

DeepRank: A New Deep Architecture for Relevance Ranking in Information Retrieval

Liang Pang^{†*}, Yanyan Lan^{†*}, Jiafeng Guo^{†*}, Jun Xu^{†*}, Jingfang Xu[‡], Xueqi Cheng^{†*}

pl8787@gmail.com, {lanyanyan, guojiafeng, junxu, cxq}@ict.ac.cn, xujingfang@sogou-inc.com

[†]CAS Key Lab of Network Data Science and Technology, Institute of Computing Technology,

Chinese Academy of Sciences, Beijing, China

*University of Chinese Academy of Sciences, Beijing, China

[‡]Sogou Inc, Beijing, China

ABSTRACT

This paper concerns a deep learning approach to relevance ranking in information retrieval (IR). Existing deep IR models such as DSSM and CDSSM directly apply neural networks to generate ranking scores, without explicit understandings of the relevance. According to the human judgement process, a relevance label is generated by the following three steps: 1) relevant locations are detected; 2) local relevances are determined; 3) local relevances are aggregated to output the relevance label. In this paper we propose a new deep learning architecture, namely DeepRank, to simulate the above human judgment process. Firstly, a detection strategy is designed to extract the relevant contexts. Then, a measure network is applied to determine the local relevances by utilizing a convolutional neural network (CNN) or two-dimensional gated recurrent units (2D-GRU). Finally, an aggregation network with sequential integration and term gating mechanism is used to produce a global relevance score. DeepRank well captures important IR characteristics, including exact/semantic matching signals, proximity heuristics, query term importance, and diverse relevance requirement. Experiments on both benchmark LETOR dataset and a large scale clickthrough data show that DeepRank can significantly outperform learning to ranking methods, and existing deep learning methods.

CCS CONCEPTS

•Information systems → Retrieval models and ranking;

KEYWORDS

Deep Learning; Ranking; Text Matching; Information Retrieval

1 INTRODUCTION

Relevance ranking is a core problem of information retrieval. Given a query and a set of candidate documents, a scoring function is usually utilized to determine the relevance degree of a document with respect to the query. Then a ranking list is produced by sorting in descending order of the relevance score. Modern learning to

rank approach applies machine learning techniques to the ranking function, which combines different kinds of human knowledge (i.e. relevance features such as BM25 [27] and PageRank [22]) and therefore has achieved great improvements on the ranking performances [18]. However, a successful learning to rank algorithm usually relies on effective handcrafted features for the learning process. The feature engineering work is usually time-consuming, incomplete and over-specified, which largely hinder the further development of this approach [11].

Recently, deep learning approach [17] has shown great success in many machine learning applications such as speech recognition, computer vision, and natural language processing (NLP), owing to their ability of automatically learning the effective data representations (features). Therefore, a new direction of Neural IR is proposed to resort to deep learning for tackling the feature engineering problem of learning to rank, by directly using only automatically learned features from raw text of query and document. There have been some pioneer work, including DSSM [13], CDSSM [29], and DRMM [11]. Both DSSM and CDSSM directly apply deep neural networks to obtain the semantic representations of query and document, and the ranking score is produced by computing their cosine similarity. Guo et al. [11] argued that DSSM and CDSSM only consider the semantic matching between query and document, but ignore the more important relevance matching characteristics, such as exact matching signals, query term importance, and diverse matching requirement [27]. Thus they proposed another deep architecture, i.e. DRMM, to solve this problem. However, DRMM does not explicitly model the relevance generation process, and fails to capture important IR characteristics such as passage retrieval intrinsics [19] and proximity heuristics [31].

Inspired by the human judgement process, we propose a new deep learning architecture, namely DeepRank, to better capture the relevance intrinsics. According to the illustration in [33], the human judgement process can be divided into three steps. Human annotators first scan the whole document to detect the relevant locations. Then the local relevance for each detected location is decided. Finally, those local relevances are combined to form the global relevance of the entire document. Consequently, DeepRank contains three parts to simulate the human judgement process, by tackling the following three problems:

Where does the relevance occur? According to the query-centric assumption proposed in [33], the relevant information for a query

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.

CIKM'17, November 6–10, 2017, Singapore.

© 2017 ACM. ISBN 978-1-4503-4918-5/17/11...\$15.00

DOI: <http://dx.doi.org/10.1145/3132847.3132914>

only locates in the contexts around query terms. Therefore, the context with a query term at the center position, namely query-centric context, is recognized as the relevant location in the detection step.

How to measure the local relevance? After the detection step, a measure network is utilized to determine the local relevance between query and each query-centric context. Firstly, a tensor is constructed to incorporate both the word representations of query/query-centric context, and the interactions between them. Then a CNN or 2D-GRU is applied on the tensor to output the representation of the local relevance. In this way, important IR characteristics such as exact/semantic matching signals, passage retrieval intrinsics, and proximity heuristics can be well captured.

How to aggregate such local relevances to determine the global relevance score? As shown by [33], two factors are important for user's complex principles of aggregating local relevances, i.e. query term importance [6] and diverse relevance requirement [27]. Therefore we propose to first aggregate local relevances at query term level, and then make the combination by considering weights of different terms, via a term gating network. To obtain the term level relevance, we first group the query-centric contexts with the same central word together. Then a recurrent neural network (RNN) such as GRU [4] or LSTM [10] is utilized to aggregate such local relevances sequentially, by further considering the position information of these query-centric contexts in the whole document.

Above all, DeepRank is a new architecture composed of three components, i.e. a detection strategy, a measure network with CNN/2D-GRU, and an aggregation network with term gating and RNN inside. Therefore, DeepRank can be trained end-to-end with the pairwise ranking loss via stochastic gradient descent. We conduct experiments on both benchmark LETOR4.0 data and a large scale clickthrough data collected from a commercial search engine. The experimental results show that: 1) Existing deep IR methods such as DSSM, CDSSM, and DRMM perform much worse, if using only automatically learned features, than the pairwise and listwise learning to rank methods. 2) DeepRank significantly outperforms not only all the existing deep IR models but also all the pairwise and listwise learning to rank baseline methods. 3) If we incorporate handcrafted features into the model, as did in SQA [28], DeepRank will be further improved, and the performance is better than SQA. We also conduct a detailed experimental analysis on DeepRank to investigate the influences of different settings.

To the best of our knowledge, DeepRank is the first deep IR model to outperform existing learning to rank models.

2 RELATED WORK

We first review related work on relevance ranking for IR, including learning to rank methods and deep learning methods.

2.1 Learning to Rank Methods

In the past few decades, machine learning techniques have been applied to IR, and gained great improvements to this area. This direction is called learning to rank. Major learning to rank methods can be grouped into three categories: pointwise, pairwise and listwise approach. Different approaches define different input and output spaces, use different hypotheses, and employ different loss

functions. Pointwise approach, such as logistic regression [8], inputs a feature vector of each single document and outputs the relevance degree of each single document. Pairwise approach, such as RankSVM [14] and RankBoost [7], inputs pairs of documents, both represented by feature vectors and outputs the pairwise preference between each pair of documents. Listwise approach, such as ListNet [2], AdaRank [34] and LambdaMart [1], inputs a set of document features associated with query and outputs the ranked list. All these approaches focus on learning the optimal way of combining features through discriminative training. However, a successful learning to rank algorithm relies on effective handcrafted features for the learning process. The feature engineering work is usually time-consuming, incomplete and over-specified, which largely hinder the further development of this direction [11].

