

Machine Learning and Algorithms for Data Mining

Assessment 2*

Analysis of graph-structured data

Sebastian Borgeaud — spb61@cam.ac.uk

February 22, 2018

Abstract

I explore various algorithms for the task of transductive classification on a graph and report their accuracy on the Cora dataset [1]. First, I present a simple baseline method based on support vector machines that disregards all information of the graph structure. Next, I propose two methods that incorporate the features of neighbouring nodes by averaging over them. As expected, incorporating knowledge about the graph structure gives a significant improvement in classification accuracy, even when incorporated in a simple way. Finally, I reimplement two state-of-the-art methods: Graph Convolutional Networks [2], and Graph Attention Networks [3] which achieve the highest accuracy amongst the tested methods.

1 Main contributions, Key Concepts and Ideas

1.1 Introduction

Many real world datasets occur naturally in the form of graphs, examples from different fields include protein-protein interaction networks in biology [4], social networks in sociology [5], and scientific collaboration networks [6]. I focus on the problem of node classification for citation networks in a transductive setting and present multiple approaches to solve this problem.

*Word count: 2467 — Computed using TexCount

In the transductive case, as opposed to the inductive case, the entire structure of the network is known at training time but only some of the nodes are labelled. At test time, the task is to infer labels for the remaining nodes. This differs from supervised learning in two ways. i) The data points (the nodes) are connected to each other. These connections or edges contain further information that might be helpful for classification. ii) The features of the nodes that will be classified at test time are known in advance as the entire graph is known at test time.

1.2 Related work

It is generally possible to distinguish between two categories of approaches for graph-based semi-supervised learning: methods based on an explicit Laplacian regularisation cost and methods based on graph-embedding approaches.

In the first approach, the problem is seen as a graph-based semi-supervised learning problem, where the label information is smoothed over the graph via some form of explicit graph-based regularisation. For example, we could make the assumption that adjacent nodes are more likely to have the same class label and incorporate this in some loss

$$\mathcal{L} = \mathcal{L}_0 + \lambda \mathcal{L}_{reg}. \quad (1)$$

\mathcal{L}_0 can be any supervised loss w.r.t to labeled nodes and \mathcal{L}_{reg} is a graph Laplacian regularisation term based on our assumption (e.g. that adjacent nodes are more likely to have the same label). Examples of such approaches include the label propagation algorithm (Zhu et. al [7]) and manifold regularisation (Belkin et al. [8]).

In the second approach, the problem consists of learning embeddings for the nodes using the local graph structure. The embeddings are then used in a semi-supervised setting to learn the corresponding class labels. For example, DeepWalk (Perozzi et al. [9]) learns embeddings using sampled random walks on the graph. This method is extended with more sophisticated random walk search schemes in LINE (Tang et al. [10]) and node2vec (Grover & Leskovec [11]).

1.3 Cora Dataset

To evaluate the different models I focus on the Cora dataset¹ [1]. The dataset consists of $N = 2078$ nodes representing machine learning publications. Each publication is classified into one of $C = 7$ classes: ‘Case Based’, ‘Genetic Algorithms’, ‘Neural Networks’, ‘Probabilistic Methods’, ‘Reinforcement Learning’, ‘Rule Learning’ or ‘Theory’. The features associated to each node are a binary bag-of-word feature vector for a vocabulary of $F = 1433$ unique words. Let \mathbf{X} be the $N \times F$ matrix containing the features of the nodes in its rows, that is the i^{th} row \mathbf{X}_i contains the bag-of-word binary vector for node i . The edges are directed and represent citations, where an edge $pub_1 \rightarrow pub_2$ means that publication pub_2 is cited in publication pub_1 . The papers were selected in such a way that every paper cites or is cited by at least one other paper.

In particular, I use the data split introduced by Kipf and Welling [2]. At training time the label of only 140 nodes is given (about 6.7% of the nodes). Let \mathbf{X}_L and \mathbf{y}_L be the matrices containing respectively the features and the labels of the labeled nodes. A further 500 nodes are split into a validation set and another 1000 nodes make up the test set. Furthermore, the edge orientations are ignored by constructing a symmetric adjacency matrix \mathbf{A} , where $\mathbf{A}_{ij} = 1$ if and only if there is an edge between publications i and j .

1.3.1 Node degree distribution

The distribution of node degrees is plotted in figures 1 and 2. Most nodes have few outgoing edges, with 59.9% of the nodes having 3 outgoing edges or less and over 96.5% having 10 or fewer. There seems to be one extreme outlier publication making 168 citations, whereas the second most citing paper makes only 78 citations.

