

Hybrid Approach of Relation Network and Localized Graph Convolutional Filtering for Breast Cancer Subtype Classification

Sungmin Rhee and Seokjun Seo

Seoul National University
Computer Science and Engineering
{lars, dane2522}@snu.ac.kr

Sun Kim

Seoul National University
Computer Science and Engineering
Interdisciplinary Program in Bioinformatics
sunkim.bioinfo@snu.ac.kr

Abstract

Network biology has been successfully used to help reveal complex mechanisms of disease, especially cancer. On the other hand, network biology requires in-depth knowledge to construct disease-specific networks, but our current knowledge is very limited even with the recent advances in human cancer biology. Deep learning has shown a great potential to address the difficult situation like this. However, deep learning technologies conventionally use grid-like structured data, thus application of deep learning technologies to the classification of human disease subtypes is yet to be explored. Recently, graph based deep learning techniques have emerged, which becomes an opportunity to leverage analyses in network biology. In this paper, we proposed a hybrid model, which integrates two key components 1) graph convolution neural network (graph CNN) and 2) relation network (RN). We utilize graph CNN as a component to learn expression patterns of cooperative gene community, and RN as a component to learn associations between learned patterns. The proposed model is applied to the PAM50 breast cancer subtype classification task, the standard breast cancer subtype classification of clinical utility. In experiments of both subtype classification and patient survival analysis, our proposed method achieved significantly better performances than existing methods. We believe that this work is an important starting point to realize the upcoming personalized medicine.

1 Introduction

Breast cancer is one of the most common cancers, especially leading type of cancer in women with more than 1,300,000 cases and 450,000 deaths per year worldwide (Network and others 2012). Breast cancer has multiple risk factors for development and proliferation including genetic change (Sineshaw et al. 2014), epigenetic change (Rhee et al. 2013), and environmental factors like age (Nixon et al. 1994) or residual environments (Sineshaw et al. 2014). Also, breast cancer is a complex, multifactorial disease where the interplay between these risk factors decide the phenotype of cancer such as progression, development, or metastasis (Martin and Weber 2000). As a result, there are multiple scenarios of tumorigenesis. Thus, it is a challenging problem to determine how a cancer cell is developed and progressed.

Characterizing mechanisms of a complex disease as a whole is not possible. An effective approach is to define subtypes by dividing cancer into several categories according to

various criteria such as phenotype, molecular portraits, and histopathology. For the molecular portrait-based subtype scheme of breast cancer, numerous studies have made efforts to build gene-expression based models based on differences in molecular mechanisms among the patient cohort (Perou et al. 2000; Herschkowitz et al. 2007; Paik et al. 2004; Prat et al. 2010; Parker et al. 2009). Among the several molecular properties based breast cancer subtypes, PAM50 (Parker et al. 2009) has become a standardized model with the clinical utility to make diagnosis decisions in practice or building a treatment plan for a patient. For instance of practical usage of subtype in cancer therapy, St. Gallen international expert consensus panel introduced a system for recommending adjuvant systemic therapy based on breast cancer subtype since 2011 (Goldhirsch et al. 2013). However, despite the practical utility of breast cancer subtypes, they are still remained to be suboptimal as the complex mechanism underlying breast cancer cell is not fully investigated.

The main technical issue in elucidating biological mechanisms of breast cancer is that innate relational and cooperative characteristic of genes should be considered. In any specific biological context, multiple dysregulated genes derive phenotypic differences by mechanisms such as forming complexes, regulate each other, or affect signal transduction. To address the technical difficulties related to complex associations, many previous studies utilized biological networks (Barabasi and Oltvai 2004; Barabási, Gulbahce, and Loscalzo 2011; Cowen et al. 2017). Early biological network studies tried to discover distinct patterns based on edge information (Barabasi and Oltvai 2004; Barabási, Gulbahce, and Loscalzo 2011), and recent approaches rely on the common paradigm *network propagation*, which assumes that the information is propagated to nearby nodes through the edges (Cowen et al. 2017). However, network biology requires in-depth knowledge to construct disease-specific networks, but our current knowledge is very limited even with the recent advances in human cancer biology.

Deep learning has shown a great potential to address the difficult situation like this. However, application of deep learning technologies to the classification of human disease subtypes is not straightforward since deep learning technologies conventionally use grid-like structured data and they are not designed to handle network data. Recently, graph based deep learning techniques have emerged, which becomes an

opportunity to leverage analyses in network biology.

In this paper, we tried to advance network bioinformatics by incorporating a novel hybrid method of relation network (RN) and graph convolution neural network (graph CNN). The concept of the problem is illustrated in Fig. 1. As traditional deep learning approaches assume input data with Euclidean or grid-like structure, there have been limited attempt to incorporate deep learnings in network biology. To overcome this issue, we incorporated and extended recent deep learning methods. To sum up, given the prior knowledge of putative associating genes represented in a graph structure, our proposed method is to capture localized patterns of associating genes with graph CNN, and then learn the relation between these patterns by RN.