2.2 Deep Learning Methods

Recently, deep learning techniques have been applied to IR, for automatically learning effective ranking features. Examples include DSSM [13], CDSSM [29], and DRMM [11]. DSSM uses a deep neural network (DNN) to map both query and document to a common semantic space. Then the relevance score is calculated as the cosine similarity between these two vectors. Rather than using DNN, CDSSM is proposed to use CNN to better preserve the local word order information when capturing contextual information of query/document. Then max-pooling strategies are adopted to filter the salient semantic concepts to form a sentence level representation. However, DSSM and CDSSM view IR as a semantic matching problem, and focus on generating a good sentence level representations for query and document. They ignore the important intrinsics of relevance ranking in IR. Guo et al. [11] first point out the differences between semantic matching and relevance matching. They propose a new deep learning architecture DRMM to model IR's own characteristics, including exact matching signals, query terms importance, and diverse matching requirement. Specifically, DRMM first builds local interactions between each pair of words from query and document based on word embeddings, and then maps the local interactions to a matching histogram for each query term. Then DRMM employs DNN to learn hierarchical matching patterns. Finally, the relevance score is generated by aggregating the term level scores via a term gating network. Though DRMM has made the first step to design deep learning model specially for IR, it does not explicitly model the relevance generation process of human. It also fails to model the important IR intrinsics such as passage retrieval strategies and proximity heuristics.

Another related sort of deep models, flourishing in NLP, provide a new way of thinking if we treat IR task as a general text matching task, i.e. query matches document. These work can be mainly categorized as representation focused models and interaction focused models. The representation focused models try to build a good representation for each single text with a neural network, and then conduct matching between the two vectors. The DSSM and CDSSM mentioned above belong to this category. There are also some other ones such as ARC-I [12] model, which builds on word embeddings and makes use of convolutional layers and pooling layers to extract compositional text representation. The interaction focused models first build the local interactions between two texts, and

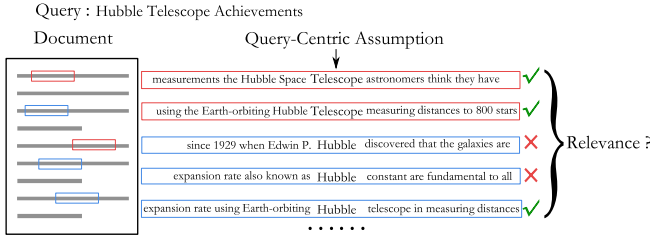


Figure 1: An example of human judgement process.

then use neural networks to learn the more complicated interaction patterns for matching. Typical examples include ARC-II [12] and MatchPyramid [23, 24], and Match-SRNN [32]. These models have been shown effective in text matching tasks such as paraphrase identification and question answering. DRMM can also be viewed as an interaction focused model.

3 MOTIVATION

The motivation of our deep architecture comes from the process of human relevance judgement. As illustrated by [33], a human annotator will first examine the whole document for local relevance information. Scanning involves iterating through every document location, and deciding for each whether local relevance information is found. As suggested by [33], if one can find relevant information in a document, the relevant information must locate around the query terms inside the document, namely query-centric assumption. Finally, the local relevance for each query-centric context is combined to form the relevance of the entire document. Figure 1 gives an example to describe the process of human relevance judgement, where user's information need is represented by the query 'Hubble Telescope Achievements'. For a given document, annotators will directly extract several contexts containing the keyword 'Telescope' and 'Hubble', and determine whether they are relevant to the query. Then these local relevance information will be considered together to determine the global relevance label of the document with respect to the query.

Therefore, the relevance label is generated following three steps in the human judgement process: 1) detection of relevant locations; 2) measurement of local relevance; 3) aggregation of local relevances. Consequently, three problems are going to be tackled if we want to well capture the relevance: (1) Where does the relevance occur? (2) How to measure the local relevance? (3) How to aggregate such local relevance to determine the final relevance label? Accordingly, we design our deep learning architecture, namely DeepRank, to model the above relevance process.

4 DEEPRANK

In this section, we demonstrate a new deep learning architecture for IR, namely DeepRank. DeepRank includes three parts to tackle the above three problems: a *Detection Strategy*, a *Measure Network*, and an *Aggregation Network*. In the detection step, the query-centric contexts are extracted to represent where the relevance occur. In the measurement step, CNN or 2D-GRU is adopted to measure the local relevance between query and each query-centric context. Finally in the aggregation step, RNN and a term gating network

are utilized to aggregate those local relevances to a global one for ranking. Figure 2 gives an illustration of DeepRank.

We will first give some mathematical notations. Each query and document are represented as a sequence of word $\mathbf{q} = (w_1, \dots, w_M)$ and $\mathbf{d} = (v_1, \dots, v_N)$, where w_i denotes the i -th word in the query and v_j denotes the j -th word in the document, respectively. Consequently, we can use $\mathbf{d}_{(k)}[p] = (v_{p-k}, \dots, v_p, \dots, v_{p+k})$ to denote a piece of continuous text within a document, centered on the p -th word with the sequence length $2k + 1$. If we use word2vec technique [20] to represent each word to a vector, each word sequence can be represented as a matrix. Taking query $\mathbf{q} = (w_1, \dots, w_M)$ for example, if the word embedding of w_i is \mathbf{x}_i , the query will be represented as $\mathbf{Q} = [\mathbf{x}_1, \dots, \mathbf{x}_M]$, where each column stands for a word vector.

4.1 Detection Strategy

According to the query-centric assumption [33], the relevance usually occurs at the locations where the query terms appear in the documents. Similar observations have also been obtained by the eye-tracking studies [5], showing that annotators focus more on the location of query terms when they are scanning the whole document for relevance judgment. Therefore, we define the query-centric context as the relevant location, which is a context with a query term at the center position. Mathematically, given a query \mathbf{q} and document \mathbf{d} , if query term w_u appears at the p -th position in the document, i.e. $w_u = v_p$, the query-centric context centered on this query term is represented as $\mathbf{s}_u^p(k) = \mathbf{d}_k[p]$. After the detection step, we can obtain a set of word sequences $\{\mathbf{s}_u^p(k)\}$, with each one represents a local relevant location.

4.2 Measure Network

The goal of the measurement step is to determine the local relevance, i.e. relevance between query and each query-centric context. As reviewed in Section 2, previous deep learning models for text matching can be divided into representation focused methods and interaction focused methods. The representation focused methods focus on extracting high level semantic representations of texts, starting from basic word representations. While the interaction focused methods turn to directly model the interactions between the two texts, starting from word-level interaction signals. In order to combine the best of both approaches, we propose a new measure network, as shown in Figure 3. Firstly a tensor is constructed as the input. Then CNN or 2D-GRU is applied on the tensor to output a vector, which stands for the representation of local relevance. In this way, the important IR characteristics such as exact/semantic matching signals, passage retrieval intrinsics, and proximity heuristics can be well captured in the measurement step.

4.2.1 Input Tensor. The key idea of this layer is to feed both the word representations of query/query-centric context and the interactions between them into the input of the measure network. Specifically for a given query \mathbf{q} and query-centric context $\mathbf{s}_u^p(k)$ with $w_u = v_p$, we denote the word-level interaction matrix used in MatchPyramid [24] and Match-SRNN [32] as \mathbf{S} , where each element S_{ij} is defined as the similarity of corresponding words w_i and v_j . For example, indicator function or cosine similarity can be used to capture the word-level exact or semantic matching signals,

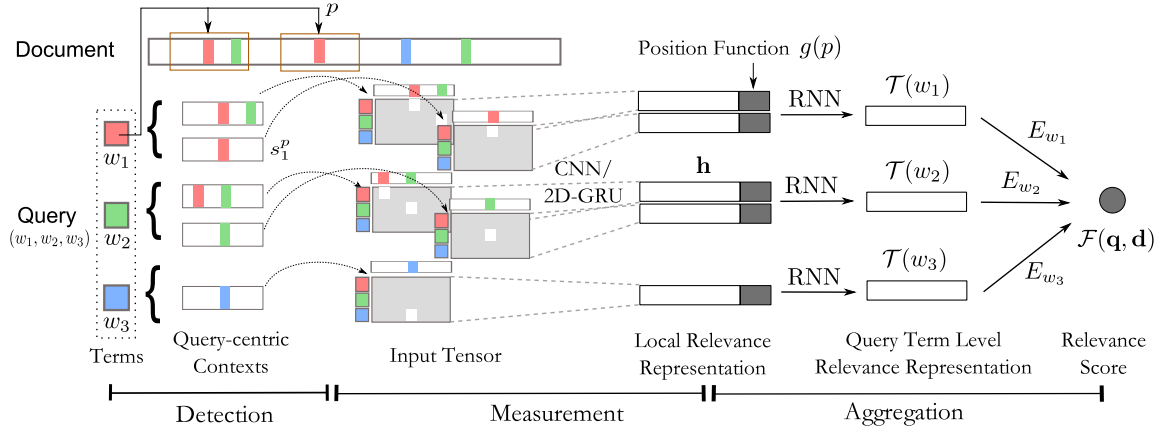


Figure 2: An illustration of DeepRank.