1.3.2 Clusters

Next, I inspect how densely clustered the network is. TODO

¹<https://linqs.soe.ucsc.edu/data>

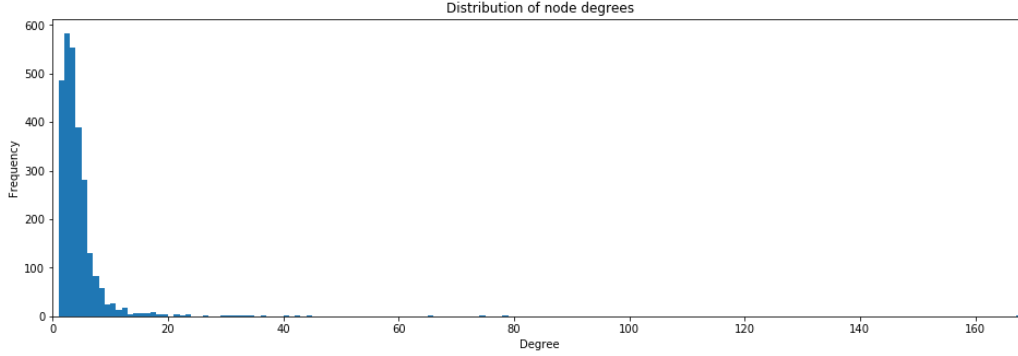


Figure 1: Distribution of node degrees in Cora dataset

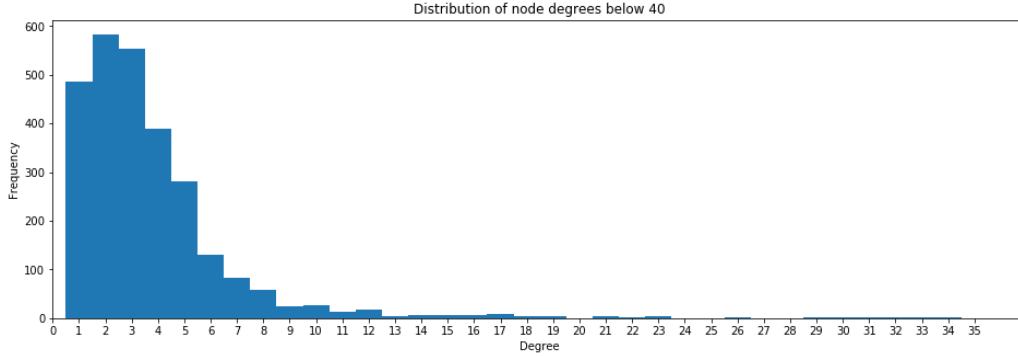


Figure 2: Distribution of node degrees for nodes with degree below 40. This includes 2701 of the 2708 publications, or above 99.7% of the nodes.

2 Methods

In this section I present the different methods I implemented to perform the classification task.

2.1 SVM on node features alone (node SVM)

The first method provides a simple baseline model and doesn't incorporate any information about the structure of the network. I regard the problem as a standard supervised learning task and train a support vector machine (SVM) on the features of the nodes only. Thus, given the binary bag-of-word vector representation of a publication, the task is to infer its publication type.

Using scikit-learn² [12] I train a SVM on the labeled training nodes. The trained model is then used to predict the class labels of the test publications.

2.2 SVM with neighbour features (average SVM)

For the second approach, I incorporate information about the graph structure into the SVM model. To do this I combine the features of the neighbours of a node by taking the average of their features and concatenating the resulting vector to the feature vector of the node itself. That is, we train the SVM on the features \mathbf{X}' where for each node i ,

$$\mathbf{X}'_i = [\mathbf{X}_i || \mathbf{X}_i^{neighbours}]$$

where

$$\mathbf{X}_i^{neighbours} = \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \mathbf{X}_j,$$

$\mathcal{N}(i)$ is the set containing the neighbour nodes of node i , and $[||\cdot]$ denotes vector concatenation. Note that \mathbf{X}' is a $N \times 2F$ matrix and thus contains twice as many entries as \mathbf{X} .

2.3 Multilayer perceptron (average MLP)

Next, I train a simple feedforward neural network on the simple concatenation of the node features with the average of the features of the neighbouring nodes. The network consists of two fully-connected layers. The first layer consists of 64 hidden units and has a ReLU activation function. The second layer consists of 7 output units, to which is applied a softmax activation to obtain a probability distribution over the 7 classes. Furthermore, as the number of training data is very small, I make heavy use of regularisation using i) Dropout [13] on both the input layer and on the output of the first hidden layer, and ii) L2 regularisation on the weights of the first hidden layer.

