Large-Scale Question Tagging via Joint Question-Topic Embedding Learning

LIQIANG NIE and YONGQI LI, Shandong University, China FULI FENG, National University of Singapore, Singapore XUEMENG SONG, Shandong University, China MENG WANG, Hefei University of Technology, China YINGLONG WANG, Qilu University of Technology (Shandong Academy of Sciences), China

Recent years have witnessed a flourishing of community-driven question answering (cQA), like Yahoo! Answers and AnswerBag, where people can seek precise information. After 2010, some novel cQA systems, including Quora and Zhihu, gained momentum. Besides interactions, the latter enables users to label the questions with topic tags that highlight the key points conveyed in the questions. In this article, we shed light on automatically annotating a newly posted question with topic tags that are predefined and preorganized into a directed acyclic graph. To accomplish this task, we present an end-to-end deep interactive embedding model to jointly learn the embeddings of questions and topics by projecting them into the same space for a similarity measure. In particular, we first learn the embeddings of questions and topic tags via fully exploring their hierarchical structures, which is able to alleviate the problem of imbalanced topic distribution. Thereafter, we interact each question embedding with the topic tag matrix, i.e., all the topic tag embeddings. Following that, a sigmoid cross-entropy loss is appended to reward the positive question-topic pairs and penalize the negative ones. To justify our model, we have conducted extensive experiments on an unprecedented largescale social QA dataset obtained from Zhihu.com, and the experimental results demonstrate that our model achieves superior performance to several state-of-the-art baselines.

CCS Concepts: • **Information systems** → **Document topic models**; *Learning to rank*;

Additional Key Words and Phrases: Question tagging, topic hierarchy, CQA, embedding learning

ACM Reference format:

Liqiang Nie, Yongqi Li, Fuli Feng, Xuemeng Song, Meng Wang, and Yinglong Wang. 2020. Large-Scale Question Tagging via Joint Question-Topic Embedding Learning. *ACM Trans. Inf. Syst.* 38, 2, Article 20 (February 2020), 23 pages.

https://doi.org/10.1145/3380954

© 2020 Association for Computing Machinery.

1046-8188/2020/02-ART20 \$15.00

https://doi.org/10.1145/3380954

This work is supported by the National Natural Science Foundation of China, No. 61772310, No. U1936203, and No. U1836216; the Shandong Provincial Natural Science and Foundation, No. ZR2019JQ23; and the Innocation Teams in Colleges and Universities in Jinan, No. 2018GXRC014.

Authors' addresses: L. Nie, Y. Li, and X. Song, Shandong University (Tsingtao Campus), 72 Binhai Road, Jimo Qingdao, Shandong Province, China, 266237; emails: {nieliqiang, liyongqi0, sxmustc}@gmail.com; F. Feng, National University of Singapore, Singapore; email: fulifeng93@gmail.com; M. Wang, Hefei University of Technology, China; email: eric.mengwang@gmail.com; Y. Wang, Qilu University of Technology (Shandong Academy of Sciences), China; email: wangyl@sdas.org.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

1 INTRODUCTION

The last three decades have witnessed the proliferation of community-based question answering (cQA). cQA sites, like Yahoo! Answers,¹ AnswerBag,² and Stack Overflow,³ are places to gain and share knowledge whereby users are encouraged to ask or answer questions and are able to connect with the contributors of unique insights and quality answers [4, 40, 42]. Users are thus empowered to learn from each other and to better understand the world. After 2010, some novel sites in the cQA family emerged, such as Zhihu⁴ and Quora.⁵ We hereafter refer to the novel cQA sites as collaborative cQA [7, 37], which refers to cQA systems that allow collaborative question answering. Compared to the conventional cOA sites, the collaborative ones have the following two prominent characteristics: (1) in the collaborative cQA sites, the spirit of crowdsourcing is highly encouraged, whereby any authenticated user is allowed to perform editing on any question, answer, and even topic, to refine the quality of community content, and (2) the collaborate cQA sites highlight more on social connections among users. For example, users in the collaborative cQA sites, like Quora and Zhihu, are enabled to follow each other, and even topics they are interested in, which is not allowed in conventional cQA sites like Yahoo! Answers and AnswerBag.⁶ These prominent features make them thriving. Considering Quora as an example, it had obtained over 300 million monthly unique visitors and archived 38 million-plus distinct questions as of May 2019.⁷

Owing to the proliferation of cQA sites, the amount of questions generated and answered by people is growing exponentially. So it becomes increasingly difficult and expensive for users to locate the questions they need and are interested in. Therefore, cQA sites organize questions by user-generated topics because tagging is a simple and efficient method to organize resources. To be more specific, each asker is strongly encouraged to select multiple topic tags from the suggested candidate list for labeling the newly posted question, which summarizes the question content in a coarse-grained but semantically meaningful level. One typical example of question tagging is demonstrated in Figure 1. Topic tags play pivotal roles in cQA sites, including but not limited to the following aspects: (1) Question routing. In addition to the unidirectional user-follow-user relations, in cQA sites, users can also follow the topics of interest. In light of this, cQA sites can put the questions into the feeds of associated topic followers to draw more attention from the potential answerers, and thus receive more quick and accurate answers. (2) Topic tags can be leveraged to benefit index, search, navigation, and organization. Therefore, question tagging in cQA sites deserves our attention.

Despite its significance and value, question tagging in cQA is nontrivial for the following reasons: (1) Topic tags are not independent. Questions are organized by topic tags (e.g., Quora, Stack Overflow, and Zhihu) or topic categories (e.g., AnswerBag, Yahoo! Answer). For the cQA systems based on topic tags, the structures of topic tags are hierarchical, like Zhihu, or flat, like Quora and Stack Overflow. In this work, we focus on cQA systems with structural topic tags. For example, in Zhihu systems, topic tags are organized into a directed acyclic graph (DAG) by experienced users and hired experts, as shown in Figure 1. The DAG is more like a tree structure except that some nodes have multiple parents. From the root-to-leaf nodes, topic tags tend to be more specific. A question can be annotated by either the leaf or the internal nodes at the same time. How to encode

¹https://answers.yahoo.com/.

²https://www.answerbag.com/.

³https://stackoverflow.com/.

⁴https://www.zhihu.com/.

⁵https://www.quora.com/.

 $^{^{6}} https://socialcompare.com/en/comparison/compare-question-answer-sites-quora-vs-yahoo-answers-vs-stackoverflow-vs-ted-conversations.$

⁷https://www.quora.com/How-many-users-does-Quora-have-in-2019.



Fig. 1. Exemplars of question tagging in cQA sites and the directed acyclic graph structure of topic tags.

the hierarchical structure of topic tags to fully explore the correlations among topic tags is the first challenge we are facing. It is worth mentioning that some cQA sites based on flat topic tags, such as Quora and Stack Overflow, do not organize topic tags into a hierarchical structure explicitly. However, their topic tag structures do exist in pairwise forms and their graph-based ontologies can be built. For example, the tags "programming language" and "python" in Stack Overflow can be treated as the father-son pair, as python is one kind of programming language. We will discuss how to model the tag hierarchy on those cQA sites in Section 6. (2) In real-world settings, the cQA sites may archive dozens of thousands, even hundreds of thousands of topic tags. The number of questions tagged by each topic tag varies significantly. Statistics tell us that the leaf nodes account for around 70% of the DAG, while they label only about 25% of the questions, whereas 30% of the DAG are the internal nodes used to label the rest of the questions. Considering such imbalanced distribution, it may be suboptimal to directly cast the tagging task as a classification problem. Instead of classification, in this article, we treat the labeling task as a ranking problem by measuring the similarity between the given question and each topic tag candidate. In view of this, how to jointly learn the embeddings of questions and imbalanced topic tags within the same space is largely untapped. Meanwhile, the knowledge is supposed to be transferred from the more frequent internal nodes to the less frequent leaf ones for alleviating the imbalance based on the given DAG structure. (3) The publicly accessible collaborative cQA dataset with explicit topic tag hierarchy is still in small scale. For example, the largest publicly available collaborative cQA dataset with topic tag hierarchy thus far has been released by the contest with fewer than 2,000 topic tags after desensitization,⁸ which is not representative enough for evaluation. How to justify the effectiveness of our model on a real-world and large-scale dataset is also a challenge we have to solve.

To address the aforementioned challenges, in this article, we present a scheme for question tagging on cQA systems with structural topic tags, as illustrated in Figure 2. This scheme consists of offline learning and online tagging. As to the offline part, we propose an end-to-end deeP inteRactive mOdel For questIon Tagging, dubbed PROFIT, which is capable of simultaneously projecting the questions and imbalanced topic tags into the same space and learning their representations for similarity measurement. In particular, we first leverage a convolutional neural network (CNN) [18] and a multilayer perception (MLP) model [51] to learn the embeddings of each question and all the imbalanced topic tags, respectively. Notably, the questions and topic tags are forced to share the same set of word embeddings that are pretrained over a very large cQA set. To learn

⁸https://zhuanlan.zhihu.com/p/26843044



Fig. 2. Schematic illustration of our proposed scheme for question tagging in the cQA setting. In the offline part, we train a PROFIT model to learn the question and topic embedding. As to the online one, we recommend relevant topic candidates based on the matching score.

the discriminative embeddings, we encode the DAG structure of topic tags into our learning model by regularizing the hierarchical relations among topic tags. Considering the less frequent problem of the leaf nodes in the DAG, we leverage the weighted linear combination of the parent nodes to represent the child one, which is able to transfer the knowledge from the top down. We thereafter perform the interaction between the question embedding and the topic tag matrix to obtain a similarity vector, whereby each element in the vector correspondingly represents the matching score between the question and the topic tag. We then normalize the similarity vector by a sigmoid layer that computes the sigmoid result for each element in the similarity matrix.

In the training phase, a cross-entropy loss is followed to reward the positive question-tag pairs and penalize the negative ones, which indeed implicitly guarantees the question and topic tag embeddings toward the same semantic space. In our work, we treat the (question_i, tag_j) as a positive pair if the tag_j labels the question_i before; otherwise, we view it as the negative. To train our proposed PROFIT model, we cooperate with the Zhihu company, China's biggest question-andanswer-style knowledge base, enabling us to run our model on an unprecedented large-scale cQA data collection. After the training, we can learn the embeddings of all the archived topic tags offline and a newly posted question online. Based on this, we can recommend users the top topic tags according to the similarity between the embeddings of each topic tag and the newly posted question. This is in fact the online part. By conducting extensive experiments, our model is demonstrated to yield superior performance to several state-of-the-art baselines.

To sum up, the contributions in this work are threefold:

- As far as we know, this is the first work on question tagging in the cQA setting that leverages the DAG structure of topic tags by regularizing their hierarchical relations to transfer the knowledge from the top-down discriminative embedding learning. This model somehow alleviates the problem of imbalanced topic tag distribution.
- We present a deep parallel scheme, which jointly learns the topic tag and question embeddings and recommends topic candidates to the given question based on their embedding interaction. Unlike pure classification models, we actually cast the question tagging task into a ranking problem, which considers the topic semantics during the interaction.
- We comparatively justify our proposed PROFIT model over a large-scale and real-world dataset. In addition, we have released the codes and the involved parameters in this article.⁹

⁹https://question-tagging.wixsite.com/anonymous.

Particularly, we have released the new big cQA data with hierarchical topic tags, which will help boost the research.

The rest of this article is structured as follows. In Section 2, we briefly review the related literature. In Section 3, we detail our proposed model, followed by experimental results and analyses in Section 4. We finally conclude the work and discuss future directions in Section 5.

2 RELATED WORK

QA systems alleviate information overload by providing users simple and accurate answers. Besides, a great many research works have been conducted in QA systems, such as answer ranking [12, 16, 41] and question routing [28, 53]. Topic tags play a key role in question organization, routing, and search, and many question tagging methods have been proposed [1, 20, 24]. On the basis of state-of-the-art reviews, we mainly review related work about the hierarchical question classification, question tagging in cQA.

2.1 Hierarchical Question Classification

In the era of community-based question answering, representative cQA communities like Yahoo! Answers and Oshiete! Goo organize their content by hierarchical ontologies. In particular, when posting a question, the user is encouraged to select a single category tag at the leaf of the given ontology to label the question. In such a context, several automatic question classification approaches have sprung up miscellaneously: big-bang and top-down approaches. The former often trains a single classifier and employs it to assign one leaf node of the category tree to the given question [3, 36]. By contrast, the latter constructs one classifier per level of the hierarchy in the training phase and classifies each given question from the higher level to lower ones until it reaches a leaf category [29, 31].

Nevertheless, questions, more often than not, convey multiple topics. Thus, it gets more difficult to find a single appropriate category label to describe a given question. What is more, the leaf categories in the predefined ontology are insufficient to handle the ever-increasing questions with various topics. Taking Yahoo! Answer as an example, it has only 1,263 leaf-level nodes distributed over 26 top-level categories. This category vocabulary is extremely limited. Nishida and Fujimura in 2010 [26] partially alleviated these phenomena by annotating a given question with different abstractions, namely category, theme, and keywords. However, automatically generated themes and keywords are still far from satisfactory as compared to the user-generated topic tags in the cQA sites.

2.2 Question Tagging in cQA

Stack Overflow is a popular Q&A site focusing on technical questions about software development, helping software engineers to strengthen their abilities in software development, maintenance, and test processes. Topic tags are popular in Stack Overflow, used to search, describe, identify, and bookmark various software objects, as well as bridge the gap between social needs and technical development. In 2013, Saha et al. [33] introduced a discriminative model to suggest tags for questions on Stack Overflow. This model consists of three main steps: converting questions into vectors with a term frequency weighting scheme, training a discriminative model with an SVM classifier, and suggesting tags with the top similarity. In the same year, Xia et al. [47] presented a composite framework, named TagCombine, to solve the tag recommendation problem from three different views with three components. One year later, Wang et al. in [45] introduced EnTagRec, an automatic tag recommender based on historical tag assignments, which improved TagCombine by 27.3% with respect to Recall@5. To further improve the quality of tags in Stack Overflow, Wang

et al. [46] in 2018 proposed EnTagRec++, an advanced version of their prior work, i.e., EnTagRec. Beyond EnTagRec, EnTagRec++ not only integrates the historical tag assignments to software objects but also leverages the user information and an initial set of tags that a user may provide for tag recommendation. In 2018, Zhang et al. [49] introduced a multitasking-like convolutional neural network to learn semantic vectors for tag recommendation.

Besides Stack Overflow, there are some works on other cQA sites. Early in the last decade, pioneers conducted an in-depth analysis of the question labeling practices by contrasting the use of community-generated tags in the Live OnA service with the use of topic categories from a fixed taxonomy in the Yahoo! Answers service [22, 32]. They found that community tagging was related to higher levels of social interactions among users. The analysis of the most frequently used community tags reveals that active users may establish strong social ties around specific tags [11]. With the rise of collaborative cQA sites, a data-driven study over Quora was presented in 2013 [43], reporting that the user-follow-topic graph appeals to users in browsing and answering general questions. Nie et al. in 2014 [25] noticed the incompleteness of question tags in cQA sites and devised a novel scheme to automatically annotate questions, which was accomplished by finding appropriate tags from similar questions via an adaptive probabilistic hypergraph. Although existing studies have achieved compelling success in the question tagging on cQA sites, they failed to take into account the valuable hierarchical structures among topic tags. We have to mention that Zhihu.com organized a contest on question annotation in 2017¹⁰ and attracted a notable amount of participants. To support this competition, Zhihu released 3 million questions and 1,999 structured topic tags after desensitization. After going through all the submitted solutions, we found that none of them took the DAG hierarchy into consideration during the learning process.

3 DATA COLLECTION

Zhihu is a Chinese socialized question-and-answer website where questions are collaboratively created, answered, edited, and organized by its user community. Its website, zhihu.com, was launched on January 26, 2011. As of January 2019, Zhihu had obtained 220 million registered users, 29 billion monthly page views, and 26 million active users surfing the website on average for 1 hour daily, as well as accumulated 30 million questions, 130 million answers, and 35 million votes in total.¹¹ In Zhihu, users present their professionalism, find high-quality information to facilitate their decision making, contact people from whom they seek help, and build collaboration or partnership. Because users in Zhihu can generate a tremendous amount of content everyday, how to understand the created content and distribute it in a highly effective way is greatly desired in all cQA sites.

Routing content to appropriate users according to user-follow-topic relationship is natural, since the followed topic tags are able to better fulfill users' needs on the desired knowledge. In light of this, automatic and accurate labeling topic tags for questions in cQA sites will play a key role in enhancing the user experience and content distribution efficiency.

Cooperating with Zhihu company, we have obtained two datasets. One is a benchmark dataset manually pruned for a public contest (Dataset I), and the other is a real-world dataset without any preprocessing (Dataset II).

3.1 Dataset I

The Zhihu algorithm team, together with the IEEE Computer Association and IEEE China office, host Zhihu Machine Learning Challenge 2017, aiming to recommend topic tags to questions.

¹⁰https://biendata.com/competition/zhihu/.

¹¹https://www.infoq.cn/article/BNcC3WccELmPw6_LFFxP.



Table 1. Statistics of the Given Hierarchical DAG Structure of Topic Tags in Dataset I

Fig. 3. Statistics of our Datasets I and II, obtained from Zhihu.com. (a,d) The length distribution of topic tags measured by the number of words. (b,e) The frequency distribution of topic tags measured by the number of questions labeled by the same topic tag. (c,f) The number of questions labeled by different numbers of topic tags.

Participants should design an automatic tagging model for untagged questions. We thereafter name this dataset as Dataset I.

Dataset I consists of 2,999,967 questions labeled by 1,999 distinct topic tags. Thereinto, 90.4% of the topic tags come with descriptions. On average, each topic tag has 3.73 words and labels 3,513 questions. As mentioned before, most of the topic tags are coordinately organized into a hierarchical DAG by users and experts hired by Zhihu. In the given DAG, the 1,999 topic tags are linked by 2,653 edges and each edge represents the parent-child relation. It is worth emphasizing that, according to our statistics, 41.6% of question tags are internal nodes within the DAG structure, and the rest are leaf nodes. We list the statistical information of the given DAG structure in Dataset I in Table 1. In addition, we display the tag length distribution, frequency distribution, and number of topic tags per question in Figures 3(a) to 3(c).

On average, each question has been labeled by 2.34 topic tags. Stepping into the topic tags of the same question, we find that 16.7% of the topic tags are siblings and 30.4% of them are ancestors and descendants. Among the ancestor-descendant topic tags of the same question, the average depth distance is 1.52. Before training our model, we first automatically label each question with extra topic tags on the path between two ancestor-descendant topic tags, as shown in Figure 1, whereinto the topic tag in green is the missing one and we complete it. Ultimately, each question has 2.57 topic tags after completion.

Considering user privacy and data security, the contest does not provide the original texts of the questions and topics but uses numbered codes and numbered segmented words to represent text messages. Meanwhile, considering the vast use of Distributed Representation [19, 27], the contest provides embedding vectors at the level of character and word. These embedding vectors

Table 2.	Statistics of	the Given	Hierarchical	DAG Structure	e of Topic	Tags in	Dataset	11

Max Depth	Min Depth	Average Depth	# Internal Node	# Leaf Node	# Edge	# Average Children
19	1	8.83	4,033 (30.5%)	9,225 (69.5%)	17,874	4.4

are obtained by conducting training with Google word2vec and taking advantage of the mega text corpora provided by Zhihu.

3.2 Dataset II

Considering that Dataset I is preprocessed for the contest usage, whereby the data distribution does not coincide with the real-world cQA and the original textual data are not available due to privacy concerns, we hence cowork with Zhihu to build a large-scale and representative dataset. To be more specific, our dataset contains 1,707,023 textual questions and 10,843,647 corresponding answers. The questions are labeled with 13,258 distinct topic tags. Thereinto, 30.1% of the topic tags come with descriptions. The length and frequency distribution of topic tags are illustrated in Figures 3(d) and 3(e), respectively. On average, each topic tag has 5.73 words and labels 229 questions. As mentioned before, most of the topic tags are coordinately organized into a hierarchical DAG by users and experts hired by Zhihu. In the DAG, there are 13,258 topic tags linked by 17,874 edges, and each edge represents the parent-child relation. The meta-information of the DAG is shown in Table 2. It is worth emphasizing that, according to our statistics, 30.5% of question tags are internal nodes within the DAG, while 69.5% are leaf nodes.

In our dataset, the number distribution of topic tags for each question is illustrated in Figure 3(f). On average, each question has been labeled by 1.78 topic tags. Going deep into the topic tags of the same question, we find that 17.2% of the topic tags are siblings and 34.7% are ancestors and descendants. Among the ancestor-descendant topic tags of the same question, the average depth distance is 1.68. Before training our model, we first automatically label each question with extra topic tags on the path between two ancestor-descendant topic tags, as shown in Figure 1. Ultimately, each question has 2.07 topic tags after completion.

It's worth mentioning that ordinary users are not allowed to create new topic tags.¹² In other words, the DAG is relatively stable. In light of this, we only consider the topic tags in the DAG and target recommending the relevant topic tags to the newly posted question.

4 OUR PROPOSED METHOD

In this section, we first define some notations and then detail our proposed PROFIT model.

4.1 Notation and Problem Formulation

For ease of problem formulation, we first declare some notations. In particular, we use bold capital letters (e.g., **X**) and bold lowercase letters (e.g., **x**) to denote matrices and vectors, respectively. We employ nonbold letters (e.g., *X*) to represent scalars, and Greek letters (e.g., λ) as parameters. If not clarified, all vectors are in the column form.

Assume that we are given a set of N questions $Q = \{q_1, q_2, ..., q_N\}$ labeled by M topic tags $\mathcal{T} = \{t_1, t_2, ..., t_M\}$. The topic tags are preorganized into a tree-like DAG ontology \mathcal{G} , whereby the leaf and internal nodes in the DAG are topic tags and the edges to represent the parent-child relationship. Notably, a node in \mathcal{G} may have more than one parent. The deeper level in the DAG

¹²Only experienced users are allowed to create new topic tags. In particular, experienced users are the people who have more than five answers with no less than five votes.



Fig. 4. Schematic illustration of our proposed PROFIT model. It consists of three components: question embedding, topic tag embedding, and question-topic interaction.

a topic tag locates, the finer-grained concept it conveys, and vice versa. Therefore, the leaf nodes represent the most specific topic tags. Our research objective is to train a deep embedding model with Q and G, toward identifying the relevant topic tags from \mathcal{T} for a newly posted question q.

4.2 Deep Interactive Embedding Model

The offline part of our developed scheme for question tagging is the PROFIT model. Given the questions and their associated topic tags, the PROFIT model targets learning their embeddings for further similarity matching. It is worth noting that the topic tags in the given DAG are relatively stable and hence the embeddings of all the topic tags can be further used during the online tagging. As demonstrated in Figure 4, the PROFIT model comprises three components, namely, the question embedding layer, topic tag embedding layer, and question-topic interaction layer. We will detail them separately.

4.2.1 *Question Embedding Layer.* On the one hand, unlike the normal documents, in cQA, questions are typically short and most of the topic tags are usually phrases containing a few words. Statistically, on average, each question and each topic tag has 12.91 and 1.35 terms in Dataset I (13.12 and 2.08 terms in Dataset II), respectively. They thus do not provide sufficient contexts for similarity matching between questions and topic tags, especially for the keyword-based matching.

On the other hand, the traditional widely used methods on text embedding mainly rely on N-gram models, like uni-gram, bi-gram, and tri-gram. Despite their popularity, N-grams have two kinds of defects: (1) Uni-gram embedding destroys the word sequence, which in turn destroys much of the semantic structure. Even though N-gram models partially consider the word order in short or local context, they have very little sense about the semantics of the words or more formal distances among the words. (2) The dimensionality of N-gram representation is proportional to the dictionary size, leading to sparse and high-dimensional embeddings.

Considering that the CNN model and its variants have been very successful in computer vision tasks [9, 10, 35] and recommender systems [5, 48], with excellent performance in natural language processing tasks [44, 50, 52], we leverage the CNN model to produce the fixed-length vector representation of questions that preserve the semantic structure. In particular, we first build a very large vocabulary of phrases by segmenting all the sentences in our cQA set using the jieba tool.¹³ In contrast with words, phrases are critical for capturing the lexical meanings. On the ground of word embedding techniques [2], we map the semantic meaning of phrases into a geometric space. This is accomplished by associating a numeric vector to every phrase in the dictionary, such that the distance (e.g., Euclidean distance or, more commonly, cosine distance) between any two vectors would capture part of the semantic relationship between the two associated phrases. The dictionary of phrases is shared between the question embedding learning and the topic tag embedding learning. We pretrain phrase embeddings instead of random initialization. Embeddings of phrases in a question concatenate together to form a matrix, representing the question. In our method, we use simple convolution layers with filters of multiple sizes on top of the learned embedding matrix. This extracts various high-level features and encodes the semantic meaning of phrases by considering the intrinsic sequential information among phrases [15]. Our method in fact addresses the word order problem of the uni-gram model while avoiding dimensionality explosion of all the N-gram models.

In the cQA sites, apart from brief questions, users often describe their queries in detail, the socalled question description, which actually encodes rich contexts and empowers the short questions. In our Dataset I, approximately 2,142,746 out of 2,999,965 questions have the associated descriptions (71.4%). In Dataset II, the ratio is around 67.4%. To make full use of the question descriptions, we devise a Siamese-style neural network whereby two disjoint CNNs are trained. These two CNNs share the same network structure but not necessarily the identical parameters [38]. It is worth emphasizing that if a question has no description, we will by default treat the question itself as its description. After two CNNs, we concatenate the embeddings of the question and its description, followed by the multilayer perception to project the ultimate question embedding into a 1,024-D space.

Because the questions and their descriptions are all of varying lengths, we pad short questions with dummy words and truncate very long questions. In this way, we can deal with questions with the same length. In our work, we set the threshold of question and question description lengths as 30 and 150, respectively.¹⁴ It must be be noted that we do not add any max-pool or min-pool layer, which is frequently utilized in visual computing to capture the spatial invariance; e.g., they detect a dog regardless of where in the given image the dog is located. Nevertheless, in the field of text understanding, the spatial location within the given question is of importance.

4.2.2 Topic Tag Embedding Layer. Analogous to the question embedding, we first form the embedding matrix for each topic tag. We then obtain a $M \times L \times D_p$ tensor, whereby M, D_p , and L

¹³https://github.com/fxsjy/jieba.

¹⁴According to our statistics, the average lengths of questions and question descriptions are 12.91 and 62.10 in Dataset I (13.12 and 56.07 in Dataset II), respectively.

respectively denote the number of topic tags, the dimension of phrase embedding, and the length threshold we set for each topic tag (L = 5 in this work). After the multilayer perception, we reach a new $M \times 1024$ matrix, in which each row denotes a topic tag embedding.

As mentioned before, topic tags are not independent but hierarchically correlated. In our work, the inherent structural relatedness among topic tags is characterized via a DAG ontology \mathcal{G} , whereby the sibling and ancestor-descendent relationships are well organized. Graph has been applied successfully in social networks [23] and recommender systems [21]. To learn a discriminant and robust embedding of each topic tag, we have to encode the ontology into our model with the expectation that learned embeddings are capable of capturing the hierarchy among topic tags, such that the topic tags can reinforce each other to alleviate the problem of imbalanced topic tag distribution.

By observing and analyzing the given DAG carefully, we found the following: (1) Nodes locating at a deeper layer of the DAG are less frequently used to label questions in the cQA sites and vice versa. For example, 70% of the DAG nodes in Dataset II are leaf nodes. They, however, only label 25% of the questions. Thereby, the internal nodes see more samples and contain richer information. Inspired by this observation, transferring knowledge from the parent nodes to the child ones is a natural way to alleviate the problem of imbalanced topic tag distribution. (2) As known, a child node in a DAG may have multiple parents. We find that a parent node usually captures only one aspect of its child node. For instance, considering "apple" as a child node, it has two parents, "electronic product" and "fruit." (3) The semantics expressed by different parents of the same child are often complementary rather than redundant or conflicting, and they together make up the overall semantic of the child node. In light of this, we assume that the embedding of a child node can be approximated by a convex combination of the embeddings and parents of itself. It is a good way to consider a node itself for remembering the historical results during the updating. Based on this assumption, we regularize the topic tag embedding learning,

$$\begin{cases} \mathbf{u}_{i} = \sum_{t_{j} \in C(i)} \alpha_{ij} \mathbf{u}_{j}, \\ s.t \sum_{t_{j} \in C(i)} \alpha_{ij} = 1, \alpha_{ij} \ge 0, \end{cases}$$
(1)

where $\mathbf{u}_i \in \mathbb{R}^D$ is a *D*-dimensional embedding vector of the node $t_i \in \mathcal{T}$, C(i) denotes a set of nodes containing node t_i and its parents, and $\alpha_{ij} \in \mathbb{R}^+$ refers to the attention weight on t_j when calculating t_i . The above formulation is transmissible; namely, beyond parents, a child node can indirectly benefit from its ancestors after a few updates. Inspired by the work in [6], we adopt the softmax function to estimate the attention weight α_{ij} :

$$\begin{cases} \alpha_{ij} = \frac{\exp(f(\mathbf{u}_i, \mathbf{u}_j))}{\sum_{k \in C(i)} \exp(f(\mathbf{u}_i, \mathbf{u}_k))}, \\ f(\mathbf{u}_i, \mathbf{u}_j) = \mathbf{z}^T \tanh(\mathbf{W}[\mathbf{u}_i, \mathbf{u}_j]^T + \mathbf{b}), \end{cases}$$
(2)

where $f(\mathbf{u}_i, \mathbf{u}_j)$ measures how much \mathbf{u}_j constitutes its child \mathbf{u}_i , which is calculated via multilayer perception with an *H*-D hidden layer. The projecting matrix $\mathbf{W} \in \mathbb{R}^{H \times 2D}$, bias vector $\mathbf{b} \in \mathbb{R}^H$, and weight vector $\mathbf{z} \in \mathbb{R}^H$ are all parameters we need to learn.

4.2.3 Question-Topic Interaction Layer. Given the embedding of a question and the embeddings of all the topic tags, we perform interactions via dot production between the embedding vector and the embedding matrix. In this way, we obtain a score vector, whereby each element indicates the similarity degree between the given question and the corresponding topic tag. It is worth

emphasizing that the element values may be out of the range of [0,1]. In light of this, we leverage a sigmoid cross-entropy loss to project the similarity score into [0,1]:

$$\Phi = -\sum_{i} (\mathbf{y}_{i} * \log(s(\widehat{\mathbf{y}}_{i})) + (1 - y)\log(1 - s(\widehat{\mathbf{y}}_{i}))),$$
(3)

where $s(\cdot)$ denotes the sigmoid function given the input vector. $\hat{\mathbf{y}}_i \in \mathbb{R}^M$ and $\mathbf{y}_i \in \mathbb{R}^M$ are multihot vectors. Thereinto, $\mathbf{y}_i \in \mathbb{R}^M$ is the ground truth for the *i*th question, whereby the *j*th column is one if the *j*th topic tag is used to label the *i*th question, otherwise zero. $\hat{\mathbf{y}}_i \in \mathbb{R}^M$ is the predicted result for our proposed PROFIT model.

4.2.4 Implementation Details. The aforementioned model is trained offline. After training, for a newly posted question, we can online calculate its embedding via the well-trained PROFIT model and tag the newly posted question by measuring the embedding similarity between it and each topic tag, wherein the embedding of each topic tag is learned offline.

During the training phase, we feed our model batch by batch. Each batch comprises 128 question-description pairs, the associated topic tags, and one parent-child relation. After the word embedding layer, the batch of questions, question descriptions, and topic tags are denoted by $128 \times 30 \times D_p$, $128 \times 150 \times D_p$, and $M \times 5 \times D_p$ tensors, respectively. Thereinto, as aforementioned, 30, 150, and 5 respectively denote the maximum length we set for each question, question description, and topic tag. In addition, D_p refers to the embedding size, which is respectively set to be 128 and 256 for Dataset I and Dataset II. It is worth noting that we input all the *M* topic tags at one time and update the embeddings of topic tags batch by batch, whereby *M* stands for the number of topic tags in our DAG.

On the basis of the word embedding layer, we separately employ two 1-layer CNNs to capture the high-level abstracts of the input question and its description. Each TextCNN is equipped with 256*5 filters in 5 distinct sizes [2, 3, 4, 5, 7]. Following the convolution, we use batch norm, relu, and maxpooling. Thereafter, we obtain a 1,280-D embedding vector for each question and its description, respectively. Next, we input the concatenation of the embeddings of a question and its description to a fully connected layer (size $2,560 \times 1,024$) followed by batch norm and Relu. We ultimately reach a 1,024-D embedding vector to represent a question. Here we set the dropout rate as 0.5 for the fully connected layer, and we will detail why we choose this rate in the experimentations. Thus far, we have represented the batch of questions as a $128 \times 1,024$ matrix.

As to the topic embedding learning, we obtain a $M \times 5 \times D_p$ tensor after the word embedding layer. We then take the transpose of this tensor and reach a new $M \times D_p \times 5$ tensor. We input this tensor into a fully connected layer (size: 5×1), followed by a Relu activation function. We next lay another fully connected layer (size: $D_p \times 1,024$; dropout rate: 0.5), followed by a batch norm. Up to now, we reach a $M \times 1,024$ topic tag matrix. In fact, we only update a parent-child relation in each batch using Equation (1). To be more specific, we parse the given DAG structure and identify all the parent-child relations. We sequentially input all the relations one by one (i.e., batch by batch) to cover all the hierarchical relations.

Once we obtain the embeddings of 128 questions and all the *M* topic tags, we multiply each question embedding with the topic tag matrix to estimate the similarities between the given question and all the topic tags.

5 EXPERIMENTS

In this section, we conducted experiments to justify the effectiveness of our proposed PROFIT model and its components. Extensive experiments have been comparatively conducted over Dataset I and Dataset II.

Large-Scale Question Tagging via Joint Question-Topic Embedding Learning

5.1 Experimental Settings

Given a question and a set of topic tags, we targeted generating an ordered list of tag candidates. This is in fact a ranking problem. To measure the performance of our model and the baselines, we adopted three metrics, the same as those in the Zhihu Machine Learning Challenge 2017: $precision_{rank}$, recall and F-measure_{rank} and the standard metric *Precision@k*.

Traditionally, precision is defined as the fraction of the documents retrieved that are relevant to the user's need. Such definition, however, does not consider the position of the ranked documents. Inspired by the metric of discounted cumulative gain (DCG) that is a measure of quality rating, we reformulated the precision metric by placing stronger emphasis on retrieving the relevant documents and their positions. Formally, it is written as

$$Precision_{rank}@k = \sum_{i=1}^{k} \frac{rel_i}{log(i+1)},$$
(4)

where the relevance values of a topic tag in our work are binary, $rel_i \in \{0, 1\}$. If the *i*th topic tag in the generated ranking list is used to label the given question, $rel_i = 1$; otherwise, $rel_i = 0$. Note that $Precision_{rank} @k$ is possible to be larger than 1 according to our definition. This definition derives from the official website of the Zhihu contest.¹⁵ The reasons we use these three "nonstandard" metrics are twofold: (1) These three metrics are more suitable for our task, as cQA sites usually recommend at most a fixed number of topic tags for each question (e.g., the number is five in Stack Overflow and Zhihu, namely, each question can be tagged with at most five topics). Accordingly, given a fixed k, it is more important to consider the ranking positions of the ground-truth topic tags. The traditional standard metrics, like precision and F-measure, are unable to capture the specific position information. (2) Existing studies that use the same dataset (i.e., Zhihu Contest) as ours, like [8], adopt the same evaluation metrics, i.e., $precision_{rank}$, recall, and F-measure_{rank}. Consequently, we used the same metrics to facilitate the comparisons.

The second metric in this work is Recall@k, which measures the fraction of the topic tags used to label the given question that is successfully ranked at the top k positions. The third metric is the *F*-measure_{rank}, which is the weighted harmonic mean of $precision_{rank}$ and recall, formulated as

$$F\text{-}measure_{rank}@k = \frac{precision_{rank}@k * Recall@k}{precision_{rank}@k + Recall@k}.$$
(5)

This is also known as the F_1 measure, because *recall* and *precision* are evenly weighted. In this work, we set k as 5 and used *Precision_r*, *Recall*, and *F-measure_r* to denote *Precision_{rank}@*5, *Recall@*5, and *F-measure_{rank}@*5, respectively. Besides, we used *P@*5 and *P@*10 to denote *Precision@*5 and *Precision@*10, respectively.

We randomly split all the questions in Datasets I and II into three chunks, respectively. To be more specific, in Dataset I, we had 2,700,000 questions for training, 149,965 for validation, and 150,000 for testing. As for Dataset II, we used 1,307,023 questions for training, 192,977 for validation, and 207,023 for testing. The training set is used to learn our model, the validation set is to tune the optimal parameter settings, and the testing one is to report the final results.

5.2 Baselines

To demonstrate the effectiveness of our proposed PROFIT model, we compared it with the following state-of-the-art methods:

¹⁵https://biendata.com/competition/zhihu/evaluation/.

- **Matching**: This baseline is straightforward. We averaged the embeddings of all the words in each question-description pair and topic tag to represent the question and topic tag, respectively. Thereafter, we calculated their cosine similarity to search the top *K* related topic tags.
- Naive Bayes: It is a classical supervised learning model for classification tasks [14]. In this baseline, we cast the question tagging problem into a regression task and each topic tag was treated as a category. We worked toward selecting the top *K* most relevant topic tags. We calculated the tf.idf of the short texts as their input features.
- **FastText**: This baseline scales a linear model to a very large corpus with a large output space in the context of text classification [13]. This is accomplished by enhancing the linear models with a rank constraint and a fast loss approximation. We could train fastText on more than one billion words in less than 10 minutes by using a standard multicore CPU, while achieving performance on par with the state-of-the-art methods.
- **TextCNN**: Indeed, this baseline is a simple CNN with one layer of convolution on top of the word vectors obtained from an unsupervised neural language model [15]. It achieves great success in natural language processing tasks like sentiment analysis [52], machine translation [44], and text summarization [50].
- **TextRNN**: This baseline integrates a recurrent structure to capture contextual information as far as possible when learning word representations [17].
- L2R: This baseline presents a convolutional neural network architecture for reranking the pairs of short texts, where it learns the optimal representation of text pairs and a similarity function to relate them in a supervised way from the available training data [34]. This network takes only words as the input, thus requiring minimal preprocessing. Specifically, we treated each question and its tag as a positive pair and constructed the negative pair by randomly sampling a tag not associated with the given question.
- **RegionEmb**: This baseline represents each word with two parts: the embedding of the word itself and a weighting matrix characterizing its interaction with the local context [30]. Besides, it outperforms existing methods in the task of text classification on several benchmark datasets.
- **PBAM**: It is a recently proposed deep neural network for question tagging, which utilizes a position-based attention mechanism to model the question text [39].
- **DGCNN**: This baseline is proposed for hierarchical text classification, which utilizes a graph convolutional network to learn the original raw text and the recursive hierarchical segmentation for relationship modeling [29]. Now we used the model for question tagging, which also has hierarchical relationships among tags.
- PROFIT: This is our proposed end-to-end deep interactive model for question tagging.

To ensure a fair comparison, all the above models were trained on the same training set and justified on the same testing one. In addition, all the models were carefully tuned to reach their optimal settings, and their best performance in terms of multiple metrics was reported.

5.3 Overall Comparison

The comparative results over two datasets between our proposed PROFIT model and several stateof-the-art baselines for question tagging are summarized in Table 3. Let us analyze the experimental results over Dataset I first: (1) As expected, the unsupervised learning method, i.e., Matching, is the worst, since it does not encode any label information. (2) All the deep learning models remarkably outperform the shallow learning one. Deep learning algorithms seek to exploit the unknown structure in the input distribution in order to discover better representations, often at multiple

Methods			Dataset I					Dataset II		
	$Precision_r$	Recall	F-measure _r	P@5	P@10	$Precision_r$	Recall	F-measure _r	P@5	P@10
Matching	0.1938	0.0980	0.0651	0.0512	0.0276	0.0242	0.0150	0.0093	0.0062	0.0021
Naive Bayes	0.2895	0.1289	0.0892	0.0777	0.0499	0.2615	0.1508	0.0957	0.0497	0.0193
L2R	0.8895	0.3876	0.2700	0.1548	0.1166	0.7384	0.4118	0.2643	0.1139	0.0974
FastText	1.4026	0.5385	0.3891	0.2472	0.1531	1.1427	0.5400	0.3667	0.1975	0.1167
RNN	1.4011	0.5372	0.3883	0.2469	0.1529	1.2522	0.5837	0.3981	0.2139	0.1264
TextCNN	1.4070	0.5410	0.3907	0.2502	0.1546	1.2549	0.5867	0.3998	0.2161	0.1269
RegionEmb	1.4033	0.5404	0.3901	0.2487	0.1544	1.2457	0.5864	0.3987	0.2146	0.1266
PBAM	1.4031	0.5390	0.3894	0.2473	0.1533	1.2359	0.5820	0.3957	0.2128	0.1256
DGCNN	1.4103	0.5423	0.3917	0.2519	0.1549	1.2620	0.5870	0.4006	0.2169	0.1273
PROFIT	1.4257	0.5490	0.3964	0.2557	0.1567	1.3059	0.6108	0.4162	0.2260	0.1312

Table 3. Performance Comparison between Our Proposed Model and Several State-of-the-ArtBaselines over Two Datasets, Measured by Four Metrics

levels, with higher-level learned features defined in terms of lower-level features. Automatically learning features at multiple levels of abstraction allows the classifier model to learn complex functions by mapping the input of short texts to the output label directly from data, which do not completely depend on the human-crafted features, whereas, the shallow learning methods should heavily rely on the human-crafted features. (3) Among the deep models, FastText, TextCNN, TextRNN, RegionEmb, PBAM, and DGCNN achieve comparable performance, which is much better than that of L2R. This is probably due to the fact that for each positive sample, we only randomly constructed a negative one when training L2R, whereas the other three deep models can sufficiently incorporate the information from the negative samples. (4) Our model is consistently better than all the baselines. This is because on the basis of deep models, we also considered the label semantics and the label hierarchies. Beyond the traditional labeling models, which understand the category labels via their associated samples, our model directly encodes the label semantics during the learning. Meanwhile, the label hierarchies convey the label correlations, enabling knowledge transfer among various categories, especially from the top down. (5) Compared with the question tagging baseline PBAM, our method achieves the better performance, which may be due to the fact that PBAM overlooks both the label semantic information and the hierarchical relationships among question topics. Moreover, compared with DGCNN, which indeed also explores the label via the recursive hierarchical segmentation, our PROFIT model has better performance. This verifies the effectiveness of our DAG mechanism to transform knowledge from parent nodes to child nodes, which contributes to overcoming the unbalance of topic tags. The standard Precision@Kmetric shows consistent results with the other three metrics, which verifies the performance of our method. Besides, we can find that *Precision*_r is higher than P@5 of the same method generally due to the difference of their formulas.

As to the performance comparison over Dataset II, we observed the following points: (1) The overall trend is almost the same as Dataset I. Specifically, the unsupervised method is the worst, the shallow learning method is the second to last, and the deep model is the best. (2) It is worth highlighting that the performance of FastText drops considerably, much worse as compared to other deep models. This is because other deep models like TextCNN and TextRNN are able to capture the sequential characteristics of texts, which is even more obvious in the larger and more sparse Dataset II. As mentioned before, Dataset II has around 13,000 labels, more than six times of Dataset I. However, the average label frequency of Dataset II is only around 229, which is 3,861 in Dataset I. (3) Our model is much more robust over Dataset II as compared to the other competitors.

Methods			Dataset I					Dataset II		
	$Precision_r$	Recall	F -measure $_r$	P@5	P@10	$Precision_r$	Recall	F -measure $_r$	P@5	P@10
No_DES	1.3697	0.5265	0.3803	0.2381	0.1495	1.2559	0.5896	0.4012	0.2173	0.1274
No_DAG	1.4203	0.5450	0.3939	0.2532	0.1556	1.2769	0.5962	0.4064	0.2201	0.1281
PROFIT	1.4257	0.5490	0.3964	0.2557	0.1567	1.3059	0.6108	0.4162	0.2260	0.1312

Table 4. Component-Wise Validation of Our Proposed PROFIT Model by Eliminating One Component Each Time

We use No_DES and No_DAG to denote new models without considering the question descriptions and the DAG structure, respectively.

This indicates the superiority of integrating the DAG structure into our model, which is able to transfer the knowledge among labels and hence boost the performance, especially for those labels with sparse samples. Statistically, 70% of the DAG nodes are the leaf ones, whereby the average frequency of leaf nodes and internal nodes is 84 and 524, respectively, in Dataset II. In other words, it is very hard, if not impossible, to learn the discrimination of the leaf nodes. With the help of the DAG structure, our model is capable of transferring the knowledge from the up ancestor nodes to down descendants. When it comes to Dataset I, the frequency of the leaf node is 2,656 on average, which is more than enough to learn a robust model for each label. That is why our model demonstrates much better performance over Dataset II.

It is worth noting that although Dataset I is released by the Zhihu Machine Learning Challenge 2017, it is intractable to compare our proposed PROFIT model with the champion solution of this challenge, as Dataset I is released as the training set of the challenge, while the exact online testing set is unavailable to us. In fact, we have studied the solutions of the top five teams and noticed that they mainly adopted the ensemble strategy with the basic models, like FastText, TextRNN, and TextCNN, to fulfill the task. Accordingly, we have introduced these basic models as the baselines in this work, and our PROFIT model shows the superiority over these baselines. Moreover, we further checked our PROFIT model with the ensemble manner and achieved superior performance (i.e., 0.4421) with respect to the F-measure on Dataset I.

5.4 On the Component-Wise Validation

In this subsection, we conducted experiments to answer two research questions: (1) Do the question descriptions add value to our model? (2) How much help does our model get from the DAG structure?

In our PROFIT model, we devised a Siamese-style neural network to make full use of the brief questions and their long descriptions, whereby two disjoint CNNs were trained. These two CNNs share the same network structure but not necessarily identical parameters. To well answer the first research question, we eliminated one branch of the Siamese-style neural network and only kept the CNN corresponding to the brief question. To answer the second research question, we eliminated the DAG-guided regularizer from our model and ensured the left unchanged.

The experimental results are summarized in Table 4. We used No_DES and No_DAG to denote new models without considering question descriptions and the DAG structure, respectively. It can be seen that No_DES is the worst model across two datasets. Such phenomenon clearly reflects that the question descriptions contain rich context information, which plays a pivotal role in enhancing the question representation learning. In addition, we noted that the difference gap between our PROFIT model and the No_DAG one is widened from Dataset I to Dataset II. This further confirms our analysis before that our PROFIT model is much more robust when applied to a large-scale and sparse dataset.

Mathad			Dataset I	_	-	Dataset II				
Method	$Precision_r$	Recall	F-measure _r	P@5	P@10	$Precision_r$	Recall	F-measure _r	P@5	P@10
Down-Top	1.4126	0.5444	0.3929	0.2526	0.1552	1.2827	0.6023	0.4098	0.2232	0.1292
Top-Down	1.4257	0.5490	0.3964	0.2557	0.1567	1.3059	0.6108	0.4162	0.2260	0.1312
Top-Down & Down-Top	1.4179	0.5462	0.3943	0.2535	0.1559	1.3003	0.6104	0.4154	0.2252	0.1304

Table 5. Performance Comparison between Different DAG Encoding Methods

For simplicity, we respectively denote Father=Sum(Children)+Itself as Down-Top, Child=Sum(Fathers)+Itself as Top-Down, and the fusion of the two methods as Top-Down&Down-Top.

5.5 On the DAG Structure

In our PROFIT model, we assumed that the embedding of a child node can be approximated by a convex combination of the embeddings of itself and its parents. Two natural questions are about how to represent a parent node by a convex combination of itself and its children and utilize the two representing methods together. In this subsection, we carried out experiments over two datasets to validate these different representations of a node.

The comparison results among various DAG encoding methods are displayed in Table 5. To facilitate the notation, we denoted Father = Sum(Children) + Itself as Down-Top, which vividly demonstrates that the knowledge is transferred from down to top. Analogously, we used Top-Down to represent Child = Sum(Fathers) + Itself. Besides, we used Top-Down and Down-Top to represent the fusion of the two methods. From this table, it is obvious that the Top-Down method is consistently superior to the Down-Top method regarding all the metrics across two datasets, especially on Dataset II. This phenomenon is caused by the different frequency distribution between the child nodes and the parent nodes in our datasets; namely, nodes locating at a deeper layer of the DAG are less frequently used to label questions in the social QA sites and vice versa. According to our statistics, although 70% of the DAG nodes are leaf nodes, they label only 25% of all the archived questions. Therefore, the nodes at higher layers have more positive samples and hence hide richer information. That is why transferring the knowledge from the parent nodes to the child ones is effective. Even so, it is worth mentioning that in some cases whereby the leaf nodes are more frequently used, the Down-Top method may be an optimal option. It is also worth mentioning that the fusion method (i.e., Down-Top + Top-Down) performs better than Down-Top but worse than Top-Down. The inferior performance of the fusion method to Top-Down suggests that transferring knowledge from down to top can interfere with top-down knowledge propagation. Such phenomenon may be caused by the different frequency distribution between the child nodes and the parent nodes in our datasets; namely, nodes locating at a deeper layer of the DAG are less frequently used to label questions in the cQA sites and vice versa. According to our statistics, although 70% of the DAG nodes are leaf nodes, they label only 25% of all the archived questions. Therefore, the nodes at higher layers have more positive samples and hence hide richer information. That is why it is an effective way to transfer knowledge from the parent nodes to the child ones.

5.6 Illustration of Attentive Weights

As analyzed before, a child node may have multiple parents. On average, a child in Dataset I and II has 1.33 and 1.35 parents, respectively. We noted that a parent node usually captures only one aspect of its child node. Considering "Library" as an example, it inherits from four parents "Public Space," "Public Building," "Book," and "Organization." Each parent usually constitutes one aspect of its child with a certain degree. That was why we leveraged the attentive weighting mechanism in learning the child node representation.

Id	Examples
1	Apple Store=0.64*Apple Store+0.36*Apple Inc.
2	SoftBank =0.43*SoftBank+0.57*The Japanese group
3	Airport=0.68*Airport+0.23*public space+0.09*public building
4	Linux=0.67*Linux+0.25*open-source software+0.08*OS
5	Larry Page=0.65*Larry Page+0.23*Google+0.09*Entrepreneur+0.03*Alphabet
6	Nokia=0.67*Nokia+0.26*Mobile phone manufacturer+0.04*Technology company+0.03* Listed company
7	Milk tea=0.42* Milk tea+0.29*Milk+0.14*Tea+0.11*Drink+0.04*Soft Drinks
8	Library=0.48*Library+0.23*public space+0.19*public building+0.09*Book+0.01*Organization
9	OneNote=0.38*OneNote+0.37*Note+0.11*Note taking app+0.08*PKM+0.05*GDT+0.01*Microsoft Office
10	Worktile=0.30*Worktile+0.32*Teamwork equipment+0.18*OA+0.09*Office software +0.06*PM+0.05*PM system

Table 6. Illustration of the Attentive Weights for Node Representation in the Given DAG

We randomly select 10 nodes with one to five parents and display their weights in the convex combination.



Fig. 5. Convergence analysis of our model over two datasets by measuring the loss decrease with respect to the number of epochs. (a) Loss curve over Dataset I. (b) Loss curve over Dataset II.

To intuitively illustrate the attention results, we randomly selected 10 nodes with one to five parents and listed their convex combination in Table 6. From the selected examples in Table 6, we had the following observations: (1) The weights of the children themselves are the largest ones in most cases. This indicates that beyond the knowledge inherited from the ancestor nodes, child nodes try their best to maintain their own information. (2) It is consistent with our assumption that different parents contribute differently to their child nodes, and some even tend to have zero contribution, like "Soft Drinks" to its child "Milk Tea."

5.7 Convergence and Parameter Analysis

To demonstrate the convergence of our proposed PROFIT model, we plotted the loss curves over two datasets with respect to the number of epochs in Figures 5(a) and 5(b). It can be seen that our PROFIT model is able to converge over two datasets within less than 20 epochs. Notably, it tends to be slower on Dataset II, as compared to Dataset I. This is because Dataset II contains many more topic tags. It is mentioned that our model takes about 55 minutes and 58 minutes every epoch in the training phase for Dataset I and Dataset II, and 0.12ms and 0.83ms for one question in the testing phase for Dataset I and Dataset II, respectively.

We also studied the performance of our PROFIT model regarding the varying dropout rates and the number of convolutional kernels (also called filters). The experimental results are demonstrated in Figure 6. From Figures 6(a) to 6(c), we can see that our PROFIT model reaches the optimal performance when dropping out half neurons during the training, no matter in Dataset I or II.



Fig. 6. Parameter analysis of our PROFIT model over two datasets. (a)–(c) Performance of our model in terms of three metrics over two datasets by varying the dropout rate. (d)–(f) Performance of our model in terms of three metrics over two datasets by varying the number of convolutional kernels.

This indicates that our model has great transportability. Surprisingly, the recall@5 of our PROFIT model over Dataset II is much higher than that over Dataset I, although Dataset II is much larger and more sparse. After analyzing the data, we found that each question on average is annotated with 2.57 and 2.07 topic tags in Datasets I and II, respectively. Also, we noted that our model is capable of ranking 1.41 and 1.26 correct topic tags on average at the top five positions over Datasets I and II, respectively, with only a slight difference. This statistic reveals that our model is robust and hence performs well over a challenging dataset.

As discussed before, on the basis of the word embedding layer, we separately employed two one-layer CNNs to capture the high-level abstracts of the input question and their description. Each TextCNN is equipped with K*5 convolutional kernels (filters) in 5 distinct sizes [2, 3, 4, 5, 7]. We varied K and recorded the performance of our model as illustrated in Figures 6(d) to 6(f). It is demonstrated that our model performs the best once K is at 256. This experimental result tells us that more kernels are not necessary to add value, since we have to learn more parameters that require more samples.

5.8 Case Study

To gain deeper insights about the performance of our proposed PROFIT model in question tagging, we listed several example questions with their predicted topic tags via different methods in Table 7. From the examples in Table 7, we had the following observations: (1) Almost all the predicted topic tags by our PROFIT model are semantically related to the question, although some of them are incorrect as compared to the ground truth. This reflects that our model can recommend semantically related labels, while the accurate question tagging is still challenging due to the fact that the correct topic tags can easily be overwhelmed by semantically related tags. For example, PBAM recommends "Microsoft Edge" for the first question but misses the ground truth tag "Browser," which may be due to the semantic correlation between "Microsoft Edge" and "Microsoft." (2) Our PROFIT model performs better than other baselines in the case study, which is consistent with the aforementioned information retrieval metrics. For example, only our PROFIT model manages to predict the "DV" tag for the second question. According to our statistics, "DV"

Questions and Associated	Predicted Topic Tags							
Topic Tags	TextCNN	PBAM	DGCNN	PROFIT				
Will Microsoft try to save Internet Explorer? What measures will Microsoft take? correct topic tags:, Microsoft, Internet Explorer, Browser	Microsoft, Internet Explorer 10, Microsoft(China), Internet Explorer, Internet Explorer 9	Microsoft, Microsoft Windows, Windows 10, Microsoft Edge, Internet Explorer	Microsoft, OS, Internet Explorer, Microsoft(China), Internet Explorer 9	Microsoft, Browser, Internet Explorer, Microsoft Services, PC Browser				
I want a powerful digital video, which can take good photos and videos that won't break my bank. What are your recommendations? correct topic tags: Digital Video, Digital Product, Digital, DV	Digital, Camera, Video Equipment, Digital Video, Vidicon	DSLR, Digital Camera, Electronic Products, Digital Product , Camera	Digital Video, Camera, Camera Company, Digital, Digital Product	Digital Video, Digital Product, Digital Camera, DV, Digital				
What are the reasons people like Grey's Anatomy? correct topic tags: Grey's Anatomy, US TV Series	Grey's Anatomy, US TV Series, Actor, TV Series Recommendation, Actress	Grey's Anatomy, Hospital, TV Play, US TV Series, Gender Relations	Grey's Anatomy, US TV Series, TV Series Recommendation, US TV Series, Japanese TV Series	Grey's Anatomy, US TV Series Recommendation, US TV Series, Doctor, TV Play				

Table 7. Case Study

The test samples of PROFIT and some best baselines on Dataset II.

is a leaf node in the DAG with limited samples. This demonstrates the advantage of taking into account the DAG structure to transform knowledge from the internodes to leaf nodes.

6 CONCLUSION AND FUTURE WORK

This article presents an end-to-end deep interactive embedding model to label questions with topic tags in social QA sites, whereby the topic tags are predefined and preorganized into a DAG structure. Instead of applying the pure classification models, we actually cast the question tagging task into a ranking problem. This deep model jointly learns the embeddings of questions and topic tags by projecting them into the same semantic space and then performs the interaction for similarity measure. During the learning, it leverages the DAG structure of topic tags to regularize their hierarchical relations for discriminative embedding learning, which is able to address the problem of imbalanced topic distribution by transferring the knowledge among topic tags. To justify our model, we conducted extensive experiments over two large-scale datasets: one is a benchmark dataset manually pruned for a public contest (Dataset I), and the other is a real-world dataset without any preprocessing (Dataset II). We notice that the experimental results support the following points: (1) our proposed model remarkably outperforms several state-of-the-art methods for the question tagging task, (2) the DAG structure adds value to the topic tag embedding learning, (3) question descriptions are able to strengthen the discriminations of the short questions, and (4) as compared to the classification problem, our PROFIT model benefits from the topic semantics by incorporating the interaction component.

In the future, we plan to deepen and widen our work from the following aspects: (1) Due to the practical concern that topic tags in many cQA cites, like Stack Overflow, are not preorganized explicitly as a hierarchical structure, we will extend our method to these more challenging cQA sites. For example, we can conduct a topic tag graph by linking the related tags and synonymous

Large-Scale Question Tagging via Joint Question-Topic Embedding Learning

tags. Besides, existing topic tag structures, like the Zhihu topic tag structure, can be transformed to the cQA cites without a preorganized hierarchical structure. (2) There still remains much room for improvement with regard to enhancing the performance of DAG modeling. One possible improvement can be made by treating the topic hierarchical structure as a graph. Considering the compelling success of Graph Convolutional Neural Networks (GCNs) in various machine learning tasks, we will explore its potential in the context of question-topic embedding learning.

REFERENCES

- [1] Fabiano Belém, Eder Martins, Tatiana Pontes, Jussara Almeida, and Marcos Gonçalves. 2011. Associative tag recommendation exploiting multiple textual features. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1033–1042.
- Yoshua Bengio, Pascal Vincent, and Christian Janvin. 2003. A neural probabilistic language model. *Journal of Machine Learning Research* 3, 6 (2003), 1137–1155.
- [3] Li Cai, Guangyou Zhou, Kang Liu, and Jun Zhao. 2011. Large-scale question classification in cQA by leveraging Wikipedia semantic knowledge. In Proceedings of the ACM International Conference on Information and Knowledge Management. 1321–1330.
- [4] Xin Cao, Gao Cong, Bin Cui, Christian S. Jensen, and Quan Yuan. 2012. Approaches to exploring category information for question retrieval in community question-answer archives. ACM Transactions on Information Systems (TOIS) 30, 2 (2012), 7.
- [5] Zhiyong Cheng, Xiaojun Chang, Lei Zhu, Rose C. Kanjirathinkal, and Mohan Kankanhalli. 2019. MMALFM: Explainable recommendation by leveraging reviews and images. ACM Transactions on Information Systems (TOIS) 37, 2 (2019), 16.
- [6] Edward Choi, Mohammad Taha Bahadori, Le Song, Walter F. Stewart, and Jimeng Sun. 2017. GRAM: Graph-based attention model for healthcare representation learning. In Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining. 787–795.
- [7] Erik Choi, Vanessa Kitzie, and Chirag Shah. 2012. Developing a typology of online Q&A models and recommending the right model for each question type. *Proceedings of the American Society for Information Science and Technology* 49, 1 (2012), 1–4.
- [8] Cunxiao Du, Zhaozheng Chen, Fuli Feng, Lei Zhu, Tian Gan, and Liqiang Nie. 2019. Explicit interaction model towards text classification. In Proceedings of the AAAI Conference on Artificial Intelligence. 6359–6366.
- [9] Zan Gao, Deyu Wang, Xiangnan He, and Hua Zhang. 2018. Group-pair convolutional neural networks for multi-view based 3D object retrieval. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [10] Junwei Han, Le Yang, Dingwen Zhang, Xiaojun Chang, and Xiaodan Liang. 2018. Reinforcement cutting-agent learning for video object segmentation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 9080–9089.
- [11] Paul Heymann, Daniel Ramage, and Hector Garcia-Molina. 2008. Social tag prediction. In Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval. 531–538.
- [12] Heyan Huang, Xiaochi Wei, Liqiang Nie, Xianling Mao, and Xin-Shun Xu. 2018. From question to text: Questionoriented feature attention for answer selection. ACM Transactions on Information Systems (TOIS) 37, 1 (2018), 6.
- [13] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In Proceedings of the Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers. 427–431.
- [14] Sang Bum Kim, Kyoung Soo Han, Hae Chang Rim, and Sung Hyon Myaeng. 2006. Some effective techniques for naive Bayes text classification. *IEEE Transactions on Knowledge and Data Engineering* 18, 11 (2006), 1457–1466.
- [15] Yoon Kim. 2014. Convolutional neural networks for sentence classification. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP'14). 1746–1751.
- [16] Jeongwoo Ko, Luo Si, Eric Nyberg, and Teruko Mitamura. 2010. Probabilistic models for answer-ranking in multilingual question-answering. *ACM Transactions on Information Systems (TOIS)* 28, 3 (2010), 16.
- [17] Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. 2015. Recurrent convolutional neural networks for text classification. In 29th AAAI Conference on Artificial Intelligence.
- [18] Yann Lecun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86, 11 (1998), 2278–2324.
- [19] Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In Proceedings of the International Conference on International Conference on Machine Learning. II–1188.

