

Capturing Semantic Relevance with DSSM

TAL Future Camp Project Write-up

Qing Lyu
Tsinghua University
lvq15@mails.tsinghua.edu.cn

Abstract

This is my project write-up in TAL Future Camp 2018¹, an AI boot camp organized by a Education and Technology Enterprise in China, Tomorrow Advancing Life Education Group (NYSE: TAL)². The camp provided an intense 2-week project-based training for aspiring AI + Education practitioners, and the division I joined was NLP. Specifically, this project focuses on a semantic relatedness capturing task under an educational setting. Given a pair of elementary school maths exercise problems, it aims at building a system capable of automatically determining if they are conceptually related. As a baseline, I re-implemented an adapted version of the Deep Structured Semantic Models (DSSM) proposed by Microsoft Research in 2013. On top of that, I experimented with several extensions (from the perspective of pre-processing, regularization, architecture design, and parameter tuning) in order to refine the baseline model. In addition to quantitative evaluation, a text retrieval task is carried out to demonstrate the model performance qualitatively. A section of further thoughts is presented in the end, showing current limitations and future possibilities of improvement. All data and codes can be found here³.

1 Introduction

1.1 Background

The technique of semantic relatedness computation has multiple applications, such as search engines (given a query, retrieving relevant Web documents), question and answering systems (selecting pertinent answers for a certain question) and so on. In order to judge semantic relevance, a prerequisite is to learn good semantic representations for sentences. Distributional semantics provides traditional approaches such as Bag of Words, Word2Vec, and Doc2Vec. In addition,

¹<https://future.ailab.100tal.com/camp/main>

²<http://en.100tal.com>

³Project repository: <https://github.com/veronica320/Capturing-Semantic-Relevance-with-DSSM>

deep learning can be another way to learn sentence representations and perform semantic relatedness calculation. In the educational context, we will specifically focus on the task of determining the relevance of a pair of exercise problems, which have potential downstream applications such as exercise recommendation in adaptive learning. We will try to adapt a classical system, Deep Structured Semantic Models (DSSM) (Huang et. al., 2013[1]), and its variation, Convolutional Deep Structured Semantic Models (CDSSM) (Shen et. al., 2014a[2], Shen et. al., 2014b[3]), proposed by Microsoft Research, and apply them under our task settings.

1.2 Task Description

Given a corpus of elementary school maths problems (provided by TAL), our goal is to determine if there exists relevance between any two problems. Instead of other subjects, only maths is chosen, because it is believed to be more complex in text form (with Chinese characters, English and Greek letters, L^AT_EX notations, equations, and so on), and thus more challenging as well as interesting.

In the dataset, each problem is paired with multiple "related" problems (annotated as "1") and the same number of "unrelated" problems (annotated as "0") at the same time. In detail, our tasks are as follows:

- Build a deep semantic model with the training data, and learn vector representations for the problem texts.
- Use the learned model to predict the relatedness of each pair in the test set.
- Based on the learned semantic representations, draw a few sample problems, and retrieve the "10 most related problems" with each of them according to vector distance.

2 Methodology

2.1 Dataset

The dataset is a Chinese corpus of elementary school maths exercise problems, provided by TAL. As mentioned above, each problem is paired with an equal number of conceptually related problems and unrelated problems. There are 20286 pairs in total, 88.82% for training, and 11.18% for testing. The problems are all written in uncleaned texts, i.e., they contain Chinese characters, Arabic numerals, English and Greek letters, L^AT_EX notations, equations, and so on. A sample text is:

两箱 水果 共重 \$ \$ 92 \$ \$ 千克 , 第一 箱比 第二 箱重 \$ \$ 6 \$ \$ 千克 , 两箱 水果 各重
 多少 千克 ? \$ \$ left (92 + 6 right) div 2 = 49 \$ \$ (千克) \$ \$ cdotsdots
 \$ \$ 第一 箱 , \$ \$ 92 - 49 = 43 \$ \$ (千克) \$ \$ cdotsdots \$ \$ 第二 箱 .

