

# Modelling Gene Expression by Integrating GRNs and HMs using GCNs

Shalin Patel<sup>1,2</sup>, TBD<sup>1,2,3</sup>, and Ritambhara Singh<sup>2,3</sup>

<sup>1</sup>*Division of Applied Mathematics, Brown University*

<sup>2</sup>*Center for Computational Molecular Biology, Brown University*

<sup>3</sup>*Department of Computer Science, Brown University*

## Abstract

## 1 Introduction

## 2 Related Work

## 3 Method

### 3.1 Formulation for Task

In this paper, we use the same inputs and outputs as Attentive and DeepChrome while also adding a gene expression matrix for each cell line. Using the same formulation as Cheng *et al.* the task is formulated as measuring the gene expression as either up (1) or down (0) regulated. First, per cell line, a GRN is precomputed which utilizes a matrix of size  $S \times G$  where  $S$  denotes the number of samples in the expression matrix while  $G$  represents the number of genes that were recorded.

Hence, for a sample gene, two pieces of information are fed. The first is  $\mathcal{G}$  which is a graph describing gene-gene interactions for a particular cell line. In the case of this paper  $\mathcal{G}$  was an adjacency list representation of the graph. Second, per gene, a matrix of size  $M \times T$  was utilized where  $M$  denotes the number of histone marks utilized while  $T$  is the total number of bin positions taken into account around the TSS site of a gene.

Overall, for the training data of the GCN, we utilized  $\mathcal{G}$  and a  $N \times M \times T$  sized matrix where  $N$  is the number of gene samples. The output, accordingly, is a  $N$ -sized vector with either 0 or 1.

### 3.2 Workflow

The workflow utilized in this paper follows three key steps to generate accurate gene expression modelling. First, the gene expression matrices are passed through a random forest based learner to determine importance scores between all the different genes in a particular cell line. Second, the HM data is fed through a 2D convolutional neural network to help capture the combinatorial interactions that can occur across HM lines as well as spatially within an HM line. This network compresses the original  $M \times T$

matrix into a 1D vector. Finally, these two pieces of information are fed through a graph convolutional network to present a final two element array which, when fed through a softmax activation, determines the classification for a particular gene.

### 3.2.1 Construction of GRNs

In order to capture the effects that genes have on one another, a Gene Regulatory Network (GRN) was constructed using a gene expression matrix. This network not only captures the complex interaction between genes, but also it is able to determine the weight of these influences. The gene regulatory networks for this paper are built using the standard method of random forests. Utilizing the grnboost2 algorithm from the arboreto package on `pypi.org`, a regression task was defined for each gene in a cell line.

Let  $E_i$  denote the row vector containing all samples for gene  $i$  in the expression matrix. Then, specifically, for gene  $i$ , a random forest model  $R_i$  was defined where,  $L(R_i(E_{1:G \setminus i}), E_i)$  is minimized with  $L$  denoting the mean square error. Once this task is completed, denote the set

$$\mathfrak{N}_i = \{imp(R_i, j) \mid j \in \{1 : G\} \setminus i\}.$$

Here,  $imp(R_i, j)$  refers to the feature importance of  $j$  in random forest model  $R_i$ . Then the final graph  $\mathcal{G}$  is constructed with the neighbor list of a node  $i$  being

$$\mathcal{I}_i = \{j \mid imp(R_i, j) > \bar{\mathfrak{N}}_i + s(\mathfrak{N}_i)\}$$

in which  $s(\mathfrak{N}_i)$  denotes the sample standard deviation. Hence,  $\mathcal{G} = \{\mathcal{I}_i\}, \forall i$ .

### 3.2.2 HM Encoding

Due to the current requirements imposed on the inputs of GCNs, the input into the GCN layers along with a graph must be a one dimensional vector. As the HM data is provided as an  $M \times T$  sized matrix, a CNN architecture is utilized to encode the information stored in the matrix to a one dimensional representation. The utilization of a CNN helps capture the combinatorial interactions that take place across HM lines as well as across multiple bins.