In summary, the main contributions of this work are as follows:

- We proposed a novel hybrid approach composed of graph CNN and RN. Our method is motivated by the fact that relations between entities are traditionally modeled with a graph structure. To the best of our knowledge, this work is the first of its kind.
- We designed a model for biological networks. Since the dimension of the conventional biological network is too large, we applied fast graph convolution filtering method that can scale up to the very large dimension. In addition, we modified the relation network to fit in the task.
- We demonstrated the effectiveness of our approach by classifying PAM50 breast cancer subtype on TCGA dataset. Our model was able to achieve good performance in both classification evaluation and capturing biological characteristics such as survival hazard and subtype prognosis.

The article is organized as follows. In the next section, we review previous studies related to our work. Then, our model is described in Section 3. In Section 4, we demonstrate the experimental result in both quantitative and qualitative perspectives.

2 Related Work

2.1 Deep learnings on graphs

Recent survey papers (Niepert, Ahmed, and Kutzkov 2016; Bronstein et al. 2017) present comprehensive surveys on graph deep learnings that recently emerge. In this section, we review a representative selection of the previous studies related to this work.

In the case of recurrent neural networks (RNN), there has been an attempt (Scarselli et al. 2009) to combine the graph structure with the neural network earlier than CNN. Graph neural network (GNN) is one of such study, which is an extension of recursive neural network and random walk. The representation of each node propagates through edges until it reaches a stable equilibrium. Then it is used as the features in classification and regression problems. This approach was further extended to a method named gated graph sequence neural network (GGs-NN), introducing gated recurrent unit and modifying to output sequence (Li et al. 2015). Recently,

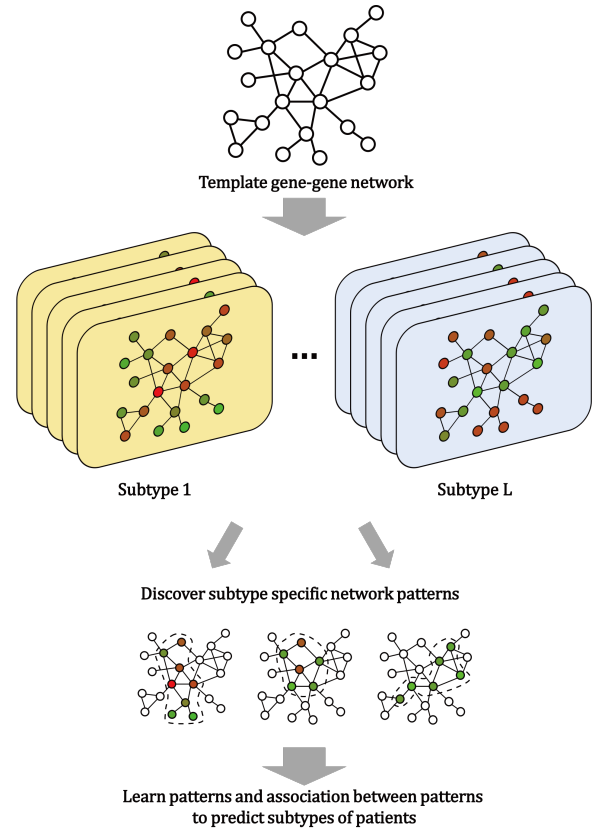


Figure 1: Conceptual illustration of our task.

Johnson proposed a method built upon GGS-NN by allowing graph-structured intermediate representations, as well as graph-structured outputs (Johnson 2016).

CNN has been successful on domains with underlying grid-like structured data such as computer vision, natural language processing, audio analysis, and DNA sequences. Recently, several works extended CNN to more general topologies like manifolds or graphs (Bruna et al. 2013; Henaff, Bruna, and LeCun 2015; Niepert, Ahmed, and Kutzkov 2016; Defferrard, Bresson, and Vandergheynst 2016). Bruna et al. introduced a spectral formulation of graph CNN (SCNN), which draws on the properties of convolutions in the Fourier domain (Bruna et al. 2013). They extended the ideas to large-scale classification problems (Henaff, Bruna, and LeCun 2015). Nevertheless, the method still does not scale up well due to the computational cost of $O(n^2)$ for matrix multiplication.

The method proposed by Defferrard, Bresson, and Vandergheynst leveraged on the spectral graph CNN (SCNN), which is a key component of our approach. Computational efficiency was improved by using Chebyshev approximation technique. As a result, their method outperformed the existing SCNNs in terms of accuracy in their experiments.

2.2 Relation reasoning

As Google DeepMind’s team mentioned in their recent study (Santoro et al. 2017), deductive reasoning methods innately reason about relations of entities from training data, which is also represented as relations (Quinlan 1990). However, these approaches lack the ability to deal with fuzzy and variational input data robustly (Harnad 1990). Meanwhile in the statistical learning domain, DeepMind recently proposed a method named relational network as a general solution to relational reasoning in neural networks (Santoro et al. 2017). The method focused on relational reasoning with an easy-to-understand and generalizable network structure, which makes it easy to be modified or combined with other methods. In addition, despite its simplicity in structure, it has demonstrated super-human performance in visual question answering problem. Owing to the simplicity of RN structure, we could easily modify and combine it with graph CNN for the first time, and shows that the relational reasoning indeed helps to improve the performance of the proposed breast cancer subtype classification task.

3 Methods

In this section, we describe the proposed method. Fig. 2 illustrates the overall workflow of the proposed method. The first step (Fig. 2 A) of the method is the graph convolution step to represent and capture localized patterns of the graph nodes (genes). The second step (Fig. 2 B) is the relational reasoning step. In this step, we built the model that can learn the complex association between graph node groups (gene sets) from the learned localized patterns of graph nodes (genes) in the previous step. The next step is to merge the representation of graph convolution layer and relation reasoning layer.

