# An Intelligent Question Answering Platform for Graduate Enrollment

Mengyuan Zhang, Yuting Wang, Jianxia Chen and Yu Cheng

School of Computer Science, Hubei University of Technology, Wuhan, China

## ABSTRACT

*To enhance the competitiveness of colleges and universities in the graduate enrollment and reduce the pressure on candidates for examination and consultation, it is necessary and practically significant to develop an intelligent Q&A platform, which can understand and analyze users' semantics and accurately return the information they need. However, there are problems such as the low volume and low quality of the corpus in the graduate enrollment, this paper develops a question answering platform based on a novel retrieval model including density-based logistic regression and the combination of convolutional neural networks and bi-directional long short-term memory. The experimental results show that the proposed model can effectively alleviate the problem of data sparseness and greatly improve the accuracy of the retrieval performance for the graduate enrollment.*

## KEYWORDS

*Question Answering System, Graduate Enrollment, Deep Learning, Sentence Semantic Similarity.*

## 1. INTRODUCTION

With the rapid development of NLP (natural language processing) technologies, the Q&A (question answering) system [1] is widely utilized in the real life. Q&A system is an human-machine dialogue service integrating knowledge base, information retrieval, machine learning, natural language understanding and other technologies [2]. It can effectively solve the problem of information overload and improve the efficiency of users' use of the system.

Recently, the number of graduate students across the country has been increasing greatly. To enhance the competitiveness of colleges and universities in the graduate enrollment and reduce the pressure on candidates for examination and consultation, it is necessary and practically significant to develop an intelligent Q&A platform for examination and enrollment in educational areas using NLP technologies, which can understand and analyze users' semantics and accurately return the information they need.

There are many traditional similarity algorithms such as SVM (support vector machine) [3], LR (logistic regression), KLR (kernel logistic regression), DT (decision tree), and NB (naive bayes) classification models. However, most of them are only suitable for specific types of data or output [4][5][6][7][8], and, have some problems such as ignoring the semantics of words or relying too much on semantic dictionaries.

With the in-depth study of deep learning, [9][10] utilized word embedding to construct word vectors to characterize the correlation between statements with vector similarity via some neural

network such as CNN [11], RNN(recurrent neural network) [12], LSTM(long short-term memory) [13], and other improved models to make much better performance in the field of the text similarity. However, these models often result in the slower response speed in the real Q&A system due to their complicated computation.

Inspired by the above approaches, this paper combines the advantages of both traditional method and DNN to propose a novel Intelligent Q&A platform for Examination and Enrollment in educational areas, IQ&AEE for short. In particular, IQ&AEE system proposes a crawler manager to crawl data in time to solve the problems of outdated information, develops a Q&A retrieval model based on DLR(density-based logistic regression) and CNN-BiLSTM (the combination of convolutional neural networks and bi-directional long short-term memory), and builds an intelligent Q&A robot, provides students with college introductions from all aspects through the web design of school information.

## 2. RELATED WORK

### 2.1. Introduction to the Question Answering System

Intelligent Q&A system refers to a system that accurately provides the knowledge required by users in the form of one question and one answer, and realizes interactive and personalized services for users.

The general architecture of the intelligent Q&A system can be divided into three parts: how to express natural language in the computer so that the computer understands its semantics, how to select the best answer, and how to map the answer to the natural language to express, namely question understanding, intelligent search and answer extraction.
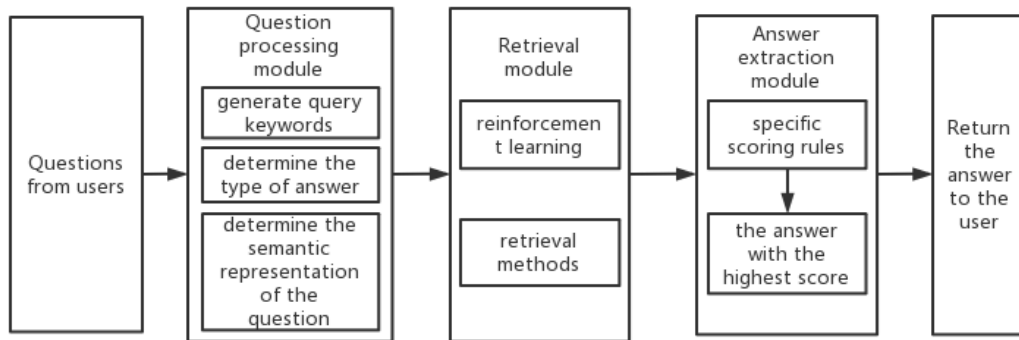


Figure 1. Q&A system flow chart

### 2.2. Sentence Semantic Matching Based on Deep Learning

Semantic matching based on deep learning is to model sentences directly, generate vectors of two sentence sequences, perform feature extraction and similarity calculation on the two sentences through the neural network mode, sort them according to the similarity score, select the most relevant pair return as the sequence pair with the highest semantic.