This paper utilized the 2DConv layers from pytorch wherein the input  $X$ , is a  $N \times 1 \times M \times T$  sized matrix with the dimension of size one taking the role of the number of channels in the feature space. Our model utilizes three of these 2DConv layers chained together eventually giving a  $N \times S$  output with  $E$  representing the number of features in the encoded space.

### 3.2.3 Activation and Dropout.

Throughout these processes, nonlinear activation and dropout is applied. These help regularize the model and help avoid overfitting during the training process. The activation utilized was an elementwise  $\tanh$ . It is useful to utilize this transformation because  $\forall x \in \mathbb{R}, -1 < \tanh(x) < 1$  which kept the model coherent through layers. The dropout utilized randomly set values in the feature map to zero with probability 0.5. This assisted in simulating dead signals and increasing the robustness of the model itself.

### 3.2.4 GCN Layer

The input into the GCN layers is, thus, a graph  $\mathcal{G}$  and the output from the encoding layer. Call this  $H$ . Let  $H^{(l)}$  denote the  $l$ th hidden layer while  $h_i^{(l)}$  is the feature matrix at node  $i$  at the  $l$ th layer.

This paper utilizes a series of six GCN combined. These layers, namely, are the GraphSAGE, ChebConv, and TAGConv layers found in the DGL package. The first layer utilized is the GraphSAGE layer which updates the hidden layer as follows,

$$\begin{aligned} h_{\mathcal{N}(i)}^{(l+1)} &= \text{aggregate}(\{h_j^l \mid j \in \mathcal{N}(i)\}) \\ h_i^{(l+1)} &= \sigma \left( W \cdot \text{concat}(h_i^l, h_{\mathcal{N}(i)}^{(l+1)} + b) \right) \\ h_i^{(l+1)} &= \text{norm}(h_i^{(l+1)}) \end{aligned}$$

where aggregate represents an aggregation that takes a  $x \times L$  feature input where  $x$  is variable and converts it to an  $L$  element vector.  $W$  and  $b$  represent learnable parameters while  $\sigma$  and norm are activation and normalization functions, respectively. The ChebConv and TAGConv layers follow similar methodologies where they first aggregate the hidden layer, apply a learnable function to the aggregated features and then provides normalization and activation. The following details the update to the hidden layer for ChebConv.

$$\begin{aligned} h_i^{(l+1)} &= \sum_{k=0}^{K-1} W^{k,l} z_i^{k,l} \\ Z^{0,l} &= H^l \\ Z^{1,l} &= \hat{L} \cdot H^l \\ Z^{k,l} &= 2\hat{L} \cdot Z^{k-1,l} - Z^{k-2,l} \\ \hat{L} &= 2 \frac{I - \hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2}}{\lambda_{max}} - I \end{aligned}$$

In this case,  $W$  represents the learnable weight matrix which changes based on layer  $l$  and hops  $k$ .  $Z$  represents the aggregated node features while which also varies according to  $l$  and  $k$ . Finally,  $\hat{L}$  denotes the graph laplacian while  $\hat{D}$  and  $\hat{A}$  denote the diagonal degree and adjacency matrices, respectively. Similarly, the update applied to the hidden layer using TAGConv is

$$H^{l+1} = \sum_{k=0}^K (\hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2}) H^l W^{k,l}$$

These layers were all chained together eventually creating an  $N \times 2$  matrix.

## 3.3 Training

As a whole the entire network can be described as

$$\hat{y}_i = \text{softmax}(f_{gcn}(f_{conv}(X_i), \mathcal{G}))$$

Suppose that all of the learnable parameters in the model, which are initially random, are in the vector  $\Theta$ . Then using the standard cross entropy loss, we can calculate the loss with respect to  $\Theta$

$$L(\Theta) = \sum_{i=0}^N \text{cross\_entropy}(y_i, \hat{y}_i)$$

|           | E003          | E066          | E071          | E096          | E114          | E116          | E118          | <i>Average</i> |
|-----------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|----------------|
| GCN       | <b>0.7146</b> | <b>0.7503</b> | <b>0.6512</b> | <b>0.6883</b> | <b>0.7966</b> | 0.9168        | <b>0.8220</b> | <b>0.7628</b>  |
| Random    | 0.4996        | 0.4968        | 0.6066        | 0.5000        | 0.6143        | 0.5032        | 0.5000        |                |
| Permute   | 0.6899        | 0.6999        | 0.6125        | 0.6389        | 0.7527        | 0.8624        | 0.7649        | 0.7094         |
| Attentive | 0.6759        | 0.7468        | 0.4906        | 0.5319        | 0.7930        | <b>0.9225</b> | 0.8181        | 0.7113         |
| MLP       | 0.6894        | 0.7035        | 0.6220        | 0.6420        | 0.7477        | 0.8606        | 0.7661        | 0.7188         |