For the explanation of our method, we will denote data elements of each sample p as $x_p \in R^n$, and weighted graph topology as $G = (V, E, A)$, where V and E represent the sets of vertices and edges, respectively. Also, we will use A to denote the weighted adjacency matrix and N to denote the number of vertices, i.e. $|V|$.

3.1 Localized pattern representation by graph convolution neural network

For capturing localized patterns of data (gene expression profile), we first mapped input data in the graph structure and used graph CNN technique to find localized patterns of the graph signal. Let $x_p \in R^n$ be the graph signal (or gene expression) in the sample p . Then graph Laplacian matrix L is used to find spectral localized patterns of x_p under the graph structure G . Laplacian matrix L of graph G is defined as $L = D - A$ where D is a weighted degree matrix of graph G and A is a weighted adjacency matrix of graph G . As graph Laplacian L is a kind of graph shift operator, it can be used to represent diverse patterns of the graph in graph signal processing. One of the usages of graph Laplacian is that we can get the difference of signals between each node and its neighboring nodes simply by multiplying L to the graph signal x .

In addition to this basic usage, the graph convolution of signal x is also defined with graph Laplacian. Let’s assume that $L = U\Lambda U^T$ is an eigenvalue decomposition of graph Laplacian L where $U = [u_1, \dots, u_n]$ is a matrix composed of eigenvectors $\{u_l\}_{l=1}^n$ and Λ is a diagonal matrix $diag([\lambda_1, \dots, \lambda_n])$ composed of eigenvalues $\{\lambda_l\}_{l=1}^n$. We can say that $\{u_l\}_{l=1}^n$ is a complete set of orthonormal eigenvectors as L is a symmetric positive semidefinite matrix (Defferrard, Bresson, and Vandergheynst 2016). Then the graph Fourier transform is defined as $\hat{x} = U^T x$ and inverse graph Fourier transform is defined as $x = U\hat{x}$ (Shuman et al. 2013).

Unlike in the classical signal processing domain, it is not straightforward to define the convolution of two signals in spectral graph domain. Thus, convolution theorem in equation 1 is borrowed from signal processing domain to define graph convolution

$$(x * y)(t) = F^{-1}(F(f)F(g)) \quad (1)$$

where F and F^{-1} in the equation denotes Fourier and inverse Fourier transform for each, and x, y denotes two input signals. As Fourier transform in graph spectral domain is already defined, we can induce graph convolution by combining convolution theorem and graph Fourier transform as following equation

$$\begin{aligned} x *_G y &= U((U^T x) \odot (U^T y)) \\ &= U((U^T y) \odot (U^T x)) \\ &= Uy(\Lambda)U^T x \end{aligned} \quad (2)$$

where \odot is the element-wise Hadamard product and $y(\Lambda) \in R^n$ is a diagonal matrix $diag([\hat{y}(\lambda_1), \dots, \hat{y}(\lambda_n)])$. Since the matrix U is determined by the topology of input graph and it is invariant, only the matrix $y(\Lambda)$ determines various forms of convolution filters and acts as learnable parameters in graph convolution. Among several graph convolution filters, the proposed method used polynomial parametrized filter $y_\theta(\Lambda) = \sum_{k=0}^{K-1} \theta_k \Lambda^k$ that can express localized signal patterns in K -hop neighboring nodes. However, evaluating polynomial parametrized filter requires very expensive computational complexity $O(n^2)$. To deal with this circumstance, previous study (Hammond, Vandergheynst, and Gribonval 2011) proposed an approximated polynomial named Chebyshev expansion. The Chebyshev polynomial $T_k(x)$ of order k is recursively defined as $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$ with $T_0 = 1$ and $T_1 = x$. Then the filter can be approximated as $y_{\theta'}(\Lambda) = \sum_{k=0}^{K-1} \theta'_k T_k(\tilde{\Lambda})$ with $\tilde{\Lambda} = 2\Lambda/\lambda_{max} - I_n$.

Going back to the graph convolution in equation 2, we can now define the final graph convolution filter as

$$x *_G y_{\theta'} = \sum_{k=0}^{K-1} \theta'_k T_k(\tilde{L}) \quad (3)$$

where $\tilde{L} = 2L/\lambda_{max} - I_n$ is a rescaled graph Laplacian. The equation 3 can be easily induced from the observation