2.4 Graph Convolution Network (GCN)

The third algorithm I implemented on this dataset is the Graph Convolution Network, presented by Kipf and Welling [2]. Here, the graph structure is

²<http://scikit-learn.org/>

encoded directly using a neural network $f(\mathbf{X}, \mathbf{A})$ where \mathbf{X} is the usual feature matrix and \mathbf{A} is the adjacency matrix. The model is trained on a single supervised loss \mathcal{L}_0 . However, as the adjacency matrix is also inputed to the network, the gradient information is distributed back from \mathcal{L}_0 to all nodes and therefore representations are learnt for nodes with and without labels [2].

A Graph Convolutional Network (GCN) is a neural network with the following layer-wise propagation rule, where the activations in each layer $H^{(l)}$ are computed as

$$\mathbf{H}^{(l+1)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}^{(l)})$$

Here, $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_N$ is the adjacency matrix with added self-connections, hence adding the $N \times N$ identity matrix \mathbf{I}_N . $\tilde{\mathbf{D}}$ is the diagonal degree matrix of $\tilde{\mathbf{D}}$, that is $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$. $\mathbf{W}^{(l)}$ are the trainable neural network weights in layer l . Finally, $\sigma(\cdot)$ is an activation function.

For the Cora dataset, the authors present a 2 layer GCN with a ReLU activation in the first layer and a softmax activation in the second layer. The computation of the forward pass of the neural network can therefore be described concisely as

$$f(\mathbf{X}, \mathbf{A}) = \text{softmax}(\hat{\mathbf{A}} \text{ReLU}(\hat{\mathbf{A}} \mathbf{X} \mathbf{W}^{(0)}) \mathbf{W}^{(1)})$$

where $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ and can be precomputed in a pre-processing step. $\mathbf{W}^{(0)}$ is an $N \times H$ weight matrix of reals, where H is the number of nodes in the hidden layer and $\mathbf{W}^{(1)}$ is an $H \times C$ weight matrix, where $C = 7$ is the number of classes. The ReLU activation function just clips the input to 0,

$$\text{ReLU}(x) = \max(0, x)$$

and the softmax activation computes for each component x_i

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

During training the cross-entropy loss

$$\mathcal{L} = - \sum_{l \in \mathbf{y}_L} \sum_{c=1}^C \mathbb{1}(\mathbf{y}_L^{(l)} = c) \ln f(X, A)_{lc}$$

is minimised over all labeled examples \mathbf{y}_L using gradient descent. As the entire dataset fits in memory, a single batch containing the entire network is used. Furthermore, Dropout [13] is applied to the input features and to the output of the first layer to reduce overfitting.

2.5 Graph Attention Networks

The final method I explore is the Graph Attention Network (GAT) presented by Veličković et al. [3]. Attention mechanism have become widely popular and have achieved state-of-the-art results in sequence-based tasks, see for example Bahdanau et al. (2014) [14] or Vaswani et al. (2017) [15] as they have the advantage of giving the neural network the ability to learn on which part of the input to focus. More specifically, in the GAT case we talk about self-attention because the attention mechanism is used to compute a representation of a single sequence. The main idea behind the algorithm is to compute a hidden representation (also called an **embedding**) for each node in the graph by attending over all its neighbours using self-attention. That is, the network computes an embedding for each node using the nodes in its neighbourhood and by choosing how much attention to pay to each individual neighbour using the self-attention mechanism.

The GAT architecture is best described in terms of its layers, called **graph attention layers**. In each layer, the input features \mathbf{X} are first linearly transformed into $\mathbf{X}' = \mathbf{X}\mathbf{W}$ where \mathbf{W} is a $F \times F'$ weight matrix. Thus each node feature vector $\mathbf{h}_i = \mathbf{X}_i$ is linearly transformed to a feature vector $\mathbf{h}'_i = \mathbf{X}'_i = \mathbf{h}_i\mathbf{W}$ of length F' . Then, self-attention is performed on the nodes by applying a shared attentional mechanism defined as a function $a : \mathbb{R}^{F'} \times \mathbb{R}^{F'} \rightarrow \mathbb{R}$. This function maps two transformed input feature vectors to a real attention coefficient:

$$e_{ij} = a(\mathbf{h}'_i, \mathbf{h}'_j) = a(\mathbf{h}_i\mathbf{W}, \mathbf{h}_j\mathbf{W})$$

Intuitively, e_{ij} indicates how important the features of node j are to node i . To inject graph structural information, we only attend over the some neighbourhood $\mathcal{N}(i)$ for a given node i , defined to be the direct neighbours of i , including i itself. Hence, we have that

$$e_{ij} = \begin{cases} a(\mathbf{h}_i\mathbf{W}, \mathbf{h}_j\mathbf{W}) & \text{if } A_{ij} = 1 \text{ or } i = j \\ 0 & \text{otherwise} \end{cases}$$

Then, the coefficients are normalised using the softmax function:

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}(i)} \exp(e_{ik})}$$

We can do this for every pair of nodes i and j to obtain an $N \times N$ matrix α_{ij} of attention coefficients.