Figure 1: A sample problem text from the dataset

The text is already tokenized. Conspicuously, it needs further pre-processing before being fed to models. The pre-processing strategies will be introduced in Section 2.1 and Section 3.2.1.

2.2 Pre-processing

Dimensionality thresholding. In the corpus, since texts contain different number of words, there will be a problem of unequal vector lengths if directly fed into neural networks. In the papers on DSSM and CDSSM by Microsoft Research, the authors address the problem by either word hashing, or bag of words plus max pooling. In Chinese, however, there is no “curse of dimensionality”, and thus the motivation for word hashing wouldn’t be as intuitive as in English. Thus, I choose to tackle the problem of unequal lengths in another way: by manually setting a maximum length parameter (default=20), longer texts would be cut short, and shorter texts would be padded with a sequence of null characters.

Word vector filtering. I use pre-trained Chinese word vectors⁴ in the model. To accelerate the loading process during training, the word vectors are filtered beforehand, i.e. only words that occur in the training or test sets remain.

2.3 Deep Structured Semantic Models (DSSM)

As mentioned above, the architecture that I will adapt, DSSM, is first published by Microsoft Research in 2013. It is a deep semantic model, initially designed for information retrieval in search engines, i.e. given a query, the model retrieves the most relevant documents on the Web. The core concept of DSSM is that the network projects a query and candidate documents into a common low-dimensional vector space, where the semantic relatedness is computed as vector similarity scores. The model later evolves in other semantic mapping tasks, including graphic matching, captioning, Sent2Vec, etc. The last layer of the original network was to compute the cosine distance between two vectors, because it was originally for similarity computation. But here we have a binary classification task, so the objective function and the loss function should be adjusted.

⁴<https://github.com/Embedding/Chinese-Word-Vectors>

2.4 Convolutional Deep Structured Semantic Models (CDSSM)

Also brought forth by Microsoft Research, CDSSM has the following motivation: in the vanilla DSSM, texts are represented as bag of words, so the positional information within sentences are discarded. In CDSSM, the authors intend to capture the relation between words in a sentence by introducing convolution. The initial layers are the same as DSSM, but later a convolutional layer is followed by a max pooling layer, in order to discover salient contextual features. At last, a similar non-linear transformation is applied. It is shown that CDSSM significantly outperforms DSSM in information retrieval. Thus, I will explore both structures in the following section.

3 Experimental Results

3.1 Baseline

3.1.1 Model Architecture

Based on Huang (2013)[1], I re-implemented a slightly different DSSM as baseline, whose architecture diagram is as follows:

