# Short Text Classification via Term Graph

**Wei Pang**

School of Artificial Intelligence,
Beijing University of Posts and Telecommunications
pangweitf@{bupt.edu.cn,163.com}

## Abstract

Short text classification is a method for classifying short sentence with predefined labels. However, short text is limited in shortness in text length that leads to a challenging problem of sparse features. Most of existing methods treat each short sentences as independently and identically distributed (IID), local context only in the sentence itself is focused and the relational information between sentences are lost. To overcome these limitations, we propose a PathWalk model that combine the strength of graph networks and short sentences to solve the sparseness of short text. Experimental results on four different available datasets show that our PathWalk method achieves the state-of-the-art results, demonstrating the efficiency and robustness of graph networks for short text classification.

## 1 Introduction

Short textual sentences are produced in an explosive way in recent years, such as user comments(Tang et al. 2017) on shopping website, movie reviews, search query(Tang et al. 2017) for web search engine, and rapidly growing bullet screen (Djamasbi et al. 2016), or named *danmu*(He et al. 2018) message on video website. Understanding these short text has become an important problem for a variety of applications. However, short text is, as the name implies, shortness in text length, a typical sentence-level textual data. Compared with document-level text data, short text lacks of context and sufficient word occurrences, since its shortness, and hence it often leads to their feature space is very sparse. Besides, it can also lead to a challenge problem of robustness, since in real world application, small perturbations of short sentence features are more likely to result in misclassification. In this paper, we aim to alleviate the problem of sparsity and robustness for short text classification by using graph network at the word level.

Traditional classification methods for document-level text heavily rely on rich feature space (Wang et al. 2018; Wang et al. 2013). Commonly used models are bag-of-words (BoW) (Hu et al. 2009) or N-grams
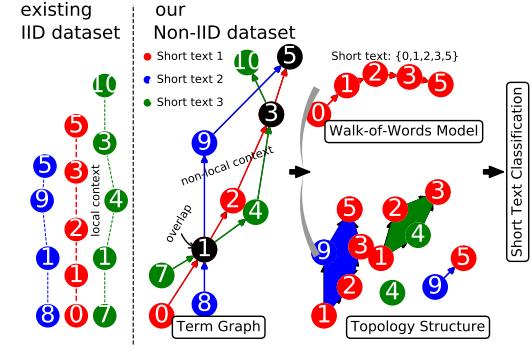
Figure 1: An idea of short text classification via term graph. The left is IID data, local context only in the sentence is focused. The right is Non-IID data, context is captured not only in the sentence itself but also in the graph networks. Richer contextual information can be used to short text classification.

(Wang and Manning 2012) and tf-idf term weighting scheme. On the one hand, owing to short text has very few words, existing methods often fail to provide sufficient features (Hu et al. 2009; Wang et al. 2013). On the other hand, conventional machine learning methods always treat dataset as independently and identically distributed (IID), as well as deep learning methods, leading to the BoW model only capture local contextual information, as the left of Fig. 1 illustrated, this IID assumption neglect the underlying correlations among sentences. A more serious problem is that existing classifiers are vulnerable to adversarial perturbations (Zugner, Akbarnejad, and Gunnemann 2018), for example, given an negative *danmu* message, if we add one or more trivial words that can't change its sentiment, it is possible to misclassify as positive message.

Enriching short text is an effective approach by introducing external knowledge(Hu et al. 2009; Wang et al. 2013; Ma et al. 2015; Li et al. 2016). Given short sentences, their external information can be mined from the re-

turned snippet of search engine(Sun 2012), or knowledge bases(Hu et al. 2009; Wang et al. 2013), such as Wikipedia and WordNet(Hu et al. 2009; Sun 2012). The challenge to expend short text is tend to topic drift(Li and Xu 2014; Tang et al. 2017), or the introduced features might be ambiguous(Hu et al. 2009; Wang et al. 2013; Li and Xu 2014; Tang et al. 2017).

However, in our view, there exists considerable redundancy information between the sentences. For example in task of sentiment analysis, the sentences that have the same sentiment, may share similar paraphrase or synonymy words in common. In particular, suppose that there are three short sentences are crossed at the word level, as shown in Fig. 1, where the circled number denote a word, they have the same context words while only few keyword or phrase is different, it may be redundant information, we can utilize these information to argument each other, and these couplings between sentences tend to reflect intrinsic characteristics in datasets. Based on these observations, we treat text datasets as non-IID, aiming to make full use of internal relations in training corpus.

Underlying the non-IID assumption, we propose a PathWalk algorithm for short text classification. It consists of two textual representation models and one classification method.

First, for representing the sentences corpus, the whole corpus is converted into a directed graph, called term graph. Where the node correspond to a word, directed edge indicate the order of words in a sentence. Moreover, edge with direction can be used to model a special word order in pair of words, since inversion of a pair of words may lead to semantic changes. For example a term graph in Fig. 2 that show complex relations between sentences, the number of edges is much more than the number of nodes, and intuitively, these edges could provide rich non-local contextual information for short text, which is the basis of our proposed PathWalk algorithm.

Second, for representing a short sentence, we present a walk-of-words model that use a sequence of nodes and directed edges in term graph to denote a sentence.

Third, we propose a method (i.e. PathWalk) that sample some non-local contextual information in term graph and use them to classify short text. In term graph, sentences that are crossed with each other form many closely connections, such as three nodes forms a triangular, four nodes forms a polygon, as shown in Fig 1. To distinguish the types of non-local contextual information in term graph, we adopt the topological feature, see more details in later.