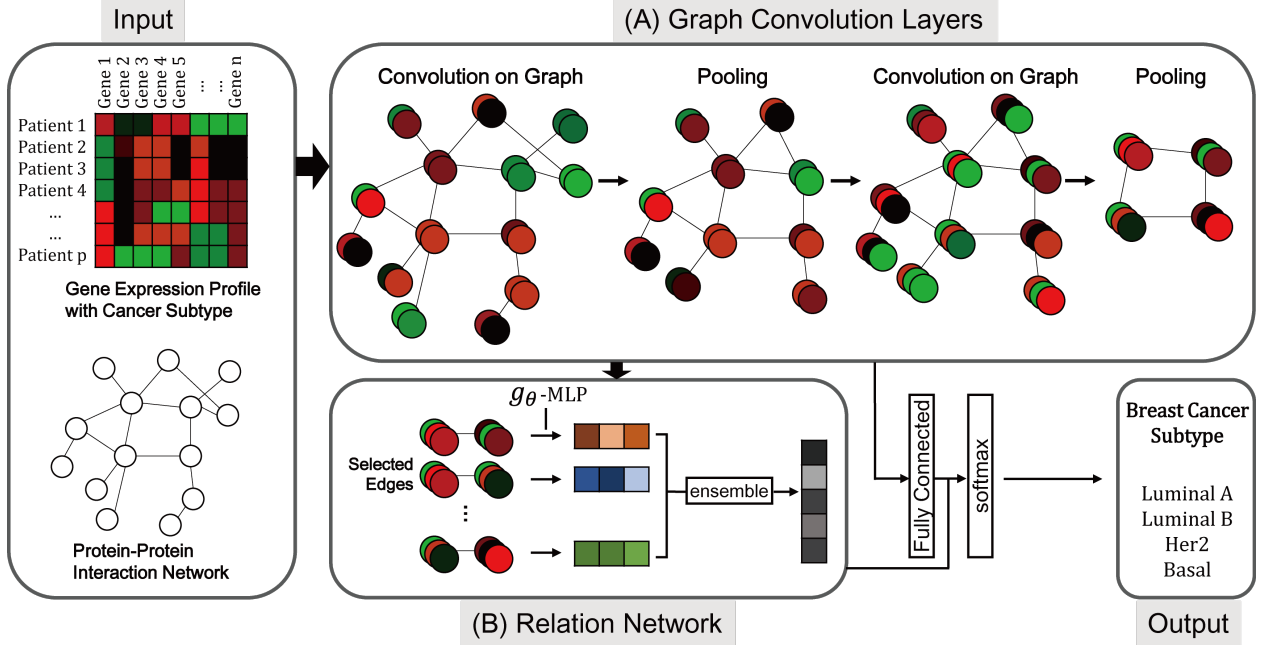


Figure 2: Overview of the proposed method.

$(U\Lambda U^T)^k = U\Lambda^k U^T$. Note that the entire filtering operation only requires $O(K|E|)$ and is fully differentiable, which makes it feasible to be learned by backpropagation algorithm (Defferrard, Bresson, and Vandergheynst 2016).

Convolutioned graph signal is further pooled with neighboring nodes identified by Graclus algorithm (Dhillon, Guan, and Kulis 2007) as proposed in the previous study (Defferrard, Bresson, and Vandergheynst 2016). There are several pooling strategies in neural network such as max pooling and average pooling. Empirically, the average pooling performed best in our experiments. Therefore, we used average pooling in the proposed method.

3.2 Learning relation between graph entities using relation network

To reason about association between graph nodes, we used relation network (RN), originally defined as

$$RN(O) = f_\phi(\sum_{i,j} g_\theta(o_i, o_j)) \quad (4)$$

in the previous study (Santoro et al. 2017) where the input is a set of objects $O = \{o_1, o_2, \dots, o_n\}, o_i \in R^m$. In the original study by Santoro et al., multi-layer perceptron (MLP) is used for function of f and g , and the parameters θ and ϕ are synaptic weights of perceptrons. Also, unlike the task in this paper, there exists a query for each of the input samples, and each query is embedded to the vector q by an LSTM. Then, Santoro et al. re-defined the RN architecture as $RN(O) = f_\phi(\sum_{i,j} g_\theta(o_i, o_j, q))$ so that it can process the query in the neural network. This query embedding q can work similarly as an attention mechanism. In other words, query embedding has an ability to select object pairs that are

important for the classification task. Also, coordinate values of objects are used to define object pairs in the work of Santoro et al. as it takes an image input, which has innate coordinate information. Therefore, even if all pairs of objects are considered in the original RN, it was able to show a good performance.

However, in our task only the information of object vectors are available in our task and the number of considering objects is larger than the original work. This lead to two technical problems for relation reasoning, 1) object pairs that are not relevant to solve the problem can interfere with learning, and 2) considering all pairs is not feasible as the number of the objects is too large. Thus, we modified the relation network to fit in our task. First, we sort the edges in the descending order of edge weights. Then top κ number of edges are selected as input object pairs. Relations between two objects in each of the selected pairs are then inferred with g function. Also, we used different functions for g and f . We used MLP with separated parameters for each of object pairs as g function and linear ensemble summation as f function. Then the final object relation network is inferred as

$$RN(O) = \sum_{i,j} \epsilon_{i,j} g_{\theta_{i,j}}(o_i, o_j) \quad (5)$$

where $\epsilon_{i,j}$ is a learnable parameter that can be interpreted as an attention of each object pairs for the task, $o_i \in R^m$ is object embedding for each of graph nodes, and m is the number of convolution filter in the last graph convolution layer.

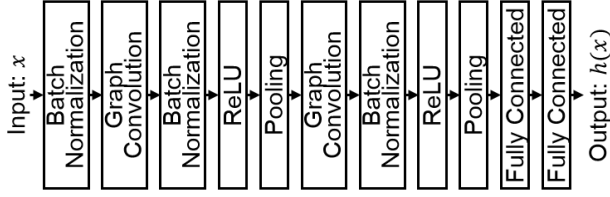


Figure 3: Illustration of h function.

