

Prediction of Protein Protein interactions using Graph Convolutional Networks

Akshat Sharma
sharma06@ads.uni-passau.de
Universität Passau

Lovesh Bishnoi
bishno01@ads.uni-passau.de
Universität Passau

Amit Manbansh
manban01@ads.uni-passau.de
Universität Passau

Mihir Shah
shah02@ads.uni-passau.de
Universität Passau

ABSTRACT

Protein-Protein Interaction (PPI) is an important subject as it is related to various fields like Microbiology, Medicine etc. The data concerning PPIs is generally in the form of a set of protein pairs which can be represented as an undirected graph and the prediction of an interaction between two proteins computationally is a much effective method which can be further verified experimentally for the various aforementioned domains. As an undirected graph the proteins can be considered as nodes and the interaction between them can be considered as an edge, therefore machine learning techniques such as link prediction can be applied in the prediction of an interaction between two proteins. In this project, Graphical Convolution Network is used for the prediction of a PPI by link prediction problem on undirected yeast *s. cerevisiae* graphical data. The Average Precision was obtained to be about 0.86309 which outperforms various computational methods to predict PPIs using link prediction.

1 INTRODUCTION

Proteins are the building blocks of the cells, interactions between these proteins determine the various cellular functions. Proteins and their interactions amongst themselves are remarkably very complex phenomena and are termed as the Protein-Protein Interactions or PPIs. These Protein-Protein Interactions can be observed via the conventional experimental methods or via computational methods. The usual experimental approaches available to study Protein-Protein Interaction are bit time consuming and as well as expensive. Moreover, large scale experiments always involve high rate of false positives. On the other hand, computational algorithms developed due to recent developments in Machine Learning, can prove to be handy in identifying many undiscovered PPIs, that require high throughput experiments. Furthermore, such computationally discovered PPIs can be verified experimentally and thus it can save time as well. One such technique is Link Prediction [6]. In general, in Link Prediction [6] the links or the edges between two nodes of a graph are predicted, given the attribute information and the observed existing link information. In case of the prediction of the Protein-Protein Interactions, Proteins are considered as the nodes and the interactions between them are considered as the links and using the technique of link prediction, interactions amongst proteins can be predicted. Knowledge of protein interaction can help pharmacists and microbiologists understand the function and

behaviour of protein, to assign a new function. Adding to this, a cluster of proteins with the same function can be generated. Biologists and pharmacists can study protein interaction so that they can characterise protein complexes [3].

1.1 Graph Convolution Network

Graph Convolution Networks (GCNs) [8] are the neural network models that share the filter parameters over all locations in the graph. The goal of these neural network models is to learn a function of features on a graph. Let $G = (V, E)$ be the graph where V is the set of the Nodes and E is the set of the Edges. The inputs to the Graph Convolution Network are a feature description x_i for every node i , where $x_i \in X^{N \times D}$ and X is the feature matrix and N is the number of the nodes and D is the number of the input features, and A is the Adjacency Matrix which represents the description of the graph structure in matrix form. The Graph Convolution Network, i.e., G produces a node level output $Z^{N \times F}$, where F is the output features per node.

Every Graphical Neural Network layer can be represented mathematically as a non linear function [8]

$$H^{(l+1)} = f(H^{(l)}, A). \quad (1)$$

where $H^{(0)} = X$ and $H^{(l)} = Z$ and ' l ' is the number of layers and ' f ' is the function chosen for a specific task.

A simple layer wise Propagation rule [8] is given by

$$f(H^{(l)}, A) = \sigma(AH^{(l)}W^{(l)}) \quad (2)$$

where $W^{(l)}$ is a weight matrix for the l -th neural network and $\sigma(\cdot)$ is a non linear activation function which could be ReLU, softmax, etc depending on the choice. Mathematically, these functions are defined as respectively:

$$ReLU(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}, \quad (3)$$

$$softmax(x) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}. \quad (4)$$

Where ' x ' is a scalar value in case of ReLU activation function and is a vector value in case of softmax activation function respectively.

There are two main limitations of this model and their possible adjustments are:

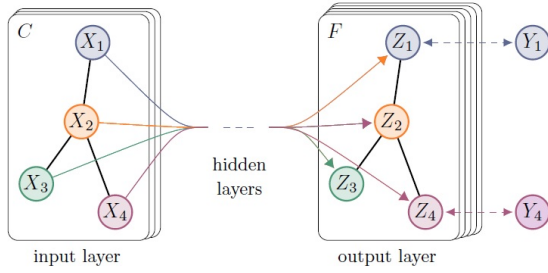


Figure 1: Graph Convolutional Network [8]

- (1) Directly multiplication with A means that for every node all the feature vectors of all the neighbouring nodes are summed but not of that node under consideration. This limitation is overcome by adding the Identity Matrix, i.e., I [8] to the Adjacency matrix, i.e., A , mathematically:

$$\tilde{A} = A + I \quad (5)$$

where \tilde{A} is the Adjacency Matrix with self loop and it is done so that each node also includes its own features at its next representation and it also helps with the numerical stability.