- [20] Chenliang Li, Shiqian Chen, Jian Xing, Aixin Sun, and Zongyang Ma. 2018. Seed-guided topic model for document filtering and classification. ACM Transactions on Information Systems (TOIS) 37, 1 (2018), 9.
- [21] Yongqi Li, Meng Liu, Jianhua Yin, Chaoran Cui, Xin-Shun Xu, and Liqiang Nie. 2019. Routing micro-videos via a temporal graph-guided recommendation system. In Proceedings of the 27th ACM International Conference on Multimedia. 1464–1472.
- [22] Yandong Liu and Eugene Agichtein. 2008. On the evolution of the yahoo! answers QA community. In Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval. 737–738.
- [23] Liqiang Nie, Xuemeng Song, and Tat-Seng Chua. 2016. Learning from multiple social networks. Synthesis Lectures on Information Concepts, Retrieval, and Services 8, 2 (2016), 1–118.
- [24] Liqiang Nie, Meng Wang, Luming Zhang, Shuicheng Yan, Bo Zhang, and Tat-Seng Chua. 2015. Disease inference from health-related questions via sparse deep learning. *IEEE Transactions on knowledge and Data Engineering* 27, 8 (2015), 2107–2119.
- [25] Liqiang Nie, Yi Liang Zhao, Xiangyu Wang, Jialie Shen, and Tat Seng Chua. 2014. Learning to recommend descriptive tags for questions in social forums. ACM Transactions on Information Systems (TOIS) 32, 1 (2014), 1–23.
- [26] Kyosuke Nishida and Ko Fujimura. 2010. Hierarchical auto-tagging: Organizing Q&A knowledge for everyone. In Proceedings of the ACM International Conference on Information and Knowledge Management. 1657–1660.
- [27] Alberto Paccanaro and Geoffrey E. Hinton. 2002. Learning distributed representations of concepts using linear relational embedding. *IEEE Transactions on Knowledge and Data Engineering* 13, 2 (2002), 232–244.
- [28] Aditya Pal. 2015. Metrics and algorithms for routing questions to user communities. ACM Transactions on Information Systems (TOIS) 33, 3 (2015), 14.
- [29] Hao Peng, Jianxin Li, Yu He, Yaopeng Liu, Mengjiao Bao, Lihong Wang, Yangqiu Song, and Qiang Yang. 2018. Largescale hierarchical text classification with recursively regularized deep graph-cnn. In Proceedings of the Conference on World Wide Web Conference on World Wide Web. 1063–1072.
- [30] Chao Qiao, Bo Huang, Guocheng Niu, Daren Li, Daxiang Dong, Wei He, Dianhai Yu, and Hua Wu. 2018. A new method of region embedding for text classification. In *International Conference on Learning Representations*.
- [31] Bo Qu, Gao Cong, Cuiping Li, Aixin Sun, and Hong Chen. 2014. An evaluation of classification models for question topic categorization. *Journal of the Association for Information Science and Technology* 63, 5 (2014), 889–903.
- [32] E. M. Rodrigues, N. Milic-Frayling, and B. Fortuna. 2008. Social tagging behaviour in community-driven question answering. In Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology. 112–119.
- [33] Avigit K. Saha, Ripon K. Saha, and Kevin A. Schneider. 2013. A discriminative model approach for suggesting tags automatically for Stack Overflow questions. In *Mining Software Repositories*. 73–76.
- [34] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to rank short text pairs with convolutional deep neural networks. In Proceedings of the ACM International Conference on Research and Development in Information Retrieval. 373–382.
- [35] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. 2014. CNN features off-the-shelf: An astounding baseline for recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. 806–813.
- [36] Amit Singh and Karthik Visweswariah. 2011. CQC: Classifying questions in CQA websites. In Proceedings of the ACM International Conference on Information and Knowledge Management. 2033–2036.
- [37] Ivan Srba and Maria Bielikova. 2016. A comprehensive survey and classification of approaches for community question answering. *ACM Transactions on the Web (TWEB)* 10, 3 (2016), 18.
- [38] Ting Su, Craig Macdonald, and Iadh Ounis. 2019. Ensembles of recurrent networks for classifying the relationship of fake news titles. In Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval. 893–896.
- [39] Bo Sun, Yunzong Zhu, Yongkang Xiao, Rong Xiao, and Yungang Wei. 2018. Automatic question tagging with deep neural networks. *IEEE Transactions on Learning Technologies* 12, 1 (2018), 29–43.
- [40] Saori Suzuki, Shin'ichi Nakayama, and Hideo Joho. 2011. Formulating effective questions for community-based question answering. In Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval. 1261–1262.
- [41] Nam Khanh Tran and Claudia Niedereée. 2018. Multihop attention networks for question answer matching. In Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval. ACM, 325–334.
- [42] Kateryna Tymoshenko and Alessandro Moschitti. 2018. Shallow and deep syntactic/semantic structures for passage reranking in question-answering systems. ACM Transactions on Information Systems (TOIS) 37, 1 (2018), 8.
- [43] Gang Wang, Konark Gill, Manish Mohanlal, Haitao Zheng, and Ben Y. Zhao. 2013. Wisdom in the social crowd: An analysis of Quora. In Proceedings of the International Conference on World Wide Web. 1341–1352.

Large-Scale Question Tagging via Joint Question-Topic Embedding Learning

- [44] Hongbin Wang, Hongxu Hou, Jing Wu, Jinting Li, and Wenting Fan. 2017. Exploring different granularity in Mongolian-Chinese machine translation based on CNN. In *International Conference on Parallel and Distributed Computing, Applications and Technologies.* 112–116.
- [45] Shaowei Wang, D. Lo, B. Vasilescu, and A. Serebrenik. 2014. EnTagRec: An enhanced tag recommendation system for software information sites. Proceedings of the IEEE International Conference on Software Maintenance and Evolution, 291–300.
- [46] Shaowei Wang, David Lo, Bogdan Vasilescu, and Alexander Serebrenik. 2018. EnTagRec ++: An enhanced tag recommendation system for software information sites. *Empirical Software Engineering* 23, 2 (2018), 800–832.
- [47] Xin Xia, David Lo, Xinyu Wang, and Bo Zhou. 2013. Tag recommendation in software information sites. In Proceedings of the Working Conference on Mining Software Repositories (MSR'13). 287–296.
- [48] Fajie Yuan, Alexandros Karatzoglou, Ioannis Arapakis, Joemon M. Jose, and Xiangnan He. 2019. A simple convolutional generative network for next item recommendation. In *Proceedings of the 12th ACM International Conference on Web Search and Data Mining*. 582–590.
- [49] Jian Zhang, Hailong Sun, Yanfei Tian, and Xudong Liu. 2018. Poster: Semantically enhanced tag recommendation for software CQAs via deep learning. In 2018 IEEE/ACM 40th International Conference on Software Engineering: Companion (ICSE-Companion'18). 294–295.
- [50] Y. Zhang, M. J. Er, R. Zhao, and M. Pratama. 2016. Multiview convolutional neural networks for multidocument extractive summarization. *IEEE Transactions on Cybernetics* 47, 10 (2016), 3230–3242.
- [51] Z. Zhang, M. Lyons, M. Schuster, and S. Akamatsu. 1998. Comparison between geometry-based and gabor-waveletsbased facial expression recognition using multi-layer perceptron. In *Proceedings of the IEEE International Conference* on Automatic Face and Gesture Recognition, 1998. 454–459.
- [52] Wei Zhao, Ziyu Guan, Long Chen, Xiaofei He, Deng Cai, Beidou Wang, and Quan Wang. 2018. Weakly-supervised deep embedding for product review sentiment analysis. *IEEE Transactions on Knowledge and Data Engineering* 30, 1 (2018), 185–197.
- [53] Tom Chao Zhou, Michael R. Lyu, and Irwin King. 2012. A classification-based approach to question routing in community question answering. In Proceedings of the 21st International Conference on World Wide Web. 783–790.

Received June 2019; revised January 2020; accepted January 2020