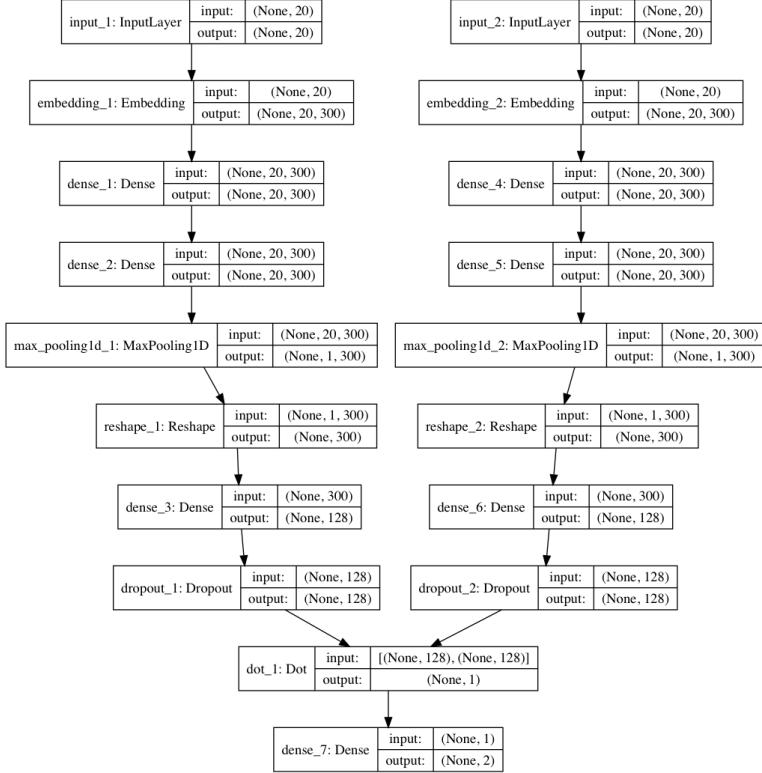


Figure 2: Baseline Model Architecture (based on DSSM)

The two texts in each pair are fed to the network as input_1 and input_2 respectively. Because the default maxlen parameter is set to 20, the input length here is limited to 20. The Input Layer is followed by the Embedding layer, where self-learned word embeddings instead of pre-trained word embeddings are used for the time being. Next are two fully-connected layers, similar to the Microsoft paper. A slight difference is that because there is no word hashing before, the dimension reduction is done by a pooling layer. Finally, the high-dimensional vectors are projected into an 128-dimensional space to produce the semantic representations for sentences. The cosine similarity between them is computed, and then we use Softmax to transform the similarity score to a binary classification probability, which is evaluated against the true label with the cross-entropy loss.

3.1.2 Several clarifications

Special character filtering. For the non-Chinese characters (e.g. mathematical symbols) in the corpus, the baseline model does not filter them out. The motivation is that if pre-trained word vectors are used, filtering is likely to

be better since the pre-trained word vectors almost do not include any mathematical symbols. But since the baseline model learns the word representations itself, the embeddings of such mathematical symbols could technically also be learned. So, I start with no filtering here.

Dropout layer. In the model structure diagram (Figure 2), there is a Dropout Layer, which does not exist in the Microsoft paper. However, for the baseline model, I temporarily set the dropout ratio parameter to 0, so the overall architecture is still mostly equivalent to the original paper.

Relevant parameters. maxlen=20, l2 regularization ratio=0.01, learning rate=0.001, embedding dim=300, train val split = 0.1, batch size=200, epoch=30, patience for early stopping=10.

3.1.3 Model performance

During the training process, 10% of the training data is used for validation. The loss and the accuracy in each epoch on the training and validation sets are plotted:

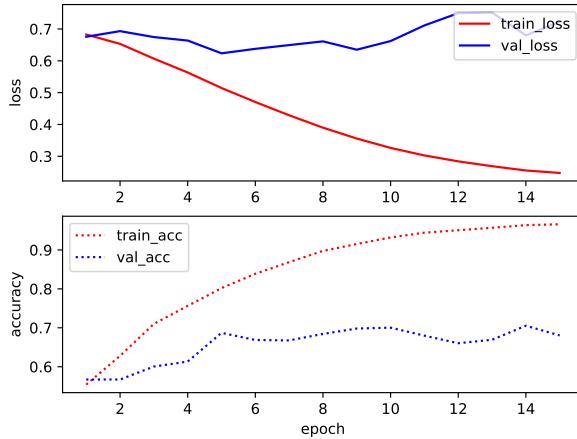


Figure 3: The performance of the baseline model during the training process

Obviously, overfitting is serious. Although the validation accuracy is slightly higher in the end than in the beginning, the overall trend fluctuates severely, and the validation loss has a tendency to rebound in the last few epochs.

On the test set, the performance of the final trained model is also evaluated (the final model is selected as the one with the lowest loss across all epochs):

Model	Accuracy	Precision	Recall	F1 score
Baseline	0.654	0.629	0.75	0.684

Table 1: The performance of the baseline model on the test set

Good news is that at least its performance is conspicuously above chance. Also, recall is much higher than precision, suggesting a tendency of over-predicting positive cases.

