

A Survey on Text Classification: From Shallow to Deep Learning

Qian Li, Hao Peng, Jianxin Li, *Member, IEEE* Congying Xia, Renyu Yang, *Member, IEEE* Lichao Sun, Philip S. Yu, *Fellow, IEEE* and Lifang He, *Member, IEEE*

Abstract—Text classification is the most fundamental and essential task in natural language processing. The last decade has seen a surge of research in this area due to the unprecedented success of deep learning. Numerous methods, datasets, and evaluation metrics have been proposed in the literature, raising the need for a comprehensive and updated survey. This paper fills the gap by reviewing the state of the art approaches from 1961 to 2020, focusing on models from shallow to deep learning. We create a taxonomy for text classification according to the text involved and the models used for feature extraction and classification. We then discuss each of these categories in detail, dealing with both the technical developments and benchmark datasets that support tests of predictions. A comprehensive comparison between different techniques, as well as identifying the pros and cons of various evaluation metrics are also provided in this survey. Finally, we conclude by summarizing key implications, future research directions, and the challenges facing the research area.

Index Terms—deep learning, shallow learning, text classification, evaluation metrics, challenges.

I. INTRODUCTION

TEXT classification – the procedure of designating pre-defined labels for text – is an essential and significant task in many Natural Language Processing (NLP) applications, such as sentiment analysis [1][2] [3], topic labeling [4] [5] [6], question answering [7] [8] and dialog act classification [9]. In the era of information explosion, it is time-consuming and challenging to process and classify large amounts of text data manually. Besides, the accuracy of manual text classification can be easily influenced by human factors, such as fatigue and expertise. It is desirable to use machine learning methods to automate the text classification procedure to yield more reliable and less subjective results. Moreover, this can also help enhance information retrieval efficiency and alleviate the problem of information overload by locating the required information.

Fig. 1 illustrates a flowchart of the procedures involved in the text classification, under the light of shallow and

Manuscript received October 11, 2020; date of current version October 10, 2020. (Corresponding author: Jianxin Li.)

Qian Li, Hao Peng and Jianxin Li are with Beijing Advanced Innovation Center for Big Data and Brain Computing, Beihang University, Beijing 100083, China. E-mail: liqian, penghao, lijx@act.buaa.edu.cn

Renyu Yang is with the School of Computing, University of Leeds, Leeds LS2 9JT, UK. E-mail: r.yang1@leeds.ac.uk.

Congying Xia and Philip S. Yu are with the Department of Computer Science, University of Illinois at Chicago, Chicago 60607, USA. E-mail: cxia8, psyu@uic.edu.

Lichao Sun and Lifang He are with the Department of Computer Science and Engineering, Lehigh University, Bethlehem, PA 18015 USA. E-mail: james.lichao.sun@gmail.com, lih319@lehigh.edu.

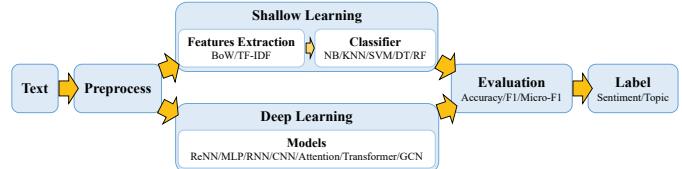


Fig. 1. Flowchart of the text classification with classic methods in each module. It is crucial to extract essential features for shallow models, but features can be extracted automatically by DNNs.

deep analysis. Text data is different from numerical, image, or signal data. It requires NLP techniques to be processed carefully. The first important step is to preprocess text data for the model. Shallow learning models usually need to obtain good sample features by artificial methods and then classify them with classic machine learning algorithms. Therefore, the effectiveness of the method is largely restricted by feature extraction. However, different from shallow models, deep learning integrates feature engineering into the model fitting process by learning a set of nonlinear transformations that serve to map features directly to outputs.

From the 1960s until the 2010s, shallow learning-based text classification models dominated. Shallow learning means statistics-based models, such as Naïve Bayes (NB) [10], K-nearest neighbor (KNN) [11], and support vector machine (SVM) [12]. Comparing with the earlier rule-based methods, this method has obvious advantages in accuracy and stability. However, these approaches still need to do feature engineering, which is time-consuming and costly. Besides, they usually disregard the natural sequential structure or contextual information in textual data, making it challenging to learn the semantic information of the words. Since the 2010s, text classification has gradually changed from shallow learning models to deep learning models. Compared with the methods based on shallow learning, deep learning methods avoid designing rules and features by humans and automatically provide semantically meaningful representations for text mining. Therefore, most of the text classification research works are based on DNNs, which are data-driven approaches with high computational complexity. Few works focus on shallow learning models to settle the limitations of computation and data.

A. Major Differences and Contributions

There have been several works reviewing text classification and its subproblems recently. Two of them are reviews of text classification. Kowsari et al. [13] surveyed different text

feature extraction, dimensionality reduction methods, basic model structure for text classification, and evaluation methods. Minaee et al. [14] reviewed recent deep learning based text classification methods, benchmark datasets, and evaluation metrics. Unlike existing text classification reviews, we conclude existing models from shallow to deep learning with works of recent years. Shallow learning models emphasize the feature extraction and classifier design. Once the text has well-designed characteristics, it can be quickly converged by training the classifier. DNNs can perform feature extraction automatically and learn well without domain knowledge. We then give the datasets and evaluation metrics for single-label and multi-label tasks and summarize future research challenges from data, models, and performance perspective. Moreover, we summarize various information in three tables, including the necessary information of classic deep learning models, primary information of main datasets, and a general benchmark of state-of-the-art methods under different applications. In summary, this study's main contributions are as follows:

- We introduce the process and development of text classification and present comprehensive analysis and research on primary models – from shallow to deep learning models – according to their model structures. We summarize the necessary information of deep learning models in terms of basic model structures in Table I, including publishing years, methods, venues, applications, evaluation metrics, datasets and code links.
- We introduce the present datasets and give the formulation of main evaluation metrics with the comparison of metrics, including single-label and multi-label text classification tasks. We summarize the necessary information of primary datasets in Table II, including the number of categories, average sentence length, the size of each dataset, related papers and data addresses.
- We summarize classification accuracy scores of models given in their articles, on benchmark datasets in Table IV and conclude the survey by discussing the main challenges facing the text classification and key implications stemming from this study.

B. Organization of the Survey

The rest of the survey is organized as follows. Section II summarizes the existing models related to text classification, including shallow learning and deep learning models, including a summary table. Section III introduces the primary datasets with a summary table and evaluation metrics on single-label and multi-label tasks. We then give quantitative results of the leading models in classic text classification datasets in Section IV. Finally, we summarize the main challenges for deep learning text classification in Section V before concluding the article in Section VI.

II. TEXT CLASSIFICATION METHODS

Text classification is referred to as extracting features from raw text data and predicting the categories of text data based on such features. Numerous models have been proposed in the

past few decades for text classification. For shallow learning models, NB [10] is the first model used for the text classification task. Whereafter, generic classification models are proposed, such as KNN, SVM [12], and RF [15], which are called classifiers, widely used for text classification. Recently, XGBoost [16] and LightGBM [17] have arguably the potential to provide excellent performance. For deep learning models, TextCNN [18] has the highest number of references in these models, wherein a CNN model has been introduced to solve the text classification problem for the first time. While not specifically designed for handling text classification tasks, BERT [19] has been widely employed when designing text classification models, considering its effectiveness on numerous text classification datasets.

A. Shallow Learning Models

Shallow learning models accelerate text classification with improved accuracy and make the application scope of shallow learning expand. The first thing is to preprocess the raw input text for training shallow learning models, which generally consists of word segmentation, data cleaning, and data statistics. Then, text representation aims to express preprocessed text in a form that is much easier for computers and minimizes information loss, such as Bag-of-words (BOW), N-gram, term frequency-inverse document frequency (TF-IDF) [20], word2vec [21] and GloVe [22]. At the core of the BOW is representing each text with a dictionary-sized vector. The individual value of the vector denotes the word frequency corresponding to its inherent position in the text. Compared to BOW, N-gram considers the information of adjacent words and builds a dictionary by considering the adjacent words. TF-IDF [20] uses the word frequency and inverses the document frequency to model the text. The word2vec [21] employs local context information to obtain word vectors. The GloVe [22] – with both the local context and global statistical features – trains on the nonzero elements in a word-word co-occurrence matrix. Finally, the represented text is fed into the classifier according to selected features. Here, we discuss some of the representative classifiers in detail:

1) *PGM-based methods*: Probabilistic graphical models (PGMs) express the conditional dependencies among features in graphs, such as the Bayesian network [23], the hidden Markov network [24]. Such models are combinations of probability theory and graph theory.

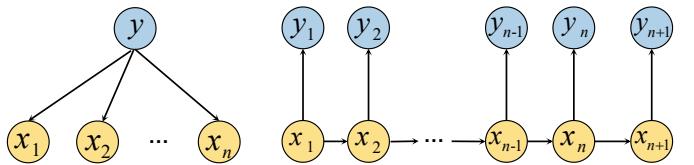


Fig. 2. The structure of NB (left) and HMM (right).

Naïve Bayes (NB) [10] is the simplest and most broadly used model based on applying Bayes' theorem. The NB algorithm has an independent assumption: when the target value has been given, the conditions between features $x = [x_1, x_2, \dots, x_n]$ are independent (see Fig. 2). The NB

algorithm primarily uses the prior probability to calculate the posterior probability. Due to its simple structure, NB is broadly used for text classification tasks. Although the assumption that the features are independent is sometimes not actual, it substantially simplifies the calculation process and performs better. To improve the performance on smaller categories, Schneider [25] proposes a feature selection score method through calculating KL-divergence [26] between the training set and corresponding categories for multinomial NB text classification. Dai et al. [27] propose a transfer learning method named Naive Bayes Transfer Classification (NBTC) to settle the different distribution between the training set and the target set. It uses the EM algorithm [28] to obtain a locally optimal posterior hypothesis on the target set.

Hidden Markov model (HMM) is a Markov model assumed to be a Markov process within hidden states [24]. It is suitable for sequential text data, effective in reducing algorithmic complexity by redesigning model structure. HMM operates under the assumption that a separate process Y exists, and its behavior depends upon X . The reachable learning goal is to learn about X by observing Y , considering the state dependencies (see Fig. 2). To consider the contextual information among pages in a text, Frasconi et al. [29] reshape a text into sequences of pages and exploit the serial order relationship among pages within a text for multi-page texts. However, these methods get no excellent performance for domain text. Motivated by this, Yi et al. [30] use prior knowledge – primarily stemming from a specialized subject vocabulary set Medical Subject Headings (MeSH) [31] – to carry out the medical text classification task.

2) *KNN-based Methods:* At the core of the K-Nearest Neighbors (KNN) algorithm [11] is to classify an unlabeled sample by finding the category with most samples on the k -nearest labeled samples. It is a simple classifier without building the model and can decrease complexity through the fast process of getting k nearest neighbors. Fig. 3 showcases the structure of KNN. We can find k training texts approaching a specific text to be classified through estimating the in-between distance. Hence, the text can be divided into the most common categories found in k training set texts. However, due to the positive correlation between model time/space complexity and the amount of data, the KNN algorithm takes an unusually long time on the large-scale datasets. To decrease the number of selected features, Soucy et al. [32] propose a KNN algorithm without feature weighting. It manages to find relevant features, building the inter-dependencies of words by using a feature selection. When the data is extremely unevenly distributed, KNN tends to classify samples with more data. The neighbor-weighted K-nearest neighbor (NWKNN) [33] is proposed to improve classification performance on the unbalanced corpora. It casts a significant weight for neighbors in a small category and a small weight for neighbors in a broad class.

3) *SVM-based Methods:* Cortes and Vapnik propose Support Vector Machine (SVM) [34] to tackle the binary classification of pattern recognition. Joachims [12], for the first time, uses the SVM method for text classification representing each text as a vector. As illustrated in Fig. 4, SVM-based

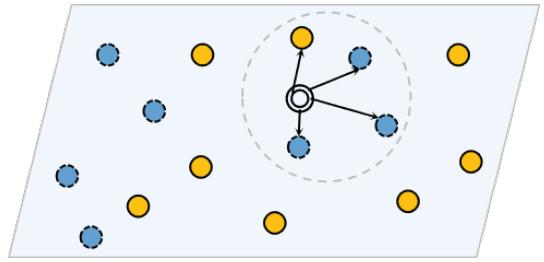


Fig. 3. The structure of KNN where $k = 4$. Each node represents a text and nodes with different contours represent different categories.

approaches turn text classification tasks into multiple binary classification tasks. In this context, SVM constructs an optimal hyperplane in the one-dimensional input space or feature space, maximizing the distance between the hyperplane and the two categories of training sets, thereby achieving the best generalization ability. The goal is to make the distance of the category boundary along the direction perpendicular to the hyperplane is the largest. Equivalently, this will result in the lowest error rate of classification. Constructing an optimal hyperplane can be transformed into a quadratic programming problem to obtain a globally optimal solution. Choosing the appropriate kernel function is of the utmost importance to ensure SVM can deal with nonlinear problems and become a robust nonlinear classifier. To analyze what the SVM algorithms learn and what tasks are suitable, Joachims [35] proposes a theoretical learning model combining the statistical traits with the generalization performance of an SVM analyzing the features and benefits using a quantitative approach. Transductive Support Vector Machine (TSVM) [36] is proposed to lessen misclassifications of the particular test collections with a general decision function considering a specific test set. It uses prior knowledge to establish a more suitable structure and study faster.

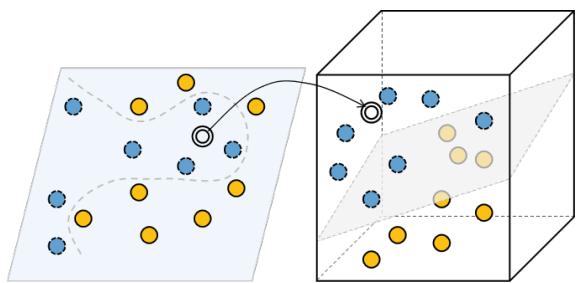


Fig. 4. The structure of SVM. Each node represents a text and nodes with different contours represent different categories.

4) *DT-based Methods:* Decision trees (DT) [37] is a supervised tree structure learning method – reflective of the idea of divide-and-conquer – and is constructed recursively. It learns disjunctive expressions and has robustness for the text with noise. As shown in Fig. 5, decision trees can be generally divided into two distinct stages: tree construction and tree pruning. It starts at the root node and tests the data samples (composed of instance sets, which have several attributes), and divides the dataset into diverse subsets according to different results. A subset of datasets constitutes a child node, and every

leaf node in the decision tree represents a category. Constructing the decision tree is to determine the correlation between classes and attributes, further exploited to predict the record categories of unknown forthcoming types. The classification rules generated by the decision tree algorithm are straightforward, and the pruning strategy can also help reduce the influence of noise. Its limitation, however, mainly derives from inefficiency in coping with explosively increasing data size. More specifically, the ID3 [38] algorithm uses information gain as the attribute selection criterion in the selection of each node – It is used to select the attribute of each branch node, and then select the attribute having the maximum information gain value to become the discriminant attribute of the current node. Based on ID3, C4.5 [39] learns to obtain a map from attributes to classes, which effectively classifies entities unknown to new categories. DT based algorithms usually need to train for each dataset, which is low efficiency. Thus, Johnson et al. [40] propose a DT-based symbolic rule system. The method represents each text as a vector calculated by the frequency of each word in the text and induces rules from the training data. The learning rules are used for classifying the other data being similar to the training data. Furthermore, to reduce the computational costs of DT algorithms, fast decision-tree (FDT) [41] uses two-pronged strategy: pre-selecting a feature set and training multiple DTs on different data subsets. Results from multiple DTs are combined through a data-fusion technique to resolve the cases of imbalanced classes.

