train_loss_{:.3f} \t epoch_time_{:.3f}'\
1239             .format(epoch, loss_avg/(i+1), epoch_end-start)
1240     else:
1241         if test_network == None:
1242             val_loss, val_acc, _ = predict(
1243                 sess,
1244                 network=(x, y, loss, outputs),
1245                 test_batch=val_batch,
1246                 bsize=config.bsize
1247             )
1248         else:
1249             # A separat test network part have been done...
1250             val_loss, val_acc, _ = predict(
1251                 sess,
1252                 network=test_network,
1253                 test_batch=val_batch,
1254                 bsize=config.bsize
1255             )
1256
1257     output_str = 'In_epoch_{:d} \t train_loss_{:.3f} \t val_loss_{:.3f} \t val_accuracy_{:.3f} % \t epoch_time_{:.3f}'\
1258         .format(epoch, loss_avg/(i+1), val_loss, val_acc*100, epoch_end-start)
1259
1260     if val_acc > best_acc:
1261         best_acc = val_acc
1262         checkpoint_path = config.model_file
1263         save_path = saver.save(sess, checkpoint_path)
1264         print('Saved_best_model_in_{:d}' .format(save_path))
1265
1266     print(output_str)
1267     if not config.screen_log_only:
1268         print(output_str, file=log_file)
1269
1270     if val_batch is None:
1271         checkpoint_path = config.model_file
1272         save_path = saver.save(sess, checkpoint_path)
1273         print('Model_at_the_last_iteration_saved_in_{:d}\n' .format(save_path))
1274         output_str = 'total_running_time_{:.3f}s' .format(total_running_time)
1275     else:
1276         output_str = 'Final_acc_{:.3f} % \t best_acc_{:.3f} % \t total_running_time_{:.3f}s'\
1277             .format(val_acc*100, best_acc*100, total_running_time)

```

```

1278
1279     print(output_str)
1280     if not config.screen_log_only:
1281         print(output_str, file=log_file)
1282         log_file.close()
1283
1284     def newton_trainer(config, sess, network, full_batch, val_batch, saver, test_network):
1285
1286         _, _, loss, outputs = network
1287         newton_solver = newton_cg(config, sess, outputs, loss)
1288         sess.run(tf.compat.v1.global_variables_initializer())
1289
1290         print('-----_initializing_network_by_methods_in_He_et_al._(2015)_-----')
1291         param = tf.compat.v1.trainable_variables()
1292         sess.run(init_model(param))
1293         newton_solver.newton(full_batch, val_batch, saver, network, test_network)
1294
1295
1296     def train_model():
1297         full_batch, num_cls, label_enum = read_data(filename=args.train_set, dim=args.dim)
1298
1299         if args.val_set is None:
1300             print('No_validation_set_is_provided._Will_output_model_at_the_last_iteration.')
1301             val_batch = None
1302         else:
1303             val_batch, _, _ = read_data(filename=args.val_set, dim=args.dim, label_enum=label_enum)
1304
1305         num_data = full_batch[0].shape[0]
1306
1307         config = ConfigClass(args, num_data, num_cls)
1308
1309         if isinstance(config.seed, int):
1310             tf.compat.v1.random.set_random_seed(config.seed)
1311             np.random.seed(config.seed)
1312
1313         if config.net in ('CNN_4layers', 'CNN_7layers', 'VGG11', 'VGG13', 'VGG16', 'VGG19'):
1314             x, y, outputs = CNN(config.net, num_cls, config.dim)
1315             test_network = None
1316         else:
1317             raise ValueError('Unrecognized_training_model')
1318
1319         if config.loss == 'MSELoss':
1320             loss = tf.reduce_sum(input_tensor=tf.pow(outputs-y, 2))
1321         else:
1322             loss = tf.reduce_sum(input_tensor=tf.nn.softmax_cross_entropy_with_logits(logits=outputs, labels=y))
1323
1324         network = (x, y, loss, outputs)
1325
1326         sess_config = tf.compat.v1.ConfigProto()
1327         sess_config.gpu_options.allow_growth = True
1328
1329         with tf.compat.v1.Session(config=sess_config) as sess:
1330
1331             full_batch[0], mean_tr = normalize_and_reshape(full_batch[0], dim=config.dim, mean_tr=None)
1332             if val_batch is not None:
1333                 val_batch[0], _ = normalize_and_reshape(val_batch[0], dim=config.dim, mean_tr=mean_tr)
1334

```

```

1335     param = tf.compat.v1.trainable_variables()
1336
1337     mean_param = tf.compat.v1.get_variable(name='mean_tr', initializer=mean_tr, trainable=False,
1338         validate_shape=True, use_resource=False)
1339     label_enum_var=tf.compat.v1.get_variable(name='label_enum', initializer=label_enum, trainable=False,
1340         validate_shape=True, use_resource=False)
1341     saver = tf.compat.v1.train.Saver(var_list=param+[mean_param])
1342
1343     if config.optim in ('SGD', 'Adam'):
1344         gradient_trainer(
1345             config, sess, network, full_batch, val_batch, saver, test_network)
1346     elif config.optim == 'NewtonCG':
1347         newton_trainer(
1348             config, sess, network, full_batch, val_batch, saver, test_network=test_network)
1349
1350     """## Train ##"""
1351
1352     train_model()
1353
1354     """## Predict ##"""
1355
1356     # Arguments for prediction PSSP dataset
1357     pred_args = ("--bsize_1024_--test_set_/" + TEST_FILE + "_--train_set_/" + TRAIN_FILE +
1358         "_--model_/saved_model/model.ckpt_--dim_" +
1359         str(WINDOW) + "_" + str(AMINO_ACID_LEN) + "_1").split()
1360
1361     test_origin = "https://gitlab.com/perf.ai/pssp_project/-/raw/master/originalData_" + dataset + "/testSet" + str(ds_num) + ".txt"
1362     train_origin = "https://gitlab.com/perf.ai/pssp_project/-/raw/master/originalData_" + dataset + "/trainSet" + str(ds_num) + ".txt"
1363     excluded_proteins = "https://gitlab.com/perf.ai/pssp_project/-/raw/master/originalData_" + dataset + "/excluded_" + dataset + ".
1364         ↪ txt"
1365
1366     train_origin, test_origin, excluded_proteins
1367
1368     import requests
1369     test_f = requests.get(test_origin)
1370     test_f = test_f.text.split('\n')[0:-1]
1371     train_f = requests.get(train_origin)
1372     train_f = train_f.text.split('\n')[0:-1]
1373     excluded_f = requests.get(excluded_proteins)
1374     excluded_f = excluded_f.text.split('\n')[0:-1]
1375
1376     excluded_f
1377
1378     TEST_PRED_FILE="pred_testSet{0}.txt".format(ds_num)
1379     TRAIN_PRED_FILE="pred_trainSet{0}.txt".format(ds_num)
1380     TEST_PRED_FILE
1381
1382     """##Declare Predict Methods##"""
1383
1384     def create_output_pred(pred_test, pred_train):
1385         pred = pred_test.astype(int)
1386         labels = ['C', 'E', 'H']
1387         counter = 0
1388         outFileName = TEST_PRED_FILE
1389         with open(outFileName, 'w') as out_file:
1390             for line in range(0, len(test_f)//3):
1391                 protein_name = test_f[line*3]
1392                 if (protein_name in excluded_f):

```

```

1391         # print(protein_name)
1392         continue
1393     primary_structure = test_f[line*3+1].replace('!', '')
1394     secondary_structure = test_f[line*3+2].replace('!', '')
1395     prediction = ""
1396     for c in secondary_structure:
1397         if (c != '!'):
1398             prediction = prediction + labels[pred[counter]]
1399             counter += 1
1400         # else:
1401         #     prediction = prediction + c
1402     # print("Protein name: " + protein_name)
1403     # print("Actual: " + secondary_structure)
1404     # print("Prediction: " + prediction)
1405     out_file.write(protein_name + "\n")
1406     out_file.write(primary_structure + "\n")
1407     out_file.write(secondary_structure + "\n")
1408     out_file.write(prediction + "\n")
1409     pred = pred_train.astype(int)
1410     counter = 0
1411     outFileNames = TRAIN_PRED_FILE
1412     with open(outFileName, 'w') as out_file:
1413         for line in range(0, len(train_f)//3):
1414             protein_name = train_f[line*3]
1415             if (protein_name in excluded_f):
1416                 # print(protein_name)
1417                 continue
1418             primary_structure = train_f[line*3+1].replace('!', '')
1419             secondary_structure = train_f[line*3+2].replace('!', '')
1420             prediction = ""
1421             for c in secondary_structure:
1422                 if (c != '!'):
1423                     prediction = prediction + labels[pred[counter]]
1424                     counter += 1
1425                 # else:
1426                 #     prediction = prediction + c
1427             # print("Protein name: " + protein_name)
1428             # print("Actual: " + secondary_structure)
1429             # print("Prediction: " + prediction)
1430             out_file.write(protein_name + "\n")
1431             out_file.write(primary_structure + "\n")
1432             out_file.write(secondary_structure + "\n")
1433             out_file.write(prediction + "\n")
1434
1435 # import tensorflow as tf
1436 # tf.compat.v1.disable_eager_execution()
1437 # from utilities import predict, read_data, normalize_and_reshape
1438 # from net.net import CNN
1439 # import numpy as np
1440 # import argparse
1441 # import pdb
1442
1443 def parse_args():
1444     parser = argparse.ArgumentParser(description='prediction')
1445     parser.add_argument('--test_set', dest='test_set',
1446                        help='provide the directory of .mat file for testing',
1447                        default=None, type=str)

```