3.2 Extensions

On the baseline model, I experimented with several extensions from the following aspects:

3.2.1 Pre-processing and word embedding

In the baseline model, the special characters (e.g. mathematical symbols) in the texts are not filtered out. Here I experimented with two other filtering schemes, respectively named as “Coarse Filtering” and “Fine Filtering” :

“Coarse Filtering”. All non-Chinese characters, including punctuation marks, are filtered out. It is expected that this scheme will be useful when using pre-trained word vectors, since they contain only Chinese characters with few exceptions. Another possibility is to use English word vectors at the same time, but the two semantic spaces then need to be projected into one shared bilingual space beforehand.

“Fine Filtering”. Observe the sample problem text again:

```
两箱 水果 共重 $ $ 92 $ $ 千克 , 第一 箱比 第二 箱重 $ $ 6 $ $ 千克 , 两箱 水果 各重
多少 千克 ? $ $ left ( 92 + 6 right ) div 2 = 49 $ $ ( 千克 ) $ $ cdots
$ $ 第一 箱 , $ $ 92 - 49 = 43 $ $ ( 千克 ) $ $ cdots $ $ 第二 箱 .
```

Figure 4: A sample problem text from the dataset

In the Fine-Filtering scheme, aside from Chinese characters, I will try to retain special characters that contribute to semantic meaning, and it is expected will work better with self-learned word vectors for the same reason as above. Detailed filtering rules include the following:

- For L^AT_EX notations, only meaningful ones such as ‘frac’, ‘sum’, and ‘pi’ remain. On the contrary, notations such as ‘end’, ‘begin’, and ‘equation’ are mainly for formatting purposes, and thus have little semantic value.
- For maths operators, some (such as *, ^, and +) are retained while others (such as -, /) are deleted. The motivation is that the first type has unambiguous semantic meaning, while the second, except as maths

operators, also appear as formatting notations (e.g. "/sum", "二面角 B-AC-D").

- Numerical values are filtered out. For example in "x + y = 5", the number "5" is deleted. This is because numerical values are believed to have little influence on the conceptual meaning of problems. If two problem texts only differ in numerical values, they should still have high semantic similarity.
- For punctuation marks, only those that have clear semantic meaning (such as ? which indicates clear questioning) are kept. Other (such as commas and periods which are only for pause) are deleted.

With the above rules, the final cleaned problem text is:

两箱 水果 共重 千克 第一 箱比 第二 箱重 千克 两箱 水果 各重 多少
千克 ? + div = 千克 cdot cdot 第一 箱 - = 千克 cdot cdot 第二 箱

Figure 5: The same sample problem text after Fine Filtering

It can be seen that the final text is free of semantically irrelevant information, and the meaning of the problem is almost completely preserved.

Using the same parameters as the baseline, the performance of the two models with new filtering schemes is as follows (Coarse Filtering is combined with pre-trained word vectors, Fine Filtering with self-learned vectors):

Model	Accuracy	Precision	Recall	F1 score
Baseline	0.654	0.629	0.75	0.684
Baseline+Coarse Filter	0.675	0.686	0.643	0.664
Baseline+Fine Filter	0.717	0.707	0.714	0.710

Table 2: The test performance of the new filtering schemes compared to baseline

Compared to the baseline model, the new filtering schemes both successfully improve test accuracy and precision. But the recall and F1 score for the Coarse Filtering scheme are lower than the baseline, which demonstrates that the overall performance of this scheme is not quite stable. It also indicates that the more characters are filtered out, the more conservative the model is when predicting positive samples. In comparison, despite a decrease in recall, the Fine Filtering scheme significantly improves F1 score along with accuracy at the same time, resulting in a well-balanced performance (all four metrics are around 0.70%). It may be safe to say that the Baseline+Fine Filtering model is the current optimal one, again suggesting the effectiveness of the carefully-designed filtering rules.