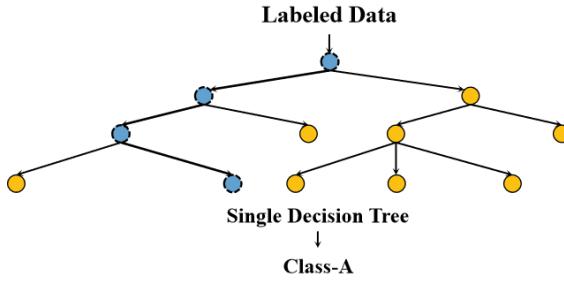


Fig. 5. The structure of the decision trees (DT). The nodes with the dotted outline represent the nodes of the decision route.

5) Integration-based Methods: Integrated algorithms aim to aggregate the results of multiple algorithms for better performance and interpretation. Conventional integrated algorithms are bootstrap aggregation, such as random forest (RF) [15], boosting such as AdaBoost [42], and XGBoost [16] and stacking. The bootstrap aggregation method trains multiple classifiers without strong dependencies and then aggregates their results. For instance, RF [15] consists of multiple tree classifiers wherein all trees depend on the value of the random vector sampled independently (depicted in Fig. 6). It is worth noting that each tree within the RF shares the same distribution. The generalization error of an RF relies on the strength of each tree and the relationship among trees and will converge to a limit with the increment of tree number in the forest. In boosting based algorithms, all labeled data are trained with the same weight to initially obtain a weaker classifier. The weights of the data will then be adjusted according to the former result of the classifier. The training procedure will continue

by repeating such steps until the termination condition is reached. Unlike bootstrap and boosting algorithms, stacking based algorithms break down the data into n parts and use n classifiers to calculate the input data in a cascade manner – Result from upstream classifier will feed into the downstream classifier as input. The training will terminate once a pre-defined iteration number is targeted. The integrated method can capture more features from multiple trees.

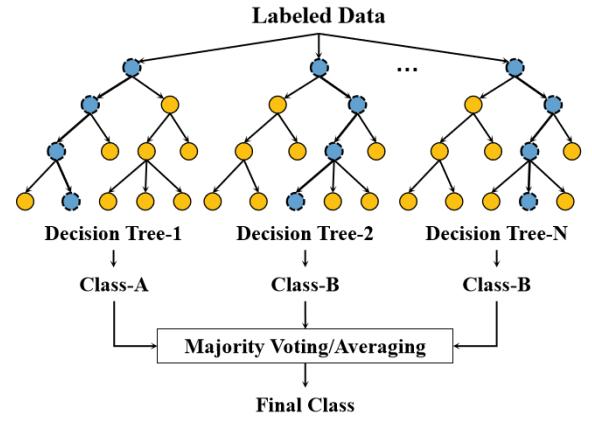


Fig. 6. The structure of the random forest (RF). The final category of input text is determined by multiple decision trees.

However, it helps little for short text. Motivated by this, Bouaziz et al. [43] combine data enrichment – with semantics in RFs for short text classification – to overcome the deficiency of sparseness and insufficiency of contextual information. In integrated algorithms, not all classifiers learn well. It is necessary to give different weights for each classifier. To differentiate contributions of trees in a forest, Islam et al. [44] exploit Semantics Aware Random Forest (SARF) classifier, choosing features similar to the features of the same class, for extracting features and producing the prediction values.

Summary. The shallow learning method is a type of machine learning. It learns from data, which are pre-defined features that are important to the performance of prediction values. However, feature engineering is tough work. Before training the classifier, we need to collect knowledge or experience to extract features from the original text. The shallow learning methods train the initial classifier based on various textual features extracted from the raw text. Toward small datasets, shallow learning models usually present better performance than deep learning models under the limitation of computational complexity. Therefore, some researchers have studied the design of shallow models for specific domains with fewer data.

B. Deep Learning Models

The DNNs consist of artificial neural networks that simulate the human brain to automatically learn high-level features from data, getting better results than shallow learning models in speech recognition, image processing, and text understanding. Input datasets should be analyzed to classify the data, such as a single-label, multi-label, unsupervised, unbalanced dataset. According to the trait of the dataset, the input word vectors are

sent into the DNN for training until the termination condition is reached. The performance of the training model is verified by the downstream task, such as sentiment classification, question answering, and event prediction. We show some DNNs over the years in Table I, including designs that are different from the corresponding basic models, evaluation metrics, and experimental datasets.

Numerous deep learning models have been proposed in the past few decades for text classification, as shown in Table I. We tabulate primary information – including publication years, venues, applications, code links, evaluation metrics, and experiment datasets – of main deep learning models for text classification. The applications in this table include sentiment analysis (SA), topic labeling (TL), news classification (NC), question answering (QA), dialog act classification (DAC), natural language inference (NLI), relation classification (RC) and event prediction (EP). The feed-forward neural network and the recursive neural network are the first two deep learning approaches used for the text classification task, which improve performance compared with shallow learning models. Then, CNNs, RNNs, and attention mechanisms are used for text classification. Many researchers advance text classification performance for different tasks by improving CNN, RNN, and attention, or model fusion and multi-task methods. The appearance of Bidirectional Encoder Representations from Transformers (BERT) [19], which can generate contextualized word vectors, is a significant turning point in the development of text classification and other NLP technologies. Many researchers have studied text classification models based on BERT, which achieves better performance than the above models in multiple NLP tasks, including text classification. Besides, some researchers study text classification technology based on GNN [6] to capture structural information in the text, which cannot be replaced by other methods. Here, we classify DNNs by structure and discuss some of the representative models in detail:

1) *ReNN-based Methods*: Shallow learning models cost lots of time on design features for each task. The recursive neural network (ReNN) can automatically learn the semantics of text recursively and the syntax tree structure without feature design, as shown in Fig. 7. We give an example of ReNN based models. First, each word of input text is taken as the leaf node of the model structure. Then all nodes are combined into parent nodes using a weight matrix. The weight matrix is shared across the whole model. Each parent node has the same dimension with all leaf nodes. Finally, all nodes are recursively aggregated into a parent node to represent the input text to predict the label.

ReNN-based models improve performance compared with shallow learning models and save on labor costs due to excluding feature designs used for different text classification tasks. The recursive autoencoder (RAE) [45] is used to predict the distribution of sentiment labels for each input sentence and learn the representations of multi-word phrases. To learn compositional vector representations for each input text, the matrix-vector recursive neural network (MV-RNN) [47] introduces a ReNN model to learn the representation of phrases and sentences. It allows that the length and type of input texts

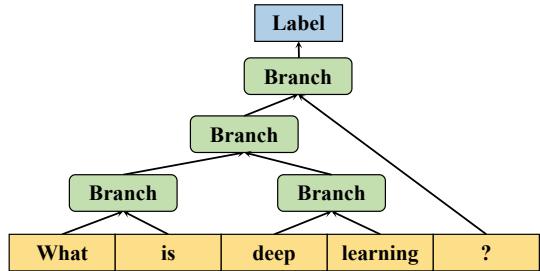


Fig. 7. The architecture of the recursive neural network (ReNN).

are inconsistent. MV-RNN allocates a matrix and a vector for each node on the constructed parse tree. Furthermore, the recursive neural tensor network (RNTN) [49] is proposed with a tree structure to capture the semantics of sentences. It inputs phrases with different length and represents the phrases by parse trees and word vectors. The vectors of higher nodes on the parse tree are estimated by the equal tensor-based composition function. For RNTN, the time complexity of building the textual tree is high, and expressing the relationship between documents is complicated within a tree structure. The performance is usually improved, with the depth being increased for DNNs. Therefore, Irsoy et al. [51] propose a deep recursive neural network (DeepReNN), which stacks multiple recursive layers. It is built by binary parse trees and learns distinct perspectives of compositionality in language.

2) *MLP-based Methods*: A multilayer perceptron (MLP) [135], sometimes colloquially called "vanilla" neural network, is a simple neural network structure that is used for capturing features automatically. As shown in Fig. 8, we show a three-layer MLP model. It contains an input layer, a hidden layer with an activation function in all nodes, and an output layer. Each node connects with a certain weight w_i . It treats each input text as a bag of words and achieves high performance on many text classification benchmarks comparing with shallow learning models.

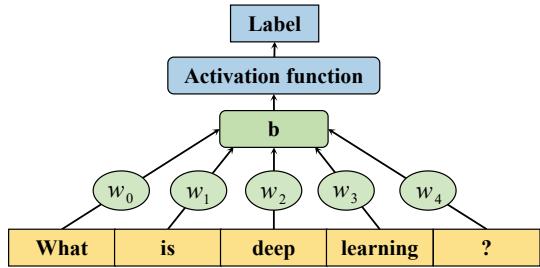


Fig. 8. The architecture of the multilayer perceptron (MLP).

There are some MLP-based methods proposed by some research groups for text classification tasks. The Paragraph Vector (Paragraph-Vec) [52] is the most popular and widely used method, which is similar to the Continuous Bag of Words (CBOW) [21]. It gets fixed-length feature representations of texts with various input lengths by employing unsupervised algorithms. Comparing with CBOW, it adds a paragraph token mapped to the paragraph vector by a matrix. The model predicts the fourth word by the connection or average of this

TABLE I
BASIC INFORMATION BASED ON DIFFERENT MODELS. TRANS: TRANSFORMER. TIME: TRAINING TIME.

| Model | Year | Method | Venue | Applications | Code Link | Metrics | Datasets |
|-----------|------|--------------------|------------|--------------|-----------|------------------------|-----------------------------|
| ReNN | 2011 | RAE [45] | EMNLP | SA, QA | [46] | Accuracy | MPQA, MR, EP |
| | 2012 | MV-RNN [47] | EMNLP | SA | [48] | Accuracy, F1 | MR |
| | 2013 | RNTN [49] | EMNLP | SA | [50] | Accuracy | SST |
| | 2014 | DeepRNN [51] | NIPS | SA;QA | - | Accuracy | SST-1;SST-2 |
| MLP | 2014 | Paragraph-Vec [52] | ICML | SA, QA | [53] | Error Rate | SST, IMDB |
| | 2015 | DAN [54] | ACL | SA, QA | [55] | Accuracy, Time | RT, SST, IMDB |
| | 2015 | Tree-LSTM [2] | ACL | SA | [56] | Accuracy | SST-1, SST-2 |
| | 2015 | S-LSTM [3] | ICML | SA | - | Accuracy | SST |
| RNN | 2015 | TextRCNN [57] | AAAI | SA, TL | [58] | <i>Macro-F1</i> , etc. | 20NG, Fudan, ACL, SST-2 |
| | 2015 | MT-LSTM [8] | EMNLP | SA,QA | [59] | Accuracy | SST-1, SST-2, QC, IMDB |
| | 2016 | oh-2LSTMp [60] | ICML | SA, TL | [61] | Error Rate | IMDB, Elec, RCV1, 20NG |
| | 2016 | BLSTM-2DCNN [62] | COLING | SA, QA, TL | [63] | Accuracy | SST-1, Subj, TREC, etc. |
| | 2016 | Multi-Task [64] | IJCAI | SA | [65] | Accuracy | SST-1, SST-2, Subj, IMDB |
| | 2017 | DeepMoji [66] | EMNLP | SA | [67] | Accuracy | SS-Twitter, SE1604, etc. |
| | 2017 | TopicRNN [68] | ICML | SA | [69] | Error Rate | IMDB |
| | 2017 | Miyato et al. [70] | ICLR | SA | [71] | Error Rate | IMDB, DBpedia, etc. |
| | 2018 | RNN-Capsule [72] | TheWebConf | SA | [73] | Accuracy | MR, SST-1, etc. |
| | 2014 | TextCNN [18] | EMNLP | SA, QA | [74] | Accuracy | MR, SST-2, Subj, etc. |
| CNN | 2014 | DCNN [7] | ACL | SA, QA | [75] | Accuracy | MR, TREC, Twitter |
| | 2015 | CharCNN [5] | NeurIPS | SA, QA, TL | [76] | Error Rate | AG, Yelp P, DBpedia, etc. |
| | 2016 | SeqTextRCNN [9] | NAACL | Dialog act | [77] | Accuracy | DSTC 4, MRDA, SwDA |
| | 2017 | XML-CNN [78] | SIGIR | NC, TL, SA | [79] | DCG@k, etc. | EUR-Lex, Wiki-30K, etc. |
| | 2017 | DPCNN [80] | ACL | SA, TL | [81] | Error Rate | AG, DBpedia, Yelp.P, etc. |
| | 2017 | KPCNN [82] | IJCAI | SA,QA,TL | - | Accuracy | Twitter, AG, Bing, etc. |
| | 2018 | TextCapsule [83] | EMNLP | SA, QA, TL | [84] | Accuracy | Subj, TREC, Reuters, etc. |
| | 2018 | HFT-CNN [85] | EMNLP | TL | [86] | <i>Micro-F1</i> , etc. | RCV1, Amazon670K |
| | 2020 | Bao et al. [87] | ICLR | TL | [88] | Accuracy | 20NG, Reuters-2157, etc. |
| | 2016 | HAN [89] | NAACL | SA, TL | [90] | Accuracy | Yelp.F, YahooA, etc. |
| Attention | 2016 | BI-Attention [91] | NAACL | SA | - | Accuracy | NLP&CC 2013 [92] |
| | 2016 | LSTMN [93] | EMNLP | SA | [94] | Accuracy | SST-1 |
| | 2017 | Lin et al. [95] | ICLR | SA | [96] | Accuracy | Yelp,SNLI Age |
| | 2018 | SGM [97] | COLING | TL | [98] | HL, <i>Micro-F1</i> | RCV1-V2, AAPD |
| | 2018 | ELMo [99] | NAACL | SA, QA, NLI | [100] | Accuracy | SQuAD, SNLI, SST-5 |
| | 2018 | BiBloSA [101] | ICLR | SA | [102] | Accuracy, Time | CR, MPQA, SUBJ, etc. |
| | 2019 | AttentionXML [103] | NeurIPS | TL | [104] | P@k, N@k, etc. | EUR-Lex, etc. |
| | 2019 | HAPN [105] | EMNLP | RC | - | Accuracy | FewRel, CSID |
| | 2019 | Proto-HATT [106] | AAAI | RC | [107] | Accuracy | FewRel |
| | 2019 | STCKA[108] | AAAI | SA;TL | [109] | Accuracy | Weibo, Product Review, etc. |
| Trans | 2019 | BERT [19] | NAACL | SA, QA | [110] | Accuracy | SST-2, QQP, QNLI, CoLA |
| | 2019 | BERT-BASE [111] | ACL | TL | [112] | P@K, R@K, etc. | EUR-LEX |
| | 2019 | Sun et al. [113] | CCL | SA, QA, TL | [114] | Error Rate | TREC, DBpedia, etc. |
| | 2019 | XLNet [115] | NeurIPS | SA, QA, NC | [116] | EM, F1, etc. | Yelp-2, AG, MNLI, etc. |
| | 2019 | RoBERTa [117] | arXiv | SA, QA | [118] | F1, Accuracy | SQuAD, MNLI-m, SST-2 |
| GNN | 2020 | ALBERT[119] | ICLR | SA, QA | [120] | F1, Accuracy | SST, MNLI, SQuAD |
| | 2018 | DGCNN [121] | TheWebConf | TL | [122] | <i>Macro-F1</i> , etc. | RCV1, NYTimes |
| | 2019 | TextGCN [6] | AAAI | SA, TL | [123] | Accuracy | 20NG, Ohsumed, R52, etc. |
| | 2019 | SGC[124] | ICML | NC, TL, SA | [125] | Accuracy, Time | 20NG, R8, Ohsumed, etc. |
| Others | 2019 | Huang et al. [126] | EMNLP | NC, TL | [127] | Accuracy | R8, R52, Ohsumed |
| | 2019 | Peng et al. [128] | arXiv | NC, TL | - | <i>Micro-F1</i> , etc. | RCV1, EUR-Lex, etc. |
| | 2020 | MAGNET [129] | ICAART | TL | [130] | <i>Micro-F1</i> , HL | Reuters, RCV1-V2, etc. |
| Others | 2017 | Miyato et al. [70] | ICLR | SA, NC | [131] | Error Rate | IMDB, RCV1, et al. |
| | 2018 | TMN [132] | EMNLP | TL | - | Accuracy, F1 | Snippets, Twitter, et al. |
| | 2019 | Zhang et al. [133] | NAACL | TL, NC | [134] | Accuracy | DBpedia, 20NG. |

vector to the three contexts of the word. Paragraph vectors can be used as a memory for paragraph themes and are used as a paragraph function and inserted into the prediction classifier.

3) *RNN-based Methods*: The recurrent neural network (RNN) is broadly used due to capturing long-range dependency through recurrent computation. The RNN language model learns historical information, considering the location information among all words suitable for text classification tasks. We show an RNN model for text classification with a simple sample, as shown in Fig. 9. Firstly, each input word is represented by a specific vector using a word embedding technology. Then, the embedding word vectors are fed into RNN cells one by one. The output of RNN cells are the same dimension with the input vector and are fed into the next hidden layer. The RNN shares parameters across different parts of the model and has the same weights of each input word. Finally, the label of input text can be predicted by the last output of the hidden layer.

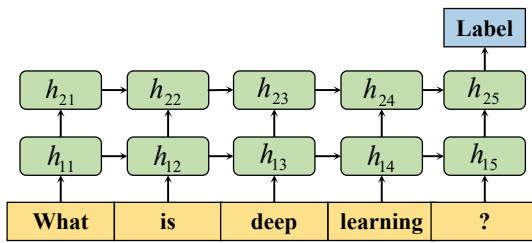


Fig. 9. The architecture of the recurrent neural network (RNN).

To diminish the time complexity of the model and capture contextual information, Liu et al. [64] introduce a model for catching the semantics of long texts. It parses the text one by one and is a biased model, making the following inputs profit over the former and decreasing the semantic efficiency of capturing the whole text. For modeling topic labeling tasks with long input sequences, TopicRNN [68] is proposed. It captures the dependencies of words in a document via latent topics and uses RNNs to capture local dependencies and latent topic models for capturing global semantic dependencies. Virtual Adversarial Training (VAT) [136] is a useful regularization method applicable to semi-supervised learning tasks. Miyato et al. [70] apply adversarial and virtual adversarial training to the text field and employ the perturbation into word embedding rather than the original input text. The model improves the quality of the word embedding and is not easy to overfit during training. Capsule network [137] captures the relationships between features using dynamic routing between capsules comprised of a group of neurons in a layer. Wang et al. [72] propose an RNN-Capsule model with a simple capsule structure for the sentiment classification task.

In the backpropagation process of RNN, the weights are adjusted by gradients, calculated by continuous multiplications of derivatives. If the derivatives are extremely small, it may cause a gradient vanishing problem by continuous multiplications. Long Short-Term Memory (LSTM) [138], the improvement of RNN, effectively alleviates the gradient vanishing problem. It is composed of a cell to remember values on arbitrary time intervals and three gate structures to control information flow.

The gate structures include input gates, forget gates, and output gates. The LSTM classification method can better capture the connection among context feature words, and use the forgotten gate structure to filter useless information, which is conducive to improving the total capturing ability of the classifier. Tree-LSTM [2] extends the sequence of LSTM models to the tree structure. The whole subtree with little influence on the result can be forgotten through the LSTM forgetting gate mechanism for the Tree-LSTM model.

Natural Language Inference (NLI) predicts whether one text's meaning can be deduced from another by measuring the semantic similarity between each pair of sentences. To consider other granular matchings and matchings in the reverse direction, Wang et al. [139] propose a model for the NLI task named Bilateral multi-perspective matching (BiMPM). It encodes input sentences by the BiLSTM encoder. Then, the encoded sentences are matched in two directions. The results are aggregated in a fixed-length matching vector by another BiLSTM layer. Finally, the result is evaluated by a fully connected layer.

4) *CNN-based Methods*: Convolutional neural networks (CNNs) are proposed for image classification with convolving filters that can extract features of pictures. Unlike RNN, CNN can simultaneously apply convolutions defined by different kernels to multiple chunks of a sequence. Therefore, CNNs are used for many NLP tasks, including text classification. For text classification, the text requires being represented as a vector similar to the image representation, and text features can be filtered from multiple angles, as shown in Fig. 10. Firstly, the word vectors of the input text are spliced into a matrix. The matrix is then fed into the convolutional layer, which contains several filters with different dimensions. Finally, the result of the convolutional layer goes through the pooling layer and concatenates the pooling result to obtain the final vector representation of the text. The category is predicted by the final vector.

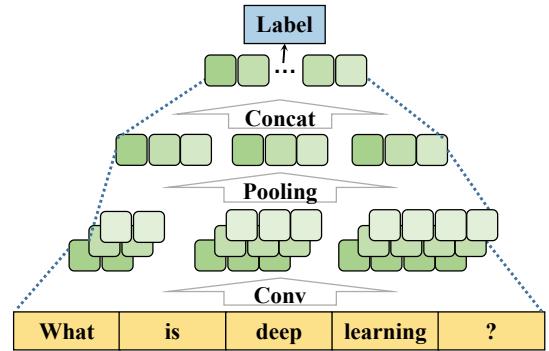


Fig. 10. The architecture of the Convolutional neural network (CNN).

To try using CNN for the text classification task, an unbiased model of convolutional neural networks is introduced by Kim, called TextCNN [18]. It can better determine discriminative phrases in the max-pooling layer with one layer of convolution and learn hyperparameters except for word vectors by keeping word vectors static. Training only on labeled data is not enough for data-driven deep models. Therefore, some

researchers consider utilizing unlabeled data. Johnson et al. [140] propose a text classification CNN model based on two-view semi-supervised learning, which first uses unlabeled data to train the embedding of text regions and then labeled data. DNNs usually have better performance, but it increases the computational complexity. Motivated by this, a deep pyramid convolutional neural network (DPCNN) [80] is proposed, with a little more computational accuracy, increasing by raising the network depth. The DPCNN is more specific than ResNet [141], as all the shortcuts are exactly simple identity mappings without any complication for dimension matching.

According to the minimum embedding unit of text, embedding methods are divided into character-level, word-level, and sentence-level embedding. Character-level embeddings can settle Out-of-Vocabulary (OOV) words. Word-level embeddings learn the syntax and semantics of the words. Moreover, sentence-level embedding can capture relationships among sentences. Motivated by these, Nguyen et al. [142] propose a deep learning method based on a dictionary, increasing information for word-level embeddings through constructing semantic rules and deep CNN for character-level embeddings. Adams et al. [143] propose a character-level CNN model, called MGTC, to classify multi-lingual texts written. TransCap [144] is proposed to encapsulate the sentence-level semantic representations into semantic capsules and transfer document-level knowledge.

RNN based models capture the sequential information to learn the dependency among input words, and CNN based models extract the relevant features from the convolution kernels. Thus some works study the fusion of the two methods. BLSTM-2DCNN [62] integrates a Bidirectional LSTM (BiLSTM) with two-dimensional max pooling. It uses a 2D convolution to sample more meaningful information of the matrix and understands the context better through BiLSTM. Moreover, Xue et al. [145] propose MTNA, a combination of BiLSTM and CNN layers, to solve aspect category classification and aspect term extraction tasks.

5) *Attention-based Methods*: CNN and RNN provide excellent results on tasks related to text classification. However, these models are not intuitive enough for poor interpretability, especially in classification errors, which cannot be explained due to the non-readability of hidden data. The attention-based methods are successfully used in the text classification. Bahdanau et al. [146] first propose an attention mechanism that can be used in machine translation. Motivated by this, Yang et al. [89] introduce the hierarchical attention network (HAN) to gain better visualization by employing the extremely informational components of a text, as shown in Fig. 11. HAN includes two encoders and two levels of attention layers. The attention mechanism lets the model pay different attention to specific inputs. It aggregates essential words into sentence vectors firstly and then aggregates vital sentence vectors into text vectors. It can learn how much contribution of each word and sentence for the classification judgment, which is beneficial for applications and analysis through the two levels of attention.

The attention mechanism can improve the performance with interpretability for text classification, which makes it popular.

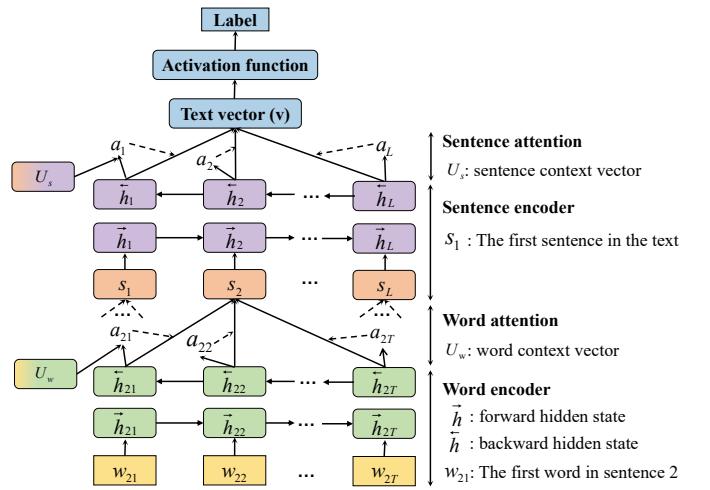


Fig. 11. The the architecture of hierarchical attention network (HAN) [89].

There are some other works based on attention. LSTMN [93] is proposed to process text step by step from left to right and does superficial reasoning through memory and attention. BI-Attention [91] is designed for cross-lingual text classification to catch bilingual long-distance dependencies. Hu et al. [147] propose an attention mechanism based on category attributes for solving the imbalance of the number of various charges which contain few-shot charges. HAPN [105] is presented for few-shot text classification.

Self-attention [148] captures the weight distribution of words in sentences by constructing K, Q and V matrices among sentences that can capture long-range dependencies on text classification. We give an example for self-attention, as shown in Fig. 12. Each input word vector a_i can be represented as three n-dimensional vectors, including q_i , k_i and v_i . After self-attention, the output vector b_i can be represented as $\sum_j softmax(a_{ij})v_j$ and $a_{ij} = q_i \cdot k_j / \sqrt{n}$. All output vectors can be parallelly computed. Lin et al. [95] used source token self-attention to explore the weight of every token to the entire sentence in the sentence representation task. To capture long-range dependencies, Bi-directional Block Self-Attention Network (Bi-BloSAN) [101] uses an intra-block self-attention network (SAN) to every block split by sequence and an inter-block SAN to the outputs.

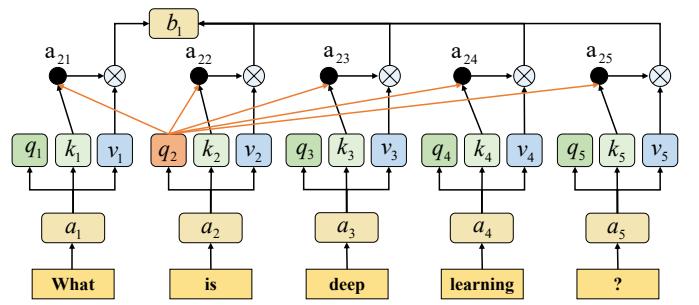


Fig. 12. An example of self-attention.

Aspect-based sentiment analysis (ABSA) breaks down a text into multiple aspects and allocates each aspect a sentiment

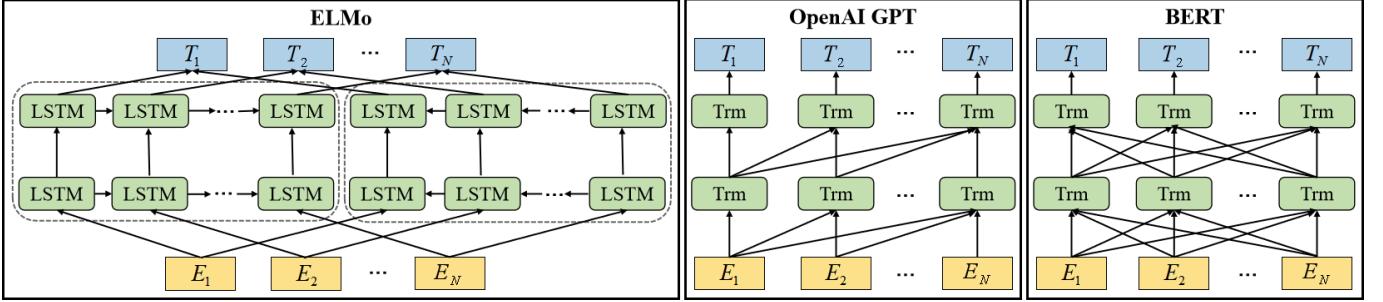


Fig. 13. Differences in pre-trained model architectures [19], including BERT, OpenAI GPT and ELMo. E_i represents embedding of i th input. Trm represents the transformer block. T_i represents predicted tag of i th input.

polarity. The sentiment polarity can be divided into two types: positive, neutral and negative. Some attention-based models are proposed to identify the fine-grained opinion polarity towards a specific aspect for aspect-based sentiment tasks. ATAE-LSTM [149] can concentrate on different parts of each sentence according to the input through the attention mechanisms. MGAN [150] presents a fine-grained attention mechanism with a coarse-grained attention mechanism to learn the word-level interaction between context and aspect.

To catch the complicated semantic relationship among each question and candidate answers for the QA task, Tan et al. [151] introduce CNN and RNN and generate answer embeddings by using a simple one-way attention mechanism affected through the question context. The attention captures the dependence among the embeddings of questions and answers. Extractive QA can be seen as the text classification task. It inputs a question and multiple candidates answers and classifies every candidate answer to recognize the correct answer. Furthermore, AP-BILSTM [152] with a two-way attention mechanism can learn the weights between the question and each candidate answer to obtain the importance of each candidate answer to the question.

6) Transformer-based Methods: Pre-trained language models effectively learn global semantic representation and significantly boost NLP tasks, including text classification. It generally uses unsupervised methods to mine semantic knowledge automatically and then construct pre-training targets so that machines can learn to understand semantics.

As shown in Fig. 13, we give differences in the model architectures among ELMo [99], OpenAI GPT [153], and BERT [19]. ELMo [99] is a deep contextualized word representation model, which is readily integrated into models. It can model complicated characteristics of words and learn different representations for various linguistic contexts. It learns each word embedding according to the context words with the bi-directional LSTM. GPT [153] employs supervised fine-tuning and unsupervised pre-training to learn general representations that transfer with limited adaptation to many NLP tasks. Furthermore, the domain of the target task does not need to be similar to the unlabeled datasets. The training procedure of the GPT algorithm usually includes two stages. Firstly, the initial parameters of a neural network model are learned by a modeling objective on the unlabeled dataset. We can then employ the corresponding supervised objective to

accommodate these parameters for the target task. To pre-train deep bidirectional representations from the unlabeled text through joint conditioning on both left and right context in every layer, BERT model [19], proposed by Google, significantly improves performance on NLP tasks, including text classification. It is fine-tuned by adding just an additional output layer to construct models for multiple NLP tasks, such as SA, QA, and machine translation. Comparing with these three models, ELMo is a feature-based method using LSTM, and BERT and OpenAI GPT are fine-tuning approaches using Transformer. Furthermore, ELMo and BERT are bidirectional training models and OpenAI GPT is training from left to right. Therefore, BERT gets a better result, which combines the advantages of ELMo and OpenAI GPT.

Transformer-based models can parallelize computation without considering the sequential information suitable for large scale datasets, making it popular for NLP tasks. Thus, some other works are used for text classification tasks and get excellent performance. RoBERTa [117] adopts the dynamic masking method that generates the masking pattern every time with a sequence to be fed into the model. It uses more data for longer pre-training and estimates the influence of various essential hyperparameters and the size of training data. ALBERT [119] uses two-parameter simplification schemes. In general, these methods adopt unsupervised objective functions for pre-training, including the next sentence prediction, masking technology, and permutation. These target functions based on the word prediction demonstrate a strong ability to learn the word dependence and semantic structure [154]. XLNet [115] is a generalized autoregressive pre-training approach. It maximizes the expected likelihood across the whole factorization order permutations to learn the bidirectional context. Furthermore, it can overcome the weaknesses of BERT by an autoregressive formulation and integrate ideas from Transformer-XL [155] into pre-training.

7) GNN-based Methods: The DNN models like CNN get great performance on regular structure, not for arbitrarily structured graphs. Some researchers study how to expand on arbitrarily structured graphs [156] [157]. With the increasing attention of graph neural networks (GNNs), GNN-based models obtain excellent performance by encoding syntactic structure of sentences on semantic role labeling task [158], relation classification task [159] and machine translation task [160]. It turns text classification into a graph node classification task.

We show a GCN model for text classification with four input texts, as shown in Fig. 14. Firstly, the four input texts and the words in the text, defined as nodes, are constructed into the graph structures. The graph nodes are connected by bold black edges, which indicates document-word edges and word-word edges. The weight of each word-word edge usually means their co-occurrence frequency in the corpus. Then, the words and texts are represented through the hidden layer. Finally, the label of all input texts can be predicted by the graph.

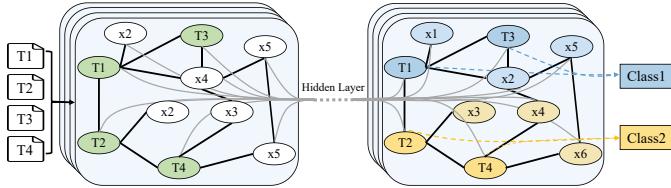


Fig. 14. The GCN based model. Black bold edges are document-word edges and word-word edges in the graph.

The GNN-based models can learn the syntactic structure of sentences making some researchers study using GNN for text classification. DGCNN [121] is a graph-CNN converting text to graph-of-words, having the advantage of learning different levels of semantics with CNN models. Yao et al. [6] propose the text graph convolutional network (TextGCN), which builds a heterogeneous word text graph for a whole dataset and captures global word co-occurrence information. To enable GNN-based models to underpin online testing, Huang et al. [126] build graphs for each text with global parameter sharing, not a corpus-level graph structure, to help preserve global information and reduce the burden. TextING [161] builds individual graphs for each document and learns text-level word interactions by GNN to effectively produce embeddings for obscure words in the new text.

Graph attention networks (GATs) [162] employ masked self-attention layers by attending over its neighbors. Thus, some GAT-based models are proposed to compute the hidden representations of each node. The heterogeneous graph attention network (HGAT) [163] with a dual-level attention mechanism learns the importance of different neighboring nodes and node types in the current node. The model propagates information on the graph and captures the relations to address the semantic sparsity for semi-supervised short text classification. MAGNET [129] is proposed to capture the correlation among the labels based on GATs, which learns the crucial dependencies between the labels and generates classifiers by a feature matrix and a correlation matrix.

Event prediction (EP) can be divided into generated event prediction and selective event prediction (also known as script event prediction). EP, referring to scripted event prediction in this review, infers the subsequent event according to the existing event context. Unlike other text classification tasks, texts in EP are composed of a series of sequential subevents. Extracting features of the relationship among such subevents is of critical importance. SGNN [164] is proposed to model event interactions and learn better event representations by constructing an event graph to utilize the event network

information better. The model makes full use of dense event connections for the EP task.

8) *Others*: In addition to all the above models, there are some other individual models. Here we introduce some exciting models.

a) *Siamese neural network*.: The siamese neural network [165] is also called a twin neural network (Twin NN). It utilizes equal weights while working in tandem using two distinct input vectors to calculate comparable output vectors. Mueller et al. [166] present a siamese adaptation of the LSTM network comprised of couples of variable-length sequences. The model is employed to estimate the semantic similarity among texts, exceeding carefully handcrafted features and proposed neural network models of higher complexity. The model further represents text employing neural networks whose inputs are word vectors learned separately from a vast dataset. To settle unbalanced data classification in the medical domain, Jayadeva et al. [167] use a Twin NN model to learn from enormous unbalanced corpora. The objective functions achieve the Twin SVM approach with non-parallel decision boundaries for the corresponding classes, and decrease the Twin NN complexity, optimizing the feature map to better discriminate among classes.

b) *Virtual adversarial training (VAT)*: Deep learning methods require many extra hyperparameters, which increase the computational complexity. VAT [168], regularization based on local distributional smoothness can be used in semi-supervised tasks, requires only a small number of hyperparameters, and can be interpreted directly as robust optimization. Miyato et al. [70] use VAT to effectively improve the robustness and generalization ability of the model and word embedding performance.

c) *Reinforcement learning (RL)*: RL learns the best action in a given environment through maximizing cumulative rewards. Zhang et al. [169] offer an RL approach to establish structured sentence representations via learning the structures related to tasks. The model has Information Distilled LSTM (ID-LSTM) and Hierarchical Structured LSTM (HS-LSTM) representation models. The ID-LSTM learns the sentence representation by choosing essential words relevant to tasks, and the HS-LSTM is a two-level LSTM for modeling sentence representation.

d) *Memory Networks*: Memory networks [170] learn to combine the inference components and the long-term memory component. Li et al. [171] use two LSTMs with extended memories and neural memory operations for jointly handling the extraction tasks of aspects and opinions via memory interactions. Topic memory networks (TMN) [132] is an end-to-end model that encodes latent topic representations indicative of class labels.

e) *QA style for the sentiment classification task*.: It is an interesting attempt to treat the sentiment classification task as a QA task. Shen et al. [172] create a high-quality annotated corpus. A three-stage hierarchical matching network was proposed to consider the matching information between questions and answers.

f) *External commonsense knowledge*.: Due to the insufficient information of the event itself to distinguish the event for

the EP task, Ding et al. [173] consider that the event extracted from the original text lacked common knowledge, such as the intention and emotion of the event participants. The model improves the effect of stock prediction, EP, and so on.

g) Quantum language model.: In the quantum language model, the words and dependencies among words are represented through fundamental quantum events. Zhang et al. [174] design a quantum-inspired sentiment representation method to learn both the semantic and the sentiment information of subjective text. By inputting density matrices to the embedding layer, the performance of the model improves.

Summary. Deep Learning consists of multiple hidden layers in a neural network with a higher level of complexity and can be trained on unstructured data. Deep learning architecture can learn feature representations directly from the input without too many manual interventions and prior knowledge. However, deep learning technology is a data-driven method, which usually needs enormous data to achieve high performance. Although self-attention based models can bring some interpretability among words for DNNs, it is not enough comparing with shallow models to explain why and how it works well.

III. DATASETS AND EVALUATION METRICS

A. Datasets

The availability of labeled datasets for text classification has become the main driving force behind the fast advancement of this research field. In this section, we summarize the characteristics of these datasets in terms of domains and give an overview in Table II, including the number of categories, average sentence length, the size of each dataset, related papers, data sources to access and applications.

Sentiment Analysis (SA). SA is the process of analyzing and reasoning the subjective text within emotional color. It is crucial to get information on whether it supports a particular point of view from the text that is distinct from the traditional text classification that analyzes the objective content of the text. SA can be binary or multi-class. Binary SA is to divide the text into two categories, including positive and negative. Multi-class SA classifies text to multi-level or fine-grained labels. The SA datasets include MR, SST, MPQA, IMDB, Yelp, AM, Subj [196], CR [198], SS-Twitter, SS-Youtube, Twitter, SE1604, EP and so on. Here we detail several datasets.

Movie Review (MR). The MR [215] [175] is a movie review dataset, each of which corresponds to a sentence. The corpus has 5,331 positive data and 5,331 negative data. 10-fold cross-validation by random splitting is commonly used to test MR.

Stanford Sentiment Treebank (SST). The SST [176] is an extension of MR. It has two categories. SST-1 with fine-grained labels with five classes. It has 8,544 training texts and 2,210 test texts, respectively. Furthermore, SST-2 has 9,613 texts with binary labels being partitioned into 6,920 training texts, 872 development texts, and 1,821 testing texts.

The Multi-Perspective Question Answering (MPQA). The MPQA [216] [178] is an opinion dataset. It has two class labels and also an MPQA dataset of opinion polarity detection sub-tasks. MPQA includes 10,606 sentences extracted from

news articles from various news sources. It should be noted that it contains 3,311 positive texts and 7,293 negative texts without labels of each text.

IMDB reviews. The IMDB review [180] is developed for binary sentiment classification of film reviews with the same amount in each class. It can be separated into training and test groups on average, by 25,000 comments per group.

Yelp reviews. The Yelp review [181] is summarized from the Yelp Dataset Challenges in 2013, 2014, and 2015. This dataset has two categories. Yelp-2 of these were used for negative and positive emotion classification tasks, including 560,000 training texts and 38,000 test texts. Yelp-5 is used to detect fine-grained affective labels with 650,000 training and 50,000 test texts in all classes.

Amazon Reviews (AM). The AM [5] is a popular corpus formed by collecting Amazon website product reviews [182]. This dataset has two categories. The Amazon-2 with two classes includes 3,600,000 training sets and 400,000 testing sets. Amazon-5, with five classes, includes 3,000,000 and 650,000 comments for training and testing.

News Classification (NC). News content is one of the most crucial information sources which has a critical influence on people. The NC system facilitates users to get vital knowledge in real-time. News classification applications mainly encompass: recognizing news topics and recommending related news according to user interest. The news classification datasets include 20NG, AG, R8, R52, Sogou, and so on. Here we detail several datasets.

20 Newsgroups (20NG). The 20NG [185] is a newsgroup text dataset. It has 20 categories with the same number of each category and includes 18,846 texts.

AG News (AG). The AG News [5] [186] is a search engine for news from academia, choosing the four largest classes. It uses the title and description fields of each news. AG contains 120,000 texts for training and 7,600 texts for testing.

R8 and R52. R8 and R52 are two subsets which are the subset of Reuters [187]. R8 [188] has 8 categories, divided into 2,189 test files and 5,485 training courses. R52 has 52 categories, split into 6,532 training files and 2,568 test files.

Sogou News (Sogou). The Sogou News [113] combines two datasets, including SogouCA and SogouCS news sets. The label of each text is the domain names in the URL.

Topic Labeling (TL). The topic analysis attempts to get the meaning of the text by defining the sophisticated text theme. The topic labeling is one of the essential components of the topic analysis technique, intending to assign one or more subjects for each document to simplify the topic analysis. The topic labeling datasets include DBPedia, Ohsmed, EUR-Lex, WOS, PubMed, and YahooA. Here we detail several datasets.

DBpedia. The DBpedia [192] is a large-scale multi-lingual knowledge base generated using Wikipedia's most ordinarily used infoboxes. It publishes DBpedia each month, adding or deleting classes and properties in every version. DBpedia's most prevalent version has 14 classes and is divided into 560,000 training data and 70,000 test data.

Ohsmed. The Ohsmed [193] belongs to the MEDLINE database. It includes 7,400 texts and has 23 cardiovascular

TABLE II
SUMMARY STATISTICS FOR THE DATASETS. C: NUMBER OF TARGET CLASSES. L: AVERAGE SENTENCE LENGTH. N: DATASET SIZE.

| Datasets | C | L | N | Related Papers | Sources | Applications |
|---------------------|-------|-------|-----------|-------------------------------|---------|--------------|
| MR | 2 | 20 | 10,662 | [18] [7] [83] [6] | [175] | SA |
| SST-1 | 5 | 18 | 11,855 | [49] [18] [2] [3][93] | [176] | SA |
| SST-2 | 2 | 19 | 9,613 | [49] [18] [8] [64] [19] | [177] | SA |
| MPQA | 2 | 3 | 10,606 | [45] [18] [101] | [178] | SA |
| Twitter | 3 | 19 | 11,209 | [7][82] | [179] | SA |
| IMDB | 2 | 294 | 50,000 | [52] [54] [8] [64] [70] [115] | [180] | SA |
| Yelp.P | 2 | 153 | 598,000 | [5] [80] | [181] | SA |
| Yelp.F | 5 | 155 | 700,000 | [5] [89] [80] | [181] | SA |
| Amz.P | 2 | 91 | 4,000,000 | [103] [5] | [182] | SA |
| Amz.F | 5 | 93 | 3,650,000 | [5] [89] [103] | [182] | SA |
| RCV1 | 103 | 240 | 807,595 | [60] [85] [121] [111] [129] | [183] | NC |
| RCV1-V2 | 103 | 124 | 804,414 | [97] [129] | [184] | NC |
| 20NG | 20 | 221 | 18,846 | [57] [60] [87] [6] [124] | [185] | NC |
| AG News | 4 | 45/7 | 127,600 | [5] [80] [82] [83] [115] | [186] | NC |
| Reuters | 90 | 1 | 10,788 | [83] [129] | [187] | NC |
| R8 | 8 | 66 | 7,674 | [6] [124] [126] | [188] | NC |
| R52 | 52 | 70 | 9,100 | [6] [124] [126] | [188] | NC |
| NYTimes | 2,318 | 629 | 1,855,659 | [121] [189] [190] | [191] | NC |
| DBpedia | 14 | 55 | 630,000 | [5] [80] [70] [113] | [192] | TL |
| YahooA | 10 | 112 | 1,460,000 | [5] [89] | [5] | TL |
| AAPD | 54 | 163 | 55,840 | [97] [129] | [98] | TL |
| Ohsumed | 23 | 136 | 7,400 | [6] [124] [126] | [193] | TL |
| Amazon670K | 670 | 244 | 643,474 | [85] [103] | [194] | TL |
| EUR-Lex | 3,956 | 1,239 | 19,314 | [78] [103] [111] [128] [111] | [195] | TL |
| Subj | 2 | 23 | 10,000 | [18] [64] [83] | [196] | QA |
| TREC | 6 | 10 | 5,952 | [18] [7] [8] [82] | [197] | QA |
| CR | 2 | 19 | 3,775 | [18] [83] | [198] | QA |
| SQuAD | - | 5,000 | 5,570 | [99] [99] [117] [119] | [199] | QA |
| WikiQA | - | 873 | 243 | [200] [152] | [200] | QA |
| DSTC 4 | 89 | - | 30,000 | [201] [9] | [201] | DAC |
| MRDA | 5 | - | 62,000 | [9] [202] [203] | [204] | DAC |
| SwDA | 43 | - | 1,022,000 | [9] [202] [203] | [205] | DAC |
| FewRel | 100 | 25 | 70,000 | [206] [105] [207] | [208] | RC |
| SemEval-2010 Task 8 | 9 | - | 10,717 | [209] [210] [211] [212] [213] | [214] | RC |

disease categories. All texts are medical abstracts and are labeled into one or more classes.

Yahoo answers (YahooA). The YahooA [5] is a topic labeling task with 10 classes. It includes 140,000 training data and 5,000 test data. All text contains three elements, being question titles, question contexts, and best answers, respectively.

