

Which is the Effective Way for Gaokao: Information Retrieval or Neural Network?

By Shangmin Guo, Xiangrong, Shizhu He, Kang Liu and Jun Zhao

Reference: [EACL 2017 Paper](#)

Abstract

National Higher Education Entrance Examination, which is commonly known as Gaokao, is designed to be difficult enough to distinguish excellent high school students in China. The paper described Gaokao History Multiple Choice Questions(GKHC) and proposed three approaches to address them. Approaches include **Information Retrieval(IR)**, **Neural Network(NN)** and the **combination of IR and NN**. Results show that IR performs better at Entity Questions(EQs), while NN performs better at Sentence Questions(SQs). The combination method achieves state-of-art performance, showing the necessity to apply hybrid method when encountering real-world scenarios.

Motivation: Difficult History Multiple Choice Questions

Answering real world questions in various subjects it increasingly getting attentions. The [Project Halo](#) was proposed to create a “digital” Aristotle which has most of the worlds’ scientific knowledge as well as solve hard problems. In terms of the history questions, there are some NLP attempts for yes-no questions: determining the correctness of the original position ([Kanayama et al., 2012](#)) and recognizing textual entailment between a description in Wikipedia and each options([Miyao et al., 2012](#)). Nevertheless, none of these approaches can solve difficult history multiple choice questions as shown in Figure 1, which require a huge amount of background knowledge.

After the World War II, U.S. and Soviet Union are fighting against each other in politics, economics and military. To promote the development of economics in Socialist Countries, Soviet Union establish The Council for Mutual Economic Assistance. This is against
A. Truman Doctrine
B. Marshall Plan
C. NATO
D. Federal Republic of Germany

Entity Question

From Qin and Han Dynasties to Ming Dynasty, businessmen are always at the bottom of hierarchy. One reason for this is that the ruling class thought the businessmen
A. are not engaged in production
B. do not respect Confucianism
C. do not respect the clan
D. do not pay tax

Sentence Question

Figure 1: Examples of history questions. EQ means all of the candidates are entities, which SQ means candidates are parts of the sentence.

This paper is not the first to solve Gaokao questions, but the former approaches based on information retrieval did not fit well and suffer from limited knowledge resources in their systems. Therefore, works introduced in this paper are mainly focus on solving difficult GKHC and proposing new approaches to improve accuracy.

Datasets: Questions and Resources

All of the questions are from Gaokao all over the country in 2011 - 2015. Questions with graphs are filtered out since solving them requires techniques beyond NLP. The remaining questions are manually tagged as EQs or SQs. Numbers of

different kinds of collected multiple-choice questions are 160 for EQs, 584 for SQs (744 in total).

Wide diversity of resources including Baidu Encyclopedia, textbooks and over 50,000 practice questions are also collected.

Approaches and Results: IR, NN and Combination

Information Retrieval (IR)

Since GKHMC questions require finding the most relevant candidate to the question stem from 4 choices, IR approach is applicable by following the pipeline below:

1. Use **Naive Bayes classifier** to classify questions.

- Features include length, entity number and verb number of candidates.
- Do 10-folder cross validation on question dataset.

2. Calculate **relevance scores** for each candidate and combine them with **specific weights** (3 different method with 7 score functions on different indices are provided for the calculation).

- **Lexical Matching Score** : for each candidate K, calculated this score by summing up score of the top i-th returned documents calculated by [Lucene's TFIDFSimilarity function](#).
- **Entity Co-Occurrence Score**: calculated by [normalized google distance](#), assuming that two entities appear at the same time are relevant.
- **Page Link Score**: inspired by [PageRank algorithm](#), calculated by finding the maximum number of links between question stem entity and candidate answer entity.
- **Training weights and loss function** For each candidate, the score can be calculated as below, which will finally be normalized.

$$score_{candidate_k} = \sum_{i=1}^7 w_i * f_i(candidate_k)$$

The loss function is:

$$loss_{questions} = -\log(1 - score_n)$$

As all operations are derivable, gradient descent algorithm can be used to train weights.

3. Candidate with highest score will be chosen as right answer.

Result: Accuracy of EQs and SQs with corresponding best weights are 49.38% and 28.60%. Obviously, IR works better over EQs than SQs.

Neural Network (NN)

Permanent-Provisional Memory Network(PPMN) is introduced in this paper as NN approach, which is designed to tackle with the joint inference between background knowledge and question stems in GKHMC. The diagram of PPMN is shown in Figure 2.

The diagram illustrates the architecture of the proposed framework, which is designed to select answer candidates from a set of material sentences based on their cultural achievements.

Input Module: This module takes material sentences as input. For example, the sentence "公元97年，东汉的班超曾派员出使欧洲强国“大秦”。东汉和“大秦”都创造了辉煌的文化。" (In 97 A.D., Ban Chao of the Eastern Han Dynasty had sent an envoy to European powerful country "DaQin". Both of the Eastern Han and "DaQin" created splendid culture.) is processed into a lead-in representation.

Encoder: The input module's output is fed into an encoder, which generates a lead-in representation. This representation is then compared with the output of the Provisional Memory Module's encoder.

Provisional Memory Module: This module contains a set of memory units (m1, m2, ..., mn) that store information from previous iterations. The output of the Provisional Memory Module's encoder is fed into the Similarity Judger.