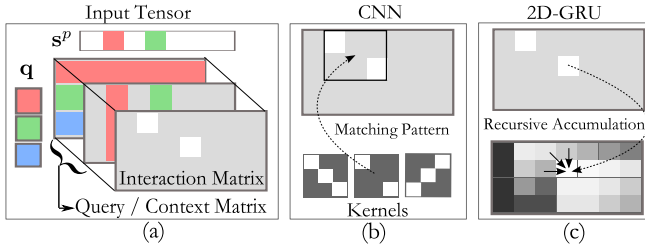


Figure 3: An illustration of measure network for a query and a query-centric context: (a) Input Tensor; (b) CNN; (c) 2D-GRU.

respectively. The mathematical formulas are shown as follows.

$$S_{ij}^{\text{ind}} = 1 \text{ if } w_i = v_j, \quad S_{ij}^{\text{ind}} = 0 \text{ otherwise,} \quad (1)$$

$$S_{ij}^{\text{cos}} = \mathbf{x}_i^T \mathbf{y}_j / (\|\mathbf{x}_i\| \cdot \|\mathbf{y}_j\|), \quad (2)$$

where \mathbf{x}_i and \mathbf{y}_j denote the word embeddings of w_i and v_j , respectively. To further incorporate the word representations of query/query-centric context to the input, we extend each element of S_{ij} to a three-dimensional vector $\tilde{S}_{ij} = [x_i, y_j, S_{ij}]^T$. Therefore, the original matrix S will become a three-order tensor, denoted as \mathcal{S} . In this way, the input tensor can be viewed as a combination of three matrices, i.e. query matrix, query-centric context matrix, and word-level interaction matrix.

Based on the input tensor \mathcal{S} , various neural networks can be directly applied to obtain the representations for the local relevance between query and query-centric context. In this paper, we choose the CNN architecture in MatchPyramid and 2D-GRU architecture in Match-SRNN, mainly because they have the ability to capture important proximity heuristics for IR.

4.2.2 Convolutional Neural Network. In this paper, we use a one-layer CNN in the measurement step, which includes a convolution operation and a max-pooling operation defined as follows. The convolution operation can extract various matching patterns from the input tensor \mathcal{S} , by using different kinds of kernels, as shown in Figure 3. Then a max-pooling operation is used to filter significant

matching patterns for further relevance determination.

$$h_{i,j}^{(\kappa)} = \sum_{l=1}^3 \sum_{s=0}^{\gamma-1} \sum_{t=0}^{\gamma-1} w_{s,t}^{(\kappa)} \cdot \mathcal{S}_{i+s,j+t}^{(l)} + b^{(\kappa)}, \quad (3)$$

$$h^{(\kappa)} = \max_{i,j} h_{i,j}^{(\kappa)}, \quad \kappa = 1, \dots, K$$

where l denotes the l -th slide of the tensor, γ denotes the fixed size of K different kernels, $\mathcal{S}_{i+s,j+t}^{(l)}$ denote the $(i+s, j+t)$ element of the l -th matrix of the input tensor \mathcal{S} , $w_{s,t}^{(\kappa)}$ and $b^{(\kappa)}$ denotes parameters.

Finally, all the significant matching patterns obtained from different kernels are concatenated to form a vector, i.e. $\mathbf{h} = [h^{(1)}, \dots, h^{(K)}]^T$, to represent the local relevance. This vector will be treated as the input to the aggregation network.

4.2.3 Two-Dimensional Gated Recurrent Units. Rather than using a hierarchical structure to capture the matching patterns, 2D-GRU in Match-SRNN [32] adopts a different sequential model to accumulate the matching signals. It is an extension of GRU [4] (a typical variant of RNN) to two-dimensional data like matrix or tensor¹. Specifically, 2D-GRU scans from top-left to bottom-right (or in a bidirectional way) recursively. At each position, the hidden representation depends on the representations of the top, left, diagonal and current positions in the matrix. The mathematical formulas are shown as follows.

$$\begin{aligned} \mathbf{c} &= [\mathbf{h}_{i-1,j}^T, \mathbf{h}_{i,j-1}^T, \mathbf{h}_{i-1,j-1}^T, \mathbf{S}_{ij}^T]^T, \\ \mathbf{r}_\theta &= \sigma(\mathbf{W}^{(r_\theta)} \mathbf{c} + \mathbf{b}^{(r_\theta)}), \quad \theta = l, t, d, \\ \mathbf{z}'_\phi &= \mathbf{W}^{(z_\phi)} \mathbf{c} + \mathbf{b}^{(z_\phi)}, \quad \phi = m, l, t, d, \\ \mathbf{r} &= [\mathbf{r}_l^T, \mathbf{r}_t^T, \mathbf{r}_d^T]^T, [\mathbf{z}_m, \mathbf{z}_l, \mathbf{z}_t, \mathbf{z}_d] = \text{RMax}([\mathbf{z}'_m, \mathbf{z}'_l, \mathbf{z}'_t, \mathbf{z}'_d]), \\ \mathbf{h}'_{ij} &= \psi(\mathbf{W} \mathbf{S}_{ij} + \mathbf{U}(\mathbf{r} \odot [\mathbf{h}_{i,j-1}^T, \mathbf{h}_{i-1,j}^T, \mathbf{h}_{i-1,j-1}^T]^T) + \mathbf{b}), \\ \mathbf{h}_{ij} &= \mathbf{z}_l \odot \mathbf{h}_{i,j-1} + \mathbf{z}_t \odot \mathbf{h}_{i-1,j} + \mathbf{z}_d \odot \mathbf{h}_{i-1,j-1} + \mathbf{z}_m \odot \mathbf{h}'_{ij}, \end{aligned} \quad (4)$$

where \mathbf{h}_{ij} stands for the hidden representation at the (i, j) -th position, $\mathbf{z}_m, \mathbf{z}_l, \mathbf{z}_t, \mathbf{z}_d$ are the four gates, \mathbf{U}, \mathbf{W} , and \mathbf{b} are parameters, σ

¹Strictly speaking, tensor is not a two-dimensional representation. However, 2D-GRU can be directly applied on tensor by treating each element of the matrix as a vector.

and ψ stand for sigmoid and tanh function respectively, and RMax is a function to conduct softmax on each dimension across gates,

$$[z_\phi]_j = \frac{e^{[z'_\phi]_j}}{e^{[z'_m]_j} + e^{[z'_l]_j} + e^{[z'_t]_j} + e^{[z'_d]_j}}, \quad \phi = m, l, t, d. \quad (5)$$

The last hidden representation of 2D-GRU will be treated as the output \mathbf{h} , which is the bottom right one at the matrix/tensor. If you use a bi-directional 2D-GRU, both the top left one $\vec{\mathbf{h}}$ and bottom right one $\overleftarrow{\mathbf{h}}$ can be concatenated together to form the output vector, i.e. $\mathbf{h} = [\vec{\mathbf{h}}^T, \overleftarrow{\mathbf{h}}^T]^T$.