3.3 Merging graph convolution layer and relation network

The final output of the proposed model architecture is defined as

$$\hat{y} = \text{softmax}(h(x_i) + \sum_{i,j} \epsilon_{ij} g_{\theta_{ij}}(o_i, o_j)) \quad (6)$$

to combine outputs from graph convolution layer and relation network, where $h(x_i)$ is a composition of functions: graph convolution, pooling, and fully connected layer.

To be more specific about h , we summarize the procedure of h as follows. First, the input signal x_i is normalized by a batch normalization method (Ioffe and Szegedy 2015) to make learning process stable since it has large absolute value and variance. Then, the normalized input signal is filtered by a graph convolution layer as defined in equation 3. Next, the convoluted signal is normalized through a batch normalization method so that the learning process is accelerated and learning parameters can have regularization effect. Then, ReLU activation function is used, defined as $\sigma(x) = \max(0, x)$. After the activation, average pooling is applied. We named the procedure from graph convolution filtering to average pooling as *graph convolution layer*. After two graph convolution layer, a final feature map is used as an input of fully connected layer. Next, the output of final fully connected layer is the output of function h . Function h is illustrated in Fig. 3.

Output at the last graph convolution layer was represented as feature maps that were directly input to the relation network. Then cross-entropy between \hat{y} in equation 6 and classification label is used as the final objective function. Xavier initialization technique (Glorot and Bengio 2010) is used for initialization of all parameters, and for parameter optimization, Adam (Kingma and Ba 2014) optimization and mini-batch training is used.

Hyperparameters for learning procedure is determined as follows. Two graph convolution layer were used. Each layer has 32 convolution filters. K of the first layer is 10 and second layer is 2. Pooling size is 2 for both of the layers. Two fully connected layers are used with 1024, 512 hidden nodes for each. For the relation network, top 1000 edges were selected. MLPs for g function have 2 layers with 128, 128 hidden nodes for each. We set hyperparameters for training procedure as mini-batch size = 100, learning rate = 1e-3, L2 regularization coefficient = 5e-4.

4 Results and Discussion

4.1 Dataset

We applied the proposed method on dataset of human breast cancer patient samples. RNA-seq based expression profiles of genes are extracted from TCGA breast cancer level 3 data (Prat Aparicio 2012). There are 57,292 genes in the original expression profile, and we excluded genes that were not expressed and further selected 4,303 genes in the cancer hall-mark gene sets (Liberzon et al. 2015) to utilize only genes that are relevant with tumor.

For the classification label of the patient, PAM50 molecular subtypes were used. PAM50 is the most commonly used breast cancer subtype scheme. The subtype includes Luminal A, Luminal B, Basal-like, and HER2. Luminal subtype cancer cells are mostly grown from inner (luminal) cells of mammary ducts and known to have better prognoses than other subtypes. Compared to Luminal A however, Luminal B subtype tumors tend to have poorer prognosis factors like higher tumor grade, larger tumor size, and lymph node involvement. Basal-like cancer cells are mostly grown from outer (basal) cells of mammary ducts and known to have worst prognoses and survival rates. HER2 subtype had its name since most HER2 subtype tumors are HER2-positive. HER2 subtype tumors tend to have poorer prognoses than luminal subtype tumors. In our study, 367 Luminal A, 295 Luminal B, 165 HER2, and 254 Basal-like patient samples were used for the experiment.

To reflect interactions among genes, gene expression values were mapped to a biological graph. For the topology of the graph, we used STRING protein-protein interaction network (Szklarczyk et al. 2014). STRING is a curated database of putatively associating genes from multiple pieces of evidence like biological experiments, text-mined literature information, computational prediction, etc. STRING is one of the most widely used networks in network bioinformatics (Han 2008).

4.2 Breast cancer subtype classification

For true breast cancer subtype label $Y = \{y_1, \dots, y_P\}$ and predicted label $\hat{Y} = \{\hat{y}_1, \dots, \hat{y}_P\}$ of P number of patients, the accuracy is calculated as

$$\text{Accuracy}(Y, \hat{Y}) = \frac{1}{P} \sum_{p=1}^P I(y_p = \hat{y}_p) \quad (7)$$

where I is an indicator function. The accuracies were measured by a monte-carlo cross validation experiment (Dubitzky, Granzow, and Berrar 2007). We repeatedly sampled 10% of patients as a validation set and used remaining 90% of patients as a training set. For each of data splits, the proposed model was fit to training set and accuracy was assessed using validation set. The accuracies were averaged over data splits.

4.3 Comparison of classification performance

Table 1 lists accuracies of the proposed model and comparing methods. Peak accuracies during learning processes

Table 1: Performance comparison of the methods on PAM50 breast cancer subtype classification. GCNN denotes graph convolution neural network that has identical hyperparameters with the proposed method. GCNN + RN denotes simple integration of GCNN and vanilla relation network.

Methods	Peak accuracy	Final accuracy
Proposed Method	87.03%	84.08%
GCNN + RN	67.05%	61.37%
GCNN	85.99%	82.93%
SVM	-	78.39%
Multinomial Naive Bayes	-	75.78%
Random Forest	-	79.60%

were listed for top 3 methods in the table, and final accuracies after learning are listed for all of the comparing methods. We can see that the proposed method performs best. Also, the simple integration of graph CNN and vanilla RN (GCNN+RN) showed the worst performance. We used identical hyperparameters and model structure with the proposed method for GCNN+RN, except that equation 5 is replaced with equation 4. We believe that GCNN+RN performs poor since, as we described in the method section, the original RN gets query encodings and coordinate values as inputs, that can work as a clue for relevant object selection. However, as there is no coordinate value or query in our task, we believe that changes in our hybrid approach is efficient to make an increase in performance.

4.4 Consistency of tSNE visualization and PAM50 subtype prognosis

To qualitatively study whether the learned representation can express the biological characteristic of the patients, tSNE plot (Maaten and Hinton 2008) of the last convolution feature map is drawn (Fig. 4). Only the representation vectors of the objects, which were inputs of relation network were used to plot. We can see distinct patterns between four subtype patients in the plot. However, the distinction between subtypes was not clear than typical examples e.g., tSNE plot of MNIST handwritten digits. We believe that this shows the complexity of the problem we are solving in the task. As we described earlier, the problem has higher input dimension and association between each feature should be considered.

More interestingly, we can see that the order of subtypes in the tSNE plot is identical to the order of prognosis of breast cancer subtypes. It is a well-known fact in the breast cancer clinical domain that Basal-like subtype has the worst prognosis, followed by HER2, Luminal B, and Luminal A. Especially, Basal-like subtype is known to have distinctive molecular characteristics from other subtypes (Bertucci, Finetti, and Birnbaum 2012), which is also represented in Fig. 4. All of these patterns do not appear in the tSNE plot with raw gene expression (Supplementary Figure). Thus, we can say that the proposed method successfully learn the la-

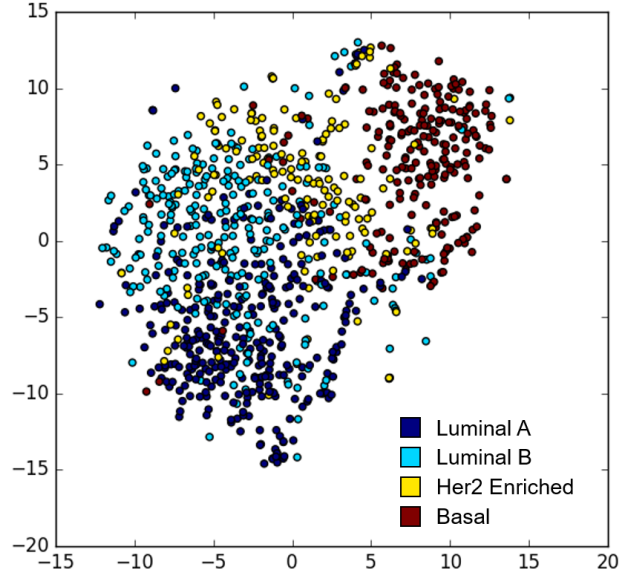


Figure 4: tSNE visualization of graph convolution feature map.

tent molecular properties in the expression profile of the patient samples.

4.5 Survival analysis

To further evaluate the ability of the model to comprehend characteristics of molecular subtypes, we performed survival analysis. We clustered the patients into two groups based on raw gene expression values and feature map data at the last graph convolution layer with dimensions reduced. Agglomerative hierarchical clustering with Ward’s criterion (Ward Jr 1963) was used for clustering and tSNE was used for dimension reduction. Then Kaplan-meier plots (KM plot) (Kaplan and Meier 1958) drawn for each of two clustering results are seen in Fig. 5. KM plot is standard analysis using non-parametric statistics to measure hazard ratio of each group. In medical science, KM plot is often used to analyze the effectiveness of treatment by comparing KM plot of treated and non-treated patient groups.

The plot generated by feature map values (Bottom of Fig. 5) shows that the patients were successfully divided into two subgroups that have distinct survival patterns with a p-value smaller than 0.05, while the plot with raw expression value (Top of Fig. 5) failed. This is an interesting result as it shows that the model could simultaneously learn the phenotypic information such as prognosis of the patient while performing the classification task, which is not directly related with the information.

5 Conclusion

In this study, we show that hybrid approach of relation network and graph convolution neural network can learn the

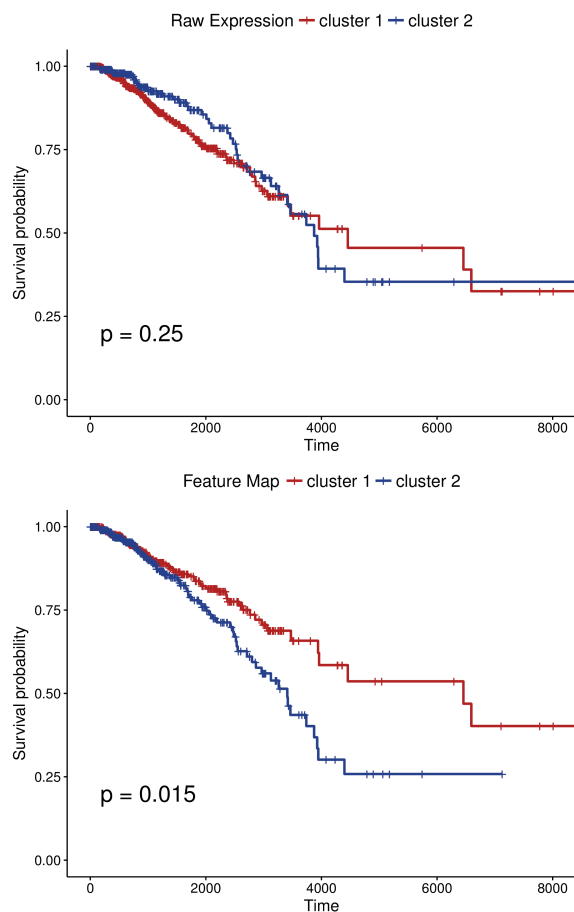


Figure 5: Kaplan meier survival plot of patients. Top plot is drawn with the patient clusters with raw gene expression and bottom plot is drawn with the patient clusters with learned convolution filter.