In the end, the key contributions of our work are:

- We consider the sentence corpus as Non-IID dataset in place of typical IID, aims to capture non-local contextual information between sentences.

- We propose a term graph that have the capacity to represent the whole sentence corpus by a directed graph. It establishes the relationship between sentence and sentence and help alleviate the problem of data sparsity.

- We propose a walk-of-words model that can denote a sentence by a sequence of nodes and directed edges, instead of the typical bag-of-words model.

- We propose a PathWalk method that enable classify short textual sentence efficiently, and our model achieves new state-of-art results.

## 2    Models

Here a short text relate to the length is limited to only dozen words(Song et al. 2014), such as comment reviews, search query and DanMu message. We start by describing some concepts. Graph networks of words is defined as following.

Term Graph, is a directed graph $G = \{V, E\}$, consist of nodes and directed edges, where $V$ is a set of nodes, correspond to the vocabulary of training corpus, a node represent a word, $E$ is the set of directed edges, which represent a pre-order relation between words, the arrow points from current word to its followed from left to right in a sentence.

As seen in Fig. 2, where the edges have direction that represent a certain relation between pairs of words, the node and edge are both embed into a low-dimensional vector space, denoted $\vec{v}$ and $\vec{e}$ respectively. Intuitively underlying term graph, we assume that the edges and nodes together play a role for expressing semantics of sentence, e.g., different word order usually convey different meanings, in this case, the edge embedding encode information on word co-occurrence. To our knowledge, we first use the edge as independent feature to capture the relationships among words.

Local context is only inside the sentence itself because of a sentence is a word sequence in a linear way. However, term graph is not linear, the non-local context can be described as:

Non-local context, based on term graph, the contextual information can be introduced via many directions from network neighborhood, not only include local context inside the sentence itself but also include external context in the graph networks.

In this work, we aim to jointly two different types of applications, network embedding and text classification, learning classification-oriented, low-dimensional vector representation for nodes and edges that can be used for short text classification.

### 2.1    Sentence Corpus Representation

Existing methods usually view the sentence corpus are IID, one of limitations is that the rich relationships among sentences are lost. We transform the textual corpus into a term graph, as illustrated in Fig.2. Particularly, the number of nodes is smaller than the number of edges, there exist a lager links in term graph, and these links provide rich internal connections among words. Therefore, different from previous works, we adopt graph networks of words as our basic data structure for short text classification.

### 2.2    Short Sentence Representation

Walk-of-words model: A short text sentence is represented as a sequence of alternating nodes and directed edges in term graph.
For example, a short sentence $S$ of length 5 is composed of: $S = \{0, (0, 1), 1, (1, 2), 2, (2, 3), 3, (3, 5), 5\}$ in Fig. 1, where nodes means a single word and word appears in succession define a directed edge, the edge with direction can
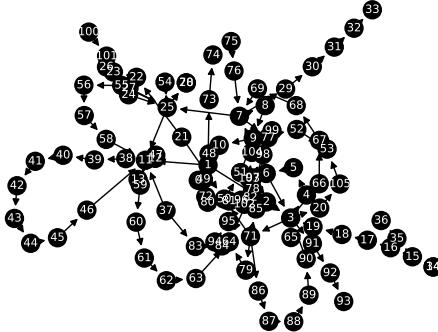
Figure 2: A real term graph, is built from the training dataset at word level, where every circled number denote a single word, the directed edge means word appears in sequential order.

also encode special information about word order, such as the order of edge $(1, 2)$.

**Network Topology**  For a given short sentence, how to draw valid topological features relative to it. Take a short sentence $S$ as example in Fig. 4(a), the nodes $\{1, 2, 3, 5\}$ are overlapped with neighbor sentences. The basic idea is that the two overlapping sentences might form a simple local triangle or higher-order polygon, these structure are frequently appears in term graph, called network topology. Moreover, the nodes and edges exist in topology that belong to external sentence might provide useful information for $S$. Even more evocatively, the semantic information can be flowed along with directed edges that make fully utilize data itself. To this end, we define four types of local network topology as contextual features, as see in Fig. 3.

Underlying Non-IID assumption, the topological structure we expect to capture the interaction among training samples that have overlapped nodes. According to Fig. 3, we define four basic types of simplical complex (network topology) as below, which can be used as contextual information for short sentence $S$ to reduce the problem of sparsity.

- 0-simplex, a set of single nodes shown in Fig.3(a), node denote a single word, which corresponding to bag of single words model.

- 1-simplex, two node with the inner directed edge between them, such as $\{9, 3\}$ in Fig.3(b), note that the ordered edge can be used to encode the word order.

- 2-simplex, a simple triangle simplicial complex see in Fig.3(c). A 2-simplex like the blue triangle $\{1, 2, 8\}$ show in Fig. 3(c).

- $3(4)$-simplex, visualized in Fig. 3(d), a 3-simplex shape such as quadrangle $\{1, 2, 3, 4\}$, a 4-simplex shape like pentagon $\{1, 2, 3, 5, 4\}$.