Please note that both CNN and 2D-GRU well capture the proximity heuristics in IR. Proximity heuristic rewards a document where the matched query terms occur close to each other, which is an important factor for a good retrieval model. For CNN, if the matched query terms occur close to each other, appropriate kernels can be utilized to extract such significant matching patterns and influence the relevance score. In this way, CNN well captures the proximity heuristics. 2D-GRU can also model proximity. When there is a document where the matched query terms occur close to each other, the representation \mathbf{h} will be strengthened by appropriately setting gates and other parameters. As a result, the relevance score of the document will be increased.

4.3 Aggregation Network

After the measurement step, we obtain a vector \mathbf{h} to represent the local relevance between query and each query-centric context. Therefore, we need a further aggregation step to output a global relevance score. In this process, two IR principles are going to be considered in our deep architecture. One is *query term importance*: query terms are critical to express user's information need and some terms are more important than others [6]. The other one is *diverse matching requirement*: the distribution of matching patterns can be quite different in a relevant document. For example, the *Verbosity Hypothesis* assumes that the relevance matching might be global. On the contrary, the *Scope Hypothesis* assumes that the relevance matching could happen in any part of a relevant document, and we do not require the document as a whole to be relevant to a query. In order to capture the two IR principles, we first conduct a query term level aggregation, in which the diverse matching requirement is taken into account. Then a term gating network is applied to capture the importance of different terms when producing the global relevance score.

4.3.1 Query Term Level Aggregation. In order to capture the principle of diverse relevance requirement, we need to consider the position of the corresponding query-centric context when conducting query term level aggregation. Therefore, we append each vector \mathbf{h} with the position indicator to encode the position information of the corresponding query-centric context. Specifically, different position functions $g(p)$ are utilized in our aggregation network:

$$\begin{aligned} \text{Constant Function: } g(p) &= C, & C &\in \mathbb{R}, \\ \text{Linear Function: } g(p) &= (L - p)/L, & L &\in \mathbb{R}, \\ \text{Reciprocal Function: } g(p) &= a/(p + b), & a, b &\in \mathbb{R}, \\ \text{Exponential Function: } g(p) &= a \cdot \exp(-p/b), & a, b &\in \mathbb{R}, \end{aligned} \quad (6)$$

where p stands for the position of the query-centric context, determined by the central word v_p . After this appending operation, the representation of local relevance for a query-centric context centered at word v_p (denoted as $\mathbf{h}(p)$) becomes $[\mathbf{h}(p)^T, g(p)]^T$.

To conduct query term level aggregation, we first group $\mathbf{h}(p)$ with the same central word together, which stands for all the local relevances with respect to a same query term. Then RNN is used to integrate such local relevances by considering position information into consideration. That is to say, we can obtain the global relevance representation $\mathcal{T}(w_u)$ for each query term w_u as follows.

$$\mathcal{T}(w_u) = \text{RNN}_{p=p_1}^{p_n} \left[\mathbf{h}^T(p), g(p) \right], \quad p_1, \dots, p_n \in \mathbb{P}(w_u), \quad (7)$$

where $\mathbb{P}(w_u)$ denotes the position set of all the query-centric contexts centered on query term w_u . For example, you can use GRU [4] to capture the sequential information, which is a typical variant of RNN. In the experimental analysis, we show the comparisons among different position functions.

4.3.2 Term Gating Network for Global Aggregation. Based on query term level global relevance representations $\mathcal{T}(w_u)$, we use a term gating network (similar to that used in DRMM) to obtain the final global relevance score, by considering importances of different query terms. Specifically, we define a weight parameter E_{w_u} for each query term, and linear combine all the query term level relevances as follows.

$$\mathcal{F}(\mathbf{q}, \mathbf{d}) = \sum_{w_u \in \mathbf{q}} (E_{w_u} \mathbb{I})^T \cdot \mathcal{T}(w_u), \quad (8)$$

where \mathbb{I} is an vector with each element set to be 1, and the dimension is set to be the same as $\mathcal{T}(w_u)$.

4.4 Model Training

DeepRank is an end-to-end deep neural network, which can be trained using stochastic gradient decent (SGD) methods, such as Adam [16]. L2 regularization and early stopping strategy [9] are also used in our implementation to deal with overfitting. More implementation details will be given in the experiments.

In our experiments, we use the following pairwise hinge loss for training, since we are considering a ranking problem. For future work, we are also willing to try other pairwise losses and listwise loss functions to conduct the training process.

$$\mathcal{L}(\mathbf{q}, \mathbf{d}^+, \mathbf{d}^-) = \max(0, 1 - \mathcal{F}(\mathbf{q}, \mathbf{d}^+) + \mathcal{F}(\mathbf{q}, \mathbf{d}^-)), \quad (9)$$

where $\mathcal{L}(\mathbf{q}, \mathbf{d}^+, \mathbf{d}^-)$ denotes the pairwise loss between a pair of positive and negative samples \mathbf{d}^+ and \mathbf{d}^- , and $\mathcal{F}(\mathbf{q}, \mathbf{d})$ denotes the relevance score produced by DeepRank.

5 EXPERIMENTS

In this section, we conduct extensive experiments to evaluate DeepRank against state-of-the-art models, including learning to rank methods, and existing deep learning methods. The experimental results on both LETOR4.0 benchmark [25] and a large scale click-through data show that our model can significantly outperform all the baselines, especially when other existing deep learning methods perform much worse than learning to rank methods. Furthermore, we give detailed experimental analysis to show more insights on our model.

5.1 Experimental Settings

We first introduce our experimental settings, including datasets, baseline methods/implementations, and evaluation measures.

5.1.1 Data Sets. Since most deep models need to learn many parameters, it is not appropriate to evaluate them with small traditional retrieval dataset, such as Robust04 and ClueWeb-09-Cat-B used in [11], which have only less than 300 queries. In our experiments, we use two datasets for evaluation, i.e. LETOR4.0 [25] and a large scale clickthrough data. The LETOR4.0 data is mainly used for comparing our model with other deep models, and the state-of-the-art learning to rank methods. The clickthrough data is larger, and we use it to compare different deep models.

LETOR4.0 dataset contains two separate data sampled from the .GOV2 corpus using the TREC 2007 and TREC 2008 Million Query track queries, denoted as MQ2007 and MQ2008, respectively. MQ2007 is a bit larger, which contains 1692 queries and 65,323 documents. While MQ2008 only contains 784 queries and 14,384 documents. Since the query number in MQ2008 is too small, which will cause the serious insufficient training problem for deep learning models, we propose to merge the training set of MQ2007 to that of MQ2008. The validation and testing set are kept unchanged. Then we form a new large data set, still denoted as MQ2008. In total, MQ2007 and MQ2008 contains 69,623 and 84,834 query-document pairs, respectively. All the baselines are conducted fairly on this new dataset for comparison. In original LETOR4.0, each query and document pair is represented as a vector containing 46 different features, which is easy for implementations of learning to rank methods. While most deep IR models (except for SQA [28]) do not use any handcrafted features, the raw text of query and document are used for implementation.

The large scale clickthrough data, namely ChineseClick, is collected from a commercial Chinese search engine. In the data collection process, the user is given the top 10 results for each proposed query. Clicked documents are viewed to be relevant, and the other ones are viewed as irrelevant. Since this is a Chinese dataset, we first conduct word segmentation for queries and documents. Then we apply some typical data preprocessing techniques, such as navigational queries filtering, stopping words and low frequency words (less than 50) removing. After these preprocessing, the final dataset contains 12,520 queries, 115,562 documents, and 118,835 query-document pairs. It is further divided into training/validation/testing set according to the proportion 3:1:1 of query numbers.

5.1.2 Baseline Methods. We adopt two types of baseline methods for comparison, including learning to rank methods and deep learning methods.