3.2.2 Dropout Layers

Based on the current optimal model in Section 3.2.1, another extension I try is to add a Dropout Layer (as shown in Figure 2, dropout ratio=0.1) before

calculating the cosine similarity, in order to alleviate overfitting. The results are presented in the table below:

Model	Accuracy	Precision	Recall	F1 score
Baseline+Fine Filter (3.2.1)	0.717	0.707	0.714	0.710
Baseline+Fine Filter+Dropout	0.719	0.698	0.740	0.718

Table 3: The test performance of the current optimal model (3.2.1) plus dropout

It is seen that after adding the Dropout Layer, precision falls by 0.09% yet recall jumps by 2.6%, which jointly improves accuracy and F1 score slightly. Generally speaking, the Baseline+Fine Filtering+Dropout model could be regarded as the current optimal.

3.2.3 Parameter tuning

Observing the training process of the above models, the problem of overfitting is still salient: in most cases, there is a bottleneck for the validation loss at around the 10th epoch, accompanied by slight fluctuations afterwards. But the training loss still drops steadily, often with early stopping. Therefore, based on the current optimal model in 3.2.2, here I try to tune certain parameters that may affect overfitting in this range:

- l2:[0.01, 0.05]
- drop out ratio:[0.0, 0.1, 0.2]
- batch size:[100, 200, 400]
- maxlen:[10, 20, 30]

Experiments show that the performance of the optimal model in the above range is as follows:

Model	Accuracy	Precision	Recall	F1 score
Baseline+Fine Filter+Dropout	0.719	0.698	0.740	0.718
Baseline+Fine Filter+Dropout+Param Tune	0.723	0.704	0.740	0.721

Table 4: The test performance of the current optimal model (3.2.2) plus parameter tuning

All metrics except recall are improved on the basis of the optimal model in 3.2.2. The optimal parameters are found to be: l2=0.05, drop out ratio=0.1, batch size=400, maxlen=10. It is seen that increasing l2, dropout ratio and batch size does help alleviate overfitting. Something worth noticing is that for the maxlen parameter (which stands for the maximum length of input texts), the baseline value is set to 20 purely for the sake of training efficiency. It is

previously thought that 20 is already a quite low value, since the average length of training texts is much longer. But surprisingly, experiments show that 10 is even better, which seems to indicate that the semantic similarity of two problem texts could be more accurately determined by only taking into consideration the first 10 characters.

3.2.4 Architecture Variation

Besides DSSM, here I also experiment with two other model architectures.

CDSSM. Based on Shen (2014a)[2] and Shen (2014b)[3], I built the following CDSSM structure:

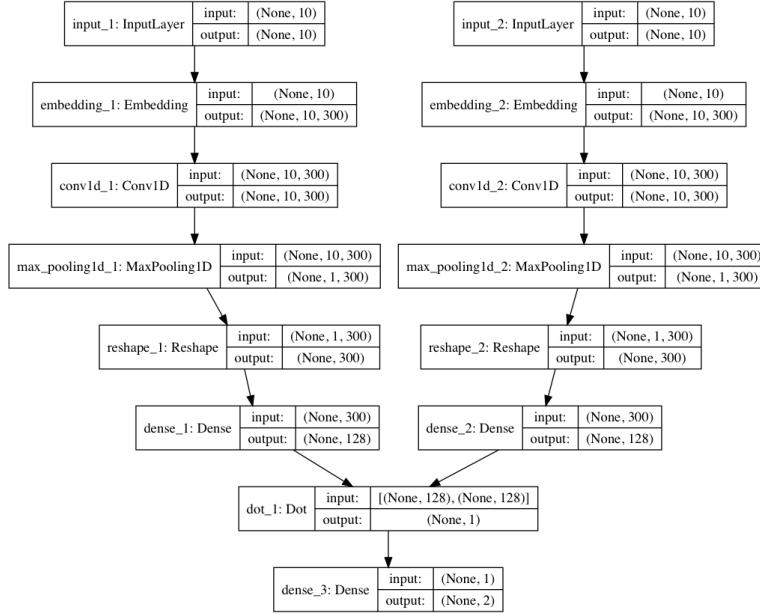


Figure 6: CDSSM Model Architecture

Similarly to the baseline model, the word hashing technique adopted by the original paper is not involved for processing Chinese texts. Additionally, since the current task is binary classification, the original parameter γ in the Softmax layer is no longer needed here.