Importantly, in order to avoid introducing unnecessary noise, we give a constraint for network topology. Take a 2-simplex $\{1, 2, 8\}$ as example in Fig. 4(c), there is only one
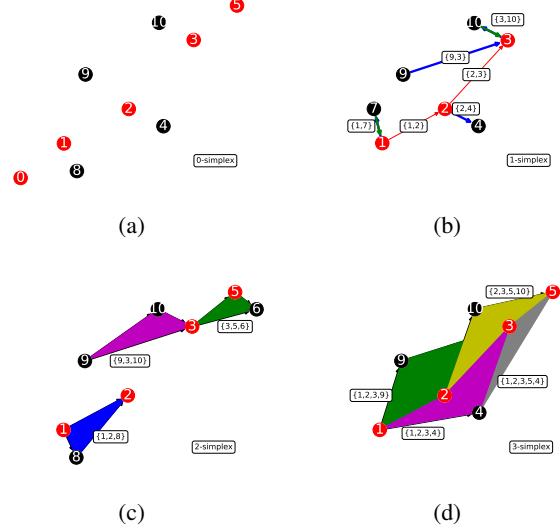


Figure 3: Four basic types of topological structure as features used for short text classification: a, $0$-simplices, such as single words in sentence. b, $1$-simplices, i.e. the directed inner-links between two single words is encoded as independent features. c, $2$-simplices, such as triangle $\{1, 2, 8\}$. d, $3(4)$-simplices, such as quadrangle or polygon show in figure.

node $\{8\}$ comes from neighbor sentence, $\{1, 2\}$ must in $S$, the node $\{1\}$ point to $\{8\}$ but immediately return to self node $\{2\}$ of $S$. We use this simple triangle as a feature w.r.t. node $\{2\}$. Ideally, node $\{8\}$ or edge $(8, 2)$ might bring contextual information, which is introduced to capture short-range context between $\{1\}$ and $\{2\}$.

The 3-simplex $\{1, 2, 3, 4\}$ seen in Fig. 4(c), is considered to have the ability to capture long-range dependencies between $\{1\}$ and $\{3\}$ in $S$. Another case is a 4-simplex $\{1, 2, 3, 5, 4\}$ in Fig. 3(d), $\{1, 2, 3, 5\}$ is a snippet of sentence $S$, although the first word $\{1\}$ is far away from the last $\{5\}$, but obviously the first have another path available to the last via external node $\{4\}$ that not in $S$. This shortcut plus node $\{4\}$ are considered as contexts of $S$, enhancement to semantic connection within $S$, in other words, they can offer richer contextual information to each other. This is why the Non-IID assumption works.

**Integrated Walk-of-Words model with Topological Structure**  In this subsection, we combine topological feature and walk-of-words model for representing short text. We describe three types of short sentence representation on top of term graph, visualized in Fig. 4.

First, Bag-of-words Model. A short sentence $S$ is only a bag of words, ignore the order of words, shown in Fig. 4(a). This representation is widely used in traditional machine learning tasks.

Second, our Walk-of-Words Model. Bag-of-words model plus with inner-links, let $E_{in}$ denote the inner-links, shown in Fig. 4(b), adding the internal edges of $S$. Different from the previous RNN-base text representation, we use the di-
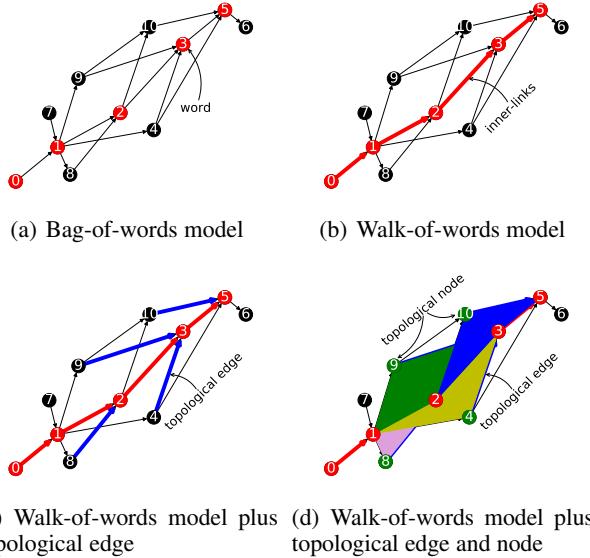
(a) Bag-of-words model     (b) Walk-of-words model

(c) Walk-of-words model plus topological edge     (d) Walk-of-words model plus topological edge and node

Figure 4: Four types of short text representation on graph, take sentence $S$ for illustration: a, Bag-of-words model ($BoW$), let $BoW = \{0,1,2,3,5\}$, then $S = BoW$. b, the proposed Walk-of-Words model, let $E_{in} = \{(0,1),(1,2),(2,3),(3,5)\}$, then $S = BoW \cup E_{in}$. c, let $E_t = \{(4,3),(8,2),(9,3),(10,5)\}$, then $S = BoW \cup E_{in} \cup E_t$. d, let $N_t = \{8,4,9,10\}$, then $S = BoW \cup E_{in} \cup E_t \cup Nt$.

rected edge as an independent feature. In details, $S = \{0,(0,1),1,(1,2),2,(2,3),3,(3,5),5\}$, is a sequence of alternating nodes and edges.

Third, Walk-of-Words Model with topology features. As see in Fig. 4(c), let $E_t$, $N_t$ denote topological edge and topological node respectively. We initially use topological edges that marked in the figure, as external features for expanding short text $S$. Furthermore, another higher-order topologies, visualized in Fig. 4(d), such as a 2-simplex $\{1,2,8\}$, 3-simplex $\{1,2,3,9\}$, $\{1,2,3,4\}$ and $\{2,3,5,10\}$, a 4-simplex $\{1,2,3,5,4\}$. The topological nodes that not in $S$, such as $\{8,9,10,4\}$, can also be used as external contexts to expand short text $S$, in order to capture non-local contextual information and provide much semantic features for short text representation.

