

Personalizing Search Results Using Hierarchical RNN with Query-aware Attention

Songwei Ge^{1,3}, Zhicheng Dou^{1,2,3}, Zhengbao Jiang^{1,2}, Jian-Yun Nie⁵, and Ji-Rong Wen^{1,2,4}

¹School of Information, Renmin University of China, ⁵DIRO, Université de Montréal

²Beijing Key Laboratory of Big Data Management and Analysis Methods

³National Engineering Laboratory of Big Data System Software (Beijing Institute of Technology)

⁴Key Laboratory of Data Engineering and Knowledge Engineering, MOE

gesongwei666@gmail.com, dou@ruc.edu.cn, rucjzb@163.com, nie@iro.umontreal.ca, jirong.wen@gmail.com

ABSTRACT

Search results personalization has become an effective way to improve the quality of search engines. Previous studies extracted information such as past clicks, user topical interests, query click entropy and so on to tailor the original ranking. However, few studies have taken into account the sequential information underlying previous queries and sessions. Intuitively, the order of issued queries is important in inferring the real user interests. And more recent sessions should provide more reliable personal signals than older sessions. In addition, the previous search history and user behaviors should influence the personalization of the current query depending on their relatedness. To implement these intuitions, in this paper we employ a hierarchical recurrent neural network to exploit such sequential information and automatically generate user profile from historical data. We propose a query-aware attention model to generate a dynamic user profile based on the input query. Significant improvement is observed in the experiment with data from a commercial search engine when compared with several traditional personalization models. Our analysis reveals that the attention model is able to attribute higher weights to more related past sessions after fine training.

KEYWORDS

search results personalization; hierarchical recurrent neural network; query-aware attention

ACM Reference Format:

Songwei Ge, Zhicheng Dou, Zhengbao Jiang, Jian-Yun Nie, and Ji-Rong Wen. 2018. Personalizing Search Results Using Hierarchical RNN with Query-aware Attention. In *The 27th ACM International Conference on Information and Knowledge Management (CIKM'18)*, October 22–26, 2018, Torino, Italy. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3269206.3271728>

1 INTRODUCTION

Users come to search engine with specific intents, but the queries they issue often fail to express accurate meanings. It is often the

case that users issue the same query to express different search intents. Therefore, returning the same results to all users is not the best strategy for a search engine. Personalization potentially solves this problem: by appropriately modeling user interest profiles, personalized search is able to alleviate the ambiguity problem by providing more precise results to individual users.

Personalizing search results to particular users based on their profiles has been a hot field in both academia [4, 6, 32, 37] and industry [15] for a long time. Among these studies, many personalization methods utilized click-through data to generate user profiles, since large-scale query logs contain strong indications on users' preferences and are also relatively cheap to acquire when compared with manual labels [4, 19, 32, 35, 41]. These studies demonstrated that through mining personalized signals from query logs, a search engine can re-rank search results and obtain more relevant rankings for individual users. In this paper, we focus on exploiting information from the sequences of past queries and sessions and build a dynamic user profile based on input query.

Previous studies have revealed that diverse features, such as click counts, user topical interests, click entropy and so on, could be extracted from user query logs to help personalize search results. [4, 11, 16, 33, 40]. However, few studies have leveraged sequential information hidden in past queries and sessions. Intuitively, a more recent query or session should contribute more than an older one to the current search. For example, a user may have searched for "cherry blossom Japan", "cherry jam" and so on, but is now interested in mechanical keyboard and intends to obtain some introductions. His query "cherry reviews" should be more related to the keyboard brand rather than the fruit or the flower. Previous studies attempted to simply apply an exponential decay to distinguish the recent and old search behaviors [41, 43]. No significant difference was observed when using such temporal weightings [4]. However, the influence of a previous query depends on more complex factors than merely the time. A highly related previous query may continue to have a strong influence on the interpretation of the current query in the same session even after a long time span. In this paper, we will use a deep learning framework to account for such factors in the sequential influence patterns.

When building the user profile, another natural intuition is that previous queries and sessions are not always useful [40]. For example, the profile constructed for previous query "peanut allergy symptoms" is probably useless when the user is searching for "JAVA book". In contrast, if the user once searched "JAVA runtime environment", then we can infer that the book is about "JAVA programming" instead of "JAVA Island". Based on such observations, we also devise

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 the author(s) 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 '18, October 22–26, 2018, Torino, Italy

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-6014-2/18/10...\$15.00

<https://doi.org/10.1145/3269206.3271728>

a query-aware attention mechanism to learn different weights for previous sessions corresponding to the input query. By doing so, we are able to build a dynamic user profile which attends to the previous sessions that are more important to the current query.