- (2) A is not normalised which can lead to the change of scale of the feature vectors. This limitation is adjusted by normalising A such that all rows sum to one, i.e., $D^{-1}A$, where D is the diagonal node degree matrix of A and $D^{-1}A$ denotes the averaging of the neighbouring node features. Practically averaging of the neighbouring node features is not sufficient therefore symmetric normalisation or spectral approximation of A is done, i.e., $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ [8].

Combining these two adjustments to the above limitations we get a propagation rule [8]:

$$f(H^{(l)}, A) = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}) \quad (6)$$

Where, \tilde{D} is the diagonal node degree matrix of \tilde{A} and it is used to normalise the nodes with large degrees. A multi layer Graph Convolutional Network can be shown with C input channels and F feature maps in the output layer as shown in Figure 1.

In the research paper "Semi-Supervised Classification with Graph Convolution Networks" by Thomas Kipf et al [8], the authors tried to solve the problem of "Citation Networks" by using the the above mentioned techniques of the GCN propagation in a two layer GCN model for semi supervised node classification. Their model was mathematically represented as:

$$Z = f(X, A) = \text{softmax}(\tilde{A} \text{ReLU}(\tilde{A} X W^{(0)}) W^{(1)}) \quad (7)$$

In the citation network the documents were considered nodes and edges were the documents cited in the published documents and on this semi-supervised classification the loss (cross entropy loss) was computed as:

$$L = - \sum_{l \in Y_L} \sum_{f=1}^F Y_{lf} \ln Z_{lf} \quad (8)$$

Where Y_L is the set of node indices with labels and F are the feature maps in the output layer

1.2 Protein-Protein Interaction

A protein is a large organic bio-molecule [2] which is formed from the combination of Amino Acids which combine with each other using a peptide bond [2] to create a protein molecule. The length of the protein depends on its genetic code which specifies the number and types of amino acids responsible in the formation of a particular protein. The genetic expression in any organism is directly related to the proteins associated with that genetic code. When these proteins interact with each other then this interaction is known as Protein-Protein Interaction (PPIs) [7]. Cellular functions are defined by the interaction between the proteins, therefore if we want to understand the basic cell functions, we need to map the protein-protein interactions. Protein interactions have a huge potential within an organism as a whole, e.g., total interactions of human PPIs [4] are estimated to be around 650,000.

There are mainly two methods to detect the protein-protein interactions [4], namely Experimental and Computational. Experimental methods include Y2H, TAP-MS, protein microarrays, mbSUS, pull-down arrays, DPI, etc. Computational methods include PTMs, gene fusion, co-expression, GO annotation, gene neighbourhood, Phylogenetic profile, topological features, sequence features, Domain interactions, protein fold, etc.

Computational techniques consider "protein-protein interactions [7]" or PPIs as the associations or links between proteins. These techniques overcome the limitations of experimental techniques, e.g., Experimental findings are often incomplete even for well-studied organisms, therefore computational methods are used to complete the missing or incomplete part of the experimental PPI data and thus help in finding the clues to map out PPI mechanisms. These methods mainly focus on individual evidence for prediction and have certain specificities and biases. The various evidence sources are integrated in a statistical learning framework, such methods are called "prediction of protein-protein interactions by evidence-combining methods". The machine learning techniques can be applied on such a statistical framework.

The machine learning algorithm on the prediction of protein-protein interactions by evidence-combining methods mainly consists of three steps:

- (1) Defining Gold Standard Datasets [4] (training datasets of interacting and non-interacting protein pairs).
- (2) Characterising the interactions between proteins by annotating the Gold Standard Datasets [4] with carefully chosen and diverse evidence.
- (3) Determining the probability of a particular interaction or interactions by individual evidence and then combining the probabilities of all the evidences.

1.2.1 Gold Standard Datasets. They are created for training or testing the PPI predictions and the datasets for training and testing are separate. These datasets could be either GSP datasets (Gold Standard Positive) or GSN datasets (Gold Standard Negative).

GSPs are PPIs with high experimental confidence or reliability or reference evidence, e.g., BioGRID, IntAct, etc, are some examples which are available in public databases. They are mainly the repositories of protein complexes and interactions are varied in terms of size and species-specificity and they contain information from both the experimental and the computational sources.

GSNs are usually not obtained by direct experimental methods or techniques. Negatome Database (2.0) provides a collection of proteins and domain pairs which are unlikely to engage in direct interaction but it wasn't able to satisfy the diverse GSP datasets of different users. GSNs can be obtained using certain methods, e.g., Negative examples can be chosen from the categories of their particular functions like annotations and subcellular localisation, etc.

1.2.2 Annotations of protein pairs with diverse evidence. Protein pairs are annotated based on their interactions with each other which could be based on cell physiology, biochemical environment, structures of the protein complexes, etc. To detect PPIs experimentally, certain conditions and criteria are met depending on the nature of protein interaction. For prediction of the PPIs by machine learning algorithms, we need to extract the protein interaction-based features. As there are several conditions for different PPIs therefore features are categorised into different categories and each category can provide a different view of protein interactions. Some of these categories and the features contained in them are: Evolutionary relationship (EVO) with features related to genes, e.g., Gene Neighbourhood (GN) etc, Functional Features (FF) and Sequence based code signatures (SEQ), Structure based signatures which are based on the structural interactions of the proteins and the features included in this category are Domain Domain Interaction (DDI), etc. DDI is represented in the Figure 2

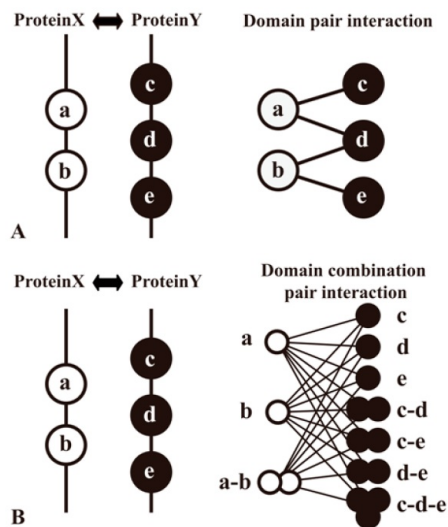


Figure 2: Types of Domain Domain Interactions [4]

Figure 2 shows two methods to predict Domain Domain Interactions from PPIs. Consider two proteins with domains {a, b} and {c, d, e} respectively, here PPIs are interpreted as the interaction amongst the domains of the two proteins.

In method A one domain of a protein will interact with one domain of the other protein.

In method B one or more domains of a protein will interact with one or more domains of the other protein.

1.2.3 Strategy for Integrative Analysis. Classification Algorithm is used to integrate the protein interaction related features, with these available features, classifiers are trained to differentiate between positive and negative examples. The usual process of PPI prediction by evidence-combining techniques [4] includes mainly three steps which are:

- (1) Step I Choose appropriate evidence.
- (2) Step II Encode protein pairs with evidence.
- (3) Step III Find strategy to merge the classifiers into the integrative datasets.

The strategies could be the use of Artificial Neural Networks, Naïve Bayes, Decision Tree, K Nearest Neighbours, Support Vector Machine, et cetera.

1.2.4 Performance evaluation of PPI Prediction. Performance evaluation is done using the ROC score and the Average precision score.

1.2.5 PPI Prediction through Domain Domain Interactions. A domain is mainly one or more submolecular parts of protein, described as a structural and functional module and usually an evolutionarily conserved unit. Recent studies[9] about domains have concluded that abnormalities of domains can cause various diseases. Therefore, studies on the protein domain can help in developing disease models and tools to diagnose them. However, experimental techniques of the domain-domain interaction for predicting PPI have often shown more false-positive and false-negative results. The researchers have started studies on PPI prediction using computational techniques, that are mainly based on different features of protein such as protein sequence, domain information, three-dimensional structure, and protein evolution. The current state of the art approach does not fully consider domain information, instead only works on important domain and domain co-occurrences. However, an overall view of the PPI network can be better understood by domain-domain interaction [9].

1.3 Related work

Our model takes inspiration from the conventional computation approach for PPI prediction and recent development on the graph convolution network for PPI. In what follows, we provide a brief description of related work in both the fields.

1.3.1 Conventional Approach for PPI. In the paper "Prediction of Protein-Protein interaction based on Domain" by Xue Li. et.al [9] proposed a novel approach based on protein domain for Protein-Protein Interaction. In this approach, the authors have used a state-of-the-art SVM model, in which physicochemical properties of domain and domain-domain interaction score are used as the features for the prediction model. The outcome of the SVM model and the domain-domain score were used to design a Protein interaction prediction model. Accuracy, sensitivity, specificity, precision, Matthews correlation coefficient, and the F1 score are used as the evaluation matrix of the prediction model. The drawback of this approach was that it was implemented on a very small scale of the dataset. Adding to this, for the unavailability of negative domain-domain pair data, the general noise was added instead.

1.3.2 Graph Convolution Network for PPI prediction. The spatial Graph convolution approach was used for protein interface prediction by Alex Fout et.al. [5] in the paper “Protein Interface Prediction using Graph Convolutional Networks”. The prediction of the protein interface was based on the graph structure of the protein where amino acid residues are nodes. A set of k residues determined by mean distance between the atoms is used as a neighborhood-based convolution, which is able to detect the node features accurately. The edge features were also derived by their network, but in limited amount and were static compared to the node feature. Thus, the model learned the latent pattern for node features only. For evaluation, the authors compared their approach with State-of-the-art SVM based protein prediction approach. The AUC score for novel approach is 0.89 whereas, AUC score for the SVM model was 0.81. The authors suggest that their approach can be improved by using a large set of protein dataset for better representation for CNN and along with node feature, the model can be improved to learn edge feature representation.

1.4 Problem Definition

In this project, we are trying to solve the problem of the time-consuming experimental process to find out the possibility of PPI for given proteins with a faster computational method. We will be using GCN to design a model which can predict the possibility of a link between two given proteins.

Mathematically, the problem can be summarized as: For a graph $G = (V, E)$, where V is the set of vertices and E is the set of edges, given two nodes $x, y \subseteq V$, where V represents proteins, we will implement a Graphical Neural Network model to predict if $\text{Edge } \{x, y\} \subseteq E \mid \text{Edge } \{x, y\} \notin E$.

1.4.1 Research Questions: We formulate our research questions as follows:

- (1) How does increase in hidden-layers in GCN impact the accuracy in a link prediction task?
- (2) What is the impact of various hyperparameters on the performance of GCN for a link prediction task?

2 DATA ACQUISITION & PRE-PROCESSING

Yeasts are the single-celled eukaryotic fungi microorganisms and they are unique because in the kingdom of Fungi as they are the only single-celled microbes whilst all other members are multicellular organisms. Here we are analysing the protein dataset of a particular yeast of Class: Saccharomycetes, Specie: *S. cerevisiae* and Genus: Saccharomyces.

2.1 Data acquisition

We are using the protein dataset¹ of yeast provided by the Deep Learning for Network Biology² by the Stanford University. The Edgelist contains the pairs of proteins which interact with each other. The yeast Edgelist is a clean and a well organised dataset. To get a general idea about our yeast dataset or the Edgelist we read the Paper by Xiaomei Wu et al [10], in which a protein protein interactions (PPI) map was derived based on Gene Ontology (GO)

annotations. It was done by measuring the similarity of two Gene Ontology terms with a semantic relation known as Relative Specificity Similarity (RSS) and using Z score analysis a positive and a negative datasets for PPI were created. A gold standard positive dataset (GSP) and a gold standard negative dataset (GSN) were also created with high levels of confidence (about 78 percent high quality) and low levels of confidence respectively.

Gene Ontology (GO) is a resource that has relative data sources integrated, which represents the knowledge of the genes in genomes, which give the particular information about certain biological roles attributed to the units associated with those genes, e.g., proteins. GO has been used in protein classification for many species, e.g., *Saccharomyces cerevisiae*, *homo sapiens* etc and has a set of three controlled vocabularies or structural ontologies, e.g., Molecular Function (MF), Biological Process (BP) and Cellular Component (CC).

Functional protein associations can be described by their shared GO terms in a structural ontology or by the semantic similarity of the protein pairs of the terms assigned to them through information or GO. Xiaomei Wu et al [10] worked on to predict a map of the yeast's PPIs by using the BP and CC ontologies or annotations. They followed mainly two particular assumptions which are:

- (1) Two proteins which mostly function in the same biological process are more likely to interact with each other than the two proteins which function in different biological processes.
- (2) For two proteins to interact they must be in a close proximity to each other.

Xiaomei Wu et al [10] used the Organelle DB which is an online resource used for the proteins in eukaryotic organisms localised to the organelles of the cell or in layman's language the subcellular structures. The proteins taken from the Organelle DB were annotated using the BP and CC ontologies from GO consortium by Xiaomei Wu et al [10].

2.2 Data preprocessing

The chosen dataset is present in Edge List³ format. A section of Yeast dataset edge list format file can be seen in Figure 3. Each line in this file represents an edge. The first string in each line represents a node from where the edge is directed. The second string of each line represents the node where the edge is directed to. The nodes here represent complex proteins. An edge between two nodes denotes that both proteins can interact with each other. For example, from line 1 of Figure, we can conclude that protein “YNL236W” can interact with protein “YGL238W”.

2.2.1 Statistics on the data: Networkx [1] library has been used to read the Yeast Edge List format file. Upon reading the Edge List format file into undirected graph⁴ format and converting it to Compressed Sparse Row matrix⁵ format, we find that there are 6526 nodes in the network. Also, there are 1062675 edges in the Compressed Sparse Row matrix⁶.

³<http://snap.stanford.edu/deepnetbio-ismb/ipynb/yeast.edgelist>

⁴https://networkx.github.io/documentation/networkx-1.9.1/_modules/networkx/classes/graph.html

⁵https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html

⁶https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html

¹<http://snap.stanford.edu/deepnetbio-ismb/ipynb/yeast.edgelist>

²<http://snap.stanford.edu/deepnetbio-ismb/>

```

YNL236W YGL238W
YNL236W YOR355W
YNL236W YJL030W
YNL236W YJL013C
YNL236W YJR034W
YNL236W YKL012W
YNL236W YFR033C
YNL236W YGR046W
YNL236W YGR117C
YGL208W YGL115W
YDR328C YLR399C
YDR328C YFL009W
YDR328C YMR094W
YDR328C YJR090C
YDR328C YIL046W
YDR328C YDR139C
YDR328C YOR057W

```

Figure 3: Edgelist

2.2.2 Handling missing data: In Figure 4, we can see node “-” is 11th most centralized node. Since “-” does not represent any protein in the real world, this missing data needs to be removed from the adjacency matrix. The node “-” is connected to 1291 proteins and it is removed from the adjacency matrix resulting in deletion of $2582(1291 * 2)$ edges from the adjacency matrix. The resulting remaining nodes are 6525 and remaining edges are 1060093.

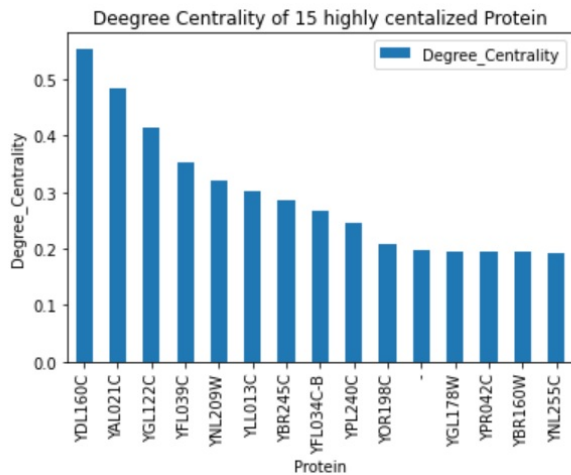


Figure 4: Degree Centrality of 15 highly centralized Proteins

2.2.3 Handling self-loop: Upon checking further we find that the matrix contains weight 1 inside 1685 elements which are present in diagonal of the matrix. Upon deletion of 1685 edges present at diagonal, there are 1058408 edges remaining in the graph. Since, this represents that A protein is can interact with itself, it doesn't carry any meaning of any significance in this project. We removed such elements from the diagonal of the matrix.

2.2.4 Handling directed and symmetrical behaviour in Adjacency matrix: Since Yeast PPI is an undirected activity, we get

rid of symmetrical data of Adjacency matrix by considering only the Upper Triangle of the matrix. We noticed that element[A, B] and element[B, A] contains the same value in the matrix. It represents that Protein A can interact with Protein B and Protein B can interact with Protein A. For example, the element at index [0,1801] and [1801, 0] has the value 1. Index 0 represents YAL008W protein and Index 1801 represents YNR020C protein. Both the cell denotes the same data and hence, we will be considering only the Upper Triangle matrix to extract interaction edges data for the training, validation and testing purpose.

After considering upper Triangle elements of the matrix, we are left with 529204 edges from the total 1060990.

2.2.5 Degree Centrality. After understanding some basic structure measures of the whole network, a good next step is identify the most important nodes in the network. In graph or network, number of edge connection of a particular node to other nodes in the network is measured by degree. Degree of a network or graph can tell about the biggest hubs, while it can not give more understanding for nodes. Adding to this, in network analysis, identifying the most important nodes in the network is measured by centrality. Moreover, the dictionaries that gives nodes as keys and centrality measures as values is called degree centrality. Thus, we can identify with numeric value, the importance of particular node in the network. Degree Centrality of some of the proteins in the dataset are shown in the Figure 4.

2.2.6 Eigenvector Centrality. After getting through the basic visualization of the network, a good next step is to identify which nodes are most important in the network, which nodes form communities/groups within the network.

In graph network, analysis of most important in the network is referred to as the degree of centrality. One of the metrics to find the most important node is a degree, which is considered as the most simple metric. The degree of a node is the number of edges extend from it. For example, a degree node with degree 3 is defined as 3 edges extending from it to other nodes in the network. Another metric for centrality measurement is eigenvector centrality, which is an extension of degree centrality. Eigenvector centrality not only looks at the edges of a particular node but it also takes into account the edges of its neighboring nodes. Centrality in PIP interaction can prove to be handy to identify which protein interacts to the most with other proteins in the network. Then it can later be helpful to study the other features of proteins interacting. In yeast dataset, Protein YDL160C has a degree centrality of 0.5546360153256705, and eigenvector centrality of 0.10615358103813279 which is highest in the network as shown in the Figure 5

2.3 Feature Engineering

2.3.1 Training/Validation/Testing set size: We aim to check the performance of the designed model on different size of training data. We plan to vary the training set size from 60% to 96% of the total data available.

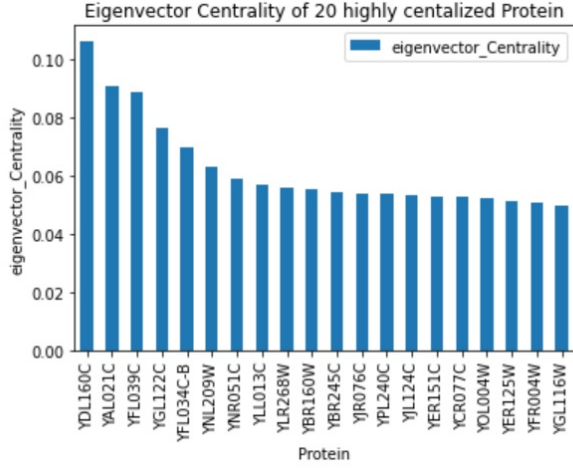


Figure 5: Eigenvector Centrality of 20 highly centralized proteins

2.3.2 Method of construction of Training/Validation/Testing set:

Suppose we are aiming to check the performance of the designed model on $x\%$ of all available data as the training set. Therefore, We pick $(100-x)/2\%$ of total edges randomly each for the validation as well as testing purpose respectively. The remaining edges are used for constructing the training edge list. We also created the same number of False edge list each for validation as well as testing purpose respectively.

For choosing Element[A, B] is eligible for being a part of the Validation false edge list, we checked if Element[A, B] satisfies 5 checks. If all the criteria are satisfied, we added the Element[A, B] to the validation false edge list. The rules are as follows:

- (1) Check whether Node A and Node B are not the same.
- (2) Check if Element[A, B] is not a member of the training edge list.
- (3) Checked if Element[B, A] is not a member of the training edge list.
- (4) Check if Element[A, B] is not a member of the validation edge list.
- (5) Checked if Element[B, A] is not a member of the validation edge list.

For choosing Element[A, B] is eligible for being a part of the testing false edge list, the checks were a little different from the checks for choosing validation false edge list. We just checked if Element [A, B] is a member of all the edges present in Upper Triangle of the Adjacency matrix or a diagonal element. If not, Element [A, B] was added to testing false edge list. The reason is to sample only those edges for the testing set which are unseen to the model.

Feature transformation of training adjacency matrix: We also normalized the Adjacency matrix created from the training edge list. We converted it into symmetric normalization, i.e., $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ as suggested in [8]. We can see the count of different edge list in the following Table 3

Table 1: Various counts from the Yeast Edgelist.

Total Nodes	6526
Total Edges	1062675
Total Nodes after removing invalid node “-”	6525
Total Edges after removing invalid node “-”	1060093
Total Edges after removing diagonal elements	1058408
Total edges in Upper Triangle of Adjacency Matrix	529204

3 MODEL IMPLEMENTATION

3.1 Methodology

3.1.1 Architecture. We have implemented a neural network which contains three hidden layers, out of which first two layers are GCN layers and the third is the Inner Product Decoder layer. The number of units in the hidden layer 1 are 32 and the number of units in the hidden layer 2 are 16. These layers are individually explained as follows:

- (1) GCN Layer 1: In this layer we are using the normalised training Adjacency matrix and the Node Feature Matrix as inputs. The normalised training adjacency matrix is the matrix representation of the graphical structure. Here, we are using the Identity Matrix as Node Feature Matrix because we do not have the node features, we take this inspiration from work done on Karate Club Network problem by Thomas Kipf et. al. in their paper Semi-Supervised Classification with Graph Convolutional Networks [8]. The normalised training Adjacency matrix is of the size $[6526 \times 6526]$ and the size of the Identity Matrix would be the same as that of the normalised training Adjacency Matrix. This GCN Layer converts its inputs into a matrix of the size $[6526 \times 32]$ as an output. The activation function used in this layer is ReLU as shown in Equation (3) and the output of this layer is then passed through this activation function and a final output is calculated which represents the low dimension node-level output of each node present in the graph.
- (2) GCN Layer 2: The functionality of this layer is same as that of the previous layer, i.e., to further reduce the dimension of the node level representation. In this layer we are using the output of the first GCN Layer and the Node feature matrix which is Identity matrix in our case as inputs. These matrices are of the sizes $[6526 \times 32]$ and $[6526 \times 6526]$ respectively. This GCN Layer converts its input node level representation into a matrix of the size $[6526 \times 16]$.
- (3) Inner Product Decoder Layer: The last layer aims to use the lower dimension node-level representation produced during the second hidden GCN layer and use that representation on the training set input data. The training set input data are positive and negative edges set selected for training purpose. Each sample has two nodes between which the link either exists or not. We select the lower dimension node-level representation of these two nodes from each sample of the training set and multiply them to get a scalar value. The value is then passed through a sigmoid activation function

in order to get the output in $[0,1]$. Here, 0 represents the prediction that there is no link present between the two nodes and 1 represents vice-versa.

3.1.2 Training. We train our model by giving the inputs to the first layer of our model as normalised training Adjacency matrix in the form of Sparse matrix and the Identity matrix as the node feature matrix as explained in 3.1.1 in the forward pass. The output of this layer is a matrix of the lower dimension which is then passed through the ReLU activation function and we obtain the low dimension node-level output of each node present in the training set. The dimension of the node-level output from the previous layer is further reduced in the second layer of the model. This is achieved by taking the output of the previous layer and the Identity matrix as the node feature matrix as inputs. This layer convolves the Second order neighbourhood of each node. In the last layer the lower dimension node-level output of the second layer representation is used on the training edge data which contains the positive edge data as well as the negative edge data as explained in 3.1.1. From the model we get the predicted edge labels and we calculate the binary cross entropy loss on comparing them with the actual edge labels. The binary cross entropy loss is shown in equation (9)

$$H_q(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i)) \quad (9)$$

where $p(y_i)$ is the predicted link label. For the back-propagation we used the Adam optimiser to update the weights on the loss calculated earlier. The default learning rate was 0.01.