## 3 PathWalk: Short Text Classification Algorithm

In this section, we present the PathWalk method for short text classification on top of term graph. The architecture is illustrated in Fig. 5, consists of four major structures.

First, Input layer, extracting topological features from term graph, including topological edges and topological nodes, those features are marked by blue color at the right of Fig. 5. Besides, Walk-of-Words representation correspond to the alternating red nodes and green edges.

Second, Embedding layer, which maps node and edge into dense vectors.

Third, BiLSTM layer, a one-layer bidirectional long short term memory (BiLSTM) encoder.

Last one is an output layer, we use SoftMax as our basic classifier.
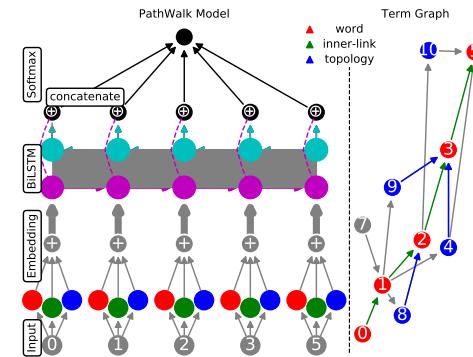


Figure 5: The model architecture of the proposed PathWalk.

According to Fig. 4, each short text representation corresponds to a feature space, obviously, from Fig. 4(a) to Fig. 4(d), the number of features used are gradually increasing. However, we would like to verify whether these contextual features can help to improve classification, Fig. 4(a) as our baseline, so that we design three variants of PathWalk algorithm for testing.

### 3.1 PathWalk-$I$ method

According to Fig. 4(b), based on walk-of-words model, a sentence $S$ of length $n$, its feature space can be defined as:

$$\overrightarrow{S}_{1:n} = \sum_{i=1}^{n} \overrightarrow{v}^{(i)} + \overrightarrow{E_{in}}^{(i)}, \tag{1}$$

where $\overrightarrow{S}_{1:n}$ denote the sentence representation of $S$, by summing of it's embedding vectors, $\overrightarrow{E_{in}}^{(i)}$ is the $i$-th inner-links embedding shown in Fig. 4(b), $\overrightarrow{v}^{(i)}$ denote the $i$-th node embedding. To simply the model, we directly use the summing of node vector and edge vector. The objective function is shown in Eq. 2:

$$\min_{\overrightarrow{v},\overrightarrow{e}} \frac{1}{M} \sum_{i=1}^{M} L(f(\overrightarrow{S}^{(i)}), y^{(i)}) + \lambda l2\_loss \tag{2}$$

where $y$ denotes the class of sentence $S$, $l2\_loss$ is a regularization that we employ half the $L2$ norm, $\lambda$ is the learning rate, $f$ denote a SoftMax function that is applied to the output of BiLSTM module and then convert it into probabilities.

### 3.2 PathWalk-$II$ method

According to Fig. 4(c), the topological edges are introduced to expand the feature space of $S$. Note that for each node of $S$, it has several topological edges, they all can be used as

contextual information, in practice, we only select the topological edge of size $K$. The sentence representation is defined as:

$$\overrightarrow{S}_{1:n} = \sum_{i=1}^{n} \{\overrightarrow{v}^{(i)} + \overrightarrow{E_{in}}^{(i)} + \sum_{j=1}^{K} \overrightarrow{E_{tj}}^{(i)}\}, \qquad (3)$$

where the symbol $\overrightarrow{E_t}$ is the embedding of topological edge.

### 3.3 PathWalk-$III$ method

According to Fig. 4(d), we calculate the higher-order topology to be used as additional features, such as the triangle or quadrangle covered with color shown in Fig. 4(d). In this case the sentence representation can be written as:

$$\overrightarrow{S}_{1:n} = \sum_{i=1}^{n} \{\overrightarrow{v}^{(i)} + \overrightarrow{E_{in}}^{(i)} + \sum_{j=1}^{K} \overrightarrow{E_{tj}}^{(i)} + \sum_{j=1}^{K} \overrightarrow{N_{tj}}^{(i)}\}, \ (4)$$

where $\overrightarrow{N_t}$ is the embedding of topological node, such as the green nodes $\{8, 4, 9, 10\}$ in figure. Similar to topological edge, we also sample $K$ topological nodes.

### 3.4 Random Sampling

The node in-degree follows a power-law-like shape distribution, that means the size of topological feature is varying in a range. To be fair, we randomly sample topologies of size $K$ as contextual features for every node of $S$.

### 3.5 The PathWalk algorithm

The pseudocode of three PathWalk variants is shown in Algorithm 1. First step, scanning each training sentences to construct term graph, where the directed edge from previous word point to current word; preprocessing each sentence using walk-of-words representation and sampling topological features; then optimizing the model using $Adam$ algorithm, these steps are performed in sequential.

## 4 Experiment and Evaluation

### 4.1 Datasets

To evaluate effectiveness of the proposed PathWalk method, we conduct experiments with four different domain datasets, detailed statistics are summarized in Table 1.

- RT-Polarity data(for short, RT)[1], a popular movie reviews dataset v1.0 for sentiment analysis with positive/negative labels.
- SST(Socher et al. 2013), a Stanford Sentiment Treebank dataset contains fine grained sentiment labels, such as very negative, negative, neutral, positive and very positive. In our work we do not use the neutral class, and merging the very negative and negative data to negative dataset, similar, merging the positive and very positive dataset to positive dataset for binary sentiment classification.