Question Answering (QA). The QA task can be divided into two types: the extractive QA and the generative QA. The extractive QA gives multiple candidate answers for each question to choose which one is the right answer. Thus, the text classification models can be used for the extractive QA task. The QA discussed in this paper is all extractive QA. The QA system can apply the text classification model to recognize the correct answer and set others as candidates. The question answering datasets include SQuAD, MS MARCO, TREC-QA, WikiQA, and Quora [217]. Here we detail several datasets.

Stanford Question Answering Dataset (SQuAD). The SQuAD [199] is a set of question and answer pairs obtained from Wikipedia articles. The SQuAD has two categories. SQuAD1.1 contains 536 pairs of 107,785 Q&A items. SQuAD2.0 combines 100,000 questions in SQuAD1.1 with more than 50,000 unanswerable questions that crowd workers face in a form similar to answerable questions [218].

MS MARCO. The MS MARCO [219] contains questions and answers. The questions and part of the answers are sampled from actual web texts by the Bing search engine. Others are generative. It is used for developing generative QA systems released by Microsoft.

TREC-QA. The TREC-QA [197] includes 5,452 training texts and 500 testing texts. It has two versions. TREC-6 contains 6 categories, and TREC-50 has 50 categories.

WikiQA. The WikiQA dataset [200] includes questions with no correct answer, which needs to evaluate the answer.

Natural Language Inference (NLI). NLI is used to predict whether the meaning of one text can be deduced from another. Paraphrasing is a generalized form of NLI. It uses the task of measuring the semantic similarity of sentence pairs to decide whether one sentence is the interpretation of another. The NLI datasets include SNLI, MNLI, SICK, STS, RTE, SciTail, MSRP, etc. Here we detail several of the primary datasets.

The Stanford Natural Language Inference (SNLI). The SNLI [220] is generally applied to NLI tasks. It contains 570,152 human-annotated sentence pairs, including training, development, and test sets, which are annotated with three categories: neutral, entailment, and contradiction.

Multi-Genre Natural Language Inference (MNLI). The Multi-NLI [221] is an expansion of SNLI, embracing a broader scope of written and spoken text genres. It includes 433,000 sentence pairs annotated by textual entailment labels.

Sentences Involving Compositional Knowledge (SICK). The SICK [222] contains almost 10,000 English sentence pairs. It consists of neutral, entailment and contradictory labels.

Microsoft Research Paraphrase (MSRP). The MSRP [223] consists of sentence pairs, usually for the text-similarity task. Each pair is annotated by a binary label to discriminate whether they are paraphrases. It respectively includes 1,725 training and 4,076 test sets.

Dialog Act Classification (DAC). A dialog act describes an utterance in a dialog based on semantic, pragmatic, and syntactic criteria. DAC labels a piece of a dialog according to its category of meaning and helps learn the speaker's intentions. It is to give a label according to dialog, such as DSTC 4 [201], MRDA [204], and SwDA [205].

Dialog State Tracking Challenge 4 (DSTC 4). The DSTC 4 [201] is used for dialog act classification. It has 89 training classes, 24,000 training texts, and 6,000 testing texts.

ICSI Meeting Recorder Dialog Act (MRDA). The MRDA [204] is used for dialog act classification. It has 5 training classes, 51,000 training texts, 11,000 testing texts, and 11,000 validation texts.

Switchboard Dialog Act (SwDA). The SwDA [205] is used for dialog act classification. It has 43 training classes, 1,003,000 training texts, 19,000 testing texts and 112,000 validation texts.

Relation Classification (RC). RC – extract relations between entities – is a crucial NLP task. It relies on information of both the sentence and the two target entities. It includes SemEval-2010 Task 8, ACE 2003-2004, TACRED, NYT-10 and FewRel. Here we detail FewRel dataset.

Few-Shot Relation Classification (FewRel). The FewRel [206] is a few-shot relation extraction dataset, which has two versions. It has 70000 samples and contains 64 relations on training data, 16 relations for validation and 20 relations on test set separately.

Multi-label datasets. In multi-label classification, an instance has multiple labels, and each label can only take one of the multiple classes. There are many datasets based on multi-label text classification. It includes Reuters, Education, Patent,

RCV1, RCV1-2K, AmazonCat-13K, BlurbGenreCollection, WOS-11967, AAPD, etc. Here we detail several datasets.

Reuters news. The Reuters [187] is a popularly used dataset for text classification from Reuters financial news services. It has 90 training classes, 7,769 training texts, and 3,019 testing texts, containing multiple labels and single labels. There are also some Reuters sub-sets of data, such as R8, BR52, RCV1, and RCV1-v2.

Patent Dataset. The Patent Dataset is obtained from USPTO ¹, which is a patent system grating U.S. patents containing textual details such title and abstract. It contains 100,000 US patents awarded in the real-world with multiple hierarchical categories.

Reuters Corpus Volume I (RCV1) and RCV1-2K. The RCV1 [183] is collected from Reuters News articles from 1996-1997, which is human-labeled with 103 categories. It consists of 23,149 training and 784,446 testing texts, respectively. The RCV1-2K dataset has the same features as the RCV1. However, the label set of RCV1-2K has been expanded with some new labels. It contains 2456 labels.

Web of Science (WOS-11967). The WOS-11967 [224] is crawled from the Web of Science, consisting of abstracts of published papers with two labels for each example. It is shallower, but significantly broader, with fewer classes in total.

Arxiv Academic Paper Dataset (AAPD). The AAPD [98] is a large dataset in the computer science field for the multi-label text classification from website ². It has 55,840 papers, including the abstract and the corresponding subjects with 54 labels in total. The aim is to predict the corresponding subjects of each paper according to the abstract.

Others. There are some datasets for other applications, such as Geonames toponyms, Twitter posts, and so on.

B. Evaluation Metrics

In terms of evaluating text classification models, accuracy and F1 score are the most used to assess the text classification methods. Later, with the increasing difficulty of classification tasks or the existence of some particular tasks, the evaluation metrics are improved. For example, evaluation metrics such as P@K and Micro-F1 are used to evaluate multi-label text classification performance, and MRR is usually used to estimate the performance of QA tasks. In Table III, we give the notations used in evaluation metrics.

1) Single-label metrics: Single-label text classification divides the text into one of the most likely categories applied in NLP tasks such as QA, SA, and dialogue systems [9]. For single-label text classification, one text belongs to just one catalog, making it possible not to consider the relations among labels. Here we introduce some evaluation metrics used for single-label text classification tasks.

Accuracy and Error Rate. Accuracy and Error Rate are the fundamental metrics for a text classification model. The Accuracy and Error Rate are respectively defined as

$$\text{Accuracy} = \frac{(TP + TN)}{N},$$

¹<https://www.uspto.gov/>

²<https://arxiv.org/>

TABLE III
THE NOTATIONS USED IN EVALUATION METRICS.

| Notations | Descriptions |
|---------------|--|
| TP | true positive |
| FP | false positive |
| TN | true negative |
| FN | false negative |
| TP_t | true positive of the t th label on a text |
| FP_t | false positive of the t th label on a text |
| TN_t | true negative of the t th label on a text |
| FN_t | false negative of the t th label on a text |
| \mathcal{S} | label set of all samples |
| Q | the number of predicted labels on each text |

$$\text{ErrorRate} = 1 - \text{Accuracy} = \frac{(FP + FN)}{N}.$$

Precision, Recall and F1. These are vital metrics utilized for unbalanced test sets regardless of the standard type and error rate. For example, most of the test samples have a class label. F1 is the harmonic average of Precision and Recall. Accuracy, Recall, and F1 as defined

$$\begin{aligned} \text{Precision} &= \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \\ F1 &= \frac{2\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}. \end{aligned}$$

The desired results will be obtained when the accuracy, F1 and recall value reach 1. On the contrary, when the values become 0, the worst result is obtained. For the multi-class classification problem, the precision and recall value of each class can be calculated separately, and then the performance of the individual and whole can be analyzed.

Exact Match (EM). The EM is a metric for QA tasks measuring the prediction that matches all the ground-truth answers precisely. It is the primary metric utilized on the SQuAD dataset.

Mean Reciprocal Rank (MRR). The MRR is usually applied for assessing the performance of ranking algorithms on QA and Information Retrieval (IR) tasks. MRR is defined as

$$\text{MRR} = \frac{1}{Q} \sum_{i=1}^Q \frac{1}{\text{rank}_i}.$$

where rank_i is the ranking of the ground-truth answer at answer i th.

Hamming-loss (HL). The HL [225] assesses the score of misclassified instance-label pairs where a related label is omitted or an unrelated is predicted.

Among these single-label evaluation metrics, the Accuracy is the earliest metric that calculates the proportion of the sample size that is predicted correctly and is not considered whether the predicted sample is a positive or a negative sample. Precision calculates how many of the positive samples are actually positive, and the Recall calculates how many of the positive examples in the sample are predicted correctly.

Furthermore, F1 is the harmonic average of them, which is the most commonly used evaluation metrics.

2) Multi-label metrics: Compared with single-label text classification, multi-label text classification divides the text into multiple category labels, and the number of category labels is variable. These metrics are designed for single label text classification, which are not suitable for multi-label tasks. Thus, there are some metrics designed for multi-label text classification.

Micro – F1. The *Micro – F1* [226] is a measure that considers the overall accuracy and recall of all labels. The *Micro – F1* is defined as:

$$\text{Micro} - F1 = \frac{2P_t \times R_t}{P + R},$$

where:

$$P = \frac{\sum_{t \in \mathcal{S}} TP_t}{\sum_{t \in \mathcal{S}} TP_t + FP_t}, \quad R = \frac{\sum_{t \in \mathcal{S}} TP_t}{\sum_{t \in \mathcal{S}} TP_t + FN_t}.$$

Macro – F1. The *Macro – F1* calculates the average *F1* of all labels. Unlike *Micro – F1*, which sets even weight to every example, *Macro – F1* sets the same weight to all labels in the average process. Formally, *Macro – F1* is defined as:

$$\text{Macro} - F1 = \frac{1}{|\mathcal{S}|} \sum_{t \in \mathcal{S}} \frac{2P_t \times R_t}{P_t + R_t},$$

where:

$$P_t = \frac{TP_t}{TP_t + FP_t}, \quad R_t = \frac{TP_t}{TP_t + FN_t}.$$

In addition to the above evaluation metrics, there are some rank-based evaluation metrics for extreme multi-label classification tasks, including P@K and NDCG@K.

Precision at Top K (P@K). The *P@K* is the precision at the top k. For *P@K*, each text has a set of \mathcal{L} ground truth labels $L_t = \{l_0, l_1, l_2, \dots, l_{\mathcal{L}-1}\}$, in order of decreasing probability $P_t = [p_0, p_1, p_2, \dots, p_{\mathcal{L}-1}]$. The precision at k is

$$\begin{aligned} \text{P}@K &= \frac{1}{k} \sum_{j=0}^{\min(\mathcal{L}, k)-1} \text{rel}_{L_i}(P_t(j)), \\ \text{rel}_L(p) &= \begin{cases} 1 & \text{if } p \in L \\ 0 & \text{otherwise} \end{cases}. \end{aligned}$$

where \mathcal{L} is the number of ground truth labels or possible answers on each text and k is the number of selected labels on extreme multi-label text classification.

Normalized Discounted Cummulated Gains (NDCG@K). The NDCG at k is

$$\text{NDCG}@K = \frac{1}{\text{IDCG}(L_i, k)} \sum_{j=0}^{n-1} \frac{\text{rel}_{L_i}(P_t(j))}{\ln(j+1)},$$

where

$$n = \min(\max(|P_i|, |L_i|), k).$$

Among these multi-label evaluation metrics, *Micro – F1* considers the number of categories, which makes it suitable for the unbalanced data distribution. *Macro – F1* does not

take into account the amount of data that treats each class equally. Thus it is easily affected by the classes with high Recall and Precision. When the number of categories is or extremely large, either P@K or NDCG@K is used.

IV. QUANTITATIVE RESULTS

In this section, we tabulate the performance of the main models given in their articles on classic datasets evaluated by classification accuracy, as shown in Table IV, including MR, SST-2, IMDB, Yelp.P, Yelp.F, Amazon.F, 20NG, AG, DBpedia and SNLI. We can see that BERT based models get better results on most datasets, which means that if you need to implement a text classification task, you can try BERT based models firstly, except MR and 20NG, which have not been experimented on BERT based models. RNN-Capsule [72] obtains the best result on MR and BLSTM-2DCNN [62] gets the best on 20NG.

V. FUTURE RESEARCH CHALLENGES

Text classification – as efficient information retrieval and mining technology – plays a vital role in managing text data. It uses NLP, data mining, machine learning, and other techniques to automatically classify and discover different text types. Text classification takes multiple types of text as input, and the text is represented as a vector by the pre-training model. Then the vector is fed into the DNN for training until the termination condition is reached, and finally, the performance of the training model is verified by the downstream task. Existing models have already shown their usefulness in text classification, but there are still many possible improvements to explore.

Although some new text classification models repeatedly brush up the accuracy index of most classification tasks, it cannot indicate whether the model “understands” the text from the semantic level like human beings. Moreover, with the emergence of the noise sample, the small sample noise may cause the decision confidence to change substantially or even lead to decision reversal. Therefore, the semantic representation ability and robustness of the model need to be proved in practice. Besides, the pre-trained semantic representation model represented by word vectors can often improve the performance of downstream NLP tasks. The existing research on the transfer strategy of context-free word vectors is still relatively preliminary. Thus, we conclude from data, models, and performance perspective, text classification mainly faces the following challenges.

A. Data

For a text classification task, data is essential to model performance, whether it is shallow learning or deep learning method. The text data mainly studied includes multi-chapter, short text, cross-language, multi-label, less sample text. For the characteristics of these data, the existing technical challenges are as follows:

Zero-shot/Few-shot learning. Zero-shot or few-shot learning for text classification aim to classify text having no or few

same labeled class data. However, the current models are too dependent on numerous labeled data. The performance of these models is significantly affected by zero-shot or few-shot learning. Thus, some works focus on tackling these problems. The main idea is to infer the features through learning kinds of semantic knowledge, such as learning relationship among classes [231] and incorporating class descriptions [133]. Furthermore, latent features generation [232] meta-Learning [233] [234] [87] and dynamic memory mechanism [235] are also efficient methods. Nevertheless, with the limitation of little unseen class data and different data distribution between seen class and unseen class, there is still a long way to go to reach the learning ability comparable to that of humans.

The external knowledge. As we all know, the more beneficial information is input into a DNN, its better performance. Therefore, adding external knowledge (knowledge base or knowledge graph) is an efficient way to promote the model’s performance. The existing knowledge includes conceptual information [82] [108] [236], commonsense knowledge [173], knowledge base information [237] [238], general knowledge graph [133] and so on, which enhances the semantic representation texts. Nevertheless, with the limitation of input scale, how and what to add for different tasks is still a challenge.

The multi-label text classification task. Multi-label text classification requires full consideration of the semantic relationship among labels, and the embedding and encoding of the model is a process of lossy compression. Therefore, how to reduce the loss of hierarchical semantics and retain rich and complex document semantic information during training is still a problem to be solved.

Special domain with many terminologies. Texts in a particular field, such as financial and medical texts, contain many specific words or domain experts intelligible slang, abbreviations, etc., which make the existing pre-trained word vectors challenging to work on.

B. Models

Most existing structures of shallow and deep learning models are tried for text classification, including integration methods. BERT learns a language representation that can be used to fine-tune for many NLP tasks. The primary method is to increase data, improve computation power, and design training procedures for getting better results. How to tradeoff between data and compute resources and prediction performance is worth studying.

C. Performance

The shallow model and the deep model can achieve good performance in most text classification tasks, but the anti-interference ability of their results needs to be improved. How to realize the interpretation of the deep model is also a technical challenge.

The semantic robustness of the model. In recent years, researchers have designed many models to enhance the accuracy of text classification models. However, when there are some adversarial samples in the datasets, the model’s

TABLE IV
ACCURACY OF DEEP LEARNING-BASED TEXT CLASSIFICATION MODELS ON PRIMARY DATASETS EVALUATED BY CLASSIFICATION ACCURACY (IN TERMS OF PUBLICATION YEAR). BOLD IS THE MOST ACCURATE.

| Model | Sentiment | | | | | | News | | Topic | NLI |
|--------------------|-------------|-----------|--------------|--------------|-------------|--------------|-------------|--------------|--------------|-------------|
| | MR | SST-2 | IMDB | Yelp.P | Yelp.F | Amz.F | 20NG | AG | DBpedia | SNLI |
| RAE [45] | 77.7 | 82.4 | - | - | - | - | - | - | - | - |
| MV-RNN [47] | 79 | 82.9 | - | - | - | - | - | - | - | - |
| RNTN [49] | 75.9 | 85.4 | - | - | - | - | - | - | - | - |
| DCNN [7] | | 86.8 | 89.4 | - | - | - | - | - | - | - |
| Paragraph-Vec [52] | | 87.8 | 92.58 | - | - | - | - | - | - | - |
| TextCNN[18] | 81.5 | 88.1 | - | - | - | - | - | - | - | - |
| TextRCNN [57] | - | - | - | - | - | - | 96.49 | - | - | - |
| DAN [54] | - | 86.3 | 89.4 | - | - | - | - | - | - | - |
| Tree-LSTM [2] | | 88 | - | - | - | - | - | - | - | - |
| CharCNN [5] | - | - | - | 95.12 | 62.05 | - | - | 90.49 | 98.45 | - |
| HAN [89] | - | - | 49.4 | - | - | 63.6 | - | - | - | - |
| SeqTextRCNN [9] | - | - | - | - | - | - | - | - | - | - |
| oh-2LSTMP [60] | - | - | 94.1 | 97.1 | 67.61 | - | 86.68 | 93.43 | 99.16 | - |
| LSTMN [93] | - | 87.3 | - | - | - | - | - | - | - | - |
| Multi-Task [64] | - | 87.9 | 91.3 | - | - | - | - | - | - | - |
| BLSTM-2DCNN [62] | 82.3 | 89.5 | - | - | - | - | 96.5 | - | - | - |
| TopicRNN [68] | - | - | 93.72 | - | - | - | - | - | - | - |
| DPCNN [80] | - | - | - | 97.36 | 69.42 | 65.19 | - | 93.13 | 99.12 | - |
| KPCNN [82] | 83.25 | - | - | - | - | - | - | 88.36 | - | - |
| RAM [227] | - | - | - | - | - | - | - | - | - | - |
| RNN-Capsule [72] | 83.8 | - | - | - | - | - | - | - | - | - |
| ULMFiT [228] | - | - | 95.4 | 97.84 | 71.02 | - | - | 94.99 | 99.2 | - |
| LEAM[229] | 76.95 | - | - | 95.31 | 64.09 | - | 81.91 | 92.45 | 99.02 | - |
| TextCapsule [83] | 82.3 | 86.8 | - | - | - | - | - | 92.6 | - | - |
| TextGCN [6] | 76.74 | - | - | - | - | - | 86.34 | 67.61 | - | - |
| BERT-base [19] | - | 93.5 | 95.63 | 98.08 | 70.58 | 61.6 | - | - | - | 91.0 |
| BERT-large [19] | - | 94.9 | 95.79 | 98.19 | 71.38 | 62.2 | - | - | - | 91.7 |
| MT-DNN[230] | - | 95.6 | 83.2 | - | - | - | - | - | - | 91.5 |
| XLNet-Large [115] | - | 96.8 | 96.21 | 98.45 | 72.2 | 67.74 | - | - | - | - |
| XLNet [115] | - | 97 | - | - | - | - | - | 95.51 | 99.38 | - |
| RoBERTa [117] | - | 96.4 | - | - | - | - | - | - | - | 92.6 |

performance decreases significantly. Consequently, how to improve the robustness of models is a current research hotspot and challenge.

The interpretability of the model. DNNs have unique advantages in feature extraction and semantic mining and have achieved excellent text classification tasks. However, deep learning is a black-box model, the training process is challenging to reproduce, and the implicit semantics and output interpretability are poor. It makes the improvement and optimization of the model, losing clear guidelines. Furthermore, we cannot accurately explain why the model improves performance.

VI. CONCLUSION

This paper principally introduces the existing models for text classification tasks from shallow learning to deep learning.

Firstly, we introduce some primary shallow learning models and deep learning models with a summary table. The shallow model improves text classification performance mainly by improving the feature extraction scheme and classifier design. In contrast, the deep learning model enhances performance by improving the presentation learning method, model structure, and additional data and knowledge. Then, we introduce the datasets with a summary table and evaluation metrics for single-label and multi-label tasks. Furthermore, we give the quantitative results of the leading models in a summary table under different applications for classic text classification datasets. Finally, we summarize the possible future research challenges of text classification.

ACKNOWLEDGMENT

This work is supported in part by the NSFC (61872022 and 61872294), NSF (III-1526499, III-1763325, III-

1909323), CNS-1930941, NSF of Guangdong Province (2017A030313339), and the UK EPSRC (EP/T01461X/1).

REFERENCES

- [1] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA*, pp. 142–150, 2011.
- [2] K. S. Tai, R. Socher, and C. D. Manning, "Improved semantic representations from tree-structured long short-term memory networks," in *Proc. ACL, 2015*, pp. 1556–1566, 2015.
- [3] X. Zhu, P. Sobhani, and H. Guo, "Long short-term memory over recursive structures," in *Proc. ICML, 2015*, pp. 1604–1612, 2015.
- [4] S. I. Wang and C. D. Manning, "Baselines and bigrams: Simple, good sentiment and topic classification," in *Proc. ACL, 2012*, pp. 90–94, 2012.
- [5] X. Zhang, J. J. Zhao, and Y. LeCun, "Character-level convolutional networks for text classification," in *Proc. NeurIPS, 2015*, pp. 649–657, 2015.
- [6] L. Yao, C. Mao, and Y. Luo, "Graph convolutional networks for text classification," in *Proc. AAAI, 2019*, pp. 7370–7377, 2019.
- [7] N. Kalchbrenner, E. Grefenstette, and P. Blunsom, "A convolutional neural network for modelling sentences," in *Proc. ACL, 2014*, pp. 655–665, 2014.
- [8] P. Liu, X. Qiu, X. Chen, S. Wu, and X. Huang, "Multi-timescale long short-term memory neural network for modelling sentences and documents," in *Proc. EMNLP, 2015*, pp. 2326–2335, 2015.
- [9] J. Y. Lee and F. Dernoncourt, "Sequential short-text classification with recurrent and convolutional neural networks," in *Proc. NAACL, 2016*, pp. 515–520, 2016.
- [10] M. E. Maron, "Automatic indexing: An experimental inquiry," *J. ACM*, vol. 8, no. 3, pp. 404–417, 1961.
- [11] T. M. Cover and P. E. Hart, "Nearest neighbor pattern classification," *IEEE Trans. Inf. Theory*, vol. 13, no. 1, pp. 21–27, 1967.
- [12] T. Joachims, "Text categorization with support vector machines: Learning with many relevant features," in *Proc. ECML, 1998*, pp. 137–142, 1998.
- [13] K. Kowsari, K. J. Meimandi, M. Heidarysafa, S. Mendu, L. E. Barnes, and D. E. Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, 2019.
- [14] 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.
- [15] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [16] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proc. ACM SIGKDD, 2016*, pp. 785–794, 2016.
- [17] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T. Liu, "Lightgbm: A highly efficient gradient boosting decision tree," in *Proc. NeurIPS, 2017*, pp. 3146–3154, 2017.
- [18] Y. Kim, "Convolutional neural networks for sentence classification," in *Proc. EMNLP, 2014*, pp. 1746–1751, 2014.
- [19] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL, 2019*, pp. 4171–4186, 2019.
- [20] "Term frequency by inverse document frequency," in *Encyclopedia of Database Systems*, p. 3035, 2009.
- [21] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *Proc. ICLR, 2013*, 2013.
- [22] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proc. EMNLP, 2014*, pp. 1532–1543, 2014.
- [23] M. Zhang and K. Zhang, "Multi-label learning by exploiting label dependency," in *Proc. ACM SIGKDD, 2010*, pp. 999–1008, 2010.
- [24] A. van den Bosch, "Hidden markov models," in *Encyclopedia of Machine Learning and Data Mining*, pp. 609–611, 2017.
- [25] K. Schneider, "A new feature selection score for multinomial naive bayes text classification based on kl-divergence," in *Proc. ACL, 2004*, 2004.
- [26] T. M. Cover and J. A. Thomas, *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. USA: Wiley-Interscience, 2006.
- [27] W. Dai, G. Xue, Q. Yang, and Y. Yu, "Transferring naive bayes classifiers for text classification," in *Proc. AAAI, 2007*, pp. 540–545, 2007.
- [28] A., P., Dempster, N., M., Laird, D., B., and Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the Royal Statistical Society*, 1977.
- [29] P. Frasconi, G. Soda, and A. Vullo, "Hidden markov models for text categorization in multi-page documents," *J. Intell. Inf. Syst.*, vol. 18, no. 2-3, pp. 195–217, 2002.
- [30] K. Yi and J. Beheshti, "A hidden markov model-based text classification of medical documents," *J. Inf. Sci.*, vol. 35, no. 1, pp. 67–81, 2009.
- [31] M. O'Donnell, "Cataloging and classification: An introduction," *Technical Services Quarterly*, vol. 26, no. 1, pp. 86–87, 2009.
- [32] P. Soucy and G. W. Mineau, "A simple KNN algorithm for text categorization," in *Proc. ICDM, 2001*, pp. 647–648, 2001.
- [33] S. Tan, "Neighbor-weighted k-nearest neighbor for unbalanced text corpus," *Expert Syst. Appl.*, vol. 28, no. 4, pp. 667–671, 2005.
- [34] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, 1995.
- [35] T. Joachims, "A statistical learning model of text classification for support vector machines," in *Proc. SIGIR, 2001*, pp. 128–136, 2001.
- [36] T. JOACHIMS, "Transductive inference for text classification using support vector machines," in *International Conference on Machine Learning, 1999*.
- [37] T. M. Mitchell, *Machine learning*. McGraw Hill series in computer science, McGraw-Hill, 1997.
- [38] R. J. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [39] J. R. Quinlan, *C4.5: Programs for Machine Learning*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1993.
- [40] D. E. Johnson, F. J. Oles, T. Zhang, and T. Götz, "A decision-tree-based symbolic rule induction system for text categorization," *IBM Syst. J.*, vol. 41, no. 3, pp. 428–437, 2002.
- [41] P. Vateekul and M. Kubat, "Fast induction of multiple decision trees in text categorization from large scale, imbalanced, and multi-label data," in *Proc. ICDM Workshops, 2009*, pp. 320–325, 2009.
- [42] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," in *Proc. EuroCOLT, 1995*, pp. 23–37, 1995.
- [43] A. Bouaziz, C. Dartigues-Pallez, C. da Costa Pereira, F. Precioso, and P. Lloret, "Short text classification using semantic random forest," in *Proc. DAWAK, 2014*, pp. 288–299, 2014.
- [44] M. Z. Islam, J. Liu, J. Li, L. Liu, and W. Kang, "A semantics aware random forest for text classification," in *Proc. CIKM, 2019*, pp. 1061–1070, 2019.
- [45] R. Socher, J. Pennington, E. H. Huang, A. Y. Ng, and C. D. Manning, "Semi-supervised recursive autoencoders for predicting sentiment distributions," in *Proc. EMNLP, 2011*, pp. 151–161, 2011.
- [46] "A MATLAB implementation of RAE." <https://github.com/vin00/Semi-Supervised-Recursive-Autoencoders-for-Predicting-Sentiment-Distributions>, 2011.
- [47] R. Socher, B. Huval, C. D. Manning, and A. Y. Ng, "Semantic compositionality through recursive matrix-vector spaces," in *Proc. EMNLP, 2012*, pp. 1201–1211, 2012.
- [48] "A Tensorflow implementation of MV_RNN." https://github.com/github-pengge/MV_RNN, 2012.
- [49] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proc. EMNLP, 2013*, pp. 1631–1642, 2013.
- [50] "A MATLAB implementation of RNTN." <https://github.com/pondruska/DeepSentiment>, 2013.
- [51] O. Irsoy and C. Cardie, "Deep recursive neural networks for compositionality in language," in *Proc. NIPS, 2014*, pp. 2096–2104, 2014.
- [52] Q. V. Le and T. Mikolov, "Distributed representations of sentences and documents," in *Proc. ICML, 2014*, pp. 1188–1196, 2014.
- [53] "A PyTorch implementation of Paragraph Vectors (doc2vec)." <https://github.com/inejc/paragraph-vectors>, 2014.
- [54] M. Iyyer, V. Manjunatha, J. L. Boyd-Graber, and H. D. III, "Deep unordered composition rivals syntactic methods for text classification," in *Proc. ACL, 2015*, pp. 1681–1691, 2015.
- [55] "A implementation of DAN." <https://github.com/miyyer/dan>, 2015.
- [56] "A PyTorch implementation of Tree-LSTM." <https://github.com/stanfordnlp/treelstm>, 2015.
- [57] S. Lai, L. Xu, K. Liu, and J. Zhao, "Recurrent convolutional neural networks for text classification," *AAAI'15*, p. 2267–2273, AAAI Press, 2015.
- [58] "A Tensorflow implementation of TextRCNN." <https://github.com/roomylee/rcnn-text-classification>, 2015.

- [59] “A implementation of MT-LSTM.” <https://github.com/AlexAntn/MTLSTM>, 2015.
- [60] R. Johnson and T. Zhang, “Supervised and semi-supervised text categorization using LSTM for region embeddings,” in *Proc. ICML*, 2016, pp. 526–534, 2016.
- [61] “A implementation of oh-2LSTMp.” http://riejohnson.com/cnn_2download.html, 2015.
- [62] P. Zhou, Z. Qi, S. Zheng, J. Xu, H. Bao, and B. Xu, “Text classification improved by integrating bidirectional LSTM with two-dimensional max pooling,” in *Proc. COLING*, 2016, pp. 3485–3495, 2016.
- [63] “A implementation of BLSTM-2DCNN.” <https://github.com/ManuelVs/NNForTextClassification>, 2016.
- [64] P. Liu, X. Qiu, and X. Huang, “Recurrent neural network for text classification with multi-task learning,” in *Proc. IJCAI*, 2016, pp. 2873–2879, 2016.
- [65] “A PyTorch implementation of Multi-Task.” https://github.com/baixl/text_classification, 2016.
- [66] B. Felbo, A. Mislove, A. Søgaard, I. Rahwan, and S. Lehmann, “Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm,” in *Proc. EMNLP*, 2017, pp. 1615–1625, 2017.
- [67] “A Keras implementation of DeepMoji.” <https://github.com/bfelbo/DeepMoji>, 2018.
- [68] A. B. Dieng, C. Wang, J. Gao, and J. W. Paisley, “Topicrnn: A recurrent neural network with long-range semantic dependency,” in *Proc. ICLR*, 2017, 2017.
- [69] “A PyTorch implementation of TopicRNN.” <https://github.com/dangitstam/topic-rnn>, 2017.
- [70] T. Miyato, A. M. Dai, and I. J. Goodfellow, “Adversarial training methods for semi-supervised text classification,” in *Proc. ICLR*, 2017, 2017.
- [71] “A Tensorflow implementation of Virtual adversarial training.” https://github.com/tensorflow/models/tree/master/adversarial_text, 2017.
- [72] Y. Wang, A. Sun, J. Han, Y. Liu, and X. Zhu, “Sentiment analysis by capsules,” in *Proc. WWW*, 2018, pp. 1165–1174, 2018.
- [73] “A PyTorch implementation of RNN-Capsule.” <https://github.com/wangjiosw/Sentiment-Analysis-by-Capsules>, 2018.
- [74] “A Keras implementation of TextCNN.” <https://github.com/alexander-rakhlina/CNN-for-Sentence-Classification-in-Keras>, 2014.
- [75] “A Tensorflow implementation of DCNN.” https://github.com/kinimod23/ATS_Project, 2014.
- [76] “A Tensorflow implementation of CharCNN.” <https://github.com/mhjabreel/CharCNN>, 2015.
- [77] “A Keras implementation of SeqTextRCNN.” <https://github.com/ilimugur/short-text-classification>, 2016.
- [78] J. Liu, W. Chang, Y. Wu, and Y. Yang, “Deep learning for extreme multi-label text classification,” in *Proc. ACM SIGIR*, 2017, pp. 115–124, 2017.
- [79] “A Pytorch implementation of XML-CNN.” <https://github.com/siddsax/XML-CNN>, 2017.
- [80] R. Johnson and T. Zhang, “Deep pyramid convolutional neural networks for text categorization,” in *Proc. ACL*, 2017, pp. 562–570, 2017.
- [81] “A PyTorch implementation of DPCNN.” <https://github.com/Cheneng/DPCNN>, 2017.
- [82] J. Wang, Z. Wang, D. Zhang, and J. Yan, “Combining knowledge with deep convolutional neural networks for short text classification,” in *Proc. IJCAI*, 2017, pp. 2915–2921, 2017.
- [83] M. Yang, W. Zhao, J. Ye, Z. Lei, Z. Zhao, and S. Zhang, “Investigating capsule networks with dynamic routing for text classification,” in *Proc. EMNLP*, 2018, pp. 3110–3119, 2018.
- [84] “A Tensorflow implementation of TextCapsule.” https://github.com/andyweizhao/capsule_text_classification, 2018.
- [85] K. Shimura, J. Li, and F. Fukumoto, “HFT-CNN: learning hierarchical category structure for multi-label short text categorization,” in *Proc. EMNLP*, 2018, pp. 811–816, 2018.
- [86] “A implementation of HFT-CNN.” <https://github.com/ShimShim46/HFT-CNN>, 2018.
- [87] Y. Bao, M. Wu, S. Chang, and R. Barzilay, “Few-shot text classification with distributional signatures,” in *Proc. ICLR*, 2020, 2020.
- [88] “A PyTorch implementation of few-shot text classification with distributional signatures.” <https://github.com/YujiaBao/Distributional-Signatures>, 2020.
- [89] Z. Yang, D. Yang, C. Dyer, X. He, A. J. Smola, and E. H. Hovy, “Hierarchical attention networks for document classification,” in *Proc. NAACL*, 2016, pp. 1480–1489, 2016.
- [90] “A Keras implementation of TextCNN.” <https://github.com/richliao/textClassifier>, 2014.
- [91] X. Zhou, X. Wan, and J. Xiao, “Attention-based LSTM network for cross-lingual sentiment classification,” in *Proc. EMNLP*, 2016, pp. 247–256, 2016.
- [92] “NLP&CC Corpus.” <http://tccf.ccf.org.cn/conference/2013/index.html>, 2013.
- [93] J. Cheng, L. Dong, and M. Lapata, “Long short-term memory-networks for machine reading,” in *Proc. EMNLP*, 2016, pp. 551–561, 2016.
- [94] “A Tensorflow implementation of LSTMN.” <https://github.com/JRC1995/Abstractive-Summarization>, 2016.
- [95] Z. Lin, M. Feng, C. N. dos Santos, M. Yu, B. Xiang, B. Zhou, and Y. Bengio, “A structured self-attentive sentence embedding,” in *Proc. ICLR*, 2017, 2017.
- [96] “A PyTorch implementation of Structured-Self-Attention.” <https://github.com/kaushalshetty/Structured-Self-Attention>, 2017.
- [97] P. Yang, X. Sun, W. Li, S. Ma, W. Wu, and H. Wang, “SGM: sequence generation model for multi-label classification,” in *Proc. COLING*, 2018, pp. 3915–3926, 2018.
- [98] “A PyTorch implementation of SGM.” <https://github.com/lancopku/SGM>, 2018.
- [99] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, “Deep contextualized word representations,” in *Proc. NAACL*, 2018, pp. 2227–2237, 2018.
- [100] “A PyTorch implementation of ELMo.” <https://github.com/flairNLP/flair>, 2018.
- [101] T. Shen, T. Zhou, G. Long, J. Jiang, and C. Zhang, “Bi-directional block self-attention for fast and memory-efficient sequence modeling,” in *Proc. ICLR*, 2018, 2018.
- [102] “A PyTorch implementation of BiBloSA.” <https://github.com/galsang/BiBloSA-pytorch>, 2018.
- [103] R. You, Z. Zhang, Z. Wang, S. Dai, H. Mamitsuka, and S. Zhu, “Attentionxml: Label tree-based attention-aware deep model for high-performance extreme multi-label text classification,” in *Proc. NeurIPS*, 2019, pp. 5812–5822, 2019.
- [104] “A PyTorch implementation of AttentionXML.” <https://github.com/yourh/AttentionXML>, 2019.
- [105] S. Sun, Q. Sun, K. Zhou, and T. Lv, “Hierarchical attention prototypical networks for few-shot text classification,” in *Proc. EMNLP*, 2019, pp. 476–485, 2019.
- [106] T. Gao, X. Han, Z. Liu, and M. Sun, “Hybrid attention-based prototypical networks for noisy few-shot relation classification,” in *Proc. AAAI*, 2019, pp. 6407–6414, 2019.
- [107] “A PyTorch implementation of HATT-Proto.” <https://github.com/thunlp/HATT-Proto>, 2019.
- [108] J. Chen, Y. Hu, J. Liu, Y. Xiao, and H. Jiang, “Deep short text classification with knowledge powered attention,” in *Proc. AAAI*, 2019, pp. 6252–6259, 2019.
- [109] “A PyTorch implementation of STCKA.” <https://github.com/AIRobotZhang/STCKA>, 2019.
- [110] “A Tensorflow implementation of BERT.” <https://github.com/google-research/bert>, 2019.
- [111] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, and I. Androutsopoulos, “Large-scale multi-label text classification on EU legislation,” in *Proc. ACL*, 2019, pp. 6314–6322, 2019.
- [112] “A Tensorflow implementation of BERT-BASE.” <https://github.com/iliaskalchididis/lmte-eurlex57k>, 2019.
- [113] C. Sun, X. Qiu, Y. Xu, and X. Huang, “How to fine-tune BERT for text classification?,” in *Proc. CCL*, 2019, pp. 194–206, 2019.
- [114] “A Tensorflow implementation of BERT4doc-Classification.” <https://github.com/xuyige/BERT4doc-Classification>, 2019.
- [115] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” in *Proc. NeurIPS*, 2019, pp. 5754–5764, 2019.
- [116] “A Tensorflow implementation of XLNet.” <https://github.com/zihangdai/xlnet>, 2019.
- [117] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized BERT pretraining approach,” *CoRR*, vol. abs/1907.11692, 2019.
- [118] “A PyTorch implementation of RoBERTa.” <https://github.com/pytorch/fairseq>, 2019.
- [119] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, “ALBERT: A lite BERT for self-supervised learning of language representations,” in *Proc. ICLR*, 2020, 2020.
- [120] “A Tensorflow implementation of ALBERT.” <https://github.com/google-research/ALBERT>, 2020.
- [121] H. Peng, J. Li, Y. He, Y. Liu, M. Bao, L. Wang, Y. Song, and Q. Yang, “Large-scale hierarchical text classification with recursively regularized deep graph-cnn,” in *Proc. WWW*, 2018, pp. 1063–1072, 2018.

- [122] "A Tensorflow implementation of DeepGraphCNNforTexts." <https://github.com/HKUST-KnowComp/DeepGraphCNNforTexts>, 2018.
- [123] "A Tensorflow implementation of TextGCN." https://github.com/yao8839836/text_gcn, 2019.
- [124] F. Wu, A. H. S. Jr., T. Zhang, C. Fifty, T. Yu, and K. Q. Weinberger, "Simplifying graph convolutional networks," in *Proc. ICML*, 2019, pp. 6861–6871, 2019.
- [125] "A implementation of SGC." <https://github.com/Tiiiger/SGC>, 2019.
- [126] L. Huang, D. Ma, S. Li, X. Zhang, and H. Wang, "Text level graph neural network for text classification," in *Proc. EMNLP*, 2019, pp. 3442–3448, 2019.
- [127] "A implementation of TextLevelGNN." <https://github.com/LidgeW/TextLevelGNN>, 2019.
- [128] H. Peng, J. Li, S. Wang, L. Wang, Q. Gong, R. Yang, B. Li, P. Yu, and L. He, "Hierarchical taxonomy-aware and attentional graph capsule rcnn for large-scale multi-label text classification," *IEEE Transactions on Knowledge and Data Engineering*, 2019.
- [129] A. Pal, M. Selvakumar, and M. Sankarasubbu, "MAGNET: multi-label text classification using attention-based graph neural network," in *Proc. ICAART*, 2020, pp. 494–505, 2020.
- [130] "A repository of MAGNET." <https://github.com/monk1337/MAGnet>, 2020.
- [131] "A Tensorflow implementation of Miyato et al.." <https://github.com/TobiasLee/Text-Classification>, 2017.
- [132] J. Zeng, J. Li, Y. Song, C. Gao, M. R. Lyu, and I. King, "Topic memory networks for short text classification," in *Proc. EMNLP*, 2018, pp. 3120–3131, 2018.
- [133] J. Zhang, P. Lertvittayakumjorn, and Y. Guo, "Integrating semantic knowledge to tackle zero-shot text classification," in *Proc. NAACL*, 2019, pp. 1031–1040, 2019.
- [134] "A Tensorflow implementation of KG4ZeroShotText." <https://github.com/JingqingZ/KG4ZeroShotText>, 2019.
- [135] M. k. Alsmadi, K. B. Omar, S. A. Noah, and I. Almarashdah, "Performance comparison of multi-layer perceptron (back propagation, delta rule and perceptron) algorithms in neural networks," in *2009 IEEE International Advance Computing Conference*, pp. 296–299, 2009.
- [136] T. Miyato, S. Maeda, M. Koyama, and S. Ishii, "Virtual adversarial training: a regularization method for supervised and semi-supervised learning," *CoRR*, vol. abs/1704.03976, 2017.
- [137] G. E. Hinton, A. Krizhevsky, and S. D. Wang, "Transforming auto-encoders," in *Proc. ICANN*, 2011 (T. Honkela, W. Duch, M. Girolami, and S. Kaski, eds.), (Berlin, Heidelberg), pp. 44–51, Springer Berlin Heidelberg, 2011.
- [138] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [139] Z. Wang, W. Hamza, and R. Florian, "Bilateral multi-perspective matching for natural language sentences," in *Proc. IJCAI*, 2017, pp. 4144–4150, 2017.
- [140] R. Johnson and T. Zhang, "Semi-supervised convolutional neural networks for text categorization via region embedding," in *Proc. NeurIPS*, 2015, pp. 919–927, 2015.
- [141] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Proc. ECCV*, 2016, pp. 630–645, 2016.
- [142] H. Nguyen and M. Nguyen, "A deep neural architecture for sentence-level sentiment classification in twitter social networking," in *Proc. PAICLING*, 2017, pp. 15–27, 2017.
- [143] B. Adams and G. McKenzie, "Crowdsourcing the character of a place: Character-level convolutional networks for multilingual geographic text classification," *Trans. GIS*, vol. 22, no. 2, pp. 394–408, 2018.
- [144] Z. Chen and T. Qian, "Transfer capsule network for aspect level sentiment classification," in *Proc. ACL*, 2019, pp. 547–556, 2019.
- [145] W. Xue, W. Zhou, T. Li, and Q. Wang, "MTNA: A neural multi-task model for aspect category classification and aspect term extraction on restaurant reviews," in *Proc. IJCNLP*, 2017, pp. 151–156, 2017.
- [146] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proc. ICLR*, 2015, 2015.
- [147] Z. Hu, X. Li, C. Tu, Z. Liu, and M. Sun, "Few-shot charge prediction with discriminative legal attributes," in *Proc. COLING*, 2018, pp. 487–498, 2018.
- [148] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NeurIPS*, 2017, pp. 5998–6008, 2017.
- [149] Y. Wang, M. Huang, X. Zhu, and L. Zhao, "Attention-based LSTM for aspect-level sentiment classification," in *Proc. EMNLP*, 2016, pp. 606–615, 2016.
- [150] F. Fan, Y. Feng, and D. Zhao, "Multi-grained attention network for aspect-level sentiment classification," in *Proc. EMNLP*, 2018, pp. 3433–3442, 2018.
- [151] M. Tan, C. dos Santos, B. Xiang, and B. Zhou, "Improved representation learning for question answer matching," in *Proc. ACL*, 2016, (Berlin, Germany), pp. 464–473, Association for Computational Linguistics, Aug. 2016.
- [152] C. N. dos Santos, M. Tan, B. Xiang, and B. Zhou, "Attentive pooling networks," *CoRR*, vol. abs/1602.03609, 2016.
- [153] A. Radford, "Improving language understanding by generative pre-training," 2018.
- [154] G. Jawahar, B. Sagot, and D. Seddah, "What does BERT learn about the structure of language?," in *Proc. ACL*, 2019, pp. 3651–3657, 2019.
- [155] Z. Dai, Z. Yang, Y. Yang, J. G. Carbonell, Q. V. Le, and R. Salakhutdinov, "Transformer-xl: Attentive language models beyond a fixed-length context," in *Proc. ACL*, 2019, pp. 2978–2988, 2019.
- [156] M. Henaff, J. Bruna, and Y. LeCun, "Deep convolutional networks on graph-structured data," *CoRR*, vol. abs/1506.05163, 2015.
- [157] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Proc. NeurIPS*, 2016, pp. 3837–3845, 2016.
- [158] D. Marcheggiani and I. Titov, "Encoding sentences with graph convolutional networks for semantic role labeling," in *Proc. EMNLP*, 2017, pp. 1506–1515, 2017.
- [159] Y. Li, R. Jin, and Y. Luo, "Classifying relations in clinical narratives using segment graph convolutional and recurrent neural networks (seg-rcnns)," *JAMIA*, vol. 26, no. 3, pp. 262–268, 2019.
- [160] J. Bastings, I. Titov, W. Aziz, D. Marcheggiani, and K. Sima'an, "Graph convolutional encoders for syntax-aware neural machine translation," in *Proc. EMNLP*, 2017, pp. 1957–1967, 2017.
- [161] Y. Zhang, X. Yu, Z. Cui, S. Wu, Z. Wen, and L. Wang, "Every document owns its structure: Inductive text classification via graph neural networks," in *Proc. ACL*, 2020, pp. 334–339, 2020.
- [162] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *Proc. ICLR*, 2018, 2018.
- [163] L. Hu, T. Yang, C. Shi, H. Ji, and X. Li, "Heterogeneous graph attention networks for semi-supervised short text classification," in *Proc. EMNLP*, 2019, pp. 4820–4829, 2019.
- [164] Z. Li, X. Ding, and T. Liu, "Constructing narrative event evolutionary graph for script event prediction," in *Proc. IJCAI*, 2018, pp. 4201–4207, 2018.
- [165] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, "Signature verification using a siamese time delay neural network," in *Proc. NeurIPS*, 1993, pp. 737–744, 1993.
- [166] J. Mueller and A. Thyagarajan, "Siamese recurrent architectures for learning sentence similarity," in *Proc. AAAI*, 2016, pp. 2786–2792, 2016.
- [167] Jayadeva, H. Pant, M. Sharma, and S. Soman, "Twin neural networks for the classification of large unbalanced datasets," *Neurocomputing*, vol. 343, pp. 34 – 49, 2019. Learning in the Presence of Class Imbalance and Concept Drift.
- [168] T. Miyato, S. ichi Maeda, M. Koyama, K. Nakae, and S. Ishii, "Distributional smoothing with virtual adversarial training," 2015.
- [169] T. Zhang, M. Huang, and L. Zhao, "Learning structured representation for text classification via reinforcement learning," in *Proc. AAAI*, 2018, pp. 6053–6060, 2018.
- [170] J. Weston, S. Chopra, and A. Bordes, "Memory networks," 2015.
- [171] X. Li and W. Lam, "Deep multi-task learning for aspect term extraction with memory interaction," in *Proc. EMNLP*, 2017, pp. 2886–2892, 2017.
- [172] C. Shen, C. Sun, J. Wang, Y. Kang, S. Li, X. Liu, L. Si, M. Zhang, and G. Zhou, "Sentiment classification towards question-answering with hierarchical matching network," in *Proc. EMNLP*, 2018, pp. 3654–3663, 2018.
- [173] X. Ding, K. Liao, T. Liu, Z. Li, and J. Duan, "Event representation learning enhanced with external commonsense knowledge," in *Proc. EMNLP*, 2019, pp. 4893–4902, 2019.
- [174] Y. Zhang, D. Song, P. Zhang, X. Li, and P. Wang, "A quantum-inspired sentiment representation model for twitter sentiment analysis," *Applied Intelligence*, 2019.
- [175] "MR Corpus." <http://www.cs.cornell.edu/people/pabo/movie-review-data/>, 2002.
- [176] "SSST Corpus." <http://nlp.stanford.edu/sentiment/>, 2013.
- [177] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proc. EMNLP*, 2013, (Seattle, Washington,