For learning to rank approach, we compare both pairwise and listwise ranking methods. The pairwise baselines include *RankSVM* [14] and *RankBoost* [7], which apply SVM and boosting techniques to the pairwise ranking problem, respectively. The listwise baselines include *AdaRank* [34] and *LambdaMart* [1], where *AdaRank* proposes to directly optimizing IR evaluation measures by boosting to obtain a ranking list, and *LambdaMart* uses gradient boosting for optimizing a listwise ranking loss, which is the winner of Yahoo Learning to Rank Challenge [3]. Though the public results on LETOR4.0 have included RankSVM, RankBoost and AdaRank as

the baselines, we are not able to conduct significant testing since the ranking list and relevance scores are missing. Therefore, we implement them on our own. Most of our results are comparable with those on LETOR4.0². For RankSVM, we directly use the implementation in SVM^{rank} [15]. RankBoost, AdaRank and LambdaMart are implemented using RankLib³, which is a widely used tool in the area of learning to rank. For LETOR4.0 dataset, we use 46 dimensional standard features provided for public to evaluate the performance of learning to rank approaches. BM25-Title which calculate BM25 score between query and document title, is one of the powerful feature among these features.

For deep learning approach, we compare three existing deep IR models, i.e. *DSSM* [13], *CDSSM* [29], and *DRMM* [11]. We also compare some popular deep methods for text matching, including one representation focused method, i.e. *ARC-I* [12], and three interaction focused methods, i.e. *ARC-II* [12], *MatchPyramid* [24], and *Match-SRNN* [32]. Implementation details are listed as follows. Firstly, all the word embeddings in these methods are learned with the Continuous Bag-of-Words (CBOW) model [20] from Wikipedia corpus, and the dimension is set to 50. In general, most deep learning baselines are applied following the original implementations. We only reduce the parameter numbers due to the relative small size of our data. For example, we use a three-layer DNN as that in the original paper of DSSM, and the node number of each layer is reduced to 100, 100, and 50. For CDSSM, we use a one-layer CNN with 50 (1×3) kernels. Therefore, a 50-dimensional vector is obtained after global pooling strategy. DRMM is directly implemented using the best configuration and the code released by [11]. For ARC-I, we use a two-layer CNN with each one containing 16 (1×3) kernels. The size of pooling in the first layer is set to 2, while the size of pooling in the last layer is set to be 2 for query representation, and 20 for document representation, respectively. For ARC-II, we use a two-layer CNN, where there are 8 kernels in each layer. The size of kernels and pooling in both layers are set to (1×3)/(3×3) and (2×2)/(2×20), respectively. For MatchPyramid, we use cosine similarity to construct the word-level interaction matrix. Then a one-layer CNN is applied on the matrix with 8 (3×3) kernels, and the size of dynamic pooling is set to (3×10). For Match-SRNN, 2D-GRU is directly applied on the same interaction matrix, and the dimension of hidden node is set to 2.

SQA [28] model combines handcraft features in the learning process, therefore it is not appropriate to directly compare it with DeepRank, which only uses automatically learned feature from raw text for ranking. For fair comparison, we delete the handcrafted features in SQA and obtain a pure deep SQA model, denoted as SQA-noFeat. Furthermore, we incorporate the handcrafted features (46 default features in LETOR4.0) into the last layer of DeepRank to obtain DeepRank-Feat, which is used to compare with SQA. For both SQA-noFeat and SQA, one-layer CNN with 50 (1×3) kernels is used in the deep architecture.

The DeepRank⁴ for performance comparison is implemented using the following settings: the window size of query-centric context is set to 15; cosine similarity is adopted to construct the input tensor; both CNN and 2D-GRU are used in the measurement step,

²<http://research.microsoft.com/en-us/um/beijing/projects/letor/letor4baseline.aspx>

³<https://sourceforge.net/p/lemur/wiki/RankLib/>

⁴The source code: <https://github.com/pl8787/textnet-release>.

therefore we have two versions of DeepRank, denoted as DeepRank-CNN and DeepRank-2DGRU; the reciprocal function is used as the positional function, and GRU is adopted in the aggregation step. We also compare different settings of DeepRank for detailed analysis.

5.1.3 Evaluation Measures. For the evaluation on LETOR4.0, we follow the data partitions on this dataset (5-fold) and the average results are reported. While for the evaluation on ChineseClick, we train the model on the training set, tune hyper-parameters on the validation set and report the results on the testing set for comparison. Three evaluation measures [26] are used in this paper, i.e. Precision, NDCG, and MAP. Furthermore, we conduct a pairwise t-test [30] for significance testing with p-value lower than 0.05 (i.e. $p\text{-value} \leq 0.05$).

5.2 Performance Comparison

The performance comparison results of DeepRank against baseline models are shown in Table 1.

5.2.1 Performance Comparison on LETOR 4.0. From the results on MQ2007 and MQ2008, we can see that: 1) None of existing deep learning models could perform comparably with learning to rank methods. Some of them are even worse than BM25. The results tell us that the automatically learned features in existing deep learning models are not better than traditional extracted ones, though they are using more complex models for training. Someone may argue that the experimental findings are inconsistent with previous studies that DSSM and CDSSM can significantly outperform traditional retrieval models, as stated in [13] and [29]. The reason lies in that LETOR4.0 is much smaller than the clickthrough data used in [13] and [29]. In the following experiments on ChineseClick, we can see that all the deep models perform better than BM25. Therefore, deep models usually need more data for optimization, which is also the reason why we do not use Robust04 and ClueWeb-09-CAT-B used in [11] for evaluation. 2) As for the comparisons between these deep models, interaction focused ones such as DRMM, ARC-II, MatchPyramid and Match-SRNN perform much better than representation focused ones such as DSSM, CDSSM, ARC-I, and SQA-noFeat. This is consistent with the understanding that interaction signals are much more important than the semantic representation of query/document in IR, as described in [11]. Furthermore, DRMM performs the best among all the deep learning baseline methods. This is because DRMM further incorporate IR characteristics into their architecture, which indicate the importance of capturing IR intrinsics in the architecture design process. 3) Our DeepRank not only significantly outperforms the deep learning baselines, but also significantly improves the results of learning to rank methods, even only use the query and document raw text data. For example, the improvement of DeepRank-CNN against the best deep learning baseline (i.e. DRMM) on MQ2007 is 16.1% w.r.t. NDCG@1, 12.9% w.r.t. P@1, and 6.4% w.r.t. MAP, respectively; while the improvement of DeepRank-CNN against the best learning to rank method (i.e. LambdaMart) on MQ2007 is 7.0% w.r.t. NDCG@1, 5.6% w.r.t. P@1, and 6.2% w.r.t. MAP, respectively. The results indicate that by appropriately modeling relevance, deep learning approach can significantly outperform learning to rank approach for IR application. 4) Though SQA has used both automatically learned features

and handcrafted features, the performance cannot compare with DeepRank by using only automatically learned features for ranking. If we incorporate handcrafted features into DeepRank, the performance will be further improved, as shown in DeepRank-CNN-Feat. The results demonstrate the superiority of our deep architecture.

5.2.2 Performance Comparison on ChineseClick. The ChineseClick data is used to compare DeepRank with other deep learning methods. We do not include the DSSM and CDSSM baselines. That is because DSSM and CDSSM are specially designed for English data, and letter-trigram is used as the input of neural network, which is not applicable for Chinese data. If we are using the whole word embedding as the input, CDSSM will become the same as ARC-I. Therefore, we omit CDSSM and directly report the results of ARC-I for comparison. The results show that deep learning baselines perform comparably with BM25, some are even better. That is because we are using a larger data, the training of deep models become more sufficient and the performances are improved. Our DeepRank still performs the best. The improvement against the BM25 is about 21.0% w.r.t. NDCG@1, and 11.5% w.r.t. MAP. While the improvement against the best deep learning baseline (i.e. Match-SRNN) is about 11.0% w.r.t. NDCG@1, and 4.3% w.r.t. MAP.