Table 1: AUC Performance Evaluation Across Cell Lines and Baselines

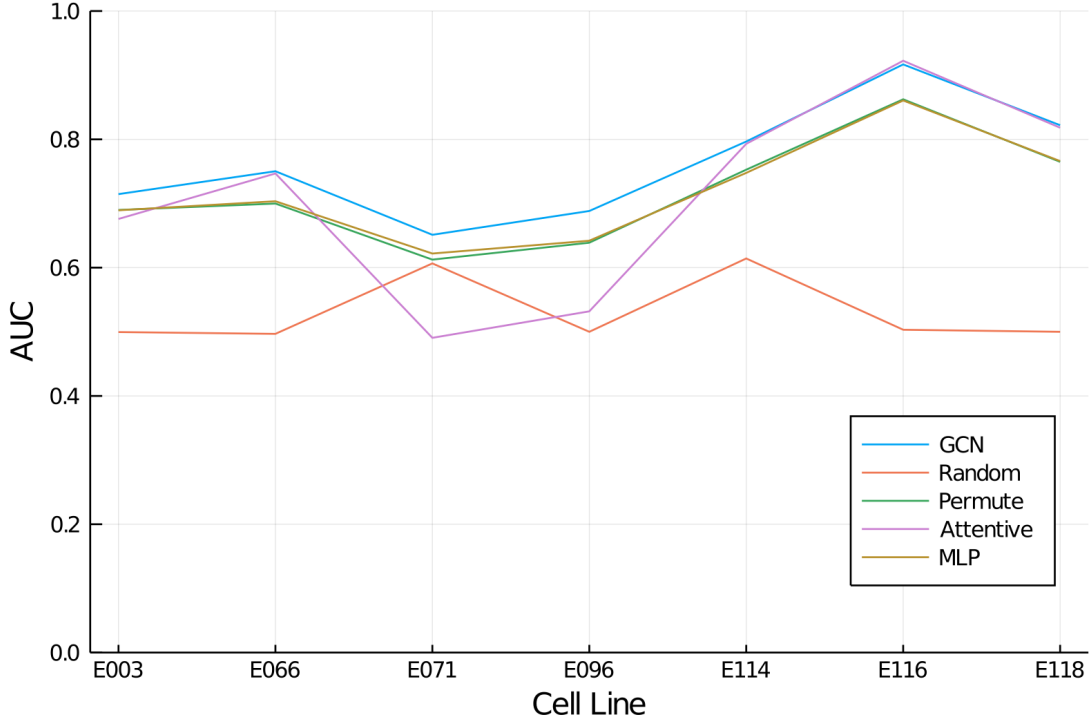


Figure 1: AUC Performance Across Cell Lines and Baselines Method

To train these parameters, and thus the entire network the ADAM optimizer is used. This is possible because all steps in the network are differentiable, and thus, the gradient of  $L$  can be calculated with respect to  $\Theta$ . Hence, ADAM can use backpropagation to train the model.

## 4 Experimental Setup

## 5 Results

### 5.1 Performance Evaluation

The performance of the model was checked against four baselines across seven different cell lines and is summarized in Table 1.

Quite clearly, the GCN model outperforms almost all of the other baseline methods, including the state of the art AttentiveChrome, across the cell lines. The only cell line where the GCN model was beaten was

in the E116 cell line which was high performing across models. It is clear, though, that the GCN model shows great consistency and proves performance gains in lower performing cell lines such as E071. This is encapsulated in the average AUC performance metric across cell lines where the GCN model has a significantly higher average than all other models and in Figure 1.

## **6 Discussion**

**Acknowledgments** .

**Funding**

## **References**