---

[1]http://www.cs.cornell.edu/people/pabo/movie-review-data/

- Video Query dataset(aka Query), Chinese queries dataset come from a video website, each query has a topic, either music or sport category. This dataset is used to binary topic classification task.
- Chinese DanMu dataset(aka DanMu), user reviews dataset, we gather from social media sites and give emotion tagging of positive or negative, and then used to binary sentiment classification problem.

### 4.2 Compared Methods and Setup

In (Kim 2014), the author proposed a shallow CNN model for sentence classification TextCNN, which achieve state-of-the-art results(Tang et al. 2017) and contains only one convolutional layer followed one max pooling layer, use three filters of size $3, 4, 5$ for capturing contextual information. fastText is another state-of-the-art(Joulin et al. 2017; ?) baseline for text representation and classification, one spot is to use subword embedding to overcome sparsity. A strong baseline method is a typical one-layer BiLSTM follows a SoftMax function(Sachan et al. 2018). We use fixed sentence length of 10 for Chinese sentence while of length 20 for English sentence, padding to fixed length with zero vector or truncating where necessary. We set the mini-batch size of 128, embedding size of 128-dimension for both node and edge, as well as the size of BiLSTM hidden states. We implement our experiments using $Tensorflow$ and use $Adam$ to optimize the model. Note that we preserve all the words and punctuation in dataset, these settings are common for all methods.

### 4.3 Parameter Sensitivity w.r.t. $K$
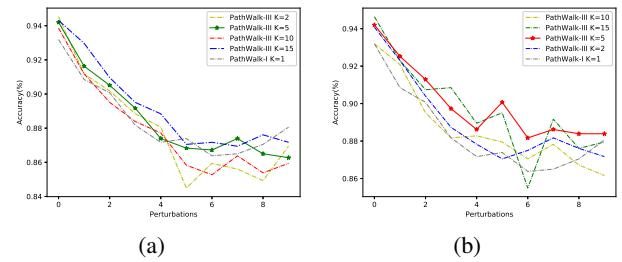


(a)          (b)

Figure 6: Parameter sensitivity with varying $K$ from 1 to 15 on $DanMu$ dataset, Perturbations means testing accuracy when adding noise to the test dataset. a, PathWalk-$II$ method. b, PathWalk-$III$ with varying $K$

The parameter $K$ is the number of topological features, it can affect the model's performance. In this subsection, when varying $K$ from 1 to 15, we perform experimental study of the influence of $K$ on the classification accuracy. Generally, a small number of topological features may lead to context sparsity, it enable us can not fully utilize information of external edges or external nodes from neighbors. Note that these features could bridge the gap on long range dependence inside $S$ as easily see from the Fig. 1. While a large $K$ might lead to semantic drift.

**Algorithm 1:** PathWalk algorithm

**Data:** Training dataset $D$
**Input:** Fixed sentence length $l$, Number of topologies $K$, Dimensions $d$
**Output:** model, Node and Edge Embedding

**1.1** $G$ = ToDirectedGraph $(D)$;
**1.2** modle = Model $(d)$;

**1.3** **for** *epoch* $\leftarrow 1$ **to** 20 **do**
**1.4**    **if** *PathWalk-I* **then**
**1.5**       batch = WalkofWord $(G, D)$;
**1.6**    **end**
**1.7**    **else if** *PathWalk-II* **then**
**1.8**       batch = WalkofWord $(G, D)$;
**1.9**       $\overrightarrow{oe}$ = SampleTopologyLinks $(G, \text{batch}, K)$;
**1.10**      batch = SumEmbedding $(\text{batch}, \overrightarrow{oe})$;
**1.11**   **end**
**1.12**   **else if** *PathWalk-III* **then**
**1.13**     batch = WalkofWord $(G, D)$;
**1.14**     $\overrightarrow{oe}$ = SampleTopologyLinks $(G, \text{batch}, K)$;
**1.15**     $\overrightarrow{\tau}$ = SampleTopologyNodes $(G, \text{batch}, K)$;
**1.16**     batch = SumEmbedding $(\text{batch}, \overrightarrow{oe}, \overrightarrow{\tau})$;
**1.17**   **end**
**1.18**   batch = PadSequences $(\text{batch}, l)$;
**1.19**   model.AdamOptimizer(batch);
**1.20** **end**

Table 1: Dataset statistics. Num., the number of datasets; Nodes, Edges, the number of nodes and edges in term graph at single word level; Lang, the language of text; AvgLen(std), average sentence length with the standard deviation.

| Dataset | Num. | Lang | AvgLen(std) | Nodes | Edges |
|---|---|---|---|---|---|
| DanMu | 50k | CH | 6.47(5.49) | 3981 | 76,321 |
| RT | 10k | EN | 22.36(**13.37**) | 18,862 | 106,752 |
| SST | 10k | EN | 19.27(9.23) | 18,302 | 88,109 |
| Query | 100k | CH | 8.40(4.12) | 5,290 | 209,085 |

Fig. 6 shows the broad trends of classification accuracy in decrease as increasing $K$.

Firstly, Fig. 6(a) shows the performance of PathWalk-$II$ is better at $K$ equal to 15, it reflects that adding more external in-degree links to walk-of-words model as contextual information can improve accuracy.

Secondly, however, the facts in Fig.6(b) are just the opposite, $K = 5$ achieve better result than the others, note that the performance of $K = 15$ is worse than $K = 5$, it means adding more external in-degree links and external nodes as contexts that lead to semantic drift. $K = 5$ is a tradeoff between contextual features and accuracy, so that we set $K$ equal to 5 both for PathWalk-$II$ and PathWalk-$III$.