```

1448 parser.add_argument('--train_set', dest='train_set',
1449                     help='provide_the_directory_of_.mat_file_for_training',
1450                     default=None, type=str)
1451 parser.add_argument('--model', dest='model_file',
1452                     help='provide_file_storing_network_parameters,i.e._/dir/model.ckpt',
1453                     default='./saved_model/model.ckpt', type=str)
1454 parser.add_argument('--bsize', dest='bsize',
1455                     help='batch_size',
1456                     default=1024, type=int)
1457 parser.add_argument('--loss', dest='loss',
1458                     help='which_loss_function_to_use:_MSELoss_or_CrossEntropy',
1459                     default='MSELoss', type=str)
1460 parser.add_argument('--dim', dest='dim', nargs='+', help='input_dimension_of_data,'\n
1461                 'shape_must_be:_height_width_num_channels',
1462                 default=[32, 32, 3], type=int)
1463
1464 args = parser.parse_args(args=pred_args)
1465 return args
1466
1467 def predict_model():
1468     args = parse_args()
1469
1470     sess_config = tf.compat.v1.ConfigProto()
1471     sess_config.gpu_options.allow_growth = True
1472
1473     with tf.compat.v1.Session(config=sess_config) as sess:
1474         graph_address = args.model_file + '.meta'
1475         imported_graph = tf.compat.v1.train.import_meta_graph(graph_address)
1476         imported_graph.restore(sess, args.model_file)
1477         mean_param = [v for v in tf.compat.v1.global_variables() if 'mean_tr:0' in v.name][0]
1478         label_enum_var = [v for v in tf.compat.v1.global_variables() if 'label_enum:0' in v.name][0]
1479
1480         sess.run(tf.compat.v1.variables_initializer([mean_param, label_enum_var]))
1481         mean_tr = sess.run(mean_param)
1482         label_enum = sess.run(label_enum_var)
1483
1484         test_batch, num_cls, _ = read_data(args.test_set, dim=args.dim, label_enum=label_enum)
1485         test_batch[0], _ = normalize_and_reshape(test_batch[0], dim=args.dim, mean_tr=mean_tr)
1486
1487         x = tf.compat.v1.get_default_graph().get_tensor_by_name('main_params/input_of_net:0')
1488         y = tf.compat.v1.get_default_graph().get_tensor_by_name('main_params/labels:0')
1489         outputs = tf.compat.v1.get_default_graph().get_tensor_by_name('output_of_net:0')
1490
1491         if args.loss == 'MSELoss':
1492             loss = tf.reduce_sum(input_tensor=tf.pow(outputs-y, 2))
1493         else:
1494             loss = tf.reduce_sum(input_tensor=tf.nn.softmax_cross_entropy_with_logits(logits=outputs, labels=tf.stop_gradient(y)))
1495
1496         network = (x, y, loss, outputs)
1497
1498         avg_loss, avg_acc, results = predict(sess, network, test_batch, args.bsize)
1499
1500         # convert results back to the original labels
1501         inverse_map = dict(zip(np.arange(num_cls), label_enum))
1502         results = np.expand_dims(results, axis=1)
1503         results = np.apply_along_axis(lambda x: inverse_map[x[0]], axis=1, arr=results)
1504

```



```

1505     train_batch, num_cls, _ = read_data(args.train_set, dim=args.dim, label_enum=label_enum)
1506     train_batch[0], _ = normalize_and_reshape(train_batch[0], dim=args.dim, mean_tr=mean_tr)
1507
1508     avg_loss_train, avg_acc_train, results_train = predict(sess, network, train_batch, args.bsize)
1509     # convert results back to the original labels
1510     inverse_map = dict(zip(np.arange(num_cls), label_enum))
1511     results_train = np.expand_dims(results_train, axis=1)
1512     results_train = np.apply_along_axis(lambda x: inverse_map[x[0]], axis=1, arr=results_train)
1513
1514     create_output_pred(results, results_train)
1515
1516     print('In_test_phase, average_loss:_{:.3f} | average_accuracy:_{:.3f}%'.\
1517         format(avg_loss, avg_acc*100))
1518
1519     print('In_train_phase, average_loss:_{:.3f} | average_accuracy:_{:.3f}%'.\
1520         format(avg_loss_train, avg_acc_train*100))
1521
1522     """##Run Predict and Display output##"""
1523
1524     predict_model()
1525
1526     # !cat "$TEST_PRED_FILE"
1527
1528     # !cat "$TRAIN_PRED_FILE"

```

Appendix G

Ensembles Program

This Python program was used to combine the results from multiple trained models using the ensembles method. It was provided by Dionysiou [24].

```
1  from numpy import *
2  import string as string
3  import sys
4
5
6  class Ensembles:
7      def run(filenamees, windowSize, ensemble, outPred, outSOV, outWeka):
8          f = open(outPred, "w")
9          files = open(filenamees, "r").readlines()
10         files = [w.replace('\n', '') for w in files]
11         files = [open(i, "r") for i in files]
12         i = 0
13         LABELS = ['C', 'E', 'H', '!']
14         if ensemble == 1:
15             for rows in zip(*files):
16                 if i == 3:
17                     for j in range(0, len(rows[0].translate(str.maketrans('', '', string.whitespace))), 1):
18                         count = [0, 0, 0, 0]
19                         for k in range(0, len(rows), 1):
20                             if rows[k][j] == 'C':
21                                 count[0] += 1
22                             elif rows[k][j] == 'E':
23                                 count[1] += 1
24                             elif rows[k][j] == 'H':
25                                 count[2] += 1
26                             else:
27                                 count[3] += 1
28                         f.write(LABELS[argmax(count)])
29                         f.write('\n')
30                         i = 0
31                 else:
32                     f.write(rows[0])
33                     i += 1
34             f.close()
35         else:
36             print('ERROR!!!_Invalid_ensemble_option.')
37
```

```

38     # count accuracy
39     f = open(outPred, "r")
40     lines = f.readlines()
41     f.close()
42     count = 0
43     countall = 0
44     for i in range(0, len(lines), 4):
45         for j in range(0, len(lines[i + 2].translate(str.maketrans('', '', string.whitespace))), 1):
46             if lines[i + 2][j] == lines[i + 3][j]:
47                 count += 1
48             countall += 1
49
50     print('Accuracy:_' + str(float(count) / float(countall) * 100) + '%')
51
52     # Confusion Matrix
53     countHH = 0
54     countHE = 0
55     countHC = 0
56     countEH = 0
57     countEE = 0
58     countEC = 0
59     countCH = 0
60     countCE = 0
61     countCC = 0
62     countH = 0
63     countE = 0
64     countC = 0
65     countHp = 0
66     countEp = 0
67     countCp = 0
68     for i in range(0, len(lines), 4):
69         for j in range(0, len(lines[i + 2].translate(str.maketrans('', '', string.whitespace))), 1):
70             if lines[i + 2][j] == 'H' and lines[i + 3][j] == 'H':
71                 countHH += 1
72             elif lines[i + 2][j] == 'H' and lines[i + 3][j] == 'E':
73                 countHE += 1
74             elif lines[i + 2][j] == 'H' and lines[i + 3][j] == 'C':
75                 countHC += 1
76             elif lines[i + 2][j] == 'E' and lines[i + 3][j] == 'H':
77                 countEH += 1
78             elif lines[i + 2][j] == 'E' and lines[i + 3][j] == 'E':
79                 countEE += 1
80             elif lines[i + 2][j] == 'E' and lines[i + 3][j] == 'C':
81                 countEC += 1
82             elif lines[i + 2][j] == 'C' and lines[i + 3][j] == 'H':
83                 countCH += 1
84             elif lines[i + 2][j] == 'C' and lines[i + 3][j] == 'E':
85                 countCE += 1
86             elif lines[i + 2][j] == 'C' and lines[i + 3][j] == 'C':
87                 countCC += 1
88
89         '''if lines[i + 2][j] == 'H':
90             countH += 1
91         elif lines[i + 2][j] == 'E':
92             countE += 1
93         elif lines[i + 2][j] == 'C':
94             countC += 1

```

```

95
96         if lines[i + 3][j] == 'H':
97             countHp += 1
98         elif lines[i + 3][j] == 'E':
99             countEp += 1
100        elif lines[i + 3][j] == 'C':
101            countCp += 1'''
102
103    print('\n\t\tCONFUSION_MATRIX\n')
104    print('{0:10}{1:10}{2:10}{3:10}'.format('_', 'H', 'E', 'C'))
105    print('{0:1}{1:10d}{2:10d}{3:10d}'.format('H', countHH, countHE, countHC))
106    print('{0:1}{1:10d}{2:10d}{3:10d}'.format('E', countEH, countEE, countEC))
107    print('{0:1}{1:10d}{2:10d}{3:10d}'.format('C', countCH, countCE, countCC))
108
109    # SOV input file
110    # f = open(outPred, "r")
111    f1 = open(outSOV, "w")
112    # lines = f.readlines()
113    # f.close()
114
115    for i in range(0, len(lines), 4):
116        f1.write('>OSEQ\n')
117        f1.write(lines[i + 2])
118        f1.write('>PSEQ\n')
119        f1.write(lines[i + 3])
120        f1.write('>AA\n')
121        f1.write(lines[i + 1])
122    f1.close()
123
124    # weka input file
125    f1 = open(outWeka, "w")
126    f1.write('@RELATION_secondary_structure\n\n')
127    for i in range(0, windowSize * 2 - 1, 1):
128        f1.write('@ATTRIBUTE_aminoacid' + str(i) + '_{C,E,H,0,0}\n')
129    f1.write('@ATTRIBUTE_output_{C,E,H}\n')
130    f1.write('\n@DATA\n')
131
132    leadingzeros = zeros((1, (windowSize - 1)))
133    for i in range(3, len(lines), 4):
134        line = leadingzeros
135        line = append(line, list(lines[i].rstrip()))
136        line = append(line, leadingzeros)
137        for j in range(0, len(lines[i].rstrip()), 1):
138            for k in range(0, windowSize * 2 - 1, 1):
139                f1.write(str(line[j + k]) + ',')
140            f1.write(lines[i - 1].rstrip()[j] + '\n')
141
142    f1.close()
143
144    files = sys.argv[1].replace(',', '')
145    run(files, int(sys.argv[2].replace(',', '')), int(sys.argv[3].replace(',', '')), sys.argv[4].replace(',', ''),
146        sys.argv[5].replace(',', ''), sys.argv[6])
147    # print('\nEnd of ensembles script\n')

```

Appendix H

External Rules Program

This Python program was used to apply the external rules filtering. It was provided by Dionysiou [24].

```
1  import sys
2
3  class externalRules:
4      def applyRules(filename, outSOV, outPred):
5          f = open(filename, "r")
6          lines = f.readlines()
7          f.close()
8          f = open(outSOV, "w")
9          fl = open(outPred, "w")
10
11         for i in range(0, len(lines), 4):
12             fl.write(lines[i])
13             fl.write(lines[i + 1])
14             fl.write(lines[i + 2])
15             f.write(">OSEQ\n")
16             f.write(lines[i + 2])
17             f.write(">PSEQ\n")
18             j = 0
19             lines[i + 3] = list(lines[i + 3].translate({ord(c):'' for c in '_\n\t\r'}))
20             # print(len(lines[i + 3]))
21             while j < len(lines[i + 3]):
22                 if len(lines[i + 3]) - j >= 4:
23                     if lines[i + 3][j] == 'H' and lines[i + 3][j + 1] == 'E' and lines[i + 3][j + 2] == 'E' and \
24                         lines[i + 3][j + 3] == 'H':
25                         lines[i + 3][j] = 'H'
26                         lines[i + 3][j + 1] = 'H'
27                         lines[i + 3][j + 2] = 'H'
28                         lines[i + 3][j + 3] = 'H'
29                         j += 4
30                     continue
31                 if lines[i + 3][j] != 'H' and lines[i + 3][j + 1] == 'H' and lines[i + 3][j + 2] == 'H' and \
32                     lines[i + 3][j + 3] != 'H':
33                     lines[i + 3][j + 1] = 'C'
34                     lines[i + 3][j + 2] = 'C'
35                     j += 4
36                 continue
37             if len(lines[i + 3]) - j >= 3:
```

```

38         if lines[i + 3][j] == 'H' and lines[i + 3][j + 1] == 'E' and lines[i + 3][j + 2] == 'H':
39             lines[i + 3][j + 1] = 'H'
40             j += 3
41             continue
42     j += 1
43
44     if lines[i + 3][0] == 'E' and lines[i + 3][1] != 'E':
45         f.write("C")
46         f1.write("C")
47     elif lines[i + 3][0] == 'H' and lines[i + 3][1] != 'H':
48         f.write("C")
49         f1.write("C")
50     else:
51         f.write(lines[i + 3][0])
52         f1.write(lines[i + 3][0])
53
54     for j in range(1, len(lines[i + 3]) - 1):
55         if lines[i + 3][j - 1] != 'E' and lines[i + 3][j] == 'E' and lines[i + 3][j + 1] != 'E':
56             f.write("C")
57             f1.write("C")
58             continue
59         elif lines[i + 3][j - 1] != 'H' and lines[i + 3][j] == 'H' and lines[i + 3][j + 1] != 'H':
60             f.write("C")
61             f1.write("C")
62             continue
63         f.write(lines[i + 3][j])
64         f1.write(lines[i + 3][j])
65
66     if lines[i + 3][len(lines[i + 3]) - 1] == 'E' and lines[i + 3][len(lines[i + 3]) - 2] != 'E':
67         f.write("C")
68         f1.write("C")
69     elif lines[i + 3][len(lines[i + 3]) - 1] == 'H' and lines[i + 3][len(lines[i + 3]) - 2] != 'H':
70         f.write("C")
71         f1.write("C")
72     else:
73         f.write(lines[i + 3][len(lines[i + 3]) - 1])
74         f1.write(lines[i + 3][len(lines[i + 3]) - 1])
75
76     f.write('\n')
77     f1.write('\n')
78     f.write(">AA\n")
79     f.write(lines[i + 1])
80
81 applyRules(sys.argv[1].replace(',', ''), sys.argv[2].replace(',', ''), sys.argv[3])
82 # print('End of external rules script\n')

```

Appendix I

SOV calculation

To calculate the SOV score the two following C programs were used. Both were provided by Dionysiou [24].

```
1  #include <stdio.h>
2  #include <stdlib.h>
3
4  int main (int argc, char* argv[]){
5      FILE *fp=fopen(argv[1], "r");
6      FILE *out;
7      char *line = NULL;
8      size_t len = 0;
9      ssize_t read;
10     fclose(fopen("resultSOV.txt", "w"));
11
12     if (fp == NULL)
13         exit(0);
14     system("cc_/q3_sov_scripts/sov.c_-o_/q3_sov_scripts/sov_-lm");
15     while ((read = getline(&line, &len, fp)) != -1) {
16         out=fopen("SOVinput.txt", "w");
17         if (out == NULL)
18             exit(0);
19         fprintf(out,"%s", line);
20         getline(&line, &len, fp);
21         fprintf(out,"%s", line);
22         getline(&line, &len, fp);
23         fprintf(out,"%s", line);
24         getline(&line, &len, fp);
25         fprintf(out,"%s", line);
26         getline(&line, &len, fp);
27         fprintf(out,"%s", line);
28         getline(&line, &len, fp);
29         fprintf(out,"%s", line);
30         fclose(out);
31
32         system("/q3_sov_scripts/sov_SOVinput.txt_>>resultSOV.txt");
33     }
34
35     free(line);
36     fclose(fp);
37     return 0;
```

```

1  /*-----
2  /
3  / Program: sov.c
4  /
5  / Secondary structure prediction accuracy evaluation
6  /
7  / SOV (Segment OVerlap) measure
8  /
9  / Copyright by Adam Zemla (11/16/1996)
10 / Email: adamz@llnl.gov
11 /
12 /-----
13 /
14 / Compile: cc sov.c -o sov -lm
15 /
16 /-----*/
17 #include <stdio.h>
18 #include <stdlib.h>
19 #include <string.h>
20 #include <math.h>
21
22 #define MAXRES 5000
23
24 typedef struct {
25     int input;
26     int order;
27     int q3_what;
28     int sov_what;
29     int sov_method;
30     float sov_delta;
31     float sov_delta_s;
32     int sov_out;
33     char fname[100];
34 } parameters;
35
36 char *letter_AA="ARNDCQEGHILKMFPSTWYV-?X"; /* 23 chars */
37
38 void default_parameters(parameters *);
39 int read_aa_osec_psec(char[MAXRES], char[MAXRES], char[MAXRES],
40                      parameters *, char*);
41 float sov(int, char[MAXRES], char[MAXRES], parameters *);
42 float q3(int, char[MAXRES], char[MAXRES], parameters *);
43 int check_aa(char, char*, int);
44
45 int main(int argc, char *argv[])
46 {
47     int i, n_aa, sov_method;
48     char c, aa[MAXRES], osec[MAXRES], psec[MAXRES];
49     parameters pdata;
50     float out0, out1, out2, out3;
51
52     if(argc<2){
53         printf("_Usage: _sov_<input_data>\n");
54         printf("_HELP: _sov_ -h\n");
55         exit(0);
56     }
57     if(!strcmp(argv[1], "-h0", 2) ||

```



```

115             out0*100.0,out1*100.0,out2*100.0,out3*100.0);
116     printf("\n");
117
118     sov_method=pdata.sov_method;
119
120     if(sov_method!=0) pdata.sov_method=1;
121
122     if(pdata.sov_method==1) {
123         pdata.sov_what=0;
124         out0=sov(n_aa,osec,psec,&pdata);
125         pdata.sov_what=1;
126         out1=sov(n_aa,osec,psec,&pdata);
127         pdata.sov_what=2;
128         out2=sov(n_aa,osec,psec,&pdata);
129         pdata.sov_what=3;
130         out3=sov(n_aa,osec,psec,&pdata);
131         printf("\n_SOV_:_____%6.1f____%6.1f____%6.1f____%6.1f",
132             out0*100.0,out1*100.0,out2*100.0,out3*100.0);
133         printf("\n");
134     }
135
136     if(sov_method!=1) pdata.sov_method=0;
137
138     if(pdata.sov_method==0) {
139         pdata.sov_delta=1.0;
140
141         pdata.sov_what=0;
142         out0=sov(n_aa,osec,psec,&pdata);
143         pdata.sov_what=1;
144         out1=sov(n_aa,osec,psec,&pdata);
145         pdata.sov_what=2;
146         out2=sov(n_aa,osec,psec,&pdata);
147         pdata.sov_what=3;
148         out3=sov(n_aa,osec,psec,&pdata);
149         printf("\n_SOV_(1994_JMB._[delta=50])_:_____%6.1f____%6.1f____%6.1f____%6.1f",
150             out0*100.0,out1*100.0,out2*100.0,out3*100.0);
151
152         pdata.sov_delta=0.0;
153
154         pdata.sov_what=0;
155         out0=sov(n_aa,osec,psec,&pdata);
156         pdata.sov_what=1;
157         out1=sov(n_aa,osec,psec,&pdata);
158         pdata.sov_what=2;
159         out2=sov(n_aa,osec,psec,&pdata);
160         pdata.sov_what=3;
161         out3=sov(n_aa,osec,psec,&pdata);
162         printf("\n_SOV_(1994_JMB._[delta=0])_:_____%6.1f____%6.1f____%6.1f____%6.1f",
163             out0*100.0,out1*100.0,out2*100.0,out3*100.0);
164
165         printf("\n");
166     }
167
168     printf("\n_-----\n");
169
170     exit(0);
171 }

```

```

172
173 /*-----
174 /
175 / check_aa - checks an amino acid
176 /
177 /-----*/
178 int check_aa(char token, char* letter, int n)
179 {
180     int i;
181
182     for(i=0;i<n;i++) {
183         if(letter[i]==token)
184             return i;
185     }
186     return n;
187 }
188
189 /*-----
190 /
191 / read_aa_osec_psec - read secondary structure segments file
192 /
193 /-----*/
194 int read_aa_osec_psec(char aa[MAXRES], char sss1[MAXRES],
195                      char sss2[MAXRES], parameters *pdata, char* letter)
196 {
197     int i, j, n_aa, n_aa_1, n_aa_2, n_aa_3, f_seq, alt_c, alt_e, alt_h;
198     float x;
199     char line[MAXRES], keyword[MAXRES], first[MAXRES], second[MAXRES], third[MAXRES], junk[MAXRES];
200     FILE *fp;
201
202     alt_c=0;
203     alt_e=0;
204     alt_h=0;
205
206     if((fp = fopen(pdata->fname, "r"))==NULL) {
207         printf("\n#_error_opening_file_%s_for_read\n\n", pdata->fname);
208         exit(0);
209     }
210
211     f_seq=0;
212     pdata->input=0;
213     n_aa=0;
214     n_aa_1=0;
215     n_aa_2=0;
216     n_aa_3=0;
217
218     while (fgets(line, MAXRES, fp) != NULL) {
219         strcpy(keyword, "    ");
220         strcpy(first, "    ");
221         strcpy(second, "    ");
222         strcpy(third, "    ");
223         strcpy(junk, "    ");
224         i=0;
225         j=0;
226         while(line[i] == ' ' && line[i] != '\n' && line[i] != '\0' && i<MAXRES) i++;
227         if(i<MAXRES) {
228             j=i;

```

```

229     while(line[i] != ' ' && line[i] != '\n' && line[i] != '\0' && i<MAXRES) i++;
230 }
231 j=i-j;
232 if(j<MAXRES && j>0) {
233     sscanf(line,"%s",keyword);
234 }
235 if(!strcmp(keyword,"#",1)) {}
236 else if(!strcmp(keyword,"-----",5)) {}
237 else if(!strcmp(keyword,"NUMBER\0",7)) {}
238 else if(!strcmp(keyword,"SECONDARY\0",10)) {}
239 else if(!strcmp(keyword,"END\0",4) && f_seq==0) {
240     fclose(fp);
241     return n_aa;
242 }
243 else if(!strcmp(keyword,"AA-OSEC-PSEC\0",13)) {
244     printf("%s", line);
245     sscanf(line,"%s_%s",keyword,first);
246     strcpy(pdata->fname,first);
247     pdata->input=1;
248 }
249 else if(line[0] == '\n' || !strcmp(keyword,"_ _ _ _\0",4)) {}
250 else if(!strcmp(keyword,"AA\0",3) && f_seq==0) {
251     sscanf(line,"%s_%s_%s",keyword,first,second);
252     if(!strcmp(keyword,"AA\0",3) &&
253         !strcmp(first,"PSEC\0",5) && !strcmp(second,"OSEC\0",5)) {
254         pdata->order=1;
255     }
256 }
257 else if(!strcmp(keyword,"SOV-DELTA\0",10)) {
258     printf("%s", line);
259     sscanf(line,"%s_%f",keyword,&x);
260     pdata->sov_delta=x;
261 }
262 else if(!strcmp(keyword,"SOV-DELTA-S\0",12)) {
263     printf("%s", line);
264     sscanf(line,"%s_%f",keyword,&x);
265     pdata->sov_delta_s=x;
266 }
267 else if(!strcmp(keyword,"SOV-METHOD\0",9)) {
268     printf("%s", line);
269     sscanf(line,"%s_%d",keyword,&i);
270     pdata->sov_method=i;
271 }
272 else if(!strcmp(keyword,"SOV-OUTPUT\0",9)) {
273     printf("%s", line);
274     sscanf(line,"%s_%d",keyword,&i);
275     pdata->sov_out=i;
276 }
277 else if(line[0]=='>') {
278     printf("%s", line);
279     if(f_seq<2) n_aa=0;
280     f_seq++;
281 }
282 else if(f_seq==0) {
283     if(j>1) {
284         if(!strcmp(keyword,"SSP\0",4)) {
285             sscanf(line,"%s_%s_%s_%s_%s",keyword,junk,first,second,third);

```

```

286     }
287     else {
288         printf("\n_ERROR!_(line:_%d)_Check_COLUMN_format_of_your_prediction!\n",n_aa+1);
289         fclose(fp);
290         exit(0);
291     }
292 }
293 else {
294     sscanf(line,"%s_%s_%s",first,second,third);
295 }
296 aa[n_aa]=first[0];
297 sss1[n_aa]=second[0];
298 sss2[n_aa]=third[0];
299 if(check_aa(aa[n_aa],letter,23)==23) {
300     printf("\n#_ERROR!\n%s",line);
301     printf("\n#_ERROR!_(line:_%d)_Check_amino_acid_code_%%c\n",n_aa+1,aa[n_aa]);
302     fclose(fp);
303     exit(0);
304 }
305 if(sss1[n_aa]=='_' || sss2[n_aa]=='_') {
306     printf("\n#_ERROR!\n%s",line);
307     printf("\n#_ERROR!_(line:_%d)_Check_secondary_structure_code\n",n_aa+1);
308     fclose(fp);
309     exit(0);
310 }
311 if(sss1[n_aa]=='L' || sss1[n_aa]=='T' || sss1[n_aa]=='S') {
312     if(alt_c==0) {
313         printf("#_WARNING!_(line:_%d)_The_%%c_characters_are_interpreted_as_'C'_(coil)\n",n_aa+1,sss1[n_aa]);
314         alt_c=1;
315     }
316     sss1[n_aa]='C';
317 }
318 if(sss1[n_aa]=='B') {
319     if(alt_e==0) {
320         printf("#_WARNING!_(line:_%d)_The_%%c_characters_are_interpreted_as_'E'_(strand)\n",n_aa+1,sss1[n_aa]);
321         alt_e=1;
322     }
323     sss1[n_aa]='E';
324 }
325 if(sss1[n_aa]=='G' || sss1[n_aa]=='I') {
326     if(alt_h==0) {
327         printf("#_WARNING!_(line:_%d)_The_%%c_characters_are_interpreted_as_'H'_(helix)\n",n_aa+1,sss1[n_aa]);
328         alt_h=1;
329     }
330     sss1[n_aa]='H';
331 }
332 if(sss2[n_aa]=='L' || sss2[n_aa]=='T' || sss2[n_aa]=='S') {
333     if(alt_c==0) {
334         printf("#_WARNING!_(line:_%d)_The_%%c_characters_are_interpreted_as_'C'_(coil)\n",n_aa+1,sss2[n_aa]);
335         alt_c=1;
336     }
337     sss2[n_aa]='C';
338 }
339 if(sss2[n_aa]=='B') {
340     if(alt_e==0) {
341         printf("#_WARNING!_(line:_%d)_The_%%c_characters_are_interpreted_as_'E'_(strand)\n",n_aa+1,sss2[n_aa]);
342         alt_e=1;

```

```

343     }
344     sss2[n_aa]='E';
345 }
346 if(sss2[n_aa]=='G' || sss2[n_aa]=='I') {
347     if(alt_h==0) {
348         printf("#_WARNING!_(line:_%d)_The_'%c'_characters_are_interpreted_as_'H'_(helix)\n",n_aa+1,sss2[n_aa]);
349         alt_h=1;
350     }
351     sss2[n_aa]='H';
352 }
353 if(sss1[n_aa]!='C' && sss1[n_aa]!='E' && sss1[n_aa]!='H') {
354     printf("\n#_ERROR!\n%s",line);
355     printf("\n#_ERROR!_(line:_%d)_Check_secondary_structure_code_%c\n",n_aa+1,sss1[n_aa]);
356     fclose(fp);
357     exit(0);
358 }
359 if(sss2[n_aa]!='C' && sss2[n_aa]!='E' && sss2[n_aa]!='H') {
360     printf("\n#_ERROR!\n%s",line);
361     printf("\n#_ERROR!_(line:_%d)_Check_secondary_structure_code_%c\n",n_aa+1,sss2[n_aa]);
362     fclose(fp);
363     exit(0);
364 }
365 n_aa++;
366 if(n_aa>=MAXRES) {
367     printf("\n#_ERROR!_Check_number_of_amino_acid_lines._(MAX=_%d)\n\n",MAXRES);
368     fclose(fp);
369     exit(0);
370 }
371 }
372 else if(f_seq==1) {
373     i=0;
374     while(line[i] != '\n') {
375         if(line[i] != '_' && line[i] != '\t' && line[i] != '\0' &&
376            line[i] != '\a' && line[i] != '\b' && line[i] != '\f' &&
377            line[i] != '\r' && line[i] != '\v' && i<MAXRES) {
378             aa[n_aa]='X';
379             sss1[n_aa]=line[i];
380             if(sss1[n_aa]=='L' || sss1[n_aa]=='T' || sss1[n_aa]=='S') {
381                 if(alt_c==0) {
382                     printf("#_WARNING!_The_'%c'_characters_are_interpreted_as_'C'_(coil)\n",sss1[n_aa]);
383                     alt_c=1;
384                 }
385                 sss1[n_aa]='C';
386             }
387             if(sss1[n_aa]=='B') {
388                 if(alt_e==0) {
389                     printf("#_WARNING!_The_'%c'_characters_are_interpreted_as_'E'_(strand)\n",sss1[n_aa]);
390                     alt_e=1;
391                 }
392                 sss1[n_aa]='E';
393             }
394             if(sss1[n_aa]=='G' || sss1[n_aa]=='I') {
395                 if(alt_h==0) {
396                     printf("#_WARNING!_The_'%c'_characters_are_interpreted_as_'H'_(helix)\n",sss1[n_aa]);
397                     alt_h=1;
398                 }
399                 sss1[n_aa]='H';

```

```

400     }
401     if(sss1[n_aa]!='C' && sss1[n_aa]!='E' && sss1[n_aa]!='H') {
402         printf("\n#_ERROR!\n%s",line);
403         printf("\n#_ERROR!_Check_secondary_structure_code:_%c\n",sss1[n_aa]);
404         fclose(fp);
405         exit(0);
406     }
407     n_aa++;
408     if(n_aa>=MAXRES) {
409         printf("\n#_ERROR!_Check_number_of_residues._(MAX_=%d)\n\n",MAXRES);
410         fclose(fp);
411         exit(0);
412     }
413 }
414 i++;
415 }
416 n_aa_1=n_aa;
417 }
418 else if(f_seq==2) {
419     i=0;
420     while(line[i] != '\n') {
421         if(line[i] != '_' && line[i] != '\t' && line[i] != '\0' &&
422            line[i] != '\a' && line[i] != '\b' && line[i] != '\f' &&
423            line[i] != '\r' && line[i] != '\v' && i<MAXRES) {
424             aa[n_aa]='X';
425             sss2[n_aa]=line[i];
426             if(sss2[n_aa]=='L' || sss2[n_aa]=='T' || sss2[n_aa]=='S') {
427                 if(alt_c==0) {
428                     printf("#_WARNING!_The_'%c'_characters_are_interpreted_as_'C'_(coil)\n",sss2[n_aa]);
429                     alt_c=1;
430                 }
431                 sss2[n_aa]='C';
432             }
433             if(sss2[n_aa]=='B') {
434                 if(alt_e==0) {
435                     printf("#_WARNING!_The_'%c'_characters_are_interpreted_as_'E'_(strand)\n",sss2[n_aa]);
436                     alt_e=1;
437                 }
438                 sss2[n_aa]='E';
439             }
440             if(sss2[n_aa]=='G' || sss2[n_aa]=='I') {
441                 if(alt_h==0) {
442                     printf("#_WARNING!_The_'%c'_characters_are_interpreted_as_'H'_(helix)\n",sss2[n_aa]);
443                     alt_h=1;
444                 }
445                 sss2[n_aa]='H';
446             }
447             if(sss2[n_aa]!='C' && sss2[n_aa]!='E' && sss2[n_aa]!='H') {
448                 printf("\n#_ERROR!\n%s",line);
449                 printf("\n#_ERROR!_Check_secondary_structure_code:_%c\n",sss2[n_aa]);
450                 fclose(fp);
451                 exit(0);
452             }
453             n_aa++;
454             if(n_aa>=MAXRES) {
455                 printf("\n#_ERROR!_Check_number_of_residues._(MAX_=%d)\n\n",MAXRES);
456                 fclose(fp);

```



```

457         exit(0);
458     }
459 }
460     i++;
461 }
462     n_aa_2=n_aa;
463 }
464     else if(f_seq==3) {
465         i=0;
466         while(line[i] != '\n') {
467             if(line[i] != '_' && line[i] != '\t' && line[i] != '\0' &&
468                 line[i] != '\a' && line[i] != '\b' && line[i] != '\f' &&
469                 line[i] != '\r' && line[i] != '\v' && i<MAXRES) {
470                 aa[n_aa_3]=line[i];
471                 if(check_aa(aa[n_aa_3],letter,23)==23) {
472                     printf("\n#_ERROR!\n%s",line);
473                     printf("\n#_ERROR!_(N_res:_%d)_Check_amino_acid_code_%%c\n",n_aa_3+1,aa[n_aa_3]);
474                     fclose(fp);
475                     exit(0);
476                 }
477                 n_aa_3++;
478                 if(n_aa_3>=MAXRES) {
479                     printf("\n#_ERROR!_Check_number_of_residues_(MAX=_%d)\n\n",MAXRES);
480                     fclose(fp);
481                     exit(0);
482                 }
483             }
484             i++;
485         }
486     }
487 }
488     if(n_aa_1!=n_aa_2) {
489         printf("\n#_ERROR!_Check_format_of_your_submission.");
490         printf("\n#_Different_length_of_observed_and_predicted_structures.\n");
491         fclose(fp);
492         exit(0);
493     }
494     return n_aa;
495 }
496
497 /*-----
498 /
499 / default_parameters - default parameters for SOV program
500 /
501 /-----*/
502 void default_parameters(parameters *pdata)
503 {
504     pdata->input=0;
505     pdata->order=0;
506     pdata->sov_method=1; // 0 - SOV definition (1994 JMB.) , 1 - SOV definition (1999 Proteins)
507     pdata->sov_delta=1.0;
508     pdata->sov_delta_s=0.5;
509     pdata->sov_out=0;
510
511     return;
512 }
513

```

```

514  /*-----
515  /
516  / sov - evaluate SSp by the Segment OVerlap quantity (SOV)
517  / Input: secondary structure segments
518  /
519  /-----*/
520  float sov(int n_aa, char sss1[MAXRES], char sss2[MAXRES], parameters *pdata)
521  {
522      int i, k, length1, length2, beg_s1, end_s1, beg_s2, end_s2;
523      int j1, j2, k1, k2, minov, maxov, d, d1, d2, n, multiple;
524      char s1, s2, sse[3];
525      float out;
526      double s, x;
527
528      sse[0]='#';
529      sse[1]='#';
530      sse[2]='#';
531
532      if(pdata->sov_what==0) {
533          sse[0]='H';
534          sse[1]='E';
535          sse[2]='C';
536      }
537      if(pdata->sov_what==1) {
538          sse[0]='H';
539          sse[1]='H';
540          sse[2]='H';
541      }
542      if(pdata->sov_what==2) {
543          sse[0]='E';
544          sse[1]='E';
545          sse[2]='E';
546      }
547      if(pdata->sov_what==3) {
548          sse[0]='C';
549          sse[1]='C';
550          sse[2]='C';
551      }
552      n=0;
553      for(i=0;i<n_aa;i++) {
554          s1=sss1[i];
555          if(s1==sse[0] || s1==sse[1] || s1==sse[2]) {
556              n++;
557          }
558      }
559      out=0.0;
560      s=0.0;
561      length1=0;
562      length2=0;
563      i=0;
564      while(i<n_aa) {
565          beg_s1=i;
566          s1=sss1[i];
567          while(sss1[i]==s1 && i<n_aa) {
568              i++;
569          }
570          end_s1=i-1;

```



```

627         minov,maxov,length1,d,beg_s1+1,end_s1+1,beg_s2+1,end_s2+1);
628     }
629 }
630 }
631 }
632 }
633 if(pdata->sov_out==2) {
634     printf("\n TEST: Number_of_considered_residues_=%d",n);
635 }
636 if(n>0) {
637     out=s/n;
638 }
639 else {
640     out=1.0;
641 }
642 return out;
643 }
644
645 /*-----
646 /
647 / Q3 - evaluate SSp by the residues predicted correctly (Q3)
648 / Input: secondary structure segments
649 /
650 /-----*/
651 float q3(int n_aa, char sss1[MAXRES], char sss2[MAXRES], parameters *pdata)
652 {
653     int i, n;
654     float out;
655     char s, sse[3];
656
657     sse[0]='#';
658     sse[1]='#';
659     sse[2]='#';
660
661     if(pdata->q3_what==0) {
662         sse[0]='H';
663         sse[1]='E';
664         sse[2]='C';
665     }
666     if(pdata->q3_what==1) {
667         sse[0]='H';
668         sse[1]='H';
669         sse[2]='H';
670     }
671     if(pdata->q3_what==2) {
672         sse[0]='E';
673         sse[1]='E';
674         sse[2]='E';
675     }
676     if(pdata->q3_what==3) {
677         sse[0]='C';
678         sse[1]='C';
679         sse[2]='C';
680     }
681
682     n=0;

```

```
683     out=0.0;
684     for(i=0;i<n_aa;i++) {
685         s=sss1[i];
686         if(s==sse[0] || s==sse[1] || s==sse[2]) {
687             n++;
688             if(sss1[i]==sss2[i]) {
689                 out=out + 1.0;
690             }
691         }
692     }
693     if(n>0) {
694         out=out/n;
695     }
696     else {
697         out=1.0;
698     }
699
700     return out;
701 }
```

Appendix J

Calculation of Q3 accuracy

The following Python program was implemented to calculate the Q3 accuracy for each class and the overall Q3 accuracy.

```
1  import sys
2  # Execute: python calc_Q3.py <pred_file>
3  import string
4  lines = None
5  labels = ['H', 'E', 'C']
6  with open(sys.argv[1]) as file:
7      lines = file.readlines()
8  if lines is None: exit(0)
9  countCor = [0, 0, 0]
10 countAll = [0, 0, 0]
11 for l in range(0, len(lines)//4):
12     protein_name = lines[4*l]
13     primary = lines[4*l + 1]
14     secondary = lines[4*l + 2]
15     prediction = lines[4*l + 3]
16     for s, p in zip(secondary, prediction):
17         if s == '\n': continue
18         if s == p:
19             countCor[labels.index(s)] += 1
20             countAll[labels.index(s)] += 1
21 total = countAll[0] + countAll[1] + countAll[2]
22 correct = countCor[0] + countCor[1] + countCor[2]
23 headers = ['Q3_All', 'Q3_C', 'Q3_E', 'Q3_H']
24 q3 = [(100*correct/total),
25        (100*countCor[0]/countAll[0]),
26        (100*countCor[1]/countAll[1]),
27        (100*countCor[2]/countAll[2])]
28 print("\n_{0:11}{1:11}{2:11}{3:11}".format('_Q3_ALL', '_Q3_H', '_Q3_E', '_Q3_C'))
29 print('{0:11.4f}{1:11.4f}{2:11.4f}{3:11.4f}\n'.format(q3[0], q3[1], q3[2], q3[3]))
```

Appendix K

Data pre-processing for filtering

This python program was used to prepare the datasets for the SVM filtering method. The same datasets were used to train the decision trees and random forests. It was provided by Dionysiou [24].

```
1 # Execute: python prepare_SVM_files.py <test_filename> <train_filename> <WINDOW> <out_test> <
  ↳ out_train>
2 import sys
3 #open TEST file to read data
4 with open(sys.argv[1],"r") as testfile:
5     lines_test = testfile.readlines()
6 #open TRAIN file to read dat
7 with open(sys.argv[2],"r") as trainfile:
8     lines_train = trainfile.readlines()
9     linenum = 1
10    window = int(sys.argv[3])
11    leftwindow = int(window/2)
12    #create train file
13    with open(sys.argv[5], "w") as svmtrain:
14        for line in lines_train:
15            if linenum == 5: linenum = 1
16            if linenum == 3:
17                target_out = line
18                # if linenum == 4:
19                for i in range(leftwindow):
20                    zeros = leftwindow - i
21                    for zer in range(zeros):
22                        svmtrain.write("0,")
23                    for rem in range(i):
24                        if line[rem] == "C": svmtrain.write("0,")
25                        if line[rem] == "E": svmtrain.write("1,")
26                        if line[rem] == "H": svmtrain.write("2,")
27                    #place right aminos
28                    for j in range(leftwindow+1):
29                        if line[i+j] == "C": svmtrain.write("0,")
30                        if line[i+j] == "E": svmtrain.write("1,")
31                        if line[i+j] == "H": svmtrain.write("2,")
32                    #place label at the end
33                    if target_out[i] == "C": svmtrain.write("0")
34                    if target_out[i] == "E": svmtrain.write("1")
35                    if target_out[i] == "H": svmtrain.write("2")
```

```

36         svmtrain.write("\n")
37     #place aminos with no boundary constraints
38     for amino in range(leftwindow,len(line)-leftwindow-1):
39         for curr in range(-leftwindow,leftwindow+1):
40             if line[amino+curr] == "C": svmtrain.write("0,")
41             if line[amino+curr] == "E": svmtrain.write("1,")
42             if line[amino+curr] == "H": svmtrain.write("2,")
43         #place label
44         if target_out[amino] == "C": svmtrain.write("0")
45         if target_out[amino] == "E": svmtrain.write("1")
46         if target_out[amino] == "H": svmtrain.write("2")
47         svmtrain.write("\n")
48     #place last aminos with padding
49     for i in range(len(line)-leftwindow-1,len(line)-1):
50         printed=0
51         for left in range(i-leftwindow-1,i):
52             if line[left] == "C": svmtrain.write("0,")
53             if line[left] == "E": svmtrain.write("1,")
54             if line[left] == "H": svmtrain.write("2,")
55         for j in range(i,len(line)-1):
56             if line[j] == "C": svmtrain.write("0,")
57             if line[j] == "E": svmtrain.write("1,")
58             if line[j] == "H": svmtrain.write("2,")
59         printed=printed+1
60         zeros = leftwindow-printed
61         for z in range(zeros):
62             svmtrain.write("0,")
63         # place label
64         if target_out[i] == "C": svmtrain.write("0")
65         if target_out[i] == "E": svmtrain.write("1")
66         if target_out[i] == "H": svmtrain.write("2")
67         svmtrain.write("\n")
68     linenum += 1
69     svmtrain.flush()
70     linenum=1
71     #create TEST file
72     with open(sys.argv[4], "w") as svmtest:
73         for line in lines_test:
74             if linenum == 5: linenum = 1
75             if linenum == 3: target_out = line
76             if linenum == 4:
77                 for i in range(leftwindow):
78                     zeros = leftwindow - i
79                     for zer in range(zeros):
80                         svmtest.write("0,")
81                 for rem in range(i):
82                     if line[rem] == "C": svmtest.write("0,")
83                     if line[rem] == "E": svmtest.write("1,")
84                     if line[rem] == "H": svmtest.write("2,")
85             #place right aminos
86             for j in range(leftwindow+1):
87                 if line[i+j] == "C": svmtest.write("0,")
88                 if line[i+j] == "E": svmtest.write("1,")
89                 if line[i+j] == "H": svmtest.write("2,")
90             #place label at the end
91             if target_out[i] == "C": svmtest.write("0")
92             if target_out[i] == "E": svmtest.write("1")

```



```

93         if target_out[i] == "H": svmtest.write("2")
94         svmtest.write("\n")
95     #place aminos with no boundary constraints
96     for amino in range(leftwindow,len(line)-leftwindow-1):
97         for curr in range(-leftwindow,leftwindow+1):
98             if line[amino+curr] == "C": svmtest.write("0,")
99             if line[amino+curr] == "E": svmtest.write("1,")
100            if line[amino+curr] == "H": svmtest.write("2,")
101        #place label
102        if target_out[amino] == "C": svmtest.write("0")
103        if target_out[amino] == "E": svmtest.write("1")
104        if target_out[amino] == "H": svmtest.write("2")
105        svmtest.write("\n")
106    #place last aminos with padding
107    for i in range(len(line)-leftwindow-1,len(line)-1):
108        printed=0
109        for left in range(i-leftwindow-1,i):
110            if line[left] == "C": svmtest.write("0,")
111            if line[left] == "E": svmtest.write("1,")
112            if line[left] == "H": svmtest.write("2,")
113        for j in range(i,len(line)-1):
114            if line[j] == "C": svmtest.write("0,")
115            if line[j] == "E": svmtest.write("1,")
116            if line[j] == "H": svmtest.write("2,")
117        printed+=1
118        zeros = leftwindow-printed
119        for z in range(zeros):
120            svmtest.write("0,")
121        # place label
122        if target_out[i] == "C": svmtest.write("0")
123        if target_out[i] == "E": svmtest.write("1")
124        if target_out[i] == "H": svmtest.write("2")
125        svmtest.write("\n")
126        linenum += 1
127    svmtest.flush()

```

Appendix L

Training Filtering Methods

The following program was implemented to train the filtering models and apply the filtering techniques on the output data of the Convolutional Neural Network.

```
1  # Execute: python train_SVM.py <test_filename> <train_filename> <WINDOW> <pred_file> <  
    ↳ out_prediction> <out_sov> <filter_opt>  
2  from __future__ import print_function  
3  import sys  
4  import string  
5  import numpy as np  
6  import numpy as np  
7  from sklearn.metrics import classification_report  
8  from sklearn.svm import SVC  
9  from sklearn import svm, pipeline  
10 from sklearn import linear_model  
11 from sklearn.metrics import confusion_matrix  
12 from sklearn.tree import DecisionTreeClassifier  
13 from sklearn.ensemble import RandomForestClassifier  
14  
15 def get_balanced_data(X_train, y_train):  
16     classH = []  
17     classE = []  
18     classC = []  
19     for i,label in enumerate(y_train):  
20         if label == 0:  
21             classH.append(i)  
22         elif label == 1:  
23             classE.append(i)  
24         else:  
25             classC.append(i)  
26     rows = min(len(classH), len(classE), len(classC))  
27  
28     # Create a balanced data set  
29     X_balanced = np.concatenate((X_train[classH][0:rows], X_train[classE][0:rows], X_train[classC][0:rows]), axis=0)  
30  
31     y_balanced = np.concatenate((y_train[classH][0:rows], y_train[classE][0:rows], y_train[classC][0:rows]), axis=0)  
32  
33     balanced = np.zeros((X_balanced.shape[0], X_balanced.shape[1]+1), dtype=int)  
34  
35     balanced[:, -1] = y_balanced  
36     balanced[:, :-1] = X_balanced
```

```

37     np.random.shuffle(balanced)
38     return balanced[:, :-1], balanced[:, -1]
39
40
41 def create_output_pred(pred, input_f, out_f, outSOV):
42     with open(input_f, "r") as pred_file:
43         pred_lines = pred_file.readlines()
44     pred = pred.astype(int)
45     labels = ['C', 'E', 'H']
46     counter = 0
47     with open(out_f, 'w') as out_file:
48         for line in range(0, len(pred_lines)//4):
49             protein_name = pred_lines[line*4][0:-1]
50             primary_structure = pred_lines[line*4+1][0:-1]
51             secondary_structure = pred_lines[line*4+2][0:-1]
52             prediction = ""
53             for c in secondary_structure:
54                 prediction = prediction + labels[pred[counter]]
55                 counter += 1
56             out_file.write(protein_name + "\n")
57             out_file.write(primary_structure + "\n")
58             out_file.write(secondary_structure + "\n")
59             out_file.write(prediction + "\n")
60
61     with open(out_f, "r") as out_file:
62         lines = out_file.readlines()
63     with open(outSOV, "w") as f1:
64         for i in range(0, len(lines), 4):
65             f1.write('>OSEQ\n')
66             f1.write(lines[i + 2])
67             f1.write('>PSEQ\n')
68             f1.write(lines[i + 3])
69             f1.write('>AA\n')
70             f1.write(lines[i + 1])
71
72     train_dataset = np.loadtxt(sys.argv[2], delimiter=",")
73     win=int(sys.argv[3])
74     X_train = train_dataset[:, 0:win]
75     y_train = train_dataset[:, [win]]
76     test_dataset = np.loadtxt(sys.argv[1], delimiter=",")
77     X_test = test_dataset[:, 0:win]
78     y_test = test_dataset[:, [win]]
79     y_train = np.reshape(y_train, len(y_train))
80     y_test = np.reshape(y_test, len(y_test))
81     X_train, y_train = get_balanced_data(X_train, y_train)
82
83     print("Training_...")
84
85     if (sys.argv[7] == '1'):
86         clf = SVC(C=10, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
87                 decision_function_shape='ovr', degree=3, gamma=0.1, kernel='rbf',
88                 max_iter=-1, probability=False, random_state=None, shrinking=True,
89                 tol=0.001, verbose=False)
90     elif (sys.argv[7] == '2'):
91         clf = DecisionTreeClassifier(max_depth=20)
92     elif (sys.argv[7] == '3'):
93         clf = RandomForestClassifier(max_depth=25, random_state=42)

```

```

94 elif (sys.argv[7] == '0'):
95     kernels = ['Polynomial', 'RBF', 'Sigmoid', 'Linear']
96     #A function which returns the corresponding SVC model
97     def getClassifier(ktype):
98         if ktype == 0:
99             # Polynomial kernel
100             return SVC(kernel='poly', degree=8, gamma="auto")
101         elif ktype == 1:
102             # Radial Basis Function kernel
103             return SVC(kernel='rbf', gamma="auto")
104         elif ktype == 2:
105             # Sigmoid kernel
106             return SVC(kernel='sigmoid', gamma="auto")
107         elif ktype == 3:
108             # Linear kernel
109             return SVC(kernel='linear', gamma="auto")
110
111     for i in range(1, 4):
112         # Train a SVC model using different kernels
113         svcclassifier = getClassifier(i)
114         svcclassifier.fit(X_train, y_train)
115         # Make prediction
116         y_pred = svcclassifier.predict(X_test)
117         # Evaluate model
118         print("Evaluation:", kernels[i], "kernel")
119         print(classification_report(y_test, y_pred))
120
121     from sklearn.model_selection import GridSearchCV
122     param_grid = {'C': [0.1, 1, 10], 'gamma': [1, 0.1, 0.01, 0.001], 'kernel': ['rbf']}
123     grid = GridSearchCV(SVC(), param_grid, refit=True, verbose=2)
124     grid.fit(X_train, y_train)
125     print(grid.best_estimator_)
126
127     y_pred = grid.predict(X_test)
128     print(confusion_matrix(y_test, y_pred))
129     print(classification_report(y_test, y_pred))
130     exit(0)
131 else:
132     print('Error!_train_SVM.py_currently_has_no_such_filtering_option.')
133     print('Please_try_again_(available_options:_0-3)')
134     exit(0)
135
136 # Predict the response for test dataset
137 clf.fit(X_train, y_train)
138 y_pred = clf.predict(X_test)
139
140 print("THE_SCORE:_", clf.score(X_test, y_test))
141 print("")
142
143 # creating a confusion matrix
144 cm = confusion_matrix(y_test, y_pred)
145 print('Confusion_Matrix')
146 print(cm)
147 print("")
148
149 create_output_pred(y_pred, sys.argv[4], sys.argv[5], sys.argv[6])

```

All filtering methods on CB513

[illegible]

```

38      cat << 'EOF'
39      _____
40      o O O | _ | / _ \ | | | \ / \
41      o | _ | ( _ ) | | | _ | | ) | _ |
42      TS_[O] _ | _ | \ _ / | _ | | _ / | _ | \ _ /
43      {=====| _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ |
44      /o--000' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0'
45      EOF
46      ::;
47      fold1)
48      cat << "EOF"
49      _____
50      o O O | _ | / _ \ | | | \ / \ | _ |
51      o | _ | ( _ ) | | | _ | | ) | _ | | _ |
52      TS_[O] _ | _ | \ _ / | _ | | _ / | _ | \ _ /
53      {=====| _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ |
54      /o--000' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0'
55      EOF
56      ::;
57      fold2)
58      cat << "EOF"
59      _____
60      o O O | _ | / _ \ | | | \ / \ _ |
61      o | _ | ( _ ) | | | _ | | ) | _ | / /
62      TS_[O] _ | _ | \ _ / | _ | | _ / | _ | / _ /
63      {=====| _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ |
64      /o--000' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0'
65      EOF
66      ::;
67      fold3)
68      cat << "EOF"
69      _____
70      o O O | _ | / _ \ | | | \ / \ | _ |
71      o | _ | ( _ ) | | | _ | | ) | _ | | _ | \
72      TS_[O] _ | _ | \ _ / | _ | | _ / | _ | \ _ /
73      {=====| _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ |
74      /o--000' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0'
75      EOF
76      ::;
77      fold4)
78      cat << "EOF"
79      _____
80      o O O | _ | / _ \ | | | \ / \ | | |
81      o | _ | ( _ ) | | | _ | | ) | _ | _ | _ |
82      TS_[O] _ | _ | \ _ / | _ | | _ / | _ | _ | _ |
83      {=====| _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ |
84      /o--000' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0'
85      EOF
86      ::;
87      fold5)
88      cat << "EOF"
89      _____
90      o O O | _ | / _ \ | | | \ / \ | _ |
91      o | _ | ( _ ) | | | _ | | ) | _ | | _ | \
92      TS_[O] _ | _ | \ _ / | _ | | _ / | _ | \ _ /
93      {=====| _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ | _ | "" _ |
94      /o--000' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0' _ | _ | --0--0'

```



```

150     filter_name="SVM"
151     ;;
152 "2")
153     filter_name="Decision_Tree"
154     ;;
155 "3")
156     filter_name="Random_Forest"
157     ;;
158 *)
159     filter_name="Unknown_Filter"
160     ;;
161 esac
162 }
163
164 get_filter_abr(){
165     case $1 in
166         "1")
167             filter_abr="svm"
168             ;;
169         "2")
170             filter_abr="dtree"
171             ;;
172         "3")
173             filter_abr="rforest"
174             ;;
175         *)
176             filter_abr="unknown"
177             ;;
178     esac
179 }
180
181 SCRIPTS="./q3_sov_scripts"
182 TEMP_FOLDER="./temp_runAll_CB513"
183 RUN_ALL_FOLDER="./CB513_runAll_out_files"
184 CROSS_VAL_FOLDER="./CB513_cross_validation"
185 [ -d "$TEMP_FOLDER" ] || mkdir "$TEMP_FOLDER"
186 [ -d "$RUN_ALL_FOLDER" ] || mkdir "$RUN_ALL_FOLDER"
187
188 echo "=====
189 echo " _>Cross_Validation_Results"
190 echo "
191     ↪ -----
192     ↪ "
191 for i in `ls "$CROSS_VAL_FOLDER"`
192 do
193     echo "$i"
194     new_folder="$RUN_ALL_FOLDER/cross_val_res"
195     [ -d "$new_folder" ] || mkdir "$new_folder"
196     out_file=("$TEMP_FOLDER/$i""_cross_val.txt")
197     for j in `ls "$CROSS_VAL_FOLDER/$i"`
198     do
199         echo "$CROSS_VAL_FOLDER/$i/$j"
200         done > "$out_file"
201         python "$SCRIPTS/ensembles.py" "$out_file" "$WINDOW" 1 "$new_folder/ens_pred.txt" "$new_folder/ens_sov.txt" "
202             ↪ $new_folder/ens_weka.txt"
202         "$SCRIPTS/runSOV" "$new_folder/ens_sov.txt"
203         print_SOV_score

```



```

204 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_pred.txt"
205 echo "
    ↪ -----
    ↪ "
206 done
207 echo "=====
208 echo ""
209 for i in `ls "$TEST_FOLDER"`
210 do
211     print_fold $i
212     new_folder="$RUN_ALL_FOLDER/$i""_results"
213     [ -d "$new_folder" ] || mkdir "$new_folder"
214     out_file=("$TEMP_FOLDER/$i""_files.txt")
215
216     for j in `ls "$TEST_FOLDER/$i"`
217     do
218         echo "$TEST_FOLDER/$i/$j"
219     done > "$out_file"
220     echo "=====
221     echo " _>Ensembles_Results"
222     echo "
    ↪ -----
    ↪ "
223 python "$SCRIPTS/ensembles.py" "$out_file" "$WINDOW" 1 "$new_folder/ensembles_pred.txt" "$new_folder/ensembles_sov.
    ↪ txt" "$new_folder/ensembles_weka.txt" > "$new_folder/ensembles_out.txt"
224 "$SCRIPTS/runSOV" "$new_folder/ensembles_sov.txt"
225 print_SOV_score
226 python "$SCRIPTS/calc_Q3.py" "$new_folder/ensembles_pred.txt"
227 echo "=====
228 echo " _>Ensembles_+_External_Rules_Results"
229 echo "
    ↪ -----
    ↪ "
230 python "$SCRIPTS/externalRules.py" "$new_folder/ensembles_pred.txt" "$new_folder/ens_rules_sov.txt" "$new_folder/
    ↪ ens_rules_pred.txt"
231 "$SCRIPTS/runSOV" "$new_folder/ens_rules_sov.txt"
232 print_SOV_score
233 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_rules_pred.txt"
234
235 for filter in "${filterOpt[@]}"
236 do
237     get_filter_name $filter
238     get_filter_abr $filter
239     echo "=====
240     echo " _>Ensembles_+_External_Rules_+_${filter_name}_Results"
241     echo "
    ↪ -----
    ↪ "
242 python "$SCRIPTS/prepare_SVM_files.py" "$new_folder/ens_rules_pred.txt" "$TRAIN_FOLDER/$i""_train_pred.txt" "
    ↪ $SVM_WIN" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt"
243 python "$SCRIPTS/train_SVM.py" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt" "$SVM_WIN" "
    ↪ $new_folder/ens_rules_pred.txt" "$new_folder/ens_rules_${filter_abr}""_pred.txt" "$new_folder/ens_rules_${filter_abr}
    ↪ ""_sov.txt" "$filter" > "$new_folder/ens_rules_${filter_abr}""_out.txt"
244 "$SCRIPTS/runSOV" "$new_folder/ens_rules_${filter_abr}""_sov.txt"
245 print_SOV_score
246 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_rules_${filter_abr}""_pred.txt"
247 echo "=====

```

```

248     echo "_>Ensembles_+_filter_name_Results"
249     echo "
        ↪ -----
        ↪ "
250     python "$SCRIPTS/prepare_SVM_files.py" "$new_folder/ensembles_pred.txt" "$TRAIN_FOLDER/$i"_train_pred.txt" "
        ↪ $SVM_WIN" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt"
251     python "$SCRIPTS/train_SVM.py" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt" "$SVM_WIN" "
        ↪ $new_folder/ensembles_pred.txt" "$new_folder/ens_$filter_abr"_pred.txt" "$new_folder/ens_$filter_abr"_sov.txt"
        ↪ "$filter" > "$new_folder/ens_$filter_abr"_out.txt"
252     "$SCRIPTS/runSOV" "$new_folder/ens_$filter_abr"_sov.txt"
253     print_SOV_score
254     python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_$filter_abr"_pred.txt"
255     echo "=====
256     echo "_>Ensembles_+_filter_name_+_External_Rules_Results"
257     echo "
        ↪ -----
        ↪ "
258     python "$SCRIPTS/externalRules.py" "$new_folder/ens_$filter_abr"_pred.txt" "$new_folder/ens_$filter_abr"_rules_sov.txt
        ↪ " "$new_folder/ens_$filter_abr"_rules_pred.txt"
259     "$SCRIPTS/runSOV" "$new_folder/ens_$filter_abr"_rules_sov.txt"
260     print_SOV_score
261     python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_$filter_abr"_rules_pred.txt"
262     done
263     echo "=====
264     echo ""
265     # exit 0
266     done
267
268     # Remove temp files
269     rm -rf "$TEMP_FOLDER"
270     rm resultSOV.txt
271     rm SOVinput.txt

```

Appendix N

View filtering results of CB513

The following bash script was implemented and used to view all the ensembles and filtering results in a table format, for the CB513 dataset.

```
1  #!/bin/bash
2
3  file="/final_results_CB513.txt"
4
5  echo "Ensembles_Results"
6  echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
7  echo "-----"
8  sed -n '/Ensembles_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e 's/^[_t]*/' | awk
    ↪ -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch=2;}else_{printf"%%.2ft%%.2
    ↪ ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,v1,v2,v3,v4;switch=1}}'
9  echo ""
10
11 echo "Ensembles_+_External_Rules_Results"
12 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
13 echo "-----"
14 sed -n '/Ensembles_+_External_Rules_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e
    ↪ 's/^[_t]*/' | awk -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch=2;}else_
    ↪ {printf"%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,v1,v2,v3,v4;switch=1}}'
15 echo ""
16
17 echo "Ensembles_+_External_Rules_+_SVM_Results"
18 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
19 echo "-----"
20 sed -n '/Ensembles_+_External_Rules_+_SVM_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_"
    ↪ " | sed -e 's/^[_t]*/' | awk -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch
    ↪ =2;}else_{printf"%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,v1,v2,v3,v4;switch
    ↪ =1}}'
21 echo ""
22
23 echo "Ensembles_+_SVM_Results"
24 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
25 echo "-----"
26 sed -n '/Ensembles_+_SVM_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e 's/^[_t]
    ↪ */' | awk -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch=2;}else_{printf_
    ↪ "%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,v1,v2,v3,v4;switch=1}}'
27 echo ""
28
```

```

29 echo "Ensembles_+_SVM_+_External_Rules_Results"
30 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
31 echo "-----"
32 sed -n '/Ensembles_+_SVM_+_External_Rules_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_"
    ↪ " | sed -e 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=$4;_switch
    ↪ =2;}_else_{printf_ "%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;_ switch
    ↪ =1}}'
33 echo ""
34
35 echo "Ensembles_+_External_Rules_+_Decision_Tree_Results"
36 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
37 echo "-----"
38 sed -n '/Ensembles_+_External_Rules_+_Decision_Tree_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' |
    ↪ tr -s "_" | sed -e 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=$4;
    ↪ _switch=2;}_else_{printf_ "%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;_
    ↪ switch=1}}'
39 echo ""
40
41 echo "Ensembles_+_Decision_Tree_Results"
42 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
43 echo "-----"
44 sed -n '/Ensembles_+_Decision_Tree_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e
    ↪ 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=$4;_switch=2;}_else_
    ↪ {printf_ "%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;_ switch=1}}'
45 echo ""
46
47 echo "Ensembles_+_Decision_Tree_+_External_Rules_Results"
48 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
49 echo "-----"
50 sed -n '/Ensembles_+_Decision_Tree_+_External_Rules_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' |
    ↪ tr -s "_" | sed -e 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=$4;
    ↪ _switch=2;}_else_{printf_ "%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;_
    ↪ switch=1}}'
51 echo ""
52
53 echo "Ensembles_+_External_Rules_+_Random_Forest_Results"
54 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
55 echo "-----"
56 sed -n '/Ensembles_+_External_Rules_+_Random_Forest_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' |
    ↪ tr -s "_" | sed -e 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=$4;
    ↪ $4;_switch=2;}_else_{printf_ "%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;_
    ↪ _switch=1}}'
57 echo ""
58
59 echo "Ensembles_+_Random_Forest_Results"
60 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
61 echo "-----"
62 sed -n '/Ensembles_+_Random_Forest_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e
    ↪ 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=$4;_switch=2;}_else_
    ↪ {printf_ "%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;_ switch=1}}'
63 echo ""
64
65 echo "Ensembles_+_Random_Forest_+_External_Rules_Results"
66 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
67 echo "-----"
68 sed -n '/Ensembles_+_Random_Forest_+_External_Rules_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]'
    ↪ ' | tr -s "_" | sed -e 's/^[^_]*//' | awk -F'_' 'BEGIN{switch=1}{if_(switch_==_1){v1=$1;_v2=$2;_v3=$3;_v4=

```

```
→ $4;_switch=2;}_else_{printf_ "%2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft%.2ft\n",_ $1,_ $2,_ $3,_ $4,_ v1,_ v2,_ v3,_ v4;  
→ _switch=1}}'
```

Appendix O

All filtering methods on PISCES

This bash script was implemented and used to apply the ensembles and the filtering methods in various orders and display the results for each fold of the PISCES dataset.

```
1  #!/bin/bash
2  # Author : Panayiotis Leontiou
3  # Since : May 2020
4  # Version: 1.0
5  # Bugs : No known bugs
6
7  TEST_FOLDER="/PISCES_test_pred"
8  TRAIN_FOLDER="/PISCES_train_pred"
9  WINDOW="15"
10 SVM_WIN="19"
11 filterOpt=( "2" "3" )
12
13 echo "
14
15
16
17 P::::::::P::::::::SS::::::::SS::::::::SS::::::::P::::::::P
18 P::::::::P_SS::::::::S_SS::::::::SP::::::::P
19 P::::::::P_S::::::::SSSSS::::::::SS::::::::SSSSS::::::::SP::::::::P
20 PP::::::::P_P::::::::P_S::::::::SSSSSSS::::::::S_P::::::::SSSSSSSP::::::::P_P::::::::P
21 P_P::::::::P_P::::::::P_S::::::::S_P::::::::S_P::::::::P_P::::::::P_P::::::::P
22 P_P::::::::P_P::::::::P_S::::::::S_P::::::::S_P::::::::P_P::::::::P_P::::::::P
23 P_P::::::::P_P::::::::P_S::::::::SSSS_P::::::::SSSS_P::::::::P_P::::::::P
24 P_P::::::::PP_SS::::::::SSSS_P::::::::SSSS_P::::::::P_P::::::::PP
25 P_P::::::::P_P::::::::SSSSS::::::::SS_P::::::::SS_P::::::::P_P::::::::P
26 P_P::::::::P_P::::::::SSSSS::::::::S_P::::::::SSSSS::::::::S_P::::::::P
27 P_P::::::::P_P::::::::S_P::::::::S_P::::::::S_P::::::::P
28 P_P::::::::P_P::::::::S_P::::::::S_P::::::::S_P::::::::P
29 PP::::::::PP_P::::::::SSSSSS_P::::::::SSSSSS_P::::::::SPP::::::::PP
30 P::::::::P_P::::::::S_P::::::::SSSSS::::::::SS_P::::::::SS_P::::::::P
31 P::::::::P_P::::::::S_P::::::::SS_P::::::::P
32 P::::::::P_P::::::::SSSSSSSSSSSS_P::::::::SSSSSSSSSSSS_P::::::::P
33
34
35
36 "
37 print_fold () {
```

```

38  case $1 in
39      fold0)
40          cat << 'EOF'
41          _____
42          o O O | _ | / _ \ | | | \ / \
43          o | _ | ( ) | | | _ | | ) | _ _ | ( ) |
44          TS_[O] _ | _ | \ _ / | _ _ | _ | / | _ | _ \ /
45          {=====| _ | "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ |
46          /o--000' _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ |
47          EOF
48          ;;
49          fold1)
50          cat << "EOF"
51          _____
52          o O O | _ | / _ \ | | | \ / \ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ |
53          o | _ | ( ) | | | _ | | ) | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ |
54          TS_[O] _ | _ | \ _ / | _ _ | _ | / | _ | _ \ /
55          _ | {=====| _ | "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ |
56          /o--000' _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ |
57          EOF
58          ;;
59          fold2)
60          cat << "EOF"
61          _____
62          o O O | _ | / _ \ | | | \ / \ _ |
63          o | _ | ( ) | | | _ | | ) | _ _ | / /
64          TS_[O] _ | _ | \ _ / | _ _ | _ | / | _ | _ \ /
65          _ | {=====| _ | "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ |
66          /o--000' _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ |
67          EOF
68          ;;
69          fold3)
70          cat << "EOF"
71          _____
72          o O O | _ | / _ \ | | | \ / \ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ |
73          o | _ | ( ) | | | _ | | ) | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ |
74          TS_[O] _ | _ | \ _ / | _ _ | _ | / | _ | _ \ /
75          _ | {=====| _ | "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ |
76          /o--000' _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ |
77          EOF
78          ;;
79          fold4)
80          cat << "EOF"
81          _____
82          o O O | _ | / _ \ | | | \ | | |
83          o | _ | ( ) | | | _ | | ) | _ _ | _ |
84          TS_[O] _ | _ | \ _ / | _ _ | _ | / | _ | _ \ /
85          _ | {=====| _ | "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ | _ | "" "" _ |
86          /o--000' _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ | _ | --0-0- _ |
87          EOF
88          ;;
89          fold5)
90          cat << "EOF"
91          _____
92          o O O | _ | / _ \ | | | \ / \ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ |
93          o | _ | ( ) | | | _ | | ) | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ | _ _ |
94          TS_[O] _ | _ | \ _ / | _ _ | _ | / | _ | _ \ /

```



```

150  case $1 in
151      "1")
152          filter_name="SVM"
153          ;;
154      "2")
155          filter_name="Decision_Tree"
156          ;;
157      "3")
158          filter_name="Random_Forest"
159          ;;
160      *)
161          filter_name="Unknown_Filter"
162          ;;
163  esac
164  }
165
166  get_filter_abr(){
167      case $1 in
168          "1")
169              filter_abr="svm"
170              ;;
171          "2")
172              filter_abr="dtree"
173              ;;
174          "3")
175              filter_abr="rforest"
176              ;;
177          *)
178              filter_abr="unknown"
179              ;;
180      esac
181  }
182
183  SCRIPTS="./q3_sov_scripts"
184  TEMP_FOLDER="./temp_runAll_PISCES"
185  RUN_ALL_FOLDER="./PISCES_runAll_out_files"
186  CROSS_VAL_FOLDER="./PISCES_cross_validation"
187  [ -d "$TEMP_FOLDER" ] || mkdir "$TEMP_FOLDER"
188  [ -d "$RUN_ALL_FOLDER" ] || mkdir "$RUN_ALL_FOLDER"
189
190  echo "=====
191  echo " _>Cross_Validation_Results"
192  echo "
193      ↪ -----
194      ↪ "
195
196  for i in `ls "$CROSS_VAL_FOLDER"`
197  do
198      echo "$i"
199      new_folder="$RUN_ALL_FOLDER/cross_val_res"
200      [ -d "$new_folder" ] || mkdir "$new_folder"
201      out_file=("$TEMP_FOLDER/$i"_cross_val.txt")
202      for j in `ls "$CROSS_VAL_FOLDER/$i"`
203      do
204          echo "$CROSS_VAL_FOLDER/$i/$j"
205          done > "$out_file"
206      python "$SCRIPTS/ensembles.py" "$out_file" "$WINDOW" 1 "$new_folder/ens_pred.txt" "$new_folder/ens_sov.txt" "
207          ↪ $new_folder/ens_weka.txt"

```

```

204 "$SCRIPTS/runSOV" "$new_folder/ens_sov.txt"
205 print_SOV_score
206 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_pred.txt"
207 echo "
    ↪ -----
    ↪ "

208 done
209 echo "=====
210 echo ""
211 for i in `ls "$TEST_FOLDER"`
212 do
213     print_fold $i
214     new_folder="$RUN_ALL_FOLDER/$i""_results"
215     [ -d "$new_folder" ] || mkdir "$new_folder"
216     out_file=("$TEMP_FOLDER/$i""_files.txt")
217
218     for j in `ls "$TEST_FOLDER/$i"`
219     do
220         echo "$TEST_FOLDER/$i/$j"
221     done > "$out_file"
222     echo "=====
223     echo "└>Ensembles_Results"
224     echo "
    ↪ -----
    ↪ "

225 python "$SCRIPTS/ensembles.py" "$out_file" "$WINDOW" 1 "$new_folder/ensembles_pred.txt" "$new_folder/ensembles_sov.
    ↪ txt" "$new_folder/ensembles_weka.txt" > "$new_folder/ensembles_out.txt"
226 "$SCRIPTS/runSOV" "$new_folder/ensembles_sov.txt"
227 print_SOV_score
228 python "$SCRIPTS/calc_Q3.py" "$new_folder/ensembles_pred.txt"
229 echo "=====
230 echo "└>Ensembles_+_External_Rules_Results"
231 echo "
    ↪ -----
    ↪ "

232 python "$SCRIPTS/externalRules.py" "$new_folder/ensembles_pred.txt" "$new_folder/ens_rules_sov.txt" "$new_folder/
    ↪ ens_rules_pred.txt"
233 "$SCRIPTS/runSOV" "$new_folder/ens_rules_sov.txt"
234 print_SOV_score
235 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_rules_pred.txt"
236 for filter in "${filterOpt[@]}"
237 do
238     get_filter_name $filter
239     get_filter_abr $filter
240     # echo "$filter $filter_name"
241     echo "=====
242     echo "└>Ensembles_+_External_Rules_+_$filter_name_Results"
243     echo "
    ↪ -----
    ↪ "

244 python "$SCRIPTS/prepare_SVM_files.py" "$new_folder/ens_rules_pred.txt" "$TRAIN_FOLDER/$i""_train_pred.txt" "
    ↪ $SVM_WIN" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt"
245 python "$SCRIPTS/train_SVM.py" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt" "$SVM_WIN" "
    ↪ $new_folder/ens_rules_pred.txt" "$new_folder/ens_rules_$filter_abr""_pred.txt" "$new_folder/ens_rules_$filter_abr
    ↪ ""_sov.txt" "$filter" > "$new_folder/ens_rules_$filter_abr""_out.txt"
246 "$SCRIPTS/runSOV" "$new_folder/ens_rules_$filter_abr""_sov.txt"
247 print_SOV_score

```

```

248 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_rules_$filter_abr""_pred.txt"
249 echo "=====
250 echo " _>Ensembles_+_filter_name_Results"
251 echo "
    ↪ -----
    ↪ "
252 python "$SCRIPTS/prepare_SVM_files.py" "$new_folder/ensembles_pred.txt" "$TRAIN_FOLDER/$i""_train_pred.txt" "
    ↪ $SVM_WIN" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt"
253 python "$SCRIPTS/train_SVM.py" "$new_folder/temp_svm_test.txt" "$new_folder/temp_svm_train.txt" "$SVM_WIN" "
    ↪ $new_folder/ensembles_pred.txt" "$new_folder/ens_$filter_abr""_pred.txt" "$new_folder/ens_$filter_abr""_sov.txt"
    ↪ "$filter" > "$new_folder/ens_$filter_abr""_out.txt"
254 "$SCRIPTS/runSOV" "$new_folder/ens_$filter_abr""_sov.txt"
255 print_SOV_score
256 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_$filter_abr""_pred.txt"
257 echo "=====
258 echo " _>Ensembles_+_filter_name_+_External_Rules_Results"
259 echo "
    ↪ -----
    ↪ "
260 python "$SCRIPTS/externalRules.py" "$new_folder/ens_$filter_abr""_pred.txt" "$new_folder/ens_$filter_abr""_rules_sov.txt
    ↪ " "$new_folder/ens_$filter_abr""_rules_pred.txt"
261 "$SCRIPTS/runSOV" "$new_folder/ens_$filter_abr""_rules_sov.txt"
262 print_SOV_score
263 python "$SCRIPTS/calc_Q3.py" "$new_folder/ens_$filter_abr""_rules_pred.txt"
264 done
265 echo "=====
266 echo ""
267 # exit 0
268 done
269
270 # Remove temp files
271 rm -rf "$TEMP_FOLDER"
272 rm resultSOV.txt
273 rm SOVinput.txt

```

Appendix P

View filtering results of PISCES

The following bash script was implemented and used to view all the ensembles and filtering results in a table format, for the PISCES dataset.

```
1  #!/bin/bash
2
3  file="/final_results_PISCES.txt"
4
5  echo "Ensembles_Results"
6  echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
7  echo "-----"
8  sed -n '/Ensembles_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e 's/^[_t]*/' | awk
    ↪ -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch=2;}else_{printf"%%.2ft%%.2
    ↪ ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,_v1,_v2,_v3,_v4;switch=1}}'
9  echo ""
10
11 echo "Ensembles_+_External_Rules_Results"
12 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
13 echo "-----"
14 sed -n '/Ensembles_+_External_Rules_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e
    ↪ 's/^[_t]*/' | awk -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch=2;}else_
    ↪ {printf"%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,_v1,_v2,_v3,_v4;switch=1}}'
15 echo ""
16
17 echo "Ensembles_+_External_Rules_+_Decision_Tree_Results"
18 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
19 echo "-----"
20 sed -n '/Ensembles_+_External_Rules_+_Decision_Tree_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' |
    ↪ tr -s "_" | sed -e 's/^[_t]*/' | awk -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;
    ↪ _switch=2;}else_{printf"%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,_v1,_v2,_v3,_v4;_
    ↪ switch=1}}'
21 echo ""
22
23 echo "Ensembles_+_Decision_Tree_Results"
24 echo -e "Q3_ALL\tQ3_H\tQ3_E\tQ3_C\tSOV_ALL\tSOV_H\tSOV_E\tSOV_C"
25 echo "-----"
26 sed -n '/Ensembles_+_Decision_Tree_Results/./====/p' "$file" | grep -E '[0-9]+' | grep -v '[a-zA-Z]' | tr -s "_" | sed -e
    ↪ 's/^[_t]*/' | awk -F'_' 'BEGIN{switch=1}{if_(switch==_1){v1=$1;v2=$2;v3=$3;v4=$4;switch=2;}else_
    ↪ {printf"%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft%%.2ft\n",_1,_2,_3,_4,_v1,_v2,_v3,_v4;switch=1}}'
27 echo ""
28
```