complex molecular mechanisms underlying breast cancer cells. The proposed method is designed to perceive cooperative patterns of genes and their associations. We observed that the method is successful to capture molecular characteristics of breast cancer in both quantitative and qualitative evaluation. We anticipate that our approach could extend the territory of both network bioinformatics and deep learnings. One important future work of the method is to extend the model to manage multiple heterogeneous data sources like sRNA sequencing, DNA methylation, as well as gene expression data. To do this, we plan to extend the model by incorporating other techniques like multi-view learning and/or transfer learning.

References

- [Barabasi and Oltvai 2004] Barabasi, A.-L., and Oltvai, Z. N. 2004. Network biology: understanding the cell's functional organization. *Nature reviews. Genetics* 5(2):101.
- [Barabási, Gulbahce, and Loscalzo 2011] Barabási, A.-L.; Gulbahce, N.; and Loscalzo, J. 2011. Network medicine: a network-based approach to human disease. *Nature reviews. Genetics* 12(1):56.
- [Bertucci, Finetti, and Birnbaum 2012] Bertucci, F.; Finetti, P.; and Birnbaum, D. 2012. Basal breast cancer: a complex and deadly molecular subtype. *Current molecular medicine* 12(1):96–110.
- [Bronstein et al. 2017] Bronstein, M. M.; Bruna, J.; LeCun, Y.; Szlam, A.; and Vandergheynst, P. 2017. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine* 34(4):18–42.
- [Bruna et al. 2013] Bruna, J.; Zaremba, W.; Szlam, A.; and LeCun, Y. 2013. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*.
- [Cowen et al. 2017] Cowen, L.; Ideker, T.; Raphael, B. J.; and Sharan, R. 2017. Network propagation: a universal amplifier of genetic associations. *Nature Reviews Genetics*.
- [Defferrard, Bresson, and Vandergheynst 2016] Defferrard, M.; Bresson, X.; and Vandergheynst, P. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*, 3844–3852.
- [Dhillon, Guan, and Kulis 2007] Dhillon, I. S.; Guan, Y.; and Kulis, B. 2007. Weighted graph cuts without eigenvectors a multilevel approach. *IEEE transactions on pattern analysis and machine intelligence* 29(11).
- [Dubitzky, Granzow, and Berrar 2007] Dubitzky, W.; Granzow, M.; and Berrar, D. P. 2007. *Fundamentals of data mining in genomics and proteomics*. Springer Science & Business Media.
- [Glorot and Bengio 2010] Glorot, X., and Bengio, Y. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 249–256.
- [Goldhirsch et al. 2013] Goldhirsch, A.; Winer, E. P.; Coates, A.; Gelber, R.; Piccart-Gebhart, M.; Thürlimann, B.; Senn, H.-J.; members, P.; Albain, K. S.; André, F.; et al. 2013. Personalizing the treatment of women with early breast cancer: highlights of the st gallen international expert consensus on the primary therapy of early breast cancer 2013. *Annals of oncology* 24(9):2206–2223.
- [Hammond, Vandergheynst, and Gribonval 2011] Hammond, D. K.; Vandergheynst, P.; and Gribonval, R. 2011. Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis* 30(2):129–150.
- [Han 2008] Han, J.-D. J. 2008. Understanding biological functions through molecular networks. *Cell research* 18(2):224.
- [Harnad 1990] Harnad, S. 1990. The symbol grounding problem. *Physica D: Nonlinear Phenomena* 42(1-3):335–346.
- [Henaff, Bruna, and LeCun 2015] Henaff, M.; Bruna, J.; and LeCun, Y. 2015. Deep convolutional networks on graph-structured data. *arXiv preprint arXiv:1506.05163*.