### 4.4 Comparison with state-of-the-art

In this subsection we compare PathWalk variants with the state-of-the-art models on four different domain datasets. All the methods are trained at the single word level. In particular, PathWalk variants are trained on term graph of single word, and the others also use bag of features at the single word level as input.

Table 2 shows the comparison result. Firstly, PathWalk variants are usually better than the baselines in all cases, and achieve a strong benchmark results.

Secondly, on $DanMu$ dataset, PathWalk variants obtain higher performance gradually when introducing richer contextual information, but this consistent trend not appears in other datasets. In details, the $RT$ and $SST$ datasets contain much longer reviews, as listed in Table1, while the $DanMu$ dataset contains much shorter reviews among the four datasets. It is worth noting that PathWalk variants are much better on shorter sentences, and shows that many overlapped shorter texts in term graph are really complementary to each other. But, to longer sentences, although network topology is still helpful to improve the overall accuracy, the fine-tuning value of $K$ play a crucial role for classification. We leave how to select the best $K$ of topological features for expanding short sentence as future work.

In summary, our PathWalk method obtains the state-of-the-art performance in all four different domain datasets. The experiment result demonstrate that the non-local context could provide much gains to improve classification. Also it shows the size $K$ of external contextual features needs to be fine-tuned during training phase.

Table 2: Classification accuracy on four different datasets at the single word level, each test datasets contains 1000 sentences.

| Method | $DanMu$ | $RT$ | $Query$ | $SST$ |
|---|---|---|---|---|
| BiLSTM | 94.08 | 61.61 | 93.08 | 74.22 |
| fastText | 92.70 | 54.10 | 92.70 | 51.85 |
| TextCNN | 92.40 | 58.70 | 91.10 | 70.09 |
| PathWalk-$I$ | 93.19 | **62.50** | **94.08** | 75.39 |
| PathWalk-$II$ | **94.20** | 60.94 | **93.53** | 77.21 |
| PathWalk-$III$ | **94.53** | **61.72** | **94.31** | 75.39 |

## 4.5 Comparison of Term Graph at Single Word vs. Word Segmentation

Here we focus on the comparison between term graph at single word and word segmentation. For Chinese sentence dataset, we would like to know how performance is influenced by Chinese word segmentation. We apply a popular word segmentation, $jieba$ tool[2] to $DanMu$ and $Query$ datasets respectively, then we build graph using the result of $jieba$. All the methods are trained on word tokens level after word segmentation.

Table. 3 shows the comparison performance.

Table 3: Compared Accuracy at the word token level after $jieba$ and at the single word level.

| Method | Word Token Level | | Single Word Level | |
|---|---|---|---|---|
| | $DanMu$ | $Query$ | $DanMu$ | $Query$ |
| BiLSTM | 96.32↑ | 94.20↑ | 94.08 | 93.08 |
| fastText | 91.00↓ | 92.00↓ | 92.70 | 92.70 |
| TextCNN | 96.10↑ | 93.00↑ | 92.40 | 91.10 |
| PathWalk-$I$ | 95.09↑ | **94.64↑** | 93.19 | 94.08 |
| PathWalk-$II$ | **96.54↑** | **95.31↑** | 94.20 | 93.53 |
| PathWalk-$III$ | **96.65↑** | 94.64↑ | 94.75 | 94.31 |

Firstly, as we can see, the overall accuracy is increased but the $fastText$ method in decrease slightly, as seen the symbol ↑ and ↓ tagging. It shows word segmentation are helpful compared with single word. The reason is that the word token is meaningful, but the drawback is their feature space become more sparse.

Secondly, our PathWalk is often better than the three baselines. Note that PathWalk variants gradually improve the accuracy in a consistent way across the two datasets, from $95.09\%$ up to $96.65\%$ and from $94.64\%$ to $95.31\%$ respectively, but PathWalk-$III$ on $Query$ dataset drops to $94.64\%$ since introducing more contexts lead to semantic drift.

## 4.6 Adversarial experiment: perturbation on test phase

The quality of robustness can be measured in adversarial experiment, in this work we adopt two types of adversarial experiment, the first is only adding noises to test dataset, while the second is attacking the model during training phase via adding noises to training dataset. In this subsection, we perform the first experiment at the single word level.

How small noises attack on the content of test sentence could affect the classification performance. We randomly sample words as noise padding to the test sentence, where the noises come from vocabulary of respective training datasets. However in this case, we assume that the perturbed sentence will preserve semantics of the original, because of a sentence size is very smaller than the size of vocabulary, so that it is less likely to change the meaning of the original sentence.

To test robustness of the pre-trained model, we gradually increase the number of noise words, with varying from 1 to 8, this is equivalent to directly attack the pre-trained model by unseen sentences, and then observe their robustness to noises in terms of classification accuracy.
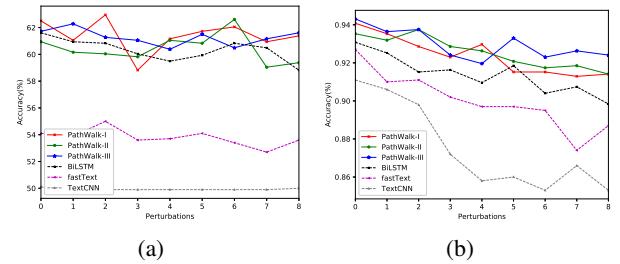


(a)                    (b)

Figure 7: Adversarial perturbation on $RT$ and $Query$ dataset with increasing number of perturbations. a, experimental result on English $RT$ dataset. b, on Chinese Video Query dataset. The noise word we randomly select from each Vocabulary.

