

EEE 586: Final Report for the Term Project

Text Classification with Distil-BERT Enhanced GNNs

Tuna Alikasıfoğlu*, Arda Can Aras[†]

^{*†}*Dept. of Electrical and Electronics Engineering*

^{*†}*Bilkent University*

{^{*}t.alikasifoglu, [†]can.aras}@bilkent.edu.tr

Abstract—In recent years, the amount of text based complex documents increased significantly, along with the importance of ability to classify texts as efficiently and accurately as possible. We had traditional algorithms to tackle the text classification problem, but with the progressive computational power, many machine learning, especially deep learning, based solutions surpassed the human capabilities in this area. Lately, BERT based models overshadow the remaining approaches in the text classification area. On the other hand, we have the trending topic of graph neural networks (GNNs). They are geometric extensions, i.e., extensions of traditional neural network architectures to the graph domain, and graph-structured data. In recent years, there has been some developments in the text classification area, especially using graph convolutional networks (GCNs). In this project, a brief overview of text classification problem and of graph neural networks are provided. Text classification overview covers the fundamental steps of a text classification task with the indication of GNN integration, where the GNN overview provides the reasons, challenges and how-tos of utilizing GNNs. Then, we provide an overview of the related work that tackles the text classification problem with the GNNs, while comparing several approaches. Finally in this project, we present a text classification approach that combines the BERT models, that provide semantic and contextual information, with the GCN models, that provide structural and global information.[§]

Index Terms—Text Classification, Graph Neural Networks (GNNs), Graph Convolutional Networks (GCNs), BERT, Distil-BERT.

I. INTRODUCTION

In this project, we have scrutinize the classic natural language processing task of text classification along with the trending approach of graph

neural networks (GNNs). In this context, we present a text classification approach that combines the BERT models, that provide semantic and contextual information, with the GCN models, that provide structural and global information. In this project, we first provide a comprehensive related work analysis in [Section II](#). Then, in [Section III](#), we present our proposed method of combining BERT based approaches with GNN based approaches in the context of text classification. Then, we present our results and compare them with the baselines in [Section IV](#). Finally, we provide a discussion and possible future directions in [Section V](#).

We have divided the related work analysis into four parts in [Section II](#). First, we provide an overview regarding the traditional text classification ([Section II-A](#)) and we provide an overview to graph neural networks ([Section II-B](#)). Both of these overviews are provided based on several survey papers for text classification [[1](#)–[3](#)] and for graph neural networks [[4](#)–[7](#)]. In [Section II-A](#), we provide a basic definition for text classification, then we disintegrate text classification process to five steps to analyse the overall process, and we point out where the integration of graph neural networks can be made. In [Section II-B](#), we provide an introduction to graph neural networks. Then, we analyse why we need graph neural networks, we provide several challenges on the route of obtaining graph based deep learning frameworks, and finally we provide bare-minimum steps in order to obtain a graph neural network based architecture, based on the aforementioned surveys, in this subsection. After this step, we are at a stage that we have provided

[§]The source code of the project is provided in its [GitHub Page](#)

background information both on the preliminary aspects of the overall term project, which are text classification and graph neural networks.

In **Section II-C**, we provide previous work directly related to our own topic of graph neural networks related to solution of the classical natural language processing task of text classification. **Section II-C** of the survey aims to investigate the related and previous work in the text classification with graph neural network field possibly with a historical order. All the papers mentioned in this section related with each other by references. Surprisingly, this structure can also be viewed as graph where each paper is node and edges of the graphs as citation. There are also studies on this topic to predict structured citation trend [8]. Most of the papers in the **Section II-C** did not directly proposed to solve text classification task. However, nearly all of the papers presented **Section II-C** used text classification bench mark data sets to evaluate model performance. Finally, in **Section II-D**, we present the similar approaches that combine BERT architecture with GNNs for text classification [9]–[13].

In **Section III**, we present our approach, and give further details that how we implemented the proposed method, by also comparing with the related work discussed in **Section II**. We provide how we are generating both BERT & GCN embeddings and how we are combining these embeddings to generate a better representation for documents to obtain a text classification model. Then, we present the results of the proposed model in **Section IV**, and compare them with the baselines discussed in **Section II**.

Finally, in **Section V**, we provide a discussion that we comment on the obtained results, issues that we had encountered, and possible future directions for the project.

II. RELATED WORK

A. Text Classification

In [2], *text classification* (text categorization) is defined as the “procedure of designating predefined labels for the text”. The task is to assign labels or tags to the text based knowledge, i.e., textual units such as sentences, paragraphs and documents, where the labels are usually defined by humans, but can also be defined by the machine. This task is

a fundamental part of Natural Language Processing (NLP), and it is significant to its applications such as sentiment analysis, question answering, text summarization, etc.. Text classification task can be partitioned into five phases as preprocessing, feature extraction, dimensionality reduction (optional), classifier selection and evaluation:

1) *Preprocessing*: Text preprocessing is a crucial prerequisite for a successful feature extraction, and summarized in [1] as follows. The input of the text classification frameworks consists of raw text data, which are in the form of a sequence of sentences. In this step, “cleaning” of the text datasets is performed to transform the data into a form that is suitable for feature extraction. The cleaning process is usually performed by tokenization, capitalization, slang and abbreviation handling, noise removal, spelling correction, stemming and lemmatization.