The framework focuses on how to efficiently use neural networks to take the semantic characteristics of sentences, such as using CNN with local perception mechanisms, or using LSTM with the memory to avoid long-term dependence problems of sentence sequences, or

using combined neural networks to learn local and contextual characteristics of sentences, construct feature vectors of sentences and learn the deep semantics of sentence sequences. Cosine similarity is calculated on the sentence characteristic vectors learned through various neural networks, and utilized to characterize the semantic matching degree of two sentence sequences. The framework is shown in Fig. 2.
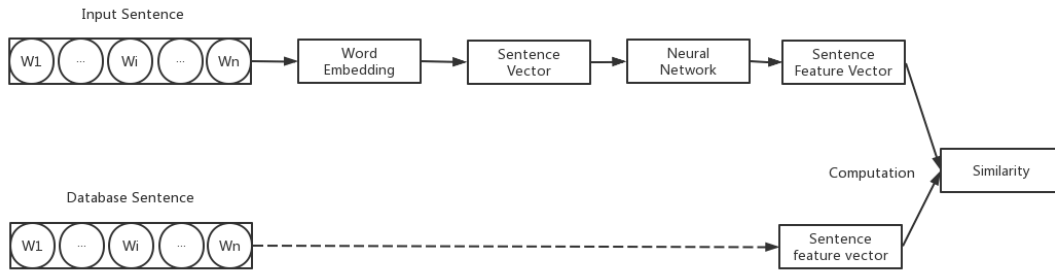


Figure 2.  Semantic matching framework based on deep learning

## 3. REQUIREMENT ANALYSIS

The proposed IQ&AEE system contains three type of users: visitors, individual users and system administrators. Visitors can ask questions and get answers instantly on the web page. Individual users can not only ask questions, but also leave messages for unanswered questions. System administrators can manage users (add or delete users), manage messages (add, delete or reply messages, update answered messages to the Q&A database), manage the Q&A database (add, edit or delete questions). The overview of the basic functional requirements of the IQ&AEE system is shown in Fig. 3 to Fig. 5.
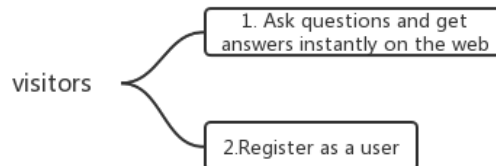


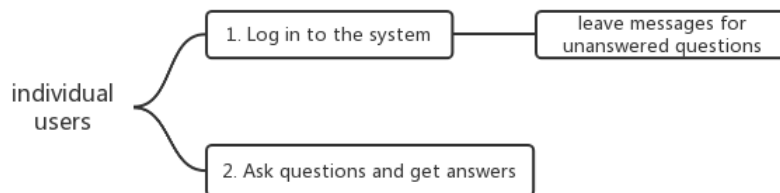Figure 3.  Functional requirements analysis diagram of visitors



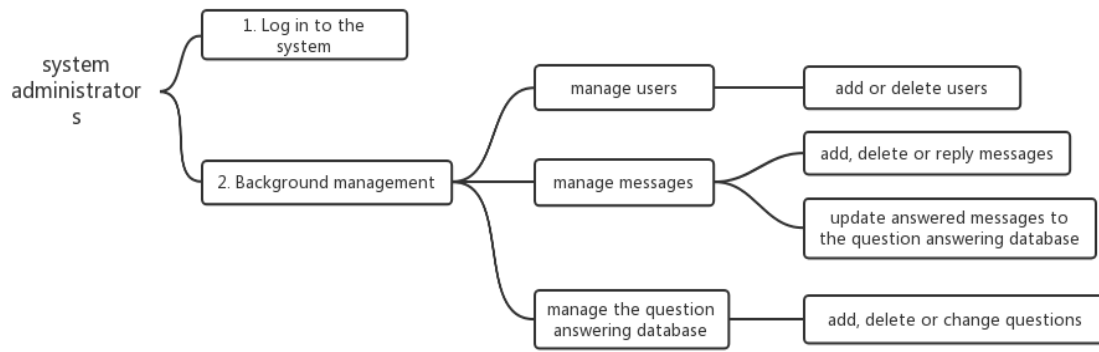Figure 4.  Functional requirements analysis diagram of individual users

Figure 5.  Functional requirements analysis diagram of system administrators

# 4. IQ&AEE DESIGN

## 4.1. IQ&AEE System Framework Design

According to the idea of software engineering MVC (model-view-controller), this paper divides the IQ&AEE system into three layers: application platform layer, functional module layer and database layer.