Similarity Judger: This module compares the lead-in representation from the input module with the output of the Provisional Memory Module's encoder to determine the similarity between them.

Encoder (Right): This module takes answer candidates as input and generates a representation. The output of this encoder is fed into the Similarity Judger.

Answer Candidates: The final output of the framework is the selected answer candidates, which are the material sentences that have the highest similarity to the lead-in representation.

Example: The example material sentence is "公元97年，东汉的班超曾派员出使欧洲强国“大秦”。东汉和“大秦”都创造了辉煌的文化。" (In 97 A.D., Ban Chao of the Eastern Han Dynasty had sent an envoy to European powerful country "DaQin". Both of the Eastern Han and "DaQin" created splendid culture.). The lead-in representation is "属于他们的文化成就分别是" (The cultural achievements of them are). The answer candidates are "A. 活字印刷术、万有引力定律" (movable type printing, the law of gravity), "B. 《九章算术》、罗马法" (Nine Chapters on the Mathematical Art, Roman law), "C. 蔡伦改进的造纸术、日心说" (papermaking technology, heliocentric theory), and "D. 《春秋繁露》、《理想国》" (Chun Qiu Fan Lu, Utopia).

Combination of IR and NN Approach

It is obvious that IR and NN approaches are complementary to some extent, which is intuitively as well. In EQs, information given by question stems is usually the description of the key entity, which is the reason why correct answer has the highest relevance score. It is more straightforward to using IR to solve EQs. However, in SQs, the key entity does not appear in any candidate, which means inference is needed. Therefore, though IR works well in EQs, it is not sufficient to find the correct choice in SQs, while NN works better in SQs.

Considering that (1) some of EQs may be more suitable to be handled as SQs and (2) both character and word embedding are more sufficient to cover the lexical meaning, a hybrid approach is proposed. IR and NN approaches can be simply combined via a weights matrix.

The performance of combined model and its comparison to IR and NN approaches are illustrated in Figure 3.

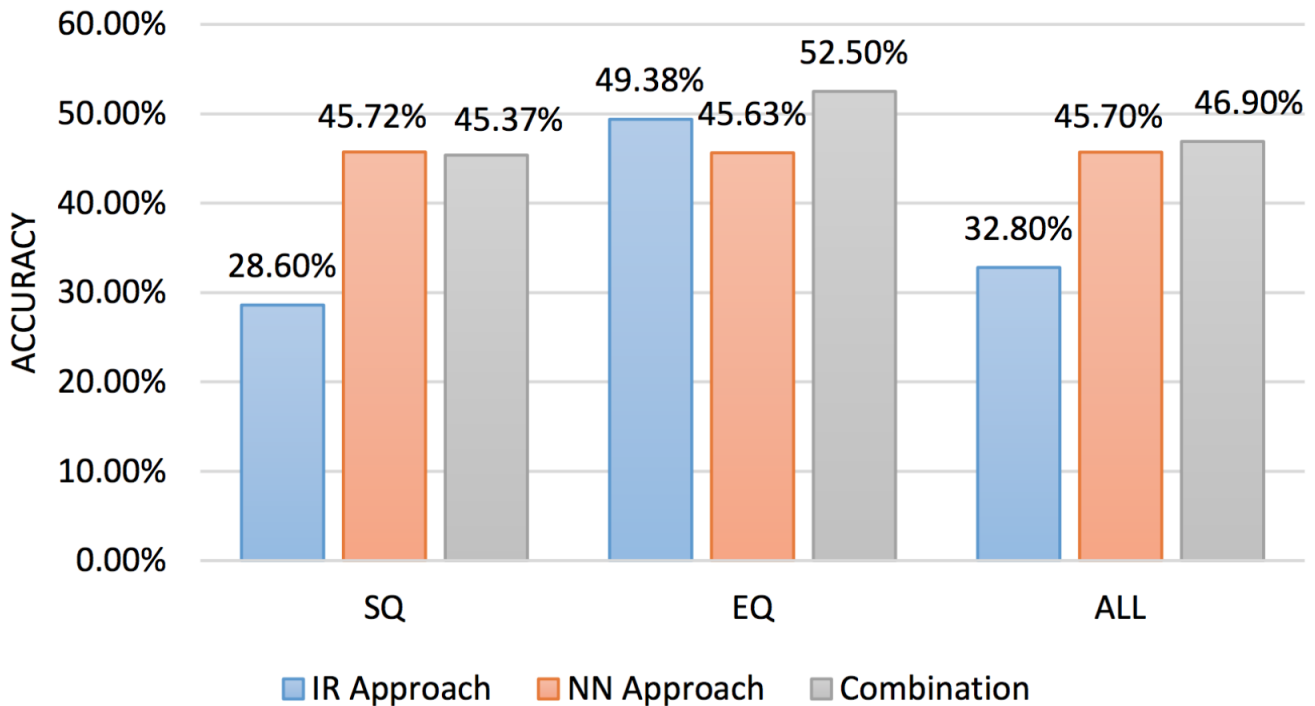


Figure 3: Result of different approaches: IR, NN and Combination

Conclusion

The paper details the GKHMC, presents different approaches to address them and compares their performances. According to the results, IR approach is more suitable for EQs while NN approach is more suitable for SQs. The combination of IR and NN has a state-of-the-art performance on GKHMC, pointing out that hybrid methods may be a better choice in real world scenarios.