- [Herschkowitz et al. 2007] Herschkowitz, J. I.; Simin, K.; Weigman, V. J.; Mikaelian, I.; Usary, J.; Hu, Z.; Rasmussen, K. E.; Jones, L. P.; Assefnia, S.; Chandrasekharan, S.; et al. 2007. Identification of conserved gene expression features between murine mammary carcinoma models and human breast tumors. *Genome biology* 8(5):R76.
- [Ioffe and Szegedy 2015] Ioffe, S., and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, 448–456.
- [Johnson 2016] Johnson, D. D. 2016. Learning graphical state transitions.
- [Kaplan and Meier 1958] Kaplan, E. L., and Meier, P. 1958. Nonparametric estimation from incomplete observations. *Journal of the American statistical association* 53(282):457–481.
- [Kingma and Ba 2014] Kingma, D., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [Li et al. 2015] Li, Y.; Tarlow, D.; Brockschmidt, M.; and Zemel, R. 2015. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493*.
- [Liberzon et al. 2015] Liberzon, A.; Birger, C.; Thorvaldsdóttir, H.; Ghandi, M.; Mesirov, J. P.; and Tamayo, P. 2015. The molecular signatures database hallmark gene set collection. *Cell systems* 1(6):417–425.
- [Maaten and Hinton 2008] Maaten, L. v. d., and Hinton, G. 2008. Visualizing data using t-sne. *Journal of Machine Learning Research* 9(Nov):2579–2605.
- [Martin and Weber 2000] Martin, A.-M., and Weber, B. L. 2000. Genetic and hormonal risk factors in breast cancer. *Journal of the National Cancer Institute* 92(14):1126–1135.
- [Network and others 2012] Network, C. G. A., et al. 2012. Comprehensive molecular portraits of human breast tumors. *Nature* 490(7418):61.
- [Niepert, Ahmed, and Kutzkov 2016] Niepert, M.; Ahmed, M.; and Kutzkov, K. 2016. Learning convolutional neural networks for graphs. In *International Conference on Machine Learning*, 2014–2023.
- [Nixon et al. 1994] Nixon, A. J.; Neuberg, D.; Hayes, D. F.; Gelman, R.; Connolly, J. L.; Schnitt, S.; Abner, A.; Recht, A.; Vicini, F.; and Harris, J. R. 1994. Relationship of patient age to pathologic features of the tumor and prognosis for patients with stage i or ii breast cancer. *Journal of Clinical Oncology* 12(5):888–894.
- [Paik et al. 2004] Paik, S.; Shak, S.; Tang, G.; Kim, C.; Baker, J.; Cronin, M.; Baehner, F. L.; Walker, M. G.; Watson, D.; Park, T.; et al. 2004. A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *New England Journal of Medicine* 351(27):2817–2826.
- [Parker et al. 2009] Parker, J. S.; Mullins, M.; Cheang, M. C.; Leung, S.; Voduc, D.; Vickery, T.; Davies, S.; Fauron, C.; He, X.; Hu, Z.; et al. 2009. Supervised risk predictor of breast cancer based on intrinsic subtypes. *Journal of clinical oncology* 27(8):1160–1167.
- [Perou et al. 2000] Perou, C. M.; Sorlie, T.; Eisen, M. B.; Van De Rijn, M.; et al. 2000. Molecular portraits of human breast tumours. *Nature* 406(6797):747.
- [Prat Aparicio 2012] Prat Aparicio, A. 2012. Comprehensive molecular portraits of human breast tumours. *Nature*, 2012, vol. 490, num. 7418, p. 61–70.
- [Prat et al. 2010] Prat, A.; Parker, J. S.; Karginova, O.; Fan, C.; Livasy, C.; Herschkowitz, J. I.; He, X.; and Perou, C. M. 2010. Phenotypic and molecular characterization of the claudin-low intrinsic subtype of breast cancer. *Breast cancer research* 12(5):R68.
- [Quinlan 1990] Quinlan, J. R. 1990. Learning logical definitions from relations. *Machine learning* 5(3):239–266.
- [Rhee et al. 2013] Rhee, J.-K.; Kim, K.; Chae, H.; Evans, J.; Yan, P.; Zhang, B.-T.; Gray, J.; Spellman, P.; Huang, T. H.-M.; Nephew, K. P.; et al. 2013. Integrated analysis of genome-wide dna methylation and gene expression profiles in molecular subtypes of breast cancer. *Nucleic acids research* 41(18):8464–8474.
- [Santoro et al. 2017] Santoro, A.; Raposo, D.; Barrett, D. G.; Malinowski, M.; Pascanu, R.; Battaglia, P.; and Lillicrap, T. 2017. A simple neural network module for relational reasoning. *arXiv preprint arXiv:1706.01427*.
- [Scarselli et al. 2009] Scarselli, F.; Gori, M.; Tsoi, A. C.; Hagenbuchner, M.; and Monfardini, G. 2009. The graph neural network model. *IEEE Transactions on Neural Networks* 20(1):61–80.
- [Shuman et al. 2013] Shuman, D. I.; Narang, S. K.; Frossard, P.; Ortega, A.; and Vandergheynst, P. 2013. The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine* 30(3):83–98.
- [Sineshaw et al. 2014] Sineshaw, H. M.; Gaudet, M.; Ward, E. M.; Flanders, W. D.; Desantis, C.; Lin, C. C.; and Jemal, A. 2014. Association of race/ethnicity, socioeconomic status, and breast cancer subtypes in the national cancer data base (2010–2011). *Breast cancer research and treatment* 145(3):753–763.
- [Szklarczyk et al. 2014] Szklarczyk, D.; Franceschini, A.; Wyder, S.; Forslund, K.; Heller, D.; Huerta-Cepas, J.; Simonovic, M.; Roth, A.; Santos, A.; Tsafou, K. P.; et al. 2014. String v10: protein–protein interaction networks, integrated over the tree of life. *Nucleic acids research* 43(D1):D447–D452.
- [Ward Jr 1963] Ward Jr, J. H. 1963. Hierarchical grouping to optimize an objective function. *Journal of the American statistical association* 58(301):236–244.