Here, the attention mechanism is defined to be a single-layer neural network, parametrised by weights \mathbf{a} and using a LeakyReLU activation:

$$\text{LeakyReLU}_\alpha(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{otherwise} \end{cases}$$

The attention coefficients can therefore be described concisely as

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{h}_i \mathbf{W} \parallel \mathbf{h}_j \mathbf{W}]))}{\sum_{k \in \mathcal{N}(i)} \exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{h}_i \mathbf{W} \parallel \mathbf{h}_k \mathbf{W}]))}$$

Next, we use the attention coefficients to compute a linear combination of the transformed features and apply a nonlinearity. This gives the final output of the layer

$$g_i = \sigma\left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} \mathbf{h}_j \mathbf{W}\right)$$

Furthermore, we can apply multi-head attention to stabilise the learning process. Instead of using a single attention mechanism, we use K_a independent attention heads, concatenating their output:

$$g_i = \left\|_{k=1}^{K_a} \sigma\left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij}^k \mathbf{h}_j \mathbf{W}^k\right)\right\|$$

Note that in this case the output feature vectors will have size $K_a F'$.

Veličković et al. propose a GAT network consisting of two graph attention layers. The first layer has $K_a = 8$ parallel attention heads and computes a hidden representation of features of size $F' = 8$. The nonlinearity used is the exponential linear unit ELU:

$$\text{ELU}_\alpha(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha(\exp(x) - 1) & \text{otherwise} \end{cases}$$

The second layer has C features and a single activation head $K_a = 1$. The nonlinearity used is the softmax activation.

To cope with the small training set and avoid overfitting L2 regularisation is applied to the layer weights and Dropout [13] is applied to the layers’ inputs and to the attention coefficients. Thus at each training epoch, a node is only exposed to a sample of its neighbourhood.

3 Results & analysis

Th results obtained with the various methods described above are shown in table 1. Next, I explain the optimisation of hyper-parameters and implementation details for each of the methods presented in more detail.

Method	Accuracy	
	validation	test
linear node SVM	0.558	0.583
RBF node SVM	0.566	0.576
linear average SVM	0.714	0.728
RBF average SVM	0.716	0.733
average MLP	0.762	0.769
GCN	0.785	0.802
GAT	0.769	0.794

Table 1: Validation accuracy and test accuracy obtained on the different models.

3.1 SVM on node features

Using the 500 validation nodes I optimise the hyper-parameters of the SVM by using a grid-search approach over a range of parameters. I consider both linear and a radial basis function (RBF) kernels. For the linear kernel, the only hyper-parameter, C , is taken at regular intervals on a logarithmic scale ranging from 2^{-15} to 2^3 . For the RBF kernel, C is similarly taken from a

logarithmic scale ranging from 2^{-15} to 2^5 . The second hyper-parameter, γ , is taken from the range 2^{-15} to 2^3 . The validation accuracy obtained for these hyper-parameters ranges are presented in figures 3 and 4. The validation accuracy is maximised at 56.6%, obtained with the RBF kernel, $C = 1.86$ and $\gamma = 2^{-7}$.

Although well below the current state-of-the-art accuracy, this provides a simple baseline accuracy on the test nodes of 57.6%. This is significantly better than a random baseline, which would have an expected accuracy of $\frac{1}{7}$ or about 14.3%. This suggests that, rather unsurprisingly, a lot of the information required for classification is contained in the nodes features, that is in the bag-of-word vectors.

3.2 Average features SVM

Using these new features \mathbf{X}' as presented in section 2.2, I train a SVM as before using the scikit-learn library. I consider both linear and a RBF kernels. The validation accuracies are presented in figures 5 and 6. Again, the RBF kernel performs slightly better with a validation accuracy of 71.6% compared to 71.4% for the linear kernel. Using the best hyper-parameters with the RBF kernel: $C = 5.411$ and $\gamma = 2^{-9}$, I obtain an accuracy of 73.3% on the test nodes.