- USA), pp. 1631–1642, Association for Computational Linguistics, Oct. 2013.
- [178] “MPQA Corpus.” <http://www.cs.pitt.edu/mpqa/>, 2005.
- [179] “Twitter Corpus.” <https://www.cs.york.ac.uk/semeval-2013/task2/>, 2013.
- [180] Q. Diao, M. Qiu, C. Wu, A. J. Smola, J. Jiang, and C. Wang, “Jointly modeling aspects, ratings and sentiments for movie recommendation (JMARS),” in *Proc. ACM SIGKDD, 2014*, pp. 193–202, 2014.
- [181] D. Tang, B. Qin, and T. Liu, “Document modeling with gated recurrent neural network for sentiment classification,” in *Proc. EMNLP, 2015*, pp. 1422–1432, 2015.
- [182] “Amazon review Corpus.” <https://www.kaggle.com/datasfiniti/consumer-reviews-of-amazon-products>, 2015.
- [183] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li, “RCV1: A new benchmark collection for text categorization research,” *J. Mach. Learn. Res.*, vol. 5, pp. 361–397, 2004.
- [184] “RCV1-V2 Corpus.” http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_README.htm, 2004.
- [185] “20NG Corpus.” <http://ana.cachopo.org/datasets-for-single-label-text-categorization>, 2007.
- [186] “AG Corpus.” http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html, 2004.
- [187] “Reuters Corpus.” <https://martin-thoma.com/nlp-reuters>, 2017.
- [188] “Reuters Corpus.” <https://www.cs.umb.edu/~smimarog/textmining/datasets/>, 2007.
- [189] J. Zhou, C. Ma, D. Long, G. Xu, N. Ding, H. Zhang, P. Xie, and G. Liu, “Hierarchy-aware global model for hierarchical text classification,” in *Proc. ACL, 2020*, pp. 1106–1117, 2020.
- [190] Y. Mao, J. Tian, J. Han, and X. Ren, “Hierarchical text classification with reinforced label assignment,” in *Proc. EMNLP, 2019*, pp. 445–455, 2019.
- [191] “NYTimes Corpus.” https://catalog.ldc.upenn.edu/docs/LDC2008T19/new_york_times_annotated_corpus.pdf, 2007.
- [192] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. N. Mendes, S. Hellmann, M. Morsey, P. van Kleef, S. Auer, and C. Bizer, “Dbpedia - A large-scale, multilingual knowledge base extracted from wikipedia,” *Semantic Web*, vol. 6, no. 2, pp. 167–195, 2015.
- [193] “Ohsumed Corpus.” <http://davis.wpi.edu/xmdv/datasets/ohsumed.html>, 2015.
- [194] “Amazon670K Corpus.” <http://manikvarma.org/downloads/XC/XMLRepository.html>, 2016.
- [195] “EUR-Lex Corpus.” <http://www.ke.tu-darmstadt.de/resources/eurlex/eurlex.html>, 2019.
- [196] B. Pang and L. Lee, “A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts,” in *Proc. ACL, 2004*, (Barcelona, Spain), pp. 271–278, July 2004.
- [197] “TREC Corpus.” <https://cogcomp.seas.upenn.edu/Data/QA/QC/>, 2002.
- [198] M. Hu and B. Liu, “Mining and summarizing customer reviews,” in *Proc. ACM SIGKDD, 2004*, pp. 168–177, 2004.
- [199] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “Squad: 100, 000+ questions for machine comprehension of text,” in *Proc. EMNLP, 2016*, pp. 2383–2392, 2016.
- [200] Y. Yang, W. Yih, and C. Meek, “Wikiqa: A challenge dataset for open-domain question answering,” in *Proc. EMNLP, 2015*, pp. 2013–2018, 2015.
- [201] S. Kim, L. F. D’Haro, R. E. Bansch, J. D. Williams, and M. Henderson, “The fourth dialog state tracking challenge,” in *Proc. IWSDS, 2016*, pp. 435–449, 2016.
- [202] Y. Wan, W. Yan, J. Gao, Z. Zhao, J. Wu, and P. S. Yu, “Improved dynamic memory network for dialogue act classification with adversarial training,” in *IEEE International Conference on Big Data, Big Data 2018, Seattle, WA, USA, December 10-13, 2018*, pp. 841–850, 2018.
- [203] V. Raheja and J. R. Tetreault, “Dialogue act classification with context-aware self-attention,” in *Proc. NAACL, 2019*, pp. 3727–3733, 2019.
- [204] J. Ang, Y. Liu, and E. Shriberg, “Automatic dialog act segmentation and classification in multiparty meetings,” in *Proc. ICASSP, 2005*, pp. 1061–1064, 2005.
- [205] D. Jurafsky and E. Shriberg, “Switchboard swbd-damsl shallow-discourse-function annotation coders manual,” 01 1997.
- [206] X. Han, H. Zhu, P. Yu, Z. Wang, Y. Yao, Z. Liu, and M. Sun, “Fewrel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation,” in *Proc. EMNLP, 2018*, pp. 4803–4809, 2018.
- [207] T. Gao, X. Han, H. Zhu, Z. Liu, P. Li, M. Sun, and J. Zhou, “Fewrel 2.0: Towards more challenging few-shot relation classification,” in *Proc. EMNLP, 2019*, pp. 6249–6254, 2019.
- [208] “FewRel Corpus.” <https://github.com/thunlp/FewRel>, 2019.
- [209] I. Hendrickx, S. N. Kim, Z. Kozareva, P. Nakov, D. Ó. Séaghdha, S. Padó, M. Pennacchiotti, L. Romano, and S. Szpakowicz, “Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals,” in *Proc. SemEval, 2010*, pp. 33–38, 2010.
- [210] D. Zeng, K. Liu, S. Lai, G. Zhou, and J. Zhao, “Relation classification via convolutional deep neural network,” in *Proc. COLING, 2014*, pp. 2335–2344, 2014.
- [211] I. Hendrickx, S. N. Kim, Z. Kozareva, P. Nakov, D. Ó. Séaghdha, S. Padó, M. Pennacchiotti, L. Romano, and S. Szpakowicz, “Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals,” *CoRR*, vol. abs/1911.10422, 2019.
- [212] F. Ren, D. Zhou, Z. Liu, Y. Li, R. Zhao, Y. Liu, and X. Liang, “Neural relation classification with text descriptions,” in *Proc. COLING, 2018*, pp. 1167–1177, 2018.
- [213] S. Wu and Y. He, “Enriching pre-trained language model with entity information for relation classification,” in *Proc. CIKM, 2019*, pp. 2361–2364, 2019.
- [214] I. Hendrickx, S. N. Kim, Z. Kozareva, P. Nakov, D. Ó. Séaghdha, S. Padó, M. Pennacchiotti, L. Romano, and S. Szpakowicz, “Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals,” in *Proc. NAACL, 2009*, pp. 94–99, 2009.
- [215] B. Pang and L. Lee, “Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales,” in *Proc. ACL, 2005*, pp. 115–124, 2005.
- [216] J. Wiebe, T. Wilson, and C. Cardie, “Annotating expressions of opinions and emotions in language,” *Language Resources and Evaluation*, vol. 39, no. 2-3, pp. 165–210, 2005.
- [217] <https://data.quora.com/First-Quora-Dataset-Release-QuestionPairs>.
- [218] P. Rajpurkar, R. Jia, and P. Liang, “Know what you don’t know: Unanswerable questions for squad,” in *Proc. ACL, 2018*, pp. 784–789, 2018.
- [219] T. Nguyen, M. Rosenberg, X. Song, J. Gao, S. Tiwary, R. Majumder, and L. Deng, “MS MARCO: A human generated machine reading comprehension dataset,” in *Proc. NeurIPS, 2016*, 2016.
- [220] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning, “A large annotated corpus for learning natural language inference,” in *Proc. EMNLP, 2015*, pp. 632–642, 2015.
- [221] A. Williams, N. Nangia, and S. R. Bowman, “A broad-coverage challenge corpus for sentence understanding through inference,” in *Proc. NAACL, 2018*, pp. 1112–1122, 2018.
- [222] M. Marelli, L. Bentivogli, M. Baroni, R. Bernardi, S. Menini, and R. Zamparelli, “Semeval-2014 task 1: Evaluation of compositional distributional semantic models on full sentences through semantic relatedness and textual entailment,” in *Proc. SemEval, 2014*, pp. 1–8, 2014.
- [223] B. Dolan, C. Quirk, and C. Brockett, “Unsupervised construction of large paraphrase corpora: Exploiting massively parallel news sources,” in *Proc. COLING, 2004*, 2004.
- [224] K. Kowsari, D. E. Brown, M. Heidarysafa, K. J. Meimandi, M. S. Gerber, and L. E. Barnes, “Hdtlx: Hierarchical deep learning for text classification,” in *Proc. IJMLA, 2017*, pp. 364–371, 2017.
- [225] R. E. Schapire and Y. Singer, “Improved boosting algorithms using confidence-rated predictions,” *Mach. Learn.*, vol. 37, no. 3, pp. 297–336, 1999.
- [226] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to information retrieval*. Cambridge University Press, 2008.
- [227] P. Chen, Z. Sun, L. Bing, and W. Yang, “Recurrent attention network on memory for aspect sentiment analysis,” in *Proc. EMNLP, 2017*, pp. 452–461, 2017.
- [228] J. Howard and S. Ruder, “Universal language model fine-tuning for text classification,” in *Proc. ACL, 2018*, pp. 328–339, 2018.
- [229] G. Wang, C. Li, W. Wang, Y. Zhang, D. Shen, X. Zhang, R. Henao, and L. Carin, “Joint embedding of words and labels for text classification,” in *Proc. ACL, 2018*, pp. 2321–2331, 2018.
- [230] X. Liu, P. He, W. Chen, and J. Gao, “Multi-task deep neural networks for natural language understanding,” in *Proc. ACL, 2019*, pp. 4487–4496, 2019.
- [231] P. K. Pushp and M. M. Srivastava, “Train once, test anywhere: Zero-shot learning for text classification,” *CoRR*, vol. abs/1712.05972, 2017.
- [232] C. Song, S. Zhang, N. Sadoughi, P. Xie, and E. P. Xing, “Generalized zero-shot text classification for ICD coding,” in *Proc. IJCAI, 2020*, pp. 4018–4024, 2020.
- [233] R. Geng, B. Li, Y. Li, X. Zhu, P. Jian, and J. Sun, “Induction networks for few-shot text classification,” in *Proc. EMNLP, 2019*, pp. 3902–3911, 2019.

- [234] S. Deng, N. Zhang, Z. Sun, J. Chen, and H. Chen, "When low resource NLP meets unsupervised language model: Meta-pretraining then meta-learning for few-shot text classification (student abstract)," in *Proc. AAAI, 2020*, pp. 13773–13774, 2020.
- [235] R. Geng, B. Li, Y. Li, J. Sun, and X. Zhu, "Dynamic memory induction networks for few-shot text classification," in *Proc. ACL, 2020*, pp. 1087–1094, 2020.
- [236] Y. Sun, S. Wang, Y. Li, S. Feng, X. Chen, H. Zhang, X. Tian, D. Zhu, H. Tian, and H. Wu, "ERNIE: enhanced representation through knowledge integration," *CoRR*, vol. abs/1904.09223, 2019.
- [237] Y. Hao, Y. Zhang, K. Liu, S. He, Z. Liu, H. Wu, and J. Zhao, "An end-to-end model for question answering over knowledge base with cross-attention combining global knowledge," in *Proc. ACL, 2017*, pp. 221–231, 2017.
- [238] R. Türker, L. Zhang, M. Koutraki, and H. Sack, "TECNE: knowledge based text classification using network embeddings," in *Proc. EKAW, 2018*, pp. 53–56, 2018.



Qian Li is currently pursuing the Ph.D. degree with the Department of Computer Science and Engineering, Beihang University (BUAA), Beijing, China. Her research interests include text mining, representation learning, and event extraction.



Renyun Yang is a Research Fellow with the University of Leeds, UK and adjunct researcher in Beijing Advanced Innovation Center for Big Data and Brain Computing in Beihang University. His research interests include reliable distributed systems, big data analytic at scale and applied machine learning.



Lichao Sun is a PhD student at University of Illinois at Chicago, US. His research interests include deep learning and data mining. He mainly focuses on security and privacy, social network and natural language processing applications. He has published more than 15 research articles in top conferences and journals like KDD, WSDM, TII, TMC, AAAI.



Hao Peng is currently an Assistant Professor at the School of Cyber Science and Technology, and Beijing Advanced Innovation Center for BigData and Brain Computing in Beihang University. His research interests include representation learning, machine learning and graph mining.



Philip S. Yu is a Distinguished Professor and the Wexler Chair in Information Technology at the Department of Computer Science, University of Illinois at Chicago. Before joining UIC, he was at the IBM Watson Research Center, where he built a world-renowned data mining and database department. He is a Fellow of the ACM and IEEE. Dr. Yu was the Editor-in-Chiefs of ACM Transactions on Knowledge Discovery from Data (2011-2017) and IEEE Transactions on Knowledge and Data Engineering (2001-2004).



Jianxin Li is currently a Professor with the State Key Laboratory of Software Development Environment, and Beijing Advanced Innovation Center for Big Data and Brain Computing in Beihang University. His current research interests include social network, machine learning, big data and trustworthy computing.



LiFang He is currently an Assistant Professor in the Department of Computer Science and Engineering at Lehigh University. Before her current position, Dr. He worked as a postdoctoral researcher in the Department of Biostatistics and Epidemiology at University of Pennsylvania. Her current research interests include machine learning, data mining, tensor analysis, with major applications in biomedical data, neuroscience, or multimodal data.



Congying Xia is currently pursing the Ph.D. degree with the Department of Computer Science at University of Illinois at Chicago. Her research interests include natural language processing, zero-shot learning and few-shot learning.