Among them, the application platform layer is the human-computer interaction interface layer, which embodies the operations provided by the system to the users, such as submitting questions, retrieving questions, obtaining answers and so on. The functional module layer implements crawler management, Q&A retrieval, knowledge base management and user management functions. The database layer is responsible for the storage of all the enrollment information of the school, and various types of databases. The framework design of the system is shown in Fig. 6.
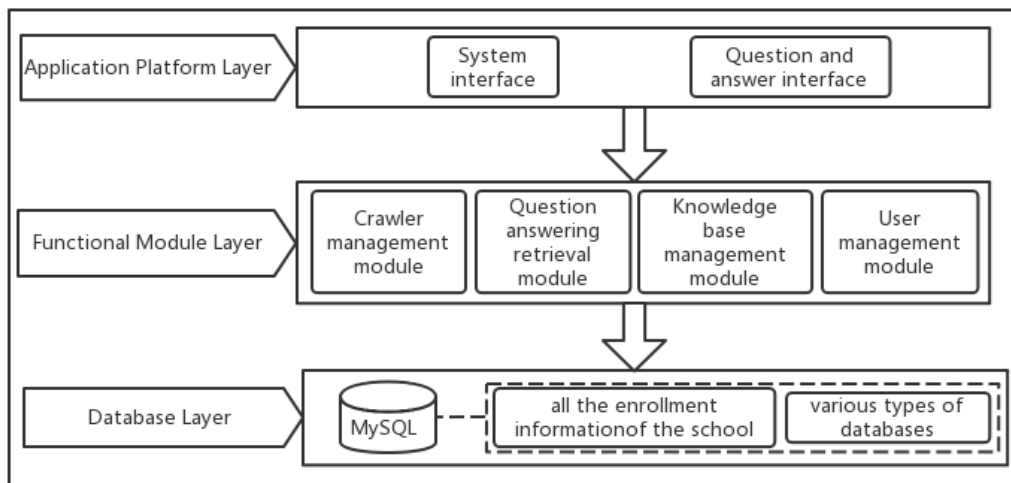


Figure 6.  IQ&AEE system framework

## 4.2. IQ&AEE System Functional Design

As shown in Fig. 7, the functional modules of IQ&AEE system are mainly divided into four main modules, namely crawler management module, Q&A retrieval module, knowledge base management mould and user management module. Each module contains several sub-modules.
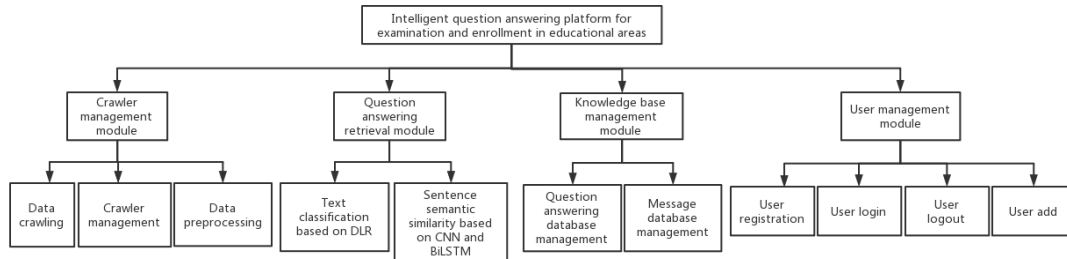


Figure 7.  IQ&AEE system function module diagram

## 4.3. Detailed Design of Each Module of the System

Crawler management module: In order to obtain rich research data, the crawler management module uses  the crawler technology to collect enrollment information data from authoritative websites of universities such as Yanzhao.com, which is utilized to build the system corpus.

Question answering retrieval module: As the critical module of the system, this module mainly provides the user with Q&A retrieval function based on frequently asked questions and knowledge, classifies the questions submitted by users, and matches the semantic similarity of the sentences in databases. The related core algorithms and experimental analysis are introduced in detail in Part 5.

Knowledge base management module: The Knowledge Base Management module provides the source of questions and answers in IQ&AEE system. Administrators can add, delete and modify questions and answers in the Q&A database, while the administrator is responsible for giving answers manually in the message database, and updating the questions and corresponding answers left by users to the Q&A database, convenient for other users to retrieve the answers to the corresponding questions.