The performance of the model on the test set is shown below:

Model	Accuracy	Precision	Recall	F1 score
CDSSM	0.685	0.681	0.694	0.687

Table 5: The performance of CDSSM on the test set

All parameters are the same as the optimal settings in Section 3.2.3, but the overall performance in terms of accuracy and F1 score is merely comparable to the baseline (3.1.3). It might be attributed to the fact that the optimal parameters are not transplantable, and another round of parameter tuning is needed. Because of the limit in time and resources, here I will skip that.

CGRU. One theoretical inadequacy of CDSSM is that the contextual information it captured is limited by the kernel size of convolution, and long-range dependencies are hardly preserved. Thus, here I adapt the previous CDSSM into a recurrent structure by adding an additional Convolutional Layer plus a Bidirectional GRU Layer, named as CGRU as a whole. The structure is shown below:

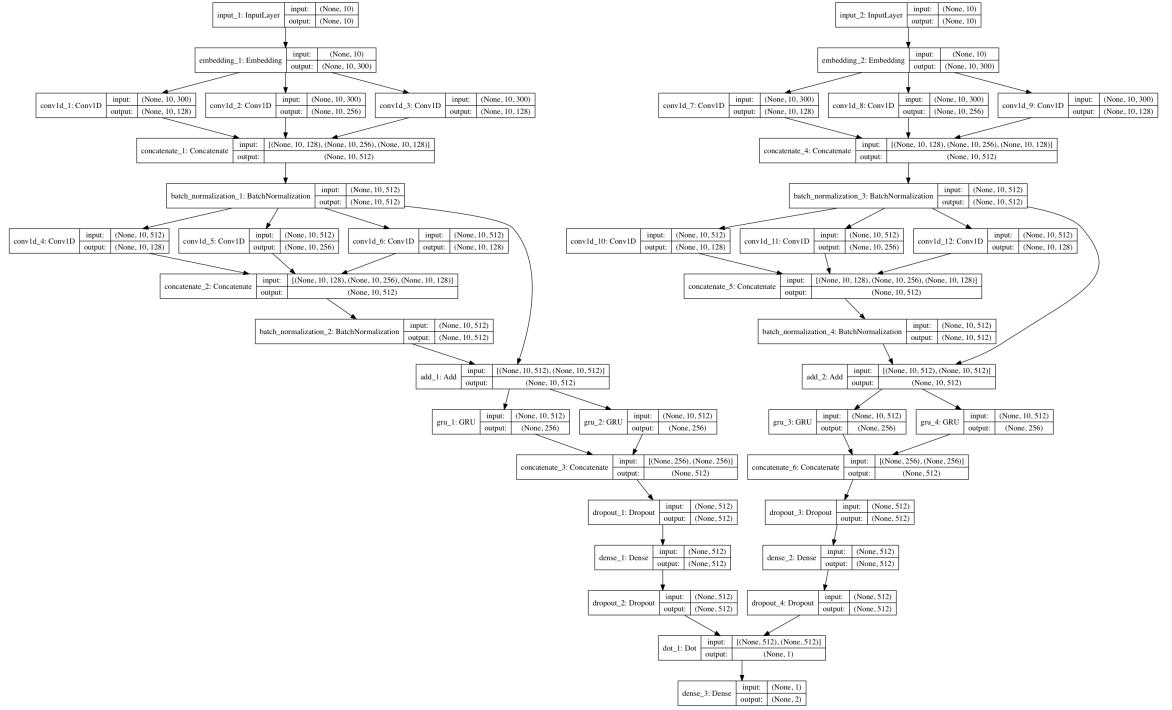


Figure 7: CGRU Model Architecture

On the test set, the performance of CGRU is:

Model	Accuracy	Precision	Recall	F1 score
CGRU	0.674	0.669	0.686	0.678

Table 6: The performance of CGRU on the test set

Unfortunately, CGRU even performs worse than simple CDSSM. Certainly, the similar lack of parameter tuning could be one possible explanation.

3.3 Text Retrieval: "10 most relevant problems"

Using the current optimal model (obtained in Section 3.2.3), here I will try to qualitatively demonstrate how the model captures semantic relatedness by drawing a few sample problems and retrieving the 10 most relevant problems from the dataset.

First, the problems in the dataset are deduplicated, resulting in a total of 19760 samples, each tagged with an ID. The three sample "queries" I select are the ones with ID 50, 100 and 150. Each of them (copied to make a list of length 19760) is fed into the model as the left input, and the rest "docs" (complied into a list) are fed as the right input. Running the prediction process yields 19760 float numbers, each standing for the relevance between the "query" and one "doc". The "10 most relevant problems" for the three sample problems are shown below:

- ID=50, problem text:

```
|小明放学回家做完作业后想出门和小朋友一起玩，妈妈允许他最多出去半个小时，否则就要挨罚，他
出门时是 5 时 14 分，那么小明最晚几时几分回到家才不会被妈妈罚？5 时 14 分+30 分=5 时 44 分，
小明最晚 5 时 44 分回到家才不会被妈妈罚。
```

10 most relevant problems:

ID	Similarity	Problem text
5809	1	小明用一个新本子写大字，第一天写了2页，第二天写了3页，第三天应从第2页开始写。小明用一个新本子写大字，第一天写了2页，第二天写了3页，第三天应从第6页开始写。
2017	1	小明家买了一箱牛奶，共15瓶，小明上个星期每个喝一瓶，还剩多少瓶？15-7=8（瓶）。
3526	1	18个同学排成一队做体操，从左边数小南是第6个。小南的右边有几个同学？18-6-1=11（个）。
7575	0.99999988	阿凡提在集市上花了600元买了一头小毛驴，转手以640元卖给了别人，随后又以650元买回了这头驴，第二天，阿凡提又以680元把驴卖了，那么阿凡提一共赚了元钱。第一天赚了40元，第二天赚了30元，一共赚了40+30=70元。
8299	0.99999973	张老师要带10个同学过河。已经有6个同学过了河，没有过的还有几人？10-6=4（人）。
1497	0.99999952	一年级有120位同学去参观动物园，每人至少带一种饮料，其中带可乐的有78人，带矿泉水的有64人，既带可乐又带矿泉水的有几人？78+64-120=22（人）。
10234	0.99999944	小亮买了一本故事书，一共20页。他第一天看了几页，第二天又看了几页后，数了数还剩6页看完，请问小亮这两天一共看了多少页？20-6=14（页）。
2012	0.99999833	一些小朋友排好队来报数，正着数小明报九，一共有多少个小朋友？7+9-1=15（个）。
8298	0.99999774	三个兄弟比个子，老大比老二高5厘米，老三比老二高2厘米，那么老大和老三相差几厘米？老大和老三相差5-2=3
4078	0.99999738	有四个人在晚上准备通过一座摇摇欲坠的小桥。此桥每次只能让2个人同时通过，否则桥会倒塌。过桥的人必须用手电筒，不然会一脚踏空，只有一个手电筒。4人的行走速度不同：小强用1分钟就可以过桥，中强要2分钟，大强要5分钟，最慢的太强需要10分钟。17分钟后桥就要倒了。请问：4个人要用分钟才能全部安全过桥？小强和中强先过桥，用2分钟，再由小强把手电筒送过去，用1分钟，现在由大强跟太强一起过桥，用10分钟，过去以后叫中强把手电筒送给小强用2分钟，最后小强与中强一起过河再用2分钟，他们一起用时间：2+1+10+2+2=17(分钟)，正好在桥倒塌的时候全部过河。（时间最短过河的原则是：时间长的一起过，时间短的来回过，这样保证总的时间是最短的）。