From the above results, we conclude that DeepRank significantly improves the results of relevance ranking, with architecture specially designed to model human's relevance generation process.

5.3 Detailed Analysis of DeepRank

DeepRank is such a flexible deep architecture that different parameter settings and neural networks can be used in the detection, measurement, and aggregation steps. Some of these settings may largely influence the final ranking performances. Therefore, we conduct a detailed analysis on MQ2007 to show the comparisons of DeepRank with different settings, with expect to give some insights for implementation. Specifically, we analyze four factors, i.e. window size of query-centric context in the detection step, input tensor in the measurement step, neural network in the measurement step, positional function in the aggregation step. We change one factor of the above DeepRank-CNN each time to conduct the comparisons.

5.3.1 Impact of Different Window Sizes of Query-Centric Context. The window size of query-centric context determines the scope of local relevance in the human judgment process. With a small window size, users would determine local relevance with less effort since contexts are short, but it is easy to introduce ambiguity due to limited context information. When window size is large, there are sufficient contexts to facilitate the precise local relevance judgment, but the cost is also increased and many noises may influence the judgment. We conduct an experiment to compare different window sizes of query-centric context, varying in the range of 1, 7, 11, 15, 19 and 23. The results listed at the top of Table 2 show that the performances of DeepRank first increase and then become stable, with the increase of window size. The best performance is obtained with window size up to 11/15 (w.r.t different evaluation measures). Therefore, with considering the computational complexity, we suggest to use a comparable medium window size in real application, and the exact number need to be tuned considering averaged document length, query length, and data size.

Table 1: Performance comparison of different models on MQ2007, MQ2008 and ChineseClick. Significant performance degradation with respect to DeepRank-CNN is denoted as (-) with p-value ≤ 0.05 .

MQ2007									
Model	NDCG@1	NDCG@3	NDCG@5	NDCG@10	P@1	P@3	P@5	P@10	MAP
BM25-TITLE	0.358 ⁻	0.372 ⁻	0.384 ⁻	0.414 ⁻	0.427 ⁻	0.404 ⁻	0.388 ⁻	0.366 ⁻	0.450 ⁻
RANKSVM	0.408 ⁻	0.405 ⁻	0.414 ⁻	0.442 ⁻	0.472 ⁻	0.432 ⁻	0.413 ⁻	0.381 ⁻	0.464 ⁻
RANKBOOST	0.401 ⁻	0.404 ⁻	0.410 ⁻	0.436 ⁻	0.462 ⁻	0.428 ⁻	0.405 ⁻	0.374 ⁻	0.457 ⁻
ADARANK	0.400 ⁻	0.410 ⁻	0.415 ⁻	0.439 ⁻	0.461 ⁻	0.431 ⁻	0.408 ⁻	0.373 ⁻	0.460 ⁻
LAMBDA MART	0.412 ⁻	0.418 ⁻	0.421 ⁻	0.446 ⁻	0.481 ⁻	0.444 ⁻	0.418 ⁻	0.384 ⁻	0.468 ⁻
DSSM	0.290 ⁻	0.319 ⁻	0.335 ⁻	0.371 ⁻	0.345 ⁻	0.359 ⁻	0.359 ⁻	0.352 ⁻	0.409 ⁻
CDSSM	0.288 ⁻	0.288 ⁻	0.297 ⁻	0.325 ⁻	0.333 ⁻	0.309 ⁻	0.301 ⁻	0.291 ⁻	0.364 ⁻
ARC-I	0.310 ⁻	0.334 ⁻	0.348 ⁻	0.386 ⁻	0.376 ⁻	0.377 ⁻	0.370 ⁻	0.364 ⁻	0.417 ⁻
SQA-NOFEAT	0.309 ⁻	0.333 ⁻	0.348 ⁻	0.386 ⁻	0.375 ⁻	0.373 ⁻	0.372 ⁻	0.364 ⁻	0.419 ⁻
DRMM	0.380 ⁻	0.396 ⁻	0.408 ⁻	0.440 ⁻	0.450 ⁻	0.430 ⁻	0.417 ⁻	0.388 ⁻	0.467 ⁻
ARC-II	0.317 ⁻	0.338 ⁻	0.354 ⁻	0.390 ⁻	0.379 ⁻	0.378 ⁻	0.377 ⁻	0.366 ⁻	0.421 ⁻
MATCHPYRAMID	0.362 ⁻	0.364 ⁻	0.379 ⁻	0.409 ⁻	0.428 ⁻	0.404 ⁻	0.397 ⁻	0.371 ⁻	0.434 ⁻
MATCH-SRNN	0.392 ⁻	0.402 ⁻	0.409 ⁻	0.435 ⁻	0.460 ⁻	0.436 ⁻	0.413 ⁻	0.384 ⁻	0.456 ⁻
DEEPRANK-2DGRU	0.439	0.439	0.447	0.473	0.513	0.467	0.443	0.405	0.489
DEEPRANK-CNN	0.441	0.447	0.457	0.482	0.508	0.474	0.452	0.412	0.497
SQA	0.423	0.432	0.442	0.466	0.491	0.463	0.443	0.404	0.479
DEEPRANK-CNN-FEAT	0.446	0.457	0.462	0.489	0.518	0.483	0.451	0.412	0.502

MQ2008									
Model	NDCG@1	NDCG@3	NDCG@5	NDCG@10	P@1	P@3	P@5	P@10	MAP
BM25-TITLE	0.344 ⁻	0.420 ⁻	0.461 ⁻	0.220 ⁻	0.408 ⁻	0.381 ⁻	0.337 ⁻	0.245 ⁻	0.465 ⁻
RANKSVM	0.375 ⁻	0.431 ⁻	0.479 ⁻	0.229	0.441 ⁻	0.390 ⁻	0.348 ⁻	0.249	0.478 ⁻
RANKBOOST	0.381	0.436 ⁻	0.477 ⁻	0.231	0.455	0.392 ⁻	0.347 ⁻	0.248	0.481 ⁻
ADARANK	0.360 ⁻	0.422 ⁻	0.462 ⁻	0.222	0.430 ⁻	0.384 ⁻	0.339 ⁻	0.247 ⁻	0.468 ⁻
LAMBDA MART	0.378	0.437 ⁻	0.477 ⁻	0.231	0.446 ⁻	0.398	0.348 ⁻	0.251	0.478 ⁻
DSSM	0.286 ⁻	0.336 ⁻	0.378 ⁻	0.178 ⁻	0.341 ⁻	0.307 ⁻	0.284 ⁻	0.221 ⁻	0.391 ⁻
CDSSM	0.283 ⁻	0.331 ⁻	0.376 ⁻	0.175 ⁻	0.335 ⁻	0.302 ⁻	0.279 ⁻	0.222 ⁻	0.395 ⁻
ARC-I	0.295 ⁻	0.363 ⁻	0.413 ⁻	0.187 ⁻	0.361 ⁻	0.336 ⁻	0.311 ⁻	0.229 ⁻	0.424 ⁻
SQA-NOFEAT	0.291 ⁻	0.350 ⁻	0.401 ⁻	0.184 ⁻	0.366 ⁻	0.332 ⁻	0.309 ⁻	0.231 ⁻	0.416 ⁻
DRMM	0.368 ⁻	0.427 ⁻	0.468 ⁻	0.220 ⁻	0.437 ⁻	0.392 ⁻	0.344 ⁻	0.245 ⁻	0.473 ⁻
ARC-II	0.299 ⁻	0.340 ⁻	0.394 ⁻	0.181 ⁻	0.366 ⁻	0.326 ⁻	0.305 ⁻	0.229 ⁻	0.413 ⁻
MATCHPYRAMID	0.351 ⁻	0.401 ⁻	0.442 ⁻	0.211 ⁻	0.408 ⁻	0.365 ⁻	0.329 ⁻	0.239 ⁻	0.449 ⁻
MATCH-SRNN	0.369 ⁻	0.426 ⁻	0.465 ⁻	0.223 ⁻	0.432 ⁻	0.383 ⁻	0.335 ⁻	0.239 ⁻	0.466 ⁻
DEEPRANK-2DGRU	0.391	0.436	0.480	0.236	0.462	0.395	0.354	0.252	0.489
DEEPRANK-CNN	0.406	0.460	0.496	0.240	0.482	0.412	0.359	0.252	0.498
SQA	0.402	0.454	0.493	0.236	0.485	0.411	0.362	0.254	0.496
DEEPRANK-CNN-FEAT	0.418	0.475	0.507	0.248	0.497	0.422	0.366	0.255	0.508