This is a significant improvement over the accuracy obtain with the SVM trained on the features of each node individually. This shows that incorporating knowledge about the graph structure does indeed improve the classification accuracy significantly even when incorporated only in a very simplistic way, i.e. by only considering the direct neighbours and averaging their features.

3.3 Feedforward neural network with neighbour features

To find the optimal hyper-parameters, I again use a grid search method, validating the accuracy using the 500 validation nodes. The network is trained using every possible combination of parameters with $\text{batch_size} \in [32, 64, 128, 256, 2708]$, $\text{num_epochs} \in [50, 100, 200, 400]$, $\text{dropout_p} \in [0.0, 0.1, 0.2, 0.4, 0.6]$ and $\lambda \in [0.05, 0.005, 0.001, 0.0001]$ where λ is the l2 regularisation weight. The best validation accuracy of 76.2% was obtained using a batch size of 2708,

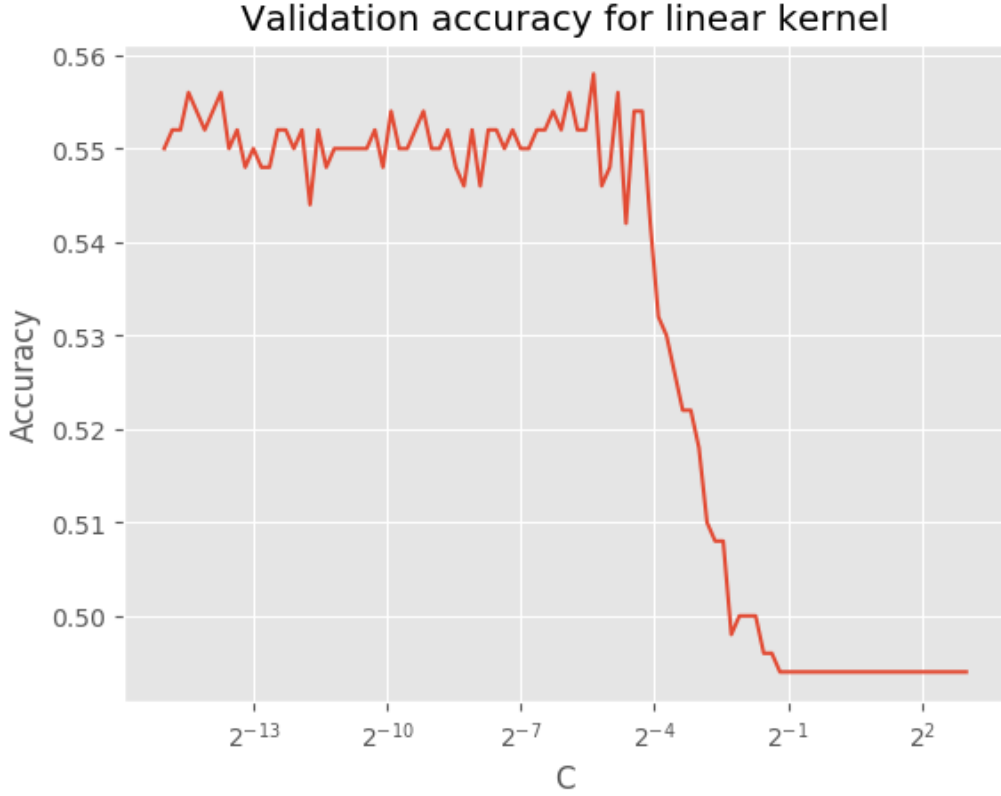


Figure 3: Validation accuracy obtained for various values of C for the SVM with a linear kernel trained on the node features alone. The best validation accuracy is obtain at $C = 0.024$, with a validation accuracy of 55.8%.

i.e. using the entire training set, training over 400 epochs with a dropout rate of 0.6 and setting $\lambda = 0.0001$. Using these hyper-parameters I obtain a test accuracy of 76.9%.

Say something about this being higher than the SVM...