- ID=100, problem text:

如图所示，将自然数有规律地填入方格表中，请问：null

10 most relevant problems:

- ID=150, problem text

甲、乙、丙、丁四人拿出同样多的钱，一起订购同样规格的若干件新年礼物，礼物买来后，甲、乙、丙分别比丁多拿了3, 7, 14件礼物，最后结算时，乙付给了丁14元钱，并且乙没有付给甲钱。那么丙应该再付给丁元钱。设丁拿了 a 件礼物，则四人花同样的钱，每人可以拿到 $a+(3+7+14)/4=a+6$ 件礼物，实际情况：丁少拿了6件，乙多拿了1件，给丁14元，则货物单价14元，丙多拿了 $14-6=8$ 件，3件给甲，5件给丁， $5 \times 14 = 70$ 元。

10 most relevant problems:

Observing the results, we have the following findings:

13

- **Problem 50.** It is a simple numerical addition and subtraction problem. Except for Problem 1497 (set intersection) and Problem 4078 (optimization), other problems are quite relevant, involving only simple addition and subtraction. The model performs fairly on this problem.
- **Problem 100.** It is a natural number inference problem. 7 of 10 problems retrieved are considered relevant, while the rest 3 (Problem 18605, Problem 1108, and Problem 11345) are less pertinent. The length of the original problem text is relatively short, which might be a reason why the retrieval performance is sub-optimal.
- **Problem 150.** It is a real-life numerical problem involving addition, subtraction, multiplication and division, and the concept of taking the average. Since keywords like "相等/倍数/一样" are salient, almost all retrieved problems are similar, except only Problem 7341, a truth-condition validation problem. But regardless of the problem type, it still involves the four mathematical operations. So, in terms of concept, it may be still seen as relevant. In all, the performance of the model is also satisfactory on this problem.

4 Further Thoughts

- Sadly, the two complex architectures that I have experimented with, CDSSM and CGRU, have not achieved expected performance. One possible reason, of course, may be the lack of parameter tuning. If time permits, it would be ideal to conduct a performance analysis on individual data points, so as to see if there are any pattern in the errors. Also, it would be possible to verify if the theoretical effect of the new architectures (e.g. CGRU capturing long-range dependencies) are actually realized.
- By reexamining the corpus, it is found that the tokenization, which is completed by the data provider, is not quite accurate. A considerable proportion of "nonsense" words like "图中缺"“可求”“开平”in the tokenized texts might have led to a significant number of Out-of-Vocabulary (OOV) words in training. It might be sensible to use character-level embeddings instead of, or combined with word-level embeddings in this case.
- Again, in terms of the difference between English and Chinese models, the original paper implements word hashing in order to reduce the dimensionality of one-hot vectors for English words, but Chinese does not have such a problem, and therefore word hashing is unnecessary. As a side effect, we have to find another way to solve the problem of unequal input lengths in the model. The method I have used, setting a maximum length plus padding, might not be careful enough. Probably it will be better to add a Pooling Layer as in CDSSM.

References

- [1] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. Learning deep structured semantic models for web search using click-through data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pages 2333–2338. ACM, 2013.
- [2] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. Learning semantic representations using convolutional neural networks for web search. In *Proceedings of the 23rd International Conference on World Wide Web*, pages 373–374. ACM, 2014.
- [3] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. A latent semantic model with convolutional-pooling structure for information retrieval. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 101–110. ACM, 2014.