2) *Feature Extraction*: After preprocessing step, another crucial step, feature extraction step is necessary. In [1], this step explained as follows. Two common methods of text based feature extraction are weighted word and word embedding techniques. In the weighted word aspect, we have old techniques like bag-of-words and term frequency-inverse document frequency (TF-IDF). In the relatively recent aspect, we have the word embedding techniques like *word2vec*, *GloVe*, *FastText*, etc.

3) *Dimensionality Reduction*: The dimensionality reduction is an optional step of a text classification task, but based on the size of the dataset, it may be a must to have a computable result. In this aspect of the task, we try to reduce the dimensionality of the feature space while preserving the information of the original features space. Some possible dimensionality reduction techniques provided in [1] include (principal / independent) component analysis, linear discriminant analysis, non-negative matrix factorization, random projection, autoencoder and stochastic neighbor embedding.

4) *Classifier Selection*: As it is stated in [2], selecting the optimal classifier is the most important aspect of a text classification task. Currently we have both traditional and deep learning oriented classifiers. The traditional classifiers are based on the statistical analysis of the training data, and the deep learning classifiers are based on the neural networks. The main distinction between the tra-

ditional and deep learning based approaches can be described as follows: Good feature extraction methodology is crucial for the traditional classifiers. They obtain sample features by artificial methods and then make classifications based on these features. Hence, the performance of the traditional classifiers are mainly restricted by feature extraction. On the other hand, by making feature mapping via nonlinear transformations a part of the learning process, deep learning based classifier selection can integrate feature extraction aspect into the model fitting process.

Examples of both traditional and deep learning based approaches are provided in [1]–[3]: Some traditional classifiers are logistic regression, (kernel) support vector machine, Naive Bayes, k -nearest neighbors, decision tree, random forest, etc. On the other hand, the deep learning classifiers are usually based on the neural networks, such as deep feed forward neural networks, convolutional neural networks (CNNs), recurrent neural networks (RNNs), lately attention and transformer based models such as BERT [14] variations and fine tuning pre-trained language models [15], and finally what we will focus on, graph neural network (GNN) based models.

5) *Evaluation*: Evaluation is step that we understand how the our model performs under the given text classification task. As it is provided in [2], [3], there are several evaluation metrics that can be used to evaluate the performance of a supervised technique. The most common metrics are accuracy, F_β -score, micro/macro-averaging. Although we also have metrics like Matthews correlation coefficient and receiver operating characteristics (ROC). In order to evaluate the performance of our model, based on the provided techniques, we need to use labeled data, i.e., we need benchmark datasets like GLUE [16], TweetEval [17], among others.

B. Graph Neural Networks

In recent years, deep learning based solutions surpassed any approach on machine learning tasks such as image classification, video processing, speech recognition and natural language processing. In these tasks, the underlying data are usually represented in the Euclidean domain. However, each day the amount of non-Euclidean data increases, which are represented as represented by graphs to

capture the underlying the complex relationships and interdependency between objects. Therefore, a need for deep learning methods that can manage graph structured data has emerged. In this context, the graph neural networks (GNNs) are born, and many of the deep learning approaches are converted to graph domain such as recurrent GNNs, convolutional GNNs (ConvGNNs) or graph convolutional networks (GCNs), graph autoencoder (GAE), graph reinforcement learning (GRL), graph adversarial methods and spatial-temporal GNNs, as summarized in [4]–[7].

1) *Reasons to use GNNs*: Hidden patterns residing under Euclidean data can be effectively obtained by traditional deep learning techniques. However, the increasing number of applications based on a non-Euclidean data structure enforces the necessity of graph based solutions. In this aspect, the following examples can be used to illustrate the benefits of having a graph based deep learning framework [5]:

- In e-commerce, highly accurate recommendation system can be achieved by using graph based deep learning techniques, since the interactions between users and products are a textbook example of graph structured data.
- For drug discovery in chemistry, we need to obtain the bioactivity of the molecules, where the molecules are modeled as graphs.
- Categorization of articles in a citation network, where the articles are linked to each other via “citationships”, i.e., forming a graph structure.

2) *Challenges to use GNNs*: In order to have a graph domain deep learning framework, we need to overcome several challenges imposed by the complexity of the graph data. Due to the nature of graphs, when they are compared with Euclidean data, they can be irregular, they can have unordered nodes with different number of neighbors. Hence, many basic operations defined in Euclidean domain can be challenging to apply to the graph domain, e.g., convolution operation. In addition, one of the fundamental assumption we have in the existing machine learning algorithms is that the instances are independent of each other, although this assumption is not valid for graph data since each instance (node) is related to others by links of various types [5]. Some of the main challenges can be categorized as follows [6]:

a) Irregular structures of graphs: We have the *geometric deep learning problem* which is the inability to define basic operations like convolution and pooling in the graph domain, which are essential aspects of traditional CNNs.

b) Heterogeneity and diversity of graphs: We have many different properties that a single graph can have: graphs can be homogenous or heterogeneous, they can be weighted or unweighted, they can be directed or undirected, and they can be signed or unsigned. Furthermore, the tasks may consist of node-level problems such as node classification, link prediction or they can consist of graph-level problems such as graph classification or graph generation. Therefore, we need a spectrum of architectures to tackle all these problems one-by-one.

c) Large-scale graphs: As in the case of e-commerce and social networks, graph structured data can have a large number of nodes and edges. However, we still need appropriate algorithms to work on the graph structure without increasing the computational and time complexity too much.

d) Incorporating interdisciplinary knowledge: Graph structured data sometimes traces back to other disciplines such as biology, chemistry and social sciences. The interdisciplinary nature helps to leverage domain knowledge to solve specific problems, but it can also complicate model designs. For the case of molecular graph generation, the chemical constraints and the generation’s objective function are often non-differentiable. Hence, gradient based training methods are out of the picture.

3) *Ways to use GNNs:* In [4], a general design to pipeline of GNNs is proposed. The following steps are necessary to obtain a graph-based deep learning framework:

a) Finding a graph structure: Based on the application in hand, we need to find out the underlying graph structure. There are two possibilities. First one is that we have an explicit graph structure, in the application such as social network, physical system or knowledge graph. The other possibility is that the underlying graph is implicit, and we need to build the graph from the task, such as obtaining a fully-connected “word” graph for text or obtaining a scene graph for an image. Then, we can obtain an

optimal GNN model for the the graph we obtained either from explicit information or from the task.

b) Design a loss function: Based on the task in hand and the training setting, a loss function needs to be determined, the loss function can be node-level, edge-level or graph-level, depending on the training setting of supervised, semi-supervised or unsupervised learning.

c) Build model using computational modules: Finally, we need computational modules to build and train our model. Based on the definition provided in [4], we need a module to conduct convolution and recurrent operations to propagate information between nodes to capture the underlying feature and topological information. We need a sampling module, and we need a pooling module. With the combination of these modules a typical architecture of GNN model can be built.

C. Text Classification with GNN

Some of the earliest success achieved on deep learning with graphs relied on finding proper ways to embed nodes into vectors using an encoder function [18]. One question arises from that definition is the “What is a good representation?”. We want these nodes embeddings to preserve interesting structures of the graphs. There are unsupervised graph representations learning algorithms like *node2vec* [19], *DeepWalk* [20] and *LINE* [21] which are trained prior to graph neural networks. These algorithms aimed to learn representative embeddings for nodes to preserve interesting structures of the graphs

Aforementioned algorithms inherently capture local similarities. Further studies find that Convolutional Graph Neural Networks (ConvGNNs) summarizes local patches of the graphs and shows that neighboring nodes tends to highly overlap [18]. Therefore, a ConvGNNs enforce similar features for neighboring nodes by its nature without needing pre-training for node embeddings. This phenomenon was also mentioned in Text Graph Convolutional Network (Text GCN) [22]. Results of this paper on multiple benchmark data sets demonstrate that a vanilla Text GCN without any external word embeddings or knowledge outperforms state-of-the-art methods for text classification. On the other

hand, Text GCN also learns predictive word and document embeddings jointly.

In [22] they evaluate Text GCN on two experimental tasks. First they seek the answer of whether their model achieve satisfactory results in text classification, even with limited labeled data and then they test whether their model learn predictive word and document embeddings. They compare their method with several state-of-art models. The suggested Text GCN may produce high text classification results and train predictive document and word embeddings, according to the experimental results. However, a major limitation of this study is that the GCN model is inherently transductive [22], in which test document nodes (without labels) are included in GCN training. Thus Text GCN could not quickly generate embeddings and make prediction for unseen test documents.

To improve the weakness of Text GCN in text classification task, [23] and [24] was mentioned in future work section of [22]. In [23] a novel algorithm Graph Attention Networks (GATs) was proposed. It was mentioned in [23] that GATs are new convolutional-style neural networks that operate on graph-structured data and use masked self-attentional layers. The graph attentional layer used in these networks is computationally efficient. Attentional layers do not require expensive matrix operations and they are parallelizable across all nodes in the graph. This structure allows for (implicitly) assigning different importance to different nodes within a neighborhood while dealing with different sized neighborhoods, and does not require knowing the entire graph structure. The experimental results yields that their attention-based models outperformed or matched state-of-the-art performance in four well-known node classification benchmarks, both transductive and inductive tasks [23].

Fast Graph Convolutional Neural Network [24] was also mentioned in the future work section of [22]. In [24] it was mentioned that, GCN in [25] represented as a useful graph model for semi-supervised learning. This model was created with the intention of being taught with both training and test data. Furthermore, for training with large, dense graphs, the recursive neighborhood expansion across layers faces time and memory issues. [24] interpret graph convolutions as integral transforms

of embedding functions under probability measures to relax the condition of simultaneous availability of test data. As a result of this interpretation, Monte Carlo techniques may be used to consistently estimate the integrals, leading to a batched training scheme like FastGCN, which is proposed [24].

After further development on top of Text GCN, Simplifying Graph Convolutional Networks [26] was proposed to overcome unnecessary complexity and redundant computation in the previous work of Fast GCN and GATs. GCNs and their variants have received a lot of attention and have become the defacto methods for learning graph representations. GCNs are primarily inspired by modern deep learning methodologies, and as a result, they may inherit extra complexity and redundant processing. In [26] they eliminate the unnecessary complexity in this paper by reducing non-linearities one by one and collapsing weight matrices between layers. The resulting linear model is theoretically analyzed and shown to correspond to a fixed low-pass filter followed by a linear classifier in. The test results in [26] shows that these simplifications have no negative influence on accuracy in a wide range of downstream applications. Furthermore, the resulting model scales to bigger datasets, is intuitively interpretable, and outperforms FastGCN by up to two orders of magnitude.

D. BERT-GNN Architectures

After coming up with the idea of combining BERT models with GNN based models to increase text classification performance, we have encountered a few similar approaches that tries to accomplish a similar task. The first significant approach in combining BERT and GNN architectures is the *VGCN-BERT* proposed in [9]. However, in this approach the authors generate a vocabulary graph to produce word embeddings using the GCN architecture, then they supply these embeddings as input to the BERT architecture. This approach is different than what we are trying to achieve, since we are proposing aggregation of these embeddings since they both embody different aspects of the information.

In addition to this approach, in [10], *BEGNN* model is proposed, which aggregates graph embeddings with BERT embeddings similar to out

approach, although they are generating graph embeddings for each document separately, so graph embeddings that are used in BEGNN are limited to the global information of that specific document. This approach can be limited, since with graph embeddings we try to convey the global and structural information of the texts, and we propose an extension to this manner by generating graph embeddings using whole training set that can embody structural and global information in a more generic sense.

After our project proposal is finalized, there has been several publications regarding the text classification with combination of GNN and BERT architectures [11]–[13]. In these publications, there are some similar approaches to our project, although we still have some differences in our methods, when it is compared to these fresh publications.

III. METHODS

A. General Structure

We propose a model to obtain structural and semantic embeddings for each of the documents and then use this information to make classification. Graph Neural Network is used to obtain structural embedding and Distil-BERT is used for retrieving semantic embeddings from text. Then these two embeddings combined with 3 different approaches that mentioned in Section III-G. Finally, the obtained document embedding passed to classification layer to get predictions for each document. The Fig. 1 is the summary of general architecture.

B. Dataset Description

The dataset of this term project is *20 News Group* (20NG) [27] from huggingface.co. Dataset contains 18,846 documents evenly categorized into 20 different categories. In total, 11,314 documents are in the training set and 7,532 documents are in the test set.

C. Graph Generation from Data

To represent documents and words in a graph structure, we identify the vocabulary for the whole dataset without explicitly dividing it into test and training. For GCN structure, PyG asks user to pass whole dataset at once and then later identify the train and test indices. Therefore, we treat whole dataset as single corpus and obtained vocabulary

out of it. Later, we put each document and words in the nodes and represent their connections with adjacency matrix \mathbf{A} . Also each node can have an embedding in arbitrary dimension. To not effect learning of the structural features, we did not initialize these node embeddings by using popular algorithms like [28], [29], since they enforce node to start with a semantic information. Therefore, we used one-hot vector representation for both document and word embeddings. It yields the matrix $\mathbf{X} \in \mathbb{R}^{(n \times n)}$, where $n = n_{doc} + |V|$. We define the entries of the adjacency matrix as follows:

$$A_{ij} = \begin{cases} \text{PMI}(i, j) & i, j \text{ are words, } \text{PMI}(i, j) > 0 \\ \text{TF-IDF}_{ij} & i \text{ is document, } j \text{ is word} \\ 1 & i = j \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

D. Learnable Adjacency Matrix

We also propose learnable adjacency matrix to have further find interesting structures in data. We define new adjacency matrix and tune the learnable part of it by α . It can be given as follows:

$$\tilde{\mathbf{A}} = \alpha \hat{\mathbf{A}} + (1 - \alpha) \mathbf{A} \quad (2)$$

where we took gradient of loss function only with respect to $\hat{\mathbf{A}}$ and $\alpha \in (0, 1)$.

E. GNN Embeddings

To implement Graph Neural Network, [PyTorch Geometric](#) library have been utilized. The aforementioned learnable adjacency matrix $\tilde{\mathbf{A}}$ have been utilized instead of default \mathbf{A} . Generic equations of the GNN structure in our model can be given as follows:

$$\mathbf{L}^{(j+1)} = \rho(\tilde{\mathbf{A}} \mathbf{L}^{(j)} \mathbf{W}_j), \quad (3)$$

Our model has 2 layers:

$$\mathbf{Z} = \text{softmax} \left(\tilde{\mathbf{A}} \text{ReLU}(\tilde{\mathbf{A}} \mathbf{X} \mathbf{W}_0) \mathbf{W}_1 \right) \quad (4)$$

and the cross-entropy error over all labeled documents:

$$\mathcal{L} = - \sum_{d \in \mathcal{Y}_D} \sum_{f=1}^F Y_{df} \ln Z_{df} \quad (5)$$

Where \mathcal{Y}_D is the set of documents indices that have labels and F is the dimension of the output features.

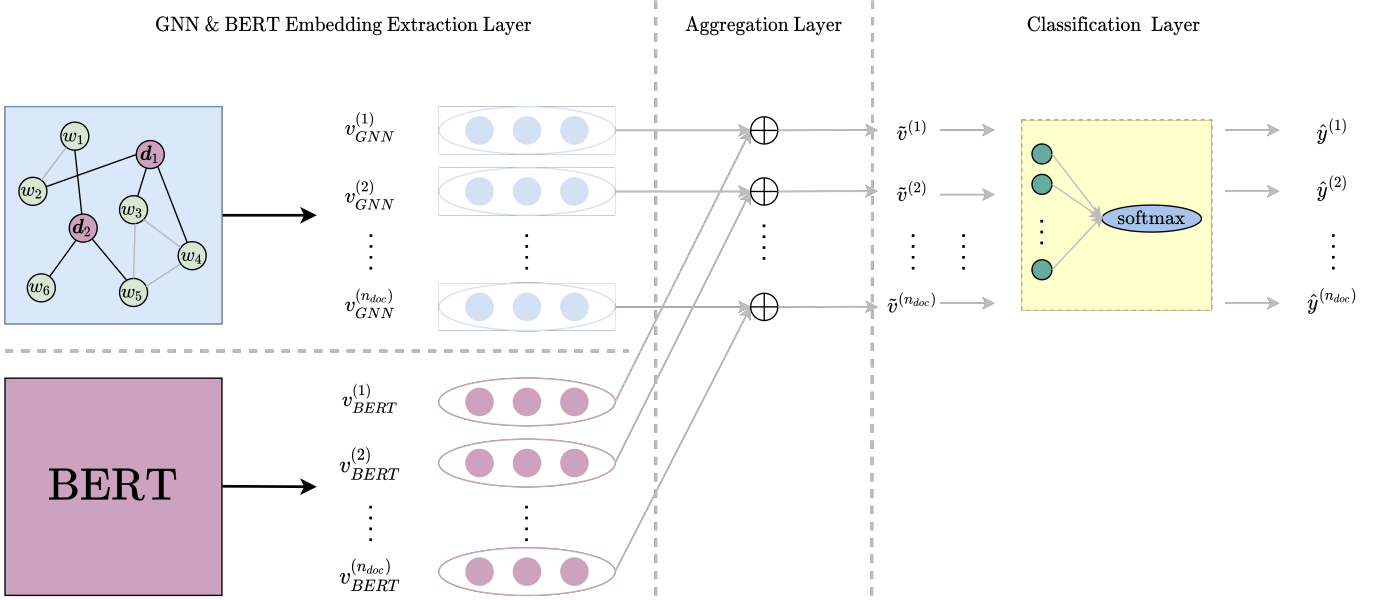


Fig. 1: Model Architecture

Finally, we can obtain embeddings from the layer equations as follows:

$$E_1 = \tilde{A}XW_0 \quad (6)$$

$$E_2 = \tilde{A} \text{ReLU}(\tilde{A}XW_0)W_1 \quad (7)$$

Where $W_0 \in \mathbb{R}^{(n \times 200)}$ and $E_1 \in \mathbb{R}^{(n \times 200)}$ which gives 200 dimensional representation for both documents and words. We used those 200 dimensional embedding vector and represented as v_{GNN} .

F. Distil-BERT Embeddings

To obtain embeddings from Distil-BERT, we fine-tuned the pre-trained model on 20NG dataset. Default parameters have 13 stacked attention layers so it has a output of dimension (13, 768) for each document. There is no pre-defined way to obtain a single 768 dimensional vector from that stacked representation. We choose the 13rd head of the attention layer as document representation. It is also possible to choose different head of the attention or one can also take element wise max of each layer etc. The obtained 768 dimensional document embeddings will be represented as v_{BERT} .

G. Aggregation of Embeddings

There are several possible ways to aggregate the embeddings retrieved from GNN and Distil-BERT. We will show the structure of each of them.

Each of the document embedding vector will be represented as \tilde{v}_{doc} . The dimension of the \tilde{v}_{doc} will be generically represented as d_{doc}

1) *Concatenation:* We simply concatenate v_{BERT} at the end of v_{GNN} to obtain 968 dimensional document representation.

$$\tilde{v}_{doc} = [v_{GNN} || v_{BERT}], \tilde{v}_{doc} \in \mathbb{R}^{(968 \times 1)} \quad (8)$$

2) *Element-wise Sum:* To sum the embeddings, they must have same dimension. Therefore, we apply PCA to v_{BERT} to reduce its dimension to 200. Then we simply sum them and obtain \tilde{v}_{doc} as follows:

$$\tilde{v}_{BERT} = \text{PCA}_{200}\{v_{BERT}\} \quad (9)$$

$$\tilde{v}_{doc} = v_{GNN} + \tilde{v}_{BERT}, \tilde{v}_{doc} \in \mathbb{R}^{(400 \times 1)} \quad (10)$$

3) *Trade-Off:* The trade-off version of the embedding approach is aimed to control the contribution of GNN and BERT embedding with trainable parameter λ . It can be done for both summation and concatenation strategies \tilde{v}_{doc} is given as follows:

$$\tilde{v}_{doc} = [\lambda v_{GNN} || (1 - \lambda)v_{BERT}] \in \mathbb{R}^{(968 \times 1)} \quad (11)$$

$$\tilde{v}_{doc} = \lambda v_{GNN} + (1 - \lambda)\tilde{v}_{BERT} \in \mathbb{R}^{(400 \times 1)} \quad (12)$$

H. Classification Layer for Document Embeddings

As a final step, we passed \tilde{v}_{doc} to classification layer to obtain accuracy results. The classification layer equation can be easily given as follows:

$$z_{doc} = \text{softmax}(\mathbf{W}\tilde{v}_{doc}), \mathbf{W} \in \mathbb{R}^{n_{doc} \times d_{doc}} \quad (13)$$

IV. RESULTS

The result section will have 4 different parts. In the first part we will investigate the performance of GNN structure by its own. The second part will include the performance of Distil-BERT only. Individual results explanations will help us to compare their performance when they are combined with aggregation layer. At the third part we will consider the classification results of obtained \tilde{v}_{doc} from the [Section III-G](#) section. We will investigate the each aggregation setting mentioned in [Section III-G](#). As a last part, we will investigate the performance of our model when it is compared with the different SOTA algorithms on the 20NG dataset.

A. GNN Results

In this part, we only train the GNN part of the algorithm and get the prediction results as a by product of embedding extraction procedure. We will use the best results from [Table I](#) to compare with SOTA results. For the GNN results with learnable adjacency matrix, we only use the best $\alpha = 0.13$ result for convenience. For convenience, we used some naming in our models. These are as follows: $\text{GCN}_{i,j,k} \triangleq \text{GCN}$ with i, j, k hidden layers and $\text{L-GCN}_{i,j,k} \triangleq \text{GCN}$ with i, j, k hidden layers using learnable adjacency matrix structure.

TABLE I: GNN Results

| Models | Train Accuracy (%) | Test Accuracy (%) |
|------------------------------|--------------------|-------------------|
| $\text{GCN}_{200,20}$ | 100 | 66.50 |
| $\text{L-GCN}_{200,20}$ | 100 | 67.50 |
| $\text{GCN}_{2000,200,20}$ | 100 | 60.80 |
| $\text{L-GCN}_{2000,200,20}$ | 100 | 60.70 |

B. BERT Results

To train Distil-BERT on our dataset, we used the [transformers](#) library. It allows us to further fine-tune our model on pre-trained Distil-BERT. Prediction results of the Distil-BERT model are in [Table II](#).

TABLE II: Distil-BERT Training Results

| Epoch | Training Loss | Validation Loss | Accuracy (%) |
|----------|---------------|-----------------|--------------|
| 1 | 1.1764 | 1.217 | 65.64 |
| 2 | 0.7607 | 1.174 | 68.12 |
| 3 | 0.5118 | 1.440 | 67.90 |
| \vdots | \vdots | \vdots | \vdots |
| 28 | 0.0684 | 3.372 | 69.98 |
| 29 | 0.0817 | 3.397 | 70.05 |
| 30 | 0.0799 | 3.404 | 69.99 |

C. Combined \tilde{v}_{doc} Results

The given combined results are find by using the method in [Section III-G1](#). This method gave the best results for all models. Therefore, only this aggregation approach results are mentioned.

TABLE III: GNN+BERT Combined Results

| Models | 20NG Test Accuracy (%) |
|--|------------------------|
| $\text{GCN}_{200,20}$ + Distil-BERT | 67.20 |
| $\text{L-GCN}_{200,20}$ + Distil-BERT | 69.70 |
| $\text{GCN}_{2000,200,20}$ + Distil-BERT | 64.90 |
| $\text{L-GCN}_{2000,200,20}$ + Distil-BERT | 63.20 |

D. Comparison with SOTA Algorithms

The results of the SOTA algorithms for the 20NG dataset along with our results are provided in [Table IV](#). The models are sorted according to their accuracies, and our results are provided in [color](#).

TABLE IV: SOTA Results

| Models | 20NG Test Accuracy (%) |
|--------------------------------------|------------------------|
| PV-DM | 51.10 |
| LSTM | 65.70 |
| $\text{L-GCN}_{200,20}$ | 67.60 |
| $\text{L-GCN}_{200,20}$ +Distil-BERT | 69.70 |
| Our Distil-BERT | 70.05 |
| RoBERTa | 83.80 |
| BERT | 85.30 |
| TextGCN | 86.30 |
| RoBERTaGAT | 86.50 |
| BertGAT | 87.40 |
| SGC | 88.50 |
| BertGCN | 89.30 |
| RoBERTaGCN | 89.50 |

V. DISCUSSION & CONCLUSION

A. Obtained Results

Unfortunately, obtained results did not surpass the SOTA algorithms that mentioned in [Section IV-D](#). However, we are aware of the problems that occurred during our implementation. These issues are mentioned in the [Section V-B](#). Our model still can pass some of the other algorithms.

B. Issues

During our implementation of the algorithm, we have discovered several problems and solved most of them. The very first problem that we faced with was the instability of the 20NG dataset in different websites. This dataset is used in both supervised and semi-supervised manner in the literature. The graph convolutional network proposed in [\[25\]](#) is for semi-supervised learning problem. [\[22\]](#) also proposed for only semi-supervised learning problem. [\[22\]](#) treats 20NG dataset as a semi-supervised learning problem. In that version of 20NG, they do not use labels of every data instance. However, in [huggingface.co](#), all of the data are labeled. In our structure, we tried to treat problem in inductive manner since we know all of the labels. However, we used the GCN architecture of the [PyTorch Geometric](#) which only accepts the semi-supervised learning tasks and it asks to pre-define the labeled data to do back propagation only on these labeled data. Since we know all of the labels, we pre-define the labeled data as whole training set.

C. Future Directions

To solve the mentioned problems in [Section V-B](#), we can treat whole problem in inductive manner and use more appropriate GNN structures. One of the candidate for inductive learning problem is [\[24\]](#). It solves the problem totally in inductive manner which also makes computation faster. Another candidate architecture for the GNN side is [\[30\]](#). It also enables us to learn more interesting features since it uses attention mechanism rather simple matrix multiplications in [Eq. \(3\)](#).

REFERENCES

- [1] Kowsari, J. Meimandi, Heidarysafa, Mendu, Barnes, and Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, Apr. 2019, ISSN: 2078-2489. DOI: [10.3390/info10040150](#). [Online]. Available: <http://dx.doi.org/10.3390/info10040150>.
- [2] Q. Li, H. Peng, J. Li, *et al.*, "A survey on text classification: From shallow to deep learning," *CoRR*, vol. abs/2008.00364, 2020. arXiv: [2008.00364](#). [Online]. Available: <https://arxiv.org/abs/2008.00364>.
- [3] S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad, M. Chenaghlu, and J. Gao, "Deep learning based text classification: A comprehensive review," *CoRR*, vol. abs/2004.03705, 2020. arXiv: [2004.03705](#). [Online]. Available: <https://arxiv.org/abs/2004.03705>.
- [4] J. Zhou, G. Cui, S. Hu, *et al.*, "Graph neural networks: A review of methods and applications," *AI Open*, vol. 1, pp. 57–81, 2020, ISSN: 2666-6510. DOI: [https://doi.org/10.1016/j.aiopen.2021.01.001](#). [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2666651021000012>.
- [5] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 4–24, Jan. 2021, ISSN: 2162-2388. DOI: [10.1109/tnnls.2020.2978386](#). [Online]. Available: <http://dx.doi.org/10.1109/TNNLS.2020.2978386>.
- [6] Z. Zhang, P. Cui, and W. Zhu, "Deep learning on graphs: A survey," *CoRR*, vol. abs/1812.04202, 2018. arXiv: [1812.04202](#). [Online]. Available: <http://arxiv.org/abs/1812.04202>.
- [7] L. Sun, J. Wang, P. S. Yu, and B. Li, "Adversarial attack and defense on graph data: A survey," *CoRR*, vol. abs/1812.10528, 2018. arXiv: [1812.10528](#). [Online]. Available: <http://arxiv.org/abs/1812.10528>.
- [8] D. Cummings and M. Nassar, "Structured citation trend prediction using graph neural networks," *CoRR*, vol. abs/2104.02562, 2021. arXiv: [2104.02562](#). [Online]. Available: <https://arxiv.org/abs/2104.02562>.
- [9] Z. Lu, P. Du, and J. Nie, "VGCN-BERT: augmenting BERT with graph embedding for text classification," *CoRR*, vol. abs/2004.05707, 2020. arXiv: [2004.05707](#). [Online]. Available: <https://arxiv.org/abs/2004.05707>.
- [10] Y. Yang and X. Cui, "Bert-enhanced text graph neural network for classification," *Entropy*, vol. 23, 2021.
- [11] Y. Lin, Y. Meng, X. Sun, *et al.*, *Bertgcn: Transductive text classification by combining gcn and bert*, Mar. 2022. DOI: [10.48550/ARXIV.2105.05727](#). [Online]. Available: <https://arxiv.org/abs/2105.05727>.
- [12] X. She, J. Chen, and G. Chen, "Joint learning with bert-gcn and multi-attention for event text classification and event assignment," *IEEE Access*, vol. 10, pp. 27 031–27 040, 2022. DOI: [10.1109/ACCESS.2022.3156918](#).