User management module: The user management module can perform user registration, user login, user addition, and user logout operations. Administrators can reset the user's password, add users, or log off users in the background system.

## 5. RETRIEVAL MODULE OF IQ&AEE SYSTEM

As shown in Fig. 8, this paper develops a Q&A retrieval model based on two algorithms including DLR and CNN-BiLSTM.

After users enter questions, the Q&A retrieval model first uses the keyword table to classify the questions roughly, and divides them into 14 categories as a whole: thirteen school categories and one none category. Afterward, the DLR is utilized to determine the specific category to which the questions belong to school categories, and CNN-BiLSTM judges the semantic similarity of questions belong to unknown category and the none category.
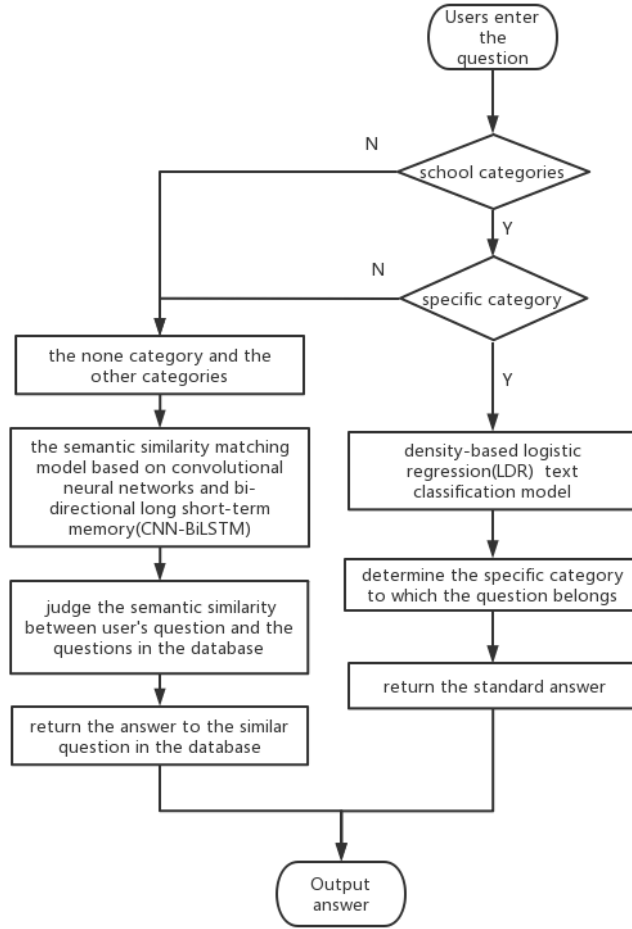
Figure 8.  IQ&AEE retrieval module flow chart

## 5.1. DLR Algorithm

DLR is a binary basic model based on LR, which is a novel type of nonlinear classifier, which is much more efficient than other nonlinear models and can naturally handle mixed data types. It also offers good interpretability and support for multiway classification [14].

The main idea of DLR is to map the training data to a specific feature space according to Nadarays-Watson density estimation algorithm, and then to build an optimization model to optimize feature weight and the width of the Nadarays-Watson density estimation algorithm. The DLR model first obtains the definition of the mapping function by calculating the probability as follows:

$$\varphi_d(x) = \ln \frac{p(y=1 \mid x^{(d)})}{p(y=0 \mid x^{(d)})} - \frac{D-1}{D} \ln \frac{p(y=1)}{p(y=0)} \tag{1}$$

Then, it is estimated by the Nadaraya-Watson estimator, and the result is obtained (2):

$$p(y = k \mid x^{(d)}) = \frac{\sum_{i \in D_k} K(\frac{x^{(d)} - x_i^{(d)}}{h_d})}{\sum_{i=1}^{N} K(\frac{x^{(d)} - x_i^{(d)}}{h_d})} \tag{2}$$

By this formula, the density of a given instance $x$ can be calculated. Then the final result can be get by substituting the standard LR. Self-adjustment of parameter is done by calculating the loss of the bias derivative.