ChineseClick									
Model	NDCG@1	NDCG@3	NDCG@5	NDCG@10	P@1	P@3	P@5	P@10	MAP
BM25	0.200 ⁻	0.320 ⁻	0.412 ⁻	0.280 ⁻	0.200 ⁻	0.174 ⁻	0.169 ⁻	0.152	0.373 ⁻
ARC-I	0.208 ⁻	0.359 ⁻	0.451 ⁻	0.286 ⁻	0.208 ⁻	0.193 ⁻	0.180 ⁻	0.153	0.393 ⁻
SQA-NOFEAT	0.232 ⁻	0.368 ⁻	0.458 ⁻	0.292 ⁻	0.232 ⁻	0.194 ⁻	0.180 ⁻	0.153 ⁻	0.403 ⁻
DRMM	0.218 ⁻	0.346 ⁻	0.442 ⁻	0.288 ⁻	0.218 ⁻	0.185 ⁻	0.177 ⁻	0.153	0.392 ⁻
ARC-II	0.190 ⁻	0.329 ⁻	0.430 ⁻	0.283 ⁻	0.190 ⁻	0.180 ⁻	0.177 ⁻	0.153	0.373 ⁻
MATCHPYRAMID	0.204 ⁻	0.342 ⁻	0.436 ⁻	0.285 ⁻	0.204 ⁻	0.184 ⁻	0.178 ⁻	0.153	0.384 ⁻
MATCH-SRNN	0.218 ⁻	0.360 ⁻	0.456 ⁻	0.295 ⁻	0.218 ⁻	0.190 ⁻	0.181 ⁻	0.153	0.399 ⁻
DEEPRANK-2DGRU	0.244	0.382	0.473	0.299	0.244	0.205	0.185	0.153	0.415
DEEPRANK-CNN	0.242	0.386	0.476	0.298	0.242	0.202	0.185	0.153	0.416

Table 2: Performance comparisons of DeepRank with different settings on MQ2007.

Model	NDCG@1	NDCG@5	MAP
DEEPRANK-W1	0.422	0.425	0.470
DEEPRANK-W7	0.426	0.444	0.490
DEEPRANK-W11	0.438	0.458	0.495
DEEPRANK-W15	0.441	0.457	0.497
DEEPRANK-W19	0.430	0.453	0.496
DEEPRANK-W23	0.441	0.455	0.496
DEEPRANK- S_I^{ind}	0.416	0.427	0.473
DEEPRANK- S_I^{cos}	0.411	0.430	0.479
DEEPRANK- S_R	0.425	0.439	0.482
DEEPRANK- S_{IR}^{cos}	0.441	0.457	0.497
DEEPRANK-DNN	0.383	0.414	0.471
DEEPRANK-2DGRU	0.440	0.447	0.489
DEEPRANK-CNN	0.441	0.457	0.497
DEEPRANK-CONST	0.384	0.419	0.473
DEEPRANK-LINEAR	0.431	0.445	0.492
DEEPRANK-EXP	0.441	0.454	0.494
DEEPRANK-RECIP	0.441	0.457	0.497

5.3.2 Impact of Different Input Tensors. In order to capture both word representations of query/query-centric context and their interactions, we propose to construct a three-order tensor \mathcal{S} as the input of the measure network. Here, we compare four different settings of tensor. S_I^{ind} and S_I^{cos} stand for the case when we use indicator or cosine function to construct the interaction matrix, and omit the other two matrices in the tensor. S_R stands for the case that only word representations of query and query-centric context is considered in the tensor, i.e. interaction matrix is ignored. S_{IR}^{cos} stands for the case when we use the three-order tensor, which is exactly the DeepRank we used in the performance comparisons. From the results listed in the second row of Table 2, we can see the performances are improved when more information is modeled in the tensor. Therefore, both word representations of query/query-centric context and word-level interactions are important to the relevance judgement.

5.3.3 Impact of Different Measure Networks. The measure network is adopted to determine the relevance between query and a detected query-centric context. In the model section, we demonstrate how to use CNN and 2D-GRU to conduct such measurement, mainly because these two kinds of neural networks have the ability to capture the proximity heuristics. Of course, you can also use other deep learning architectures, such as DNN. In this section, we conduct experiments on MQ2007 to compare the three different versions of DeepRank, denoted as DeepRank-DNN, DeepRank-CNN, and DeepRank-2DGRU. The experimental results in the third row of Table 2 show that DeepRank-CNN and DeepRank-2DGRU perform much better than DeepRank-DNN. The reason lies in that CNN and 2D-GRU both have the ability to model the proximity heuristics, while DNN cannot because it is position sensitive, which is contradict with the position independent proximity heuristic.

5.3.4 Impact of Different Position Functions. As described in the aggregation network, different kinds of position functions can be

used to model the position importance. Here we compare DeepRank with four different position functions, i.e. *Constant*, *Linear*, *Reciprocal*, and *Exponential* functions, denoted as DeepRank-Const, DeepRank-Linear, DeepRank-Recip and DeepRank-Exp, respectively. The results listed in the fourth row of Table 2 show that DeepRank-Recip is the best, while DeepRank-Const is the worst. As for the other two functions, DeepRank-Exp perform comparable with DeepRank-Recip, and DeepRank-Linear is a little worse than DeepRank-Recip and DeepRank-Exp. The results indicate that top positions are more important, which is consistent with many previous studies for relevance ranking in IR [21]. As for the reason why reciprocal and exponential function performs better than linear function, we think this is because MQ2007 is extracted from GOV data, where title and abstraction information may play a dominant role in determining the relevance. Therefore, the functions with a long tail, as that in reciprocal or exponential function, will be favored. To sum up, the position function plays an important role in DeepRank, and users should pay more attention to the choice, which need to be conducted by considering the characteristics of different applications.

5.3.5 Relations to Previous Models. We also would like to point out that DeepRank has a close relationship with previous models, such as BM25, MatchPyramid, and Match-SRNN. With some simplification, DeepRank can reduce to (or approximate) these models.

BM25 is a bag-of-words retrieval model that ranks a set of documents based on the query terms appearing in each document, regardless of the inter-relationships between the query terms within a document. The most common form is given as follows.

$$\text{BM25}(\mathbf{q}, \mathbf{d}) = \sum_{w \in \mathbf{q}} \text{IDF}(w) \cdot \frac{f(w, \mathbf{d}) \cdot (k_1 + 1)}{f(w, \mathbf{d}) + k_1 \cdot (1 - b + \frac{b|\mathbf{d}|}{\text{avgdl}})} \quad (10)$$

where k_1 and b are the hyper parameters, $f(w, \mathbf{d})$ represents term frequency of w in document \mathbf{d} , $\text{IDF}(w)$ represents inverse document frequency of w , $|\mathbf{d}|$ denotes the document length and avgdl denotes the averaged document length.

We are going to show that a simplified DeepRank has the ability to approximate BM25 function. Firstly, the window size of query-centric context is set to 1. Then the indicator function is used to construct the input tensor, therefore word representations of query/query-centric context are ignored. In this way, exact matching signals are naturally captured, like in BM25. We can see that the output of tensor will be a 0-1 vector, with only the elements at the matched positions will be 1. If CNN is used in the measurement step, we can omit the convolution layer and directly use the pooling strategy to output the value 1; while if 2D-GRU is used, we can also output the value 1 by appropriately setting gates and parameters. As a consequence, the output of the measure network will be 1 for each query-centric context. At the aggregation step, we set the position function as a constant $g(p) = 1/|\mathbf{d}|$, and term weight as the IDF value, i.e. $E_{w_u} = \text{IDF}(w_u)$. Therefore, the output of this DeepRank can be viewed as the following function:

$$\begin{aligned} \text{DeepRank}(\mathbf{q}, \mathbf{d}) &= \sum_{w \in \mathbf{q}} \text{IDF}(w) \cdot \text{RNN} \left[\begin{matrix} 1, 1/|\mathbf{d}| \end{matrix} \right] \\ &= \sum_{w \in \mathbf{q}} \text{IDF}(w) \cdot \mathcal{G}(f(w, \mathbf{d}), |\mathbf{d}|), \end{aligned} \quad (11)$$

where the second equation is obtained because RNN is an accumulative process, and the function \mathcal{G} is determined by the parameters in RNN, learned from the training data. The function \mathcal{G} has high capacities to approximate the functions in the formula of BM25 since there are many parameters. Therefore, a simplified version of DeepRank can well approximate the BM25 model.

In addition, DeepRank has closer relationships with MatchPyramid and Match-SRNN. If we set the window size of query-centric context to be $k = |d|$ and the weights of query term w_u to be $1/f(w_u, d)$, DeepRank reduces to MatchPyramid or Match-SRNN, by using CNN or 2D-GRU as the measure network, respectively.

6 CONCLUSIONS AND FUTURE WORK

In this paper, we propose a new deep learning architecture, namely DeepRank. Firstly, a detection strategy is designed to extract query-centric contexts. A measure network is then applied to determine the local relevance between query and each query-centric context, by using CNN or 2D-GRU. Finally, an aggregation network is used to produce the global relevance score, via RNN and a term gating network. DeepRank not only well simulates the relevance generation process in human judgement, but also captures important IR characteristics, i.e. exact/semantic matching signals, proximity heuristics, query term importance, and diverse relevance requirement. We conduct experiments on both benchmark LETOR4.0 data and a large clickthrough data. The results show that DeepRank significantly outperform learning to rank methods and existing deep IR models, when most existing deep IR models perform much worse than learning to rank methods. To the best of our knowledge, DeepRank is the first deep IR model to outperform existing learning to rank models. We also give a detailed analysis on DeepRank to show insights on parameter settings for implementation.

For future work, we plan to investigate the differences between the automatically learned representations of DeepRank and effective features used in learning to rank, which may introduce some insights for architecture design of more powerful deep IR models.

7 ACKNOWLEDGMENTS

This work was funded by the 973 Program of China under Grant No. 2014CB340401, the National Natural Science Foundation of China (NSFC) under Grants No. 61232010, 61433014, 61425016, 61472401, and 61203298, and the Youth Innovation Promotion Association CAS under Grants No. 20144310 and 2016102. The authors would like to thank Chengxiang Zhai (UIUC) and Yixing Fan (ICT, CAS) for their valuable suggestions on this work, and Weipeng Chen (Sogou Inc.) for providing helps on the data processing of ChineseClick.

REFERENCES

- [1] Christopher JC Burges. 2010. From ranknet to lambdarank to lambdamart: An overview. *Learning* 11 (2010), 23–581.
- [2] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *ICML*. ACM, 129–136.
- [3] Olivier Chapelle and Yi Chang. 2011. Yahoo! learning to rank challenge overview. In *Proceedings of the Learning to Rank Challenge*. 1–24.
- [4] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *EMNLP*. 1724–1734.
- [5] Carsten Eickhoff, Sebastian Dungs, and Vu Tran. 2015. An eye-tracking study of query reformulation. In *SIGIR*. ACM, 13–22.
- [6] Hui Fang, Tao Tao, and ChengXiang Zhai. 2004. A formal study of information retrieval heuristics. In *SIGIR*. ACM, 49–56.
- [7] Yoav Freund, Raj Iyer, Robert E Schapire, and Yoram Singer. 2003. An efficient boosting algorithm for combining preferences. *JMLR* 4, Nov (2003), 933–969.
- [8] Fredric C Gey. 1994. Inferring probability of relevance using the method of logistic regression. In *SIGIR*. Springer, 222–231.
- [9] Rich Caruana Steve Lawrence Lee Giles. 2001. Overfitting in Neural Nets: Back-propagation, Conjugate Gradient, and Early Stopping. In *NIPS*, Vol. 13. MIT Press, 402.
- [10] Alan Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *ICASSP*. IEEE, 6645–6649.
- [11] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. A deep relevance matching model for ad-hoc retrieval. In *CIKM*. ACM, 55–64.
- [12] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. 2014. Convolutional neural network architectures for matching natural language sentences. In *NIPS*. 2042–2050.
- [13] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *CIKM*. ACM, 2333–2338.
- [14] Thorsten Joachims. 2002. Optimizing search engines using clickthrough data. In *SIGKDD*. ACM, 133–142.
- [15] Thorsten Joachims. 2006. Training linear SVMs in linear time. In *SIGIR*. ACM, 217–226.
- [16] Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [17] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature* 521, 7553 (2015), 436–444.
- [18] Tie-Yan Liu. 2009. Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval* 3, 3 (2009), 225–331.
- [19] Yuanhua Lv and ChengXiang Zhai. 2009. Positional language models for information retrieval. In *SIGIR*. ACM, 299–306.
- [20] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS*. 3111–3119.
- [21] Shuzi Niu, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2012. Top-k learning to rank: labeling, ranking and evaluation. In *SIGIR*. ACM, 751–760.
- [22] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. *The PageRank citation ranking: Bringing order to the web*. Technical Report. Stanford InfoLab.
- [23] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, and Xueqi Cheng. 2016. A study of matchpyramid models on ad-hoc retrieval. In *Neu-IR f16 SIGIR Workshop on Neural Information Retrieval*.
- [24] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Shengxian Wan, and Xueqi Cheng. 2016. Text matching as image recognition. In *AAAI*. AAAI Press, 2793–2799.
- [25] Tao Qin, Tie-Yan Liu, Jun Xu, and Hang Li. 2010. LETOR: A benchmark collection for research on learning to rank for information retrieval. *Information Retrieval* 13, 4 (2010), 346–374.
- [26] Stephen Robertson. 2000. Evaluation in information retrieval. In *Lectures on information retrieval*. Springer, 81–92.
- [27] Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR*. Springer-Verlag New York, Inc., 232–241.
- [28] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to rank short text pairs with convolutional deep neural networks. In *Proceedings of SIGIR*. ACM, 373–382.
- [29] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. Learning semantic representations using convolutional neural networks for web search. In *WWW*. WWW, 373–374.
- [30] Mark D Smucker, James Allan, and Ben Carterette. 2007. A comparison of statistical significance tests for information retrieval evaluation. In *CIKM*. ACM, 623–632.
- [31] Tao Tao and ChengXiang Zhai. 2007. An exploration of proximity measures in information retrieval. In *SIGIR*. ACM, 295–302.
- [32] Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu, Liang Pang, and Xueqi Cheng. 2016. Match-SRNN: Modeling the Recursive Matching Structure with Spatial RNN. In *IJCAI*. 2922–2928.
- [33] Ho Chung Wu, Robert WP Luk, Kam-Fai Wong, and KL Kwok. 2007. A retrospective study of a hybrid document-context based retrieval model. *Information processing & management* 43, 5 (2007), 1308–1331.
- [34] Jun Xu and Hang Li. 2007. Adarank: a boosting algorithm for information retrieval. In *SIGIR*. ACM, 391–398.