- [13] F. Zeng, N. Chen, D. Yang, and Z. Meng, “Simplified-boosting ensemble convolutional network for text classification,” *Neural Processing Letters*, May 2022. DOI: 10.1007/s11063-022-10843-4. [Online]. Available: <https://doi.org/10.1007/s11063-022-10843-4>.
- [14] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” *CoRR*, vol. abs/1810.04805, 2018. arXiv: 1810.04805. [Online]. Available: <http://arxiv.org/abs/1810.04805>.
- [15] J. Howard and S. Ruder, “Fine-tuned language models for text classification,” *CoRR*, vol. abs/1801.06146, 2018. arXiv: 1801.06146. [Online]. Available: <http://arxiv.org/abs/1801.06146>.
- [16] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, “GLUE: A multi-task benchmark and analysis platform for natural language understanding,” arXiv preprint 1804.07461, 2018. [Online]. Available: <https://gluebenchmark.com/>.
- [17] F. Barbieri, J. Camacho-Collados, L. Espinosa-Anke, and L. Neves, “TweetEval: Unified Benchmark and Comparative Evaluation for Tweet Classification,” in *Proceedings of Findings of EMNLP*, 2020.
- [18] P. Veličković, *Theoretical foundations of graph neural networks*, 2021. [Online]. Available: <https://petar-v.com/talks/GNN-Wednesday.pdf>.
- [19] A. Grover and J. Leskovec, “Node2vec: Scalable feature learning for networks,” *CoRR*, vol. abs/1607.00653, 2016. arXiv: 1607.00653. [Online]. Available: <http://arxiv.org/abs/1607.00653>.
- [20] B. Perozzi, R. Al-Rfou, and S. Skiena, “Deepwalk,” *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, Aug. 2014. DOI: 10.1145/2623330.2623732. [Online]. Available: <http://dx.doi.org/10.1145/2623330.2623732>.
- [21] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, “Line,” *Proceedings of the 24th International Conference on World Wide Web*, May 2015. DOI: 10.1145/2736277.2741093. [Online]. Available: <http://dx.doi.org/10.1145/2736277.2741093>.
- [22] L. Yao, C. Mao, and Y. Luo, *Graph convolutional networks for text classification*, 2018. arXiv: 1809.05679 [cs.CL].
- [23] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, *Graph attention networks*, 2018. arXiv: 1710.10903 [stat.ML].
- [24] J. Chen, T. Ma, and C. Xiao, “Fastgc: Fast learning with graph convolutional networks via importance sampling,” *CoRR*, vol. abs/1801.10247, 2018. arXiv: 1801.10247. [Online]. Available: <http://arxiv.org/abs/1801.10247>.
- [25] T. N. Kipf and M. Welling, *Semi-supervised classification with graph convolutional networks*, 2017. arXiv: 1609.02907 [cs.LG].
- [26] F. Wu, T. Zhang, A. H. S. Jr., C. Fifty, T. Yu, and K. Q. Weinberger, “Simplifying graph convolutional networks,” *CoRR*, vol. abs/1902.07153, 2019. arXiv: 1902.07153. [Online]. Available: <http://arxiv.org/abs/1902.07153>.
- [27] *20 News Group DataSet*, https://huggingface.co/datasets/SetFit/20_newsgroups, Accessed: 2022-5-24.
- [28] J. Pennington, R. Socher, and C. Manning, “GloVe: Global vectors for word representation,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 1532–1543. DOI: 10.3115/v1/D14-1162. [Online]. Available: <https://aclanthology.org/D14-1162>.
- [29] T. Mikolov, K. Chen, G. Corrado, and J. Dean, *Efficient estimation of word representations in vector space*, 2013. DOI: 10.48550/ARXIV.1301.3781. [Online]. Available: <https://arxiv.org/abs/1301.3781>.
- [30] Z. Guo, Y. Zhang, and W. Lu, “Attention guided graph convolutional networks for relation extraction,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 241–251. DOI: 10.18653/v1/P19-1024. [Online]. Available: <https://aclanthology.org/P19-1024>.

APPENDIX A CONTRIBUTION OF MEMBERS

A. Tuna Alikasıfoğlu

Tuna was responsible for several aspects of the project. First, he was responsible for data pre-processing with (Distil-) BERT tokenization techniques to convert the documents to list of input IDs, which corresponds to the text representations that are both used in generation of BERT and GNN embeddings. He was also responsible for generation of (Distil-) BERT embeddings to be used in the aggregation step. Finally, he implemented an efficient way to generate graph adjacency matrix A with sparse representation.

B. Arda Can Aras

Arda was responsible for researching and understanding the novel GNN algorithms and their implementations in text classification. He has implemented Graph Neural Network architecture, with different setups and proposed a learnable adjacency matrix A as a novelty. He also implemented the three different aggregation layer strategies and classification layer. Arda was also responsible of fine-tuning the Distil-BERT on 20NG dataset.

APPENDIX B

SOURCE CODE