The calculation process is shown in Eq.(3):

$$\frac{\partial E}{\partial h_d} = \frac{\partial E}{\partial r_d} \cdot \frac{\partial r_d}{\partial h_d} = \frac{1}{h_d^3} \sum_{i=1}^{N} (b_i - y_i)\omega_d \frac{\partial \varphi_d(x_i)}{\partial r_d} \tag{3}$$

Among them:

$$r_d = -\frac{1}{2h_d^2} \tag{4}$$

$$\frac{\partial \varphi_d(x)}{\partial r_d} = \frac{\partial g_1}{\partial r_d} - \frac{\partial g_0}{\partial r_d} \tag{5}$$

$$\frac{\partial g_i}{\partial r_d} = \frac{\sum_{i \in D_j} [(x^{(d)} - x_i^{(d)})^2 \cdot \exp(r_d(x^{(d)} - x_i^{(d)})^2)]}{\sum_{i \in D_j} \exp(r_d(x^{(d)} - x_i^{(d)})^2)} \tag{6}$$

Through this series of calculations, the partial derivative of the loss can be obtained, and then the loss can be minimized through the gradient descent and adjustment.

## 5.2. CNN-BiLSTM  Algorithm

As shown in Fig. 9, since CNN is capable of extracting local features and BiLSTM can extract global features (or context features) of a sentence, the paper combine them together to generate high-quality sentence representations for measuring sentence similarity [15][16][17][18].
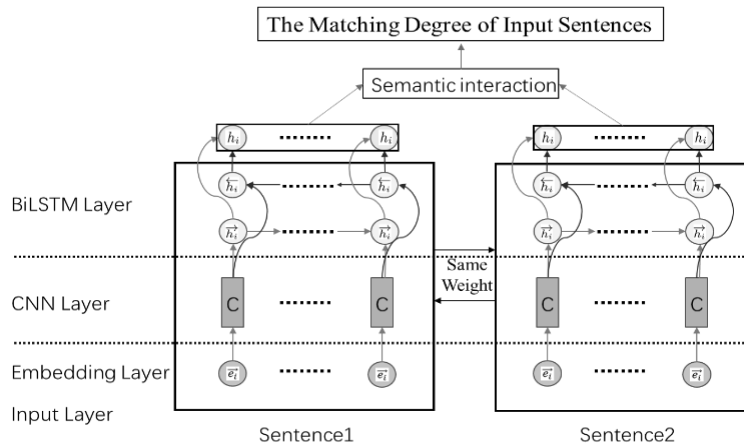


Figure 9.  CNN-BiLSTM architecture

### 5.2.1. CNN Layer

As shown in Fig. 9, CNN is utilized to extract local features of sentences and its structure. The characteristic $S_{ik}$ of the $i^{th}$ word in the sentence obtained through the convolution operation can be expressed in the Eq.(7):

$$S_{ik} = f(w_k \cdot x_{i:i+h-1} + b) \tag{7}$$

Where $w_k$ represents the weight matrix of the convolution filter $k$, and $b$ represents the deviation term. $f(.)$ is a nonlinear function, and Tanh is used in this paper.

After the convolution operation, the max-pooling method is adopted here to extract the important features of the sentence through a max pooling layer after the convolution layer, as following Eq.(8):

$$c = \max(c) \tag{8}$$

Where $c$ is the vector after the convolution operation.

### 5.2.2. Bi-LSTM Layer

In this paper, Bi-LSTM is utilized to extract the global features of the sentence and expressed in the Eq.(9-10):

$$\overrightarrow{h_i} = \overrightarrow{LSTM}(a_i, \overrightarrow{h_{i-1}}) \qquad i = 1, \cdots, N \tag{9}$$

$$\overleftarrow{h_i} = \overleftarrow{LSTM}(a_i, \overleftarrow{h_{i-1}}) \qquad i = N, \cdots, 1 \tag{10}$$

Where $\overrightarrow{h_i}$ represents the output of forwarding LSTM at $t^{th}$ time step, and $\overleftarrow{h_i}$ represents the output of backward LSTM at $t^{th}$ time step. The result of connecting $\overrightarrow{h_i}$ and $\overleftarrow{h_i}$ is the output result of Bi-LSTM at $t^{th}$ time step.

## 5.3. Experimental Results and Analysis

The system DLR text classification experiment and the CNN-BiLSTM sentence semantic similarity experiment both adopt the accuracy evaluation method. The accuracy calculation formula is as shown in Eq.(11):

$$Accuracy = \frac{n_{correct}}{n_{total}} \tag{11}$$

Among them, $n_{correct}$ represents the number of records with correct classification, and $n_{total}$ represents the number of all test data.

### 5.3.1.  Data Collection

The text classification corpus utilized in this project is shown in Table 1, which is divided into three columns: index, question and category. If the question contains "reexamination", it will be marked as "0", the "adjustment" will be marked as "1". The number of data is 540.

Table 1.  Text classification corpus.

| Index | Question | Category |
|---|---|---|
| 0 | 复试名单或者复试线多久可以看见？<br>(How long will the reexamination list or reexamination line be visible?) | 0 |
| 1 | 今年该校还会招收会计专硕（全日制）的调剂生嘛？<br>(Will the school enroll accounting adjustment students (full-time) this year?) | 1 |
| 2 | 调剂考生有可能在调剂系统打开前完成复试吗？<br>(Is it possible for the transfer candidate to complete the reexamination before the transfer system is turned on?) | 0 |
| ... | | |
| 539 | 贵校今年的复试时间在 5 月 20 日之前还是之后？<br>(Is your school's reexamination before or after May 20th this year?) | 0 |
| 540 | 今年复试是线上还是线下？<br>(Is this year's graduate reexamination online or offline?) | 0 |

The semantic similarity matching corpus utilized in this project is shown in Table 2, which is divided into three columns: question 1, question 2 and semantic match degree. The number of data is 1800.

Table 2.  The example of question pair corpus.

| Question1 | Question2 | Semantic match degree |
|---|---|---|
| 推免是否会影响到考研名额？<br>(Will the guaranteed acceptance affect the number of  the postgraduate entrance examination?) | 推免接受人数有多少？<br>(How many people can get guaranteed acceptance?) | 0 |
| 可以转专业吗？<br>(Can I switch to a major?) | 能否转专业？<br>(Can I switch to another major?) | 1 |

### 5.3.2.  The experimental results of the DLR

The number of the above text classification corpus is compared through the existing classification model, and the accuracy is utilized as the evaluation index for calculation. The results of each algorithm are shown in Table 3.

Table 3.  Text classification experiment on our corpus.

| Model | Accuracy |
|---|---|
| KNN | 75% |
| SVM | 74.2% |
| Nbayes | 75.7% |
| RanfomForest | 78.6% |
| LR | 92.13% |
| **DLR** | **93.98%** |

As can be seen from the table above, the results obtained by DLR text classification algorithm utilized in this module are significantly improved compared with the results of other traditional classification algorithms. Using DLR text classification algorithm can get better and more accurate experimental classification results.

### 5.3.3.  The experiment results of the CNN-BiLSTM

The paper selects 1500 statements in the above semantic similarity matching corpus as the training data set and 300 sentences as the test data set. The unsupervised methods and supervised methods are compared with CNN-BiLSTM, and accuracy is utilized as the evaluation index for calculation. The results of each algorithm are shown in Table 4.

Table 4.  Sentence semantic similarity matching experiment on our corpus.

| Model | Accuracy |
|---|---|
| One-hot+cos | 66.5% |
| Bert+cos | 67.5% |
| W2V+cos | 71.9% |
| W2V+tf-idf+cos | 71.2% |
| CNN | 67.3% |
| LSTM | 72.7% |
| BiLSTM | 70.6% |
| BiLSTM+CNN | 70.0% |
| **CNN+BiLSTM** | **80.7%** |

Through the above comparative experiments, it can be found that the effect of CNN-BiLSTM is better than other methods. Compared with a single CNN model and a single BiLSTM model, the fusion of these two models achieves better experimental results.

# 6. IQ&AEE SYSTEM IMPLEMENTATION

## 6.1. Experimental Environment

System hardware requirements: 64MB or more memory, 1024*768 resolution monitor, 24x speed CD-ROM or DVD-ROM drive, 50G or more available hard disk space, keyboard and mouse. The above is a test run server environment.

System web front-end development requirements: Window10, Eclipse javaee integrated development environment, apache-tomcat-9.0.16, Navicat Premium 12, MySQL 5.7, neo4j-community-3.3.5.

System background development requirements: python3.5, PyCharm2017 (includes library functions: Tensorflow 2.1.0, pymysql, sys, numpy, pandas, re, time).

## 6.2. Web Implementation

The Web is divided into a visitor part and a background administrator part. After entering the main web page, visitors can learn about the school through real-time information, popular professional introductions, topical news, popular activities and so on. As shown in Fig. 10. The real developed IQ&AEE system is shown on the website of : http://59.110.13.80:81/.

Visitors can enter the consultation page when they click on the page "school consultation" or "examination and enrollment consultation". On this web page, visitors can enter questions about the examination and enrollment, and can see the users' common questions, help, registration, login and other options. After testing, the Q&A function can be utilized normally. The consulting interface and the visitor inquiry question interface are shown in Fig. 11 to 12.

The administrator can manage the Q&A database through the knowledge base management module. On this web page, administrators can edit and delete questions in the Q&A database, or edit and delete answers in the Q&A database, and add keywords for the questions. At the same time, the administrator can edit, delete, add and do other operations to the message database and the user database through knowledge base management module and the user management module.

Figure 10.  Display of the main page of the visitors



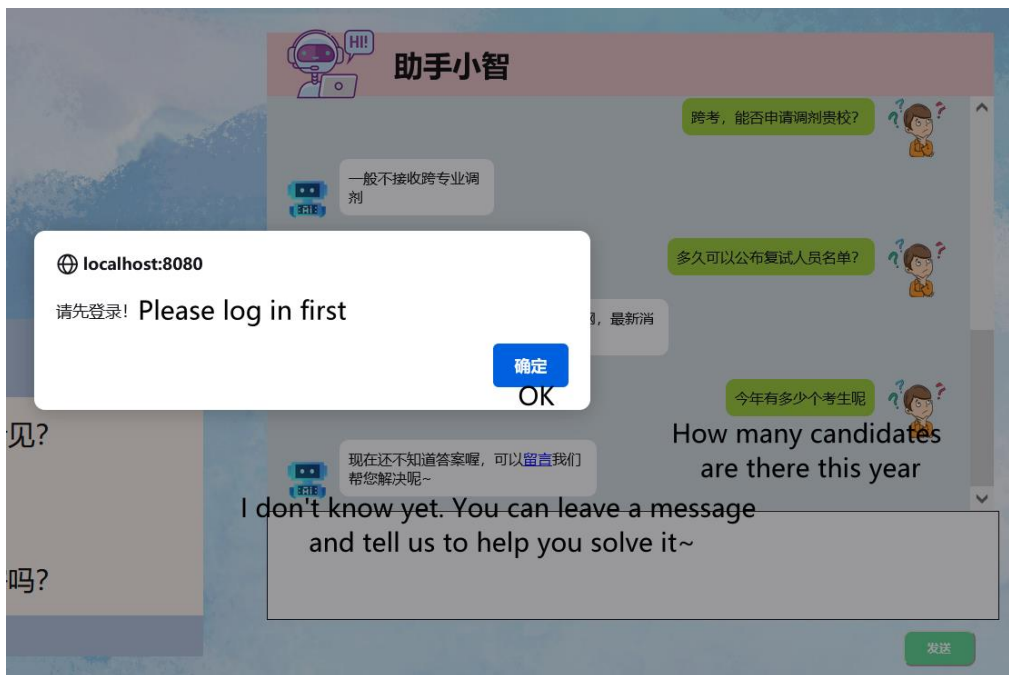Figure 11.  Display of the school consultation page

Figure 12.  Display of the school consultation message function

## 7. CONCLUSIONS

The Q&A system has experienced years of development and made some achievements. Compared with the Q&A system in the popular research domains, the Q&A system in the restricted domains such as the examination and enrollment, has the biggest problem, e.g. it is impossible to build a large corpus. Therefore, this paper proposes a novel Q&A retrieval model based on CNN-BiLSTM and DLR, which contains the ability of LR to deal with imbalances and the ability of NB to derive, solving the problem of the data sparseness, and opening up a novel way for the design of Q&A system in restricted domains.

This paper studies the technology of question classification, question component extraction and answer selection for question answering. However, due to the small data set for college examination and enrollment questions, the model CNN-BiLSTM is not adequately trained in question component extraction. On the basis of this paper, we can further study the addition of corpus automatic growth, corpus automatic error correction design, high-efficient and accurate short text similarity matching algorithm and the development of depth reasoning in Q&A systems, so that to make the Q&A system for the specific area of examination and enrollment with higher wisdom. Although there have been some achievements in the construction of the Q&A system, with booming development of big data and artificial intelligence, there are still more construction details for the Q&A system in the restricted domains, which urgently needs to be considered and explored by the developers.

**REFERENCES**

[1]    Peixuan Li, Lu Zhu and Dosheng Wu. Overview of Q&A system[J]. Digital Technology and Application, 2015(04)：69+71.(in Chinese)[李沛晏,朱露,吴多胜.问答系统综述[J].数字技术与应用,2015(04):69+71.]

[2]    Hailiang Wang, Zhuoheng Li, Xuming Lin, et al. Intelligent Q&A and deep learning[M]. Beijing: Publishing House of Electronics Industry, 2019.(in Chinese)[王海良,李卓恒,林旭鸣等著.智能问答与深度学习[M].北京:电子工业出版社,2019.]

[3]    Yao-Nan WANG and Xiao-Fang YUAN. SVM Approximate-based Internal Model Control Strategy[J]. Acta Automatica Sinica, 2008, 34(2)：172-179.

[4]    John Maindonald. Pattern Recognition and Machine Learning[J]. Journal of Statistical Software, 2007, 17(1)：1-3.

[5]    Zhi-Gang Zhao, Hui-Xian Lv, Yu-Jing Li, Jing Li. A Multi Classification SVM Based on Clustering Idea[J]. Journal of Qing dao Technological University, 2011, 32(01), 73−76.(in Chinese)[赵志刚,吕慧显,李玉景,李京.一种基于聚类思想的 SVM 多类分类方法[J].青岛理工大学学报,2011,32(01):73-76.]

[6]    Green P J, Yandell B S. Semi-parametric generalized linear models. In: Proceedings on the 2nd International GLIM Conference. New York: Springer-Verlag, 1985. 44−55.

[7]    Ji Zhu and Trevor Hastie. Kernel Logistic Regression and the Import Vector Machine[J]. Journal of Computational and Graphical Statistics, 2005, 14(1)：185-205.

[8]    Tom M. Mitchell and Jaime G. Carbonell and Ryszard S. Michalski. Machine Learning[M]. Springer, Boston, MA, 1986.

[9]    Tom Kenter and Maarten de Rijke. Short Text Similarity with Word Embeddings[C]. , 2015.

[10]   Xingjie Feng, Le Zhang, Yunze Zeng. Problem similarity calculation model based on multi-attention CNN[J]. Computer Engineering, 2019, 45(09)：284-290.(in Chinese)[冯兴杰,张乐,曾云泽.基于多注意力 CNN 的问题相似度计算模型[J].计算机工程,2019,45(09):284-290.]

[11]   Rumelhart D E . Learning Representations by Back-Propagating Errors[J]. Nature, 1986, 23.

[12]   V. Pham, T. Bluche, C. Kermorvant and J. Louradour, "Dropout Improves Recurrent Neural Networks for Handwriting Recognition," 2014 14th International Conference on Frontiers in Handwriting Recognition, 2014, pp. 285-290, doi: 10.1109/ICFHR.2014.55.

[13]   Mueller J, Thyagarajan A. Siamese recurrent architectures for learning sentence similarity[C]//Thirtieth AAAI Conference on Artificial Intelligence. 2016.

[14]   Wenlin Chen, Yixin Chen, Yi Mao, and Baolong Guo. 2013. Density-based logistic regression. In Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '13). Association for Computing Machinery, New York, NY, USA, 140–148. DOI:https://doi.org/10.1145/2487575.2487583

[15]   Y. Li, D. Zhou and W. Zhao, "Combining Local and Global Features Into a Siamese Network for Sentence Similarity," in IEEE Access, vol. 8, pp. 75437-75447, 2020, doi: 10.1109/ACCESS.2020.2988918.

[16]   Q. Chen Q. Hu J. X. Huang and L. He "CA-RNN: Using context-aligned recurrent neural networks for modeling sentence similarity" Proc. 32nd AAAI Conf. Artif. Intell. pp. 265-273 2018.

[17]   Shengxian Wan, Yanyan Lan, Jun Xu, Jiafeng Guo, Liang Pang, Xueqi Cheng. "Match-SRNN: Modeling the recursive matching structure with spatial RNN", arXiv:1604.04378. [Online]. Available: http://arxiv.org/abs/1604.04378.

[18]   Hua He, Kevin Gimpel and Jimmy Lin. "Multi-Perspective Sentence Similarity Modeling with Convolutional Neural Networks" in Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, Lisbon, Portugal, 1576–1586.2015.
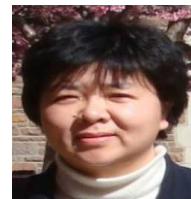
**AUTHORS**

**Mengyuan Zhang** She is currently working toward the B.S. degree in Data Science and Big Data Technology from Hubei University of Technology, Wuhan, China. She research interests include Artificial Intelligence and Natural Language Processing.

**Yuting Wang** She is currently working toward the B.S. degree in Data Science and Big Data Technology from Hubei University of Technology, Wuhan, China.

**Jianxia Chen** She is an associate professor and master tutor of the School of Computer Science, Hubei University of Technology. Research direction: cloud computing and big data, knowledge graph and natural language processing, intelligent planning research.

**Yu Cheng** She is an associate professor of the School of Computer Science, Hubei University of Technology. Research direction: database application and artificial intelligence.