Firstly, Table 4 shows perturbation result on the two different domain dataset, i.e. $DanMu$ and $SST$. When gradually increasing noise words, the overall accuracy result in decrease. Surprisingly our PathWalk variants shows more stable than others, and importantly result in a lower decrease in accuracy when increasing noises compared with the baselines.

Secondly, Fig. 7 shows the results on $RT$ and $Query$ datasets. Fig.7(b) shows the accuracy of PathWalk variants are more stable than the others on $Query$ dataset, and also consistently exceeds the baselines. In contrast, Fig.7(a) shows the accuracy of PathWalk-$II$ is worse than the other two PathWalk variant and $BiLSTM$ method. However, all the methods on $RT$ dataset achieve a lower accuracy, due to the standard deviation of $RT$ sentence length is greater than the other three datasets, in this case it indicate that our strategy of fixed-sentence length is not appropriate for dataset with high variance.

Overall, our methods outperform the strong baselines in robustness on three out of four datasets, except $RT$ since it's fixed-sentence rule is not work well.

---

[2]https://github.com/fxsjy/jieba

Table 4: Adversarial perturbation on $DanMu$ and $SST$ dataset with increasing number of perturbations, $+i$ means we add $i$ noise words to every sentence of both dataset.

| Test | PathWalk | | | Baselines | | |
|---|---|---|---|---|---|---|
| | $I$ | $II$ | $III$ | BiLSTM | fastText | TextCNN |
| $DanMu$ | 93.19 | 94.20 | 94.53 | 94.08 | 92.90 | 92.40 |
| +1 | 90.85 | 91.63 | 91.74 | 90.62 | 92.40 | 87.10 |
| +2 | 90.07 | 90.51 | 88.84 | 89.51 | 90.50 | 81.20 |
| +3 | 88.17 | 89.17 | 88.06 | 87.83 | 89.70 | 75.00 |
| +4 | 87.17 | 87.39 | 87.61 | 86.94 | 89.00 | 71.10 |
| +5 | 87.39 | 86.83 | 86.72 | 87.17 | 89.10 | 67.70 |
| +6 | 86.38 | 86.72 | 86.83 | 87.95 | 88.30 | 64.30 |
| +7 | 86.50 | 87.39 | 86.38 | 84.60 | 85.90 | 63.10 |
| +8 | 87.05 | 86.50 | 87.83 | 83.93 | 85.30 | 63.40 |
| $SST$ | 75.39 | 77.21 | 77.08 | 74.22 | 51.85 | 70.09 |
| +1 | 74.74 | 76.30 | 75.26 | 73.96 | 51.73 | 68.36 |
| +2 | 75.39 | 76.04 | 75.26 | 73.31 | 52.66 | 66.63 |
| +3 | 73.31 | 75.00 | 75.91 | 73.57 | 51.62 | 65.47 |
| +4 | 74.87 | 76.17 | 74.35 | 73.57 | - | 65.47 |
| +5 | 73.44 | 74.09 | 75.52 | 73.05 | - | 65.13 |
| +6 | 74.09 | 75.26 | 73.57 | 73.05 | - | 63.74 |
| +7 | 73.57 | 76.17 | 75.52 | 73.07 | - | 64.43 |
| +8 | 73.83 | 74.74 | 74.22 | 72.53 | - | 62.82 |

## 4.7 Adversarial experiment: perturbation during model training

Prior experiments we focus on robustness measure during test time. In this subsection we focus on the second adversarial experiment, attacking the model during training phase. Here we only conduct experiment on walk-of-words model using PathWalk-$I$ method, through adding randomly perturbations to each training sentence before fed it into model.

In details, the adversarial noises are directly added to each word of target sentence, where the noise comes from the term graph itself, including node perturbation or edge perturbations, then the sentence representation can be written as:

$$\overrightarrow{S}_{1:n} = \sum_{i=1}^{n} \{\overrightarrow{v}_i + \overrightarrow{ie}_i + \varepsilon \times \overrightarrow{o}^{(i)}\}, \qquad (5)$$

where $\varepsilon$ is a hyperparameter and $\overrightarrow{o}$ is an embedding of the sampled edge or node in term graph as noises during training. Also, the short sentence model can also be interpreted as walk-of-words model plus the adversarial perturbations for each word.

Table 5 shows the robustness performance, the last *column* as our baseline.

As seen, PathWalk-$I^{\ell}$ that attacked by only one edge perturbation, and achieves the significant improvement. Also surprisedly, PathWalk-$I^{\ell}$ outperforms PathWalk-$I$, including both of the single word level and word segmentation level. In $DanMu$ sentiment classification task, the PathWalk-$I^{\ell}$ by adversarial training improves the overall accuracy of the original PathWalk-$I$ at single word level from 93.19% to 95.42%. PathWalk-$I^{*}$ also gives better results than the baseline. While, PathWalk-$I^{\wr}$ perform worse than the baseline, because of it suffer many times perturbations.

Table 5: Adversarial training on $DanMu$ dataset using graph of single word, the different is the sentence representation is changed to Eq.(5), where $\varepsilon = 1.0$. Detailedly, the $^{\ell}$ labeled method only add one edge as noises to each word of target sentence, similarly, the $^{*}$ method add four edges as noises, the $\wr$ add four edges and nodes as noises. Where all the small perturbations are derived from term graph itself. $+i$ is the same as before.

| Test | PathWalk | | | |
|---|---|---|---|---|
| | $I^{\ell}$ | $I^{*}$ | $I^{\wr}$ | $I$ |
| $DanMu$ | 95.42↑ | 94.64↑ | 93.30 | 93.19† |
| +1 | 94.08 | 93.08 | 87.83↓ | 90.85 |
| +2 | 91.96 | 91.29 | 86.50 | 90.07 |
| +3 | 89.51 | 91.85 | 83.04 | 88.17 |
| +4 | 91.41 | 90.18 | 83.82 | 87.17 |
| +5 | 89.62 | 89.62 | 79.69 | 87.39 |
| +6 | 89.73 | 90.51 | 79.58 | 86.38 |
| +7 | 88.28 | 89.06 | 78.46 | 86.50 |
| +8 | 89.29 | 88.28 | 75.11 | 87.05 |

† the result of Table.4 as baseline.

In summary, adversarial training is an interesting problem and the result shows that a small perturbations are efficient in improving the model's performance.

## 5 Conclusions

In this work, we replace IID dataset assumption with Non-IID to perform short text classification more efficiently. Under this Non-IID, we use graph networks of words to represent the whole training corpus. In the graph, a short sen-

tence can capture more contextual information than in its own self, and in other words, it alleviate the problem of data sparsity for short text. Experiment on four different domain datasets show that term graph and network topology could improve performance in classification accuracy and robustness. When a short sentence attacked with small noises, our PathWalk method is proved to be more stable than the baselines.

In the future, we are going to explore the use of topological features based on attention methods.

## 6 Acknowledgements

## References

[Djamasbi et al. 2016] Djamasbi, S.; Hall-Phillips, A.; Liu, Z.; Li, W.; and Bian, J. 2016. Social viewing, bullet screen, user experience: A first look. In *2016 49th Hawaii International Conference on System Sciences (HICSS)*, HICSS'16, 648–657. Koloa, HI, USA: IEEE.

[He et al. 2018] He, M.; Ge, Y.; Chen, E.; Liu, Q.; and Wang, X. 2018. Exploring the emerging type of comment for online videos: Danmu. *ACM Transaction on the web(TWEB)* 12(1).

[Hu et al. 2009] Hu, X.; Sun, N.; Zhang, C.; and Chua, T.-S. 2009. Exploiting internal and external semantics for the clustering of short texts using world knowledge. In *Proceedings of the 18th ACM conference on Information and knowledge management*, CIKM'09, 919–928.

[Joulin et al. 2017] Joulin, A.; Grave, E.; Bojanowski, P.; and Mikolov, T. 2017. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, volume 2, 427–431. Association for Computational Linguistics.

[Kim 2014] Kim, Y. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, EMNLP'14, 1746–1751. Doha, Qatar: Association for Computational Linguistics.

[Li and Xu 2014] Li, H., and Xu, J. 2014. Semantic matching in search. *Foundations and Trend in Information Retrieval* 7(5):343–469.

[Li et al. 2016] Li, C.; Wang, H.; Zhang, Z.; Sun, A.; and Ma, Z. 2016. Topic modeling for short texts with auxiliary word embeddings. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, SIGIR'16, 165–174. Pisa, Italy: ACM.

[Ma et al. 2015] Ma, C.; Xu, W.; Li, P.; and Yan, Y. 2015. Distributional representations of words for short text classification. In *Proceedings of NAACL-HLT*, NAACL'15, 33–38. Denver, Colorado: Association for Computational Linguistics.

[Sachan et al. 2018] Sachan, D. S.; Zaheer, M.; Bojanowski, P.; and Salakhutdinov, R. 2018. Revisiting lstm networks for semi-supervised text classification via mixed objective function. In *KDD'18 Deep Learning Day*. London, UK: ACM.

[Socher et al. 2013] Socher, R.; Perelygin, A.; Wu, J.; Chuang, J.; Manning, C. D.; Ng, A.; and Potts, C. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Conference on Empirical Methods in Natural Language Processing*, EMNLP'13, 1631–1642. Association for Computational Linguistics.

[Song et al. 2014] Song, G.; Ye, Y.; Du, X.; Huang, X.; and Bie, S. 2014. Short text classification: A survey. *JOURNAL OF MULTIMEDIA, VOL. 9, NO. 5, MAY 2014* 9(5):635–643.

[Sun 2012] Sun, A. 2012. Short text classification using very few words. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, SIGIR'12, 1145–1146. Portland, Oregon, USA: ACM.

[Tang et al. 2017] Tang, J.; Wang, Y.; Zheng, K.; and Mei, Q. 2017. End-to-end learning for short text expansion. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD'17, 1105–1113. Halifax, NS, Canada: ACM.

[Wang and Manning 2012] Wang, S., and Manning, C. D. 2012. Baselines and bigrams: Simple, good sentiment and topic classification. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*, ACL'12, 90–94.

[Wang et al. 2013] Wang, J.; Wang, Z.; Zhang, D.; and Yan, J. 2013. Combining knowledge with deep convolutional neural networks for short text classification. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, IJCAI'17, 1631–1642. Association for Computational Linguistics.

[Wang et al. 2018] Wang, G.; Li, C.; Wang, W.; Zhang, Y.; Shen, D.; Zhang, X.; Henao, R.; and Carin, L. 2018. Joint embedding of words and labels for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, SIGIR'10, 2321–2331. Melbourne, Australia: ACM.

[Zugner, Akbarnejad, and Gunnemann 2018] Zugner, D.; Akbarnejad, A.; and Gunnemann, S. 2018. Adversarial attacks on neural networks for graph data. In *The 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD'18. New York, NY, USA: ACM.