3.4 Graph Convolution Network

I implemented the GCN model in TensorFlow and used the hyper-parameters presented by the authors which were obtained by optimising the validation accuracy on the set of 500 validation nodes. The model is trained over 200

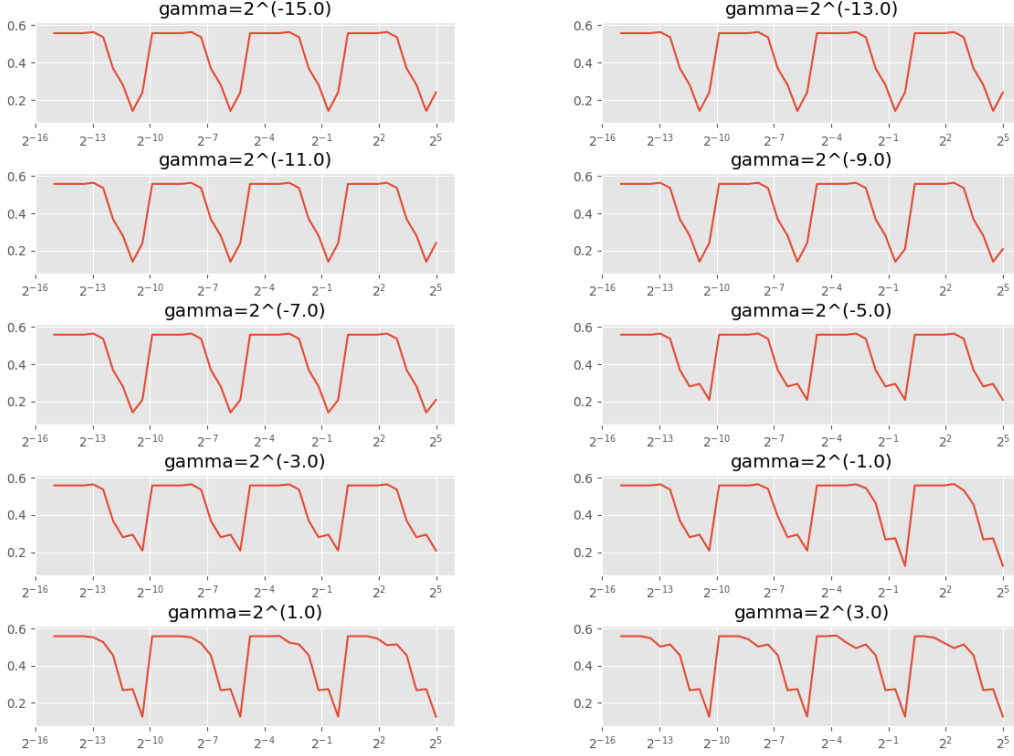


Figure 4: Validation accuracy obtained for various values of C and γ for the SVM with an RBF kernel trained on the node features alone. The best validation accuracy is obtain at $C = 1.86$ and $\gamma = 2^{-7}$, with a validation accuracy of 55.6%.

epochs using the Adam optimiser [16]. The optimal learning rate was found to be 0.01. The weights are initialised using a Glorot initialization [17]. Furthermore, the input feature vectors are row-normalized. The dropout rate was set to $p = 0.5$, the L2 regularisation weight to $\lambda = 0.0005$. The number of hidden inputs was set to 16. Averaging the result of 10 runs, I obtain an accuracy of 78.5% on the validation data and 80.3% on the test data.

The GCN is a significant improvement over the average MLP. This could be accounted for by multiple reasons. First, the graph structure is incorporated in a more principled way, using the adjacency matrix and the degree matrix of the network. Second, as the network consists of two layers, it has

Validation accuracy for linear kernel using averaged features

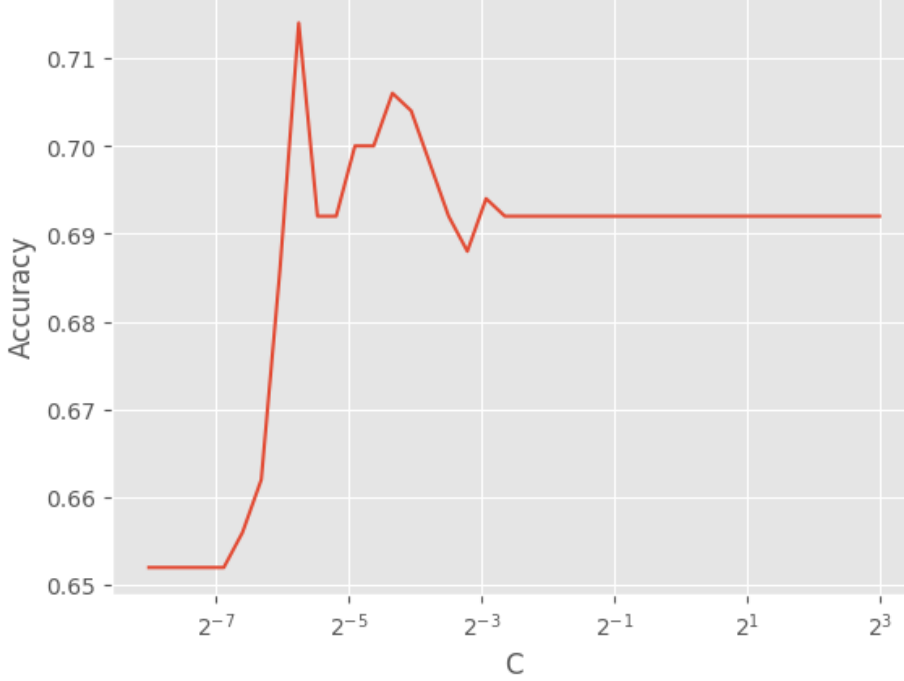


Figure 5: Validation accuracy obtained for various values of C for the SVM with a linear kernel trained on the node features and the average neighbouring features. The best validation accuracy is obtain at $C = 0.018$, with a validation accuracy of 7.14%.

the ability to incorporate information not only from the direct neighbours of a node but also from the neighbours of the neighbours of a node.

3.4.1 Hidden layer visualisation

The activations in the hidden layer of the 2 layer GCN can be visualised using a dimensionality reduction technique. Following the approach taken by Kipf and Welling, I use t-SNE [18] to reduce the 16-dimensional hidden activation to 2 dimensions. Figure 7 shows the resulting vectors for the training, validation and test datasets.

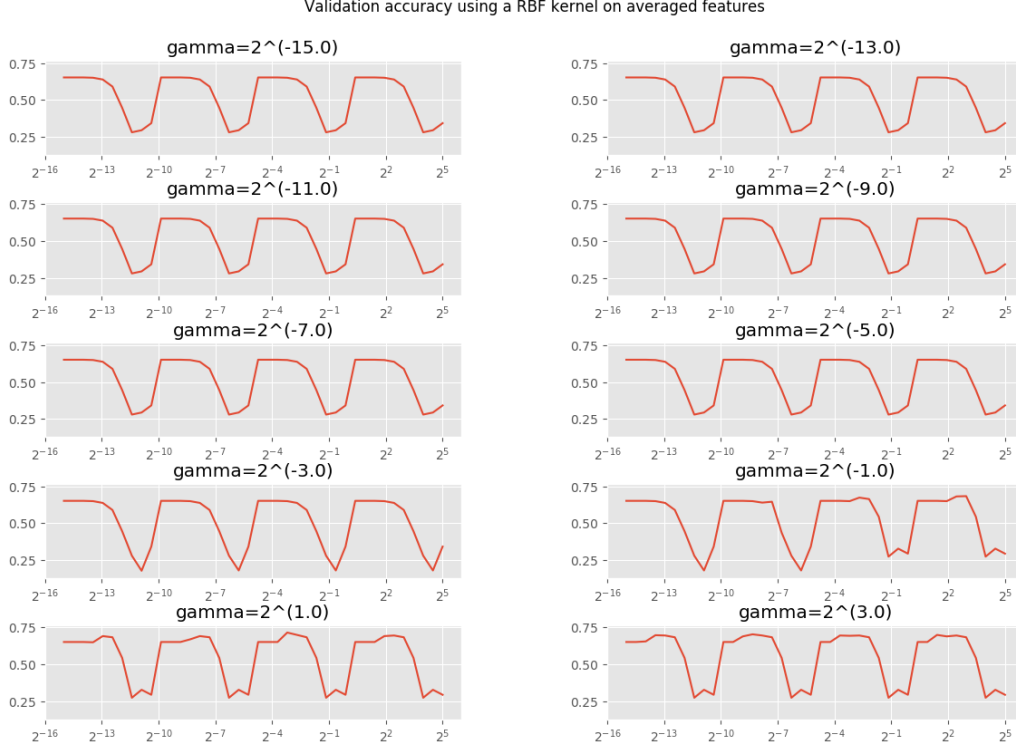


Figure 6: Validation accuracy obtained for various values of C and γ for the SVM with an RBF kernel trained on the node features and the averaged neighbouring features. The best validation accuracy is obtain at $C=5.411$ and $\gamma = 2^{-9}$, with a validation accuracy of 71.6%.

3.5 Graph Attention Network

As for the GCN, I use the hyper-parameters presented by the author but reimplement the network in TensorFlow. The regularisation weight is set to $\lambda = 0.0005$ and the dropout rate is set to $p = 0.6$. The Adam optimiser [16] is used with a learning rate of $lr = 0.005$. The model is trained over 100 epochs, each with a single batch containing all 2708 nodes.

Using these hyper-parameters I obtain a validation accuracy of 0.77% and a test accuracy of 0.79%. This is less the accuracy reported by the authors (83.0%). This discrepancy might be explained by small differences between the two implementations of the algorithm.

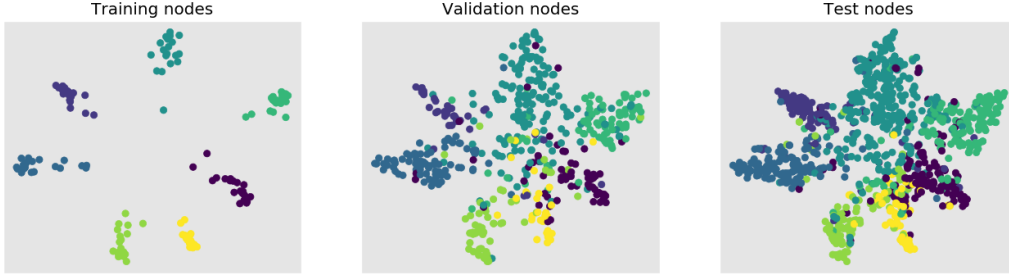


Figure 7: Activations in the hidden layer, after reducing their dimensionality using t-SNE for the training, validation and test nodes. As expected, the network is able to learn to discern the training examples well.

3.5.1 Learnt embeddings visualisation

Using t-SNE we can reduce the dimensionality of the learnt embeddings, i.e. the output of the first layer of the GAT. Reducing to 2 dimensions we can plot the embeddings in a plane, as shown in figure 8.

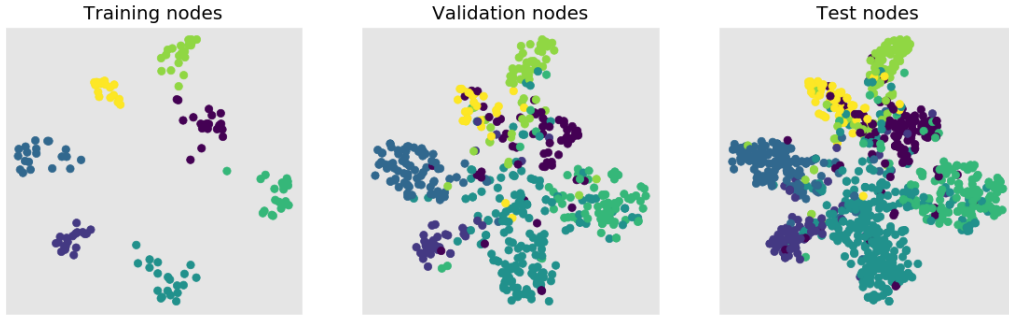


Figure 8: Visualisation of the learnt embeddings by the GAT. The dimensionality of the embeddings is reduced from 64 to 2 using t-SNE. The embeddings are shown for the training, validation and test nodes. As expected, the network is able to learn to discern the training examples well.

4 Conclusion

I implemented 5 different approaches to solve the task of transductive classification on the Cora dataset [1]. Experiments show that a simple baseline

method using the node features alone is sufficient to obtain significantly better than random accuracy. Furthermore, incorporating more information about the graph structure did increase the accuracy significantly. My implementation of the more sophisticated approaches, i.e. the GCN and the GAT obtained the best results, closely matching the current state-of-the-art results for this dataset.

References

- [1] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93, 2008.
- [2] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [3] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018. accepted as poster.
- [4] Benno Schwikowski, Peter Uetz, and Stanley Fields. A network of protein–protein interactions in yeast. *Nature biotechnology*, 18(12):1257, 2000.
- [5] Evelien Otte and Ronald Rousseau. Social network analysis: a powerful strategy, also for the information sciences. *Journal of information Science*, 28(6):441–453, 2002.
- [6] Mark EJ Newman. The structure of scientific collaboration networks. *Proceedings of the national academy of sciences*, 98(2):404–409, 2001.
- [7] Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label propagation. 2002.
- [8] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of machine learning research*, 7(Nov):2399–2434, 2006.

- [9] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710. ACM, 2014.
- [10] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *Proceedings of the 24th International Conference on World Wide Web*, pages 1067–1077. International World Wide Web Conferences Steering Committee, 2015.
- [11] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864. ACM, 2016.
- [12] 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*, 12:2825–2830, 2011.
- [13] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [14] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 6000–6010, 2017.
- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [17] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth*

International Conference on Artificial Intelligence and Statistics, pages 249–256, 2010.

- [18] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.