4 EVALUATION

The evaluation of the designed PPI Link prediction model is performed on various hyperparameters. The performance of the model is evaluated on the validation set for numerous training epochs, hidden layer 1 units, hidden layer 2 units and the learning rates. The best result producing values of the hyperparameters is then used for testing the performance of the model on the testing set.

4.1 Effect of Epochs on Area Under Curve (AUC)

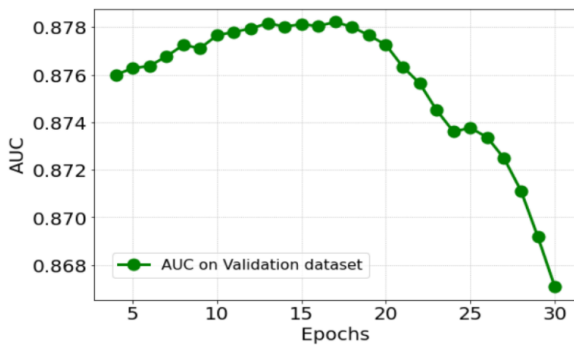


Figure 6: AUC vs Epochs

As shown in the Figure 6, we find that area under curve rises for testing set when epoch is increased from 5 to 12. However, AUC decreases between range 12 to 30. We can say that the increase in the number of epochs after epoch number 20 is making the model overfit to the training data. It eventually results in making the model unable to generalize over unseen validation data.

4.2 Effect of the Number of units in hidden layer 1 on AUC

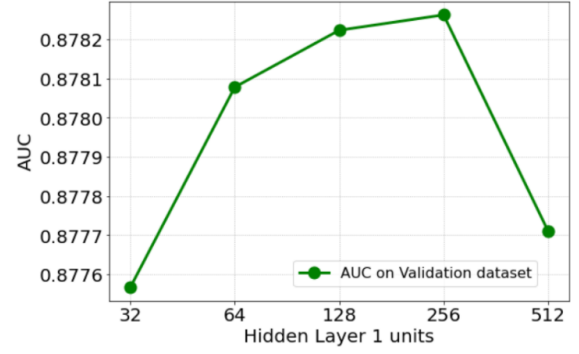


Figure 7: AUC vs No. Of Hidden layer 1 units

As shown in the Figure 7, we find that the area under the curve for testing set at increases sharply till the number of units of the hidden layer 1 is 64 but after it increases less sharply with reference to the previous increase with the increase in the the number of units of the hidden layer 1 to 256 and it decreases when the number of units are further increased to 521. We can say that the increase in the number of hidden units in hidden layer one after 256 hidden units is making the model overfit to the training data due to the increased complexity of the model. It eventually results in making the model unable to generalize over unseen validation data.

4.3 Effect of the Number of units in hidden layer 2 on AUC

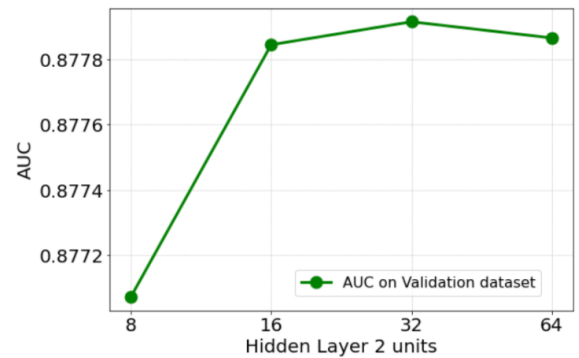


Figure 8: AUC vs No. Of Hidden layer 2 units

As shown in the figure 8, we find that the area under the curve for testing set remains increases as the number of units in hidden layer 2 increases and then remains constant when the number of units are further increased. We can say that the increase in the number of hidden units in hidden layer 2 after 32 hidden units is making the model overfit to the training data due to the increased complexity of the model. It eventually results in making the model unable to generalize over unseen validation data.

4.4 Effect of learning rate on AUC

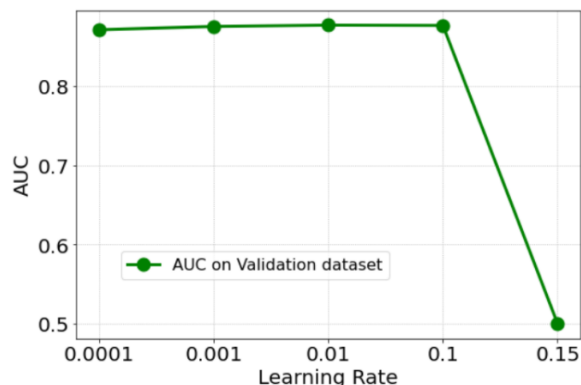


Figure 9: AUC vs Learning Rate

As shown in the Figure 9, we find that the AUC for the testing set remains constant until the learning rate 0.1 and after it, it decreases. We can say that the increase in learning after 0.1 learning rate is making the model unable to converge at global minima loss and hence it is affecting the accuracy of prediction.

4.5 Effect on the performance before and after removing the value "-" from the dataset

During the preprocessing of the data we encountered a node labeled as "-", we were not sure if it was a missing value or if a protein was interacting with itself, so we tried to run our code in two variants one with "-" and one without "-" with the number of units in the hidden layer 1 as 32 and that of the hidden layer 2 as 16, learning rate of 0.001 and for 10 epochs. The results are shown in the Table 2

Table 2: Values of ROC and Average Precision with and without "-"

	Value with "-"	Value without "-"
Validation ROC score	0.87667	0.87941
Validation Average Precision	0.86486	0.86743
Testing ROC score	0.87644	0.88135
Testing Average Precision	0.86309	0.86988

4.6 Evaluation on testing set

The model configured on the value of best hyperparameters chosen from the result of the earlier section is then used to evaluate the

performance of the model on the testing set. The result of the model is shown in the below table 3.

Table 3: Evalution on testing set.

Number of Epochs	17
Hidden Layer 1's Units	256
Hidden Layer 2's Units	32
Learning Rate	0.01
Training Loss	0.5897
AUC/ROC on Validation Set	0.873925
Average Precision score on Validation Set	0.867879
AUC/ROC on Test Set	0.870039
Average Precision score on Test Set	0.860074

5 ACKNOWLEDGEMENT

This article was written during the Data Science Lab 2020 at the University of Passau.

REFERENCES

- [1] Daniel A. Schult, Aric A. Hagberg and Pieter J. Swart. Aug 2008. "Exploring network structure, dynamics, and function using NetworkX". in *Proceedings of the 7th Python in Science Conference (SciPy2008)* 1 (Aug 2008), pp. 11–15. http://conference.scipy.org/proceedings/SciPy2008/paper_2/
- [2] Stryer L. Berg JM, Tymoczko JL. 2002. *Biochemistry*. Section 3.2: Primary Structure: Amino Acids Are Linked by Peptide Bonds to Form Polypeptide Chains, Vol. 5th edition. W H Freeman, New York. <https://www.ncbi.nlm.nih.gov/books/NBK22364/>
- [3] Norberto de Souza O et al Breda A, Valadares NF. 2006 May 1 [Updated 2007 Sep 14]. Protein Structure, Modelling and Applications. *Bioinformatics in Tropical Disease Research: A Practical and Case-Study Approach [Internet]*. 1 (2006 May 1 [Updated 2007 Sep 14]), Chapter A06. <https://www.ncbi.nlm.nih.gov/books/NBK6824/>
- [4] Ul Qamar MT, Chen LL, Ding YD, Chang JW, Zhou YQ. Nov 22, 2016. Prediction of Protein-Protein Interactions by Evidence Combining Methods. *Int J Mol Sci*. 2016;17(11):1946. 6 (Nov 22, 2016). <https://doi.org/10.3390/ijms17111946>
- [5] Alex Fout, Jonathon Byrd, Basir Shariat, and Asa Ben-Hur. 2017. Protein interface prediction using graph convolutional networks. In *Advances in neural information processing systems*. 6530–6539.
- [6] Li Yongli, Yin Han, Le Huang, Rui Mao, Guo William, Dong Liyan. 2013/09/17. "The Algorithm of Link Prediction on Social Network". *Mathematical Problems in Engineering* 2013, 1 (2013/09/17). <https://doi.org/10.1155/2013/125123>
- [7] S Jones and J M Thornton. 1996. Principles of protein-protein interactions. *Proceedings of the National Academy of Sciences* 93, 5 (1996), 13–20. <https://doi.org/10.1073/pnas.93.1.13> arXiv:https://www.pnas.org/content/93/1/13.full.pdf
- [8] Thomas N. Kipf and Max Welling. 2016. Semi-Supervised Classification with Graph Convolutional Networks. arXiv:cs.LG/1609.02907
- [9] Xue Li, Lifeng Yang, Xiaopan Zhang, and Xiong Jiao. 2019. Prediction of Protein-Protein Interactions Based on Domain. *Computational and mathematical methods in medicine* 2019 (2019).
- [10] Guo J, Zhang DY, Lin K, Wu X, Zhu L. 2006 Apr 26. "Prediction of yeast protein-protein interaction network: insights from the Gene Ontology and annotations". *Nucleic Acids Resource* vol. 34, 7 2137–50., 1 (2006 Apr 26). <https://doi.org/10.1093/nar/gkl219>