User preference evolves over time, and some previous studies attempted to capture this variation by distinguishing short- and long-term user preferences [4, 21, 43]. A short-term user profile is built from recent interactions within the same search session and is useful to predict the following intents in the current search task [43, 44]. A long-term user profile describes more long-standing user characteristics and is less sparse than a short-term profile [23], and it is built from previous sessions. Motivated by the hierarchical structures contained in the click-through data, we propose to use hierarchical recurrent neural network as our fundamental framework to model a longer dependency in the sequential data.

More specifically, a short-term user profile is constructed from the user behaviors in the same session to reflect the current search interest. This is implemented by a RNN at low-level. In addition, the long-term user profiles are built from the past sessions, which are designed to reflect the user's long-term interests. This is created using a higher-level RNN in our model. In order to use the long-term user profiles, we attribute weights to each long-term interest profile using an attention model based on the current query. Then the linear combination of all hidden state vectors is produced as a more precise long-term interest vector. Finally, we compute the similarity score between interest vectors and documents as our personalized scores to tailor the original ranking. The whole framework is trained through a learning-to-rank approach, LambdaRank algorithm [5]. Taking advantage of the two interest vectors, we can further enhance search results personalization.

In sum, our main contributions are twofold: First, we account for the sequential information contained in click-through data and generate a better user profile through a hierarchical recurrent neural network; second, we apply attention mechanism to scrutinize all previous sessions and highlight the more important sessions dynamically according to present information need.

The rest of paper is organized as follows. Section 2 summarizes previous studies that are related to our paper. The proposed framework to mine sequential information from query log is described in Section 3. We discuss experimental setups in Section 4, and analyze the results in Section 5. We conclude the work in Section 6.

2 RELATED WORK

The related work to this paper principally concerns three fields: (1) Search Results Personalization, (2) Deep Learning in Information Retrieval and (3) The Applications of Hierarchical Recurrent Neural Network.

Search Results Personalization. Users' past search interactions with search engine have been revealed beneficial to web search [6]. However, the improvement of ranking quality personalization brings to users mostly depends on the richness of user profiles. Among all intriguing research on personalizing search results, great majority focused on building user profile based on click-through data, since it is both accessible and informative. The personalized features are extracted from click-through data and can be simply categorized into two groups: click features and topical features.

Users often use search engine to navigate for the same document which was satisfactory under the previously issued queries, and this fact can be used as a safe and low-risk approach in personalization. Teevan et al. [37] have demonstrated that there is a rich opportunity to personalize search results through recognizing personal navigations. Many subsequent studies [4, 22, 38, 45] regarded click-based features as their basic features and investigated more from other angles. It is widely known that personalization sometimes plays against the goal of improving results. Therefore, click entropy is often considered to ensure the expected utility of personalization on certain queries [11, 36]. More studies were focused on building an appropriate latent topical user profiles. Early studies [3, 31, 45] attempted to build user profiles with topics of clicked documents which are learned from a manual on-line ontology, such as the Open Directory Project (ODP)¹. The potential problem of this approach is that some classes of documents may not appear in the on-line ontology, which could limit its application to dealing with new documents [7, 16]. Recent studies applied a latent topic model to determining these topics [7, 16, 39, 41]. Besides, topic entropy is also proposed to measure the topical ambiguity of different queries [33]. Li et al. [22] utilized semantic features produced by a deep learning model as semantic features to improve in-session contextual ranking. In addition to extracting features from query logs, other information can be utilized to improve search results as well. Bennett et al. [2] incorporated user location-based features into personalization model. Collins-Thompson et al. [9] studied and evaluated personalization by different reading levels. Different from the previous studies which simply aggregate the historical user behaviors as user profiles, our intuition is that the past behaviors should be differently conducive to current search according to, for example, the time span and the current information need.

Strategies to implement personalization also varied considerably from model to model. Wang et al. [42] and Song et al. [32] investigated adapting a generic ranking model for personalized search through updating parameters for individual users. The model proposed by Matthijs et al. [23] built user profiles with terms extracted from browsing history, and personalized the rankings with the matching between search snippets and these terms. As for the models focusing on topical user profile with Latent Dirichlet Allocation (LDA) [7, 16, 39, 41], the authors directly calculated the similarity between each document and established user profile, and used the similarities as personalized scores to re-rank the results. Another common approach is to train a ranking model using the LambdaMART [46] learning algorithm as in [4, 38, 45]. LambdaMART is evolved from LambdaRank [5] and can distinguish the importance of features automatically based on boosted regression tree. Different from aforementioned approaches, we employ a deep learning framework to personalize search results. Instead of highly depending on the manual features, the framework does not need specific preparation and only takes raw queries and documents with one hot representation of words as inputs.

Deep Learning in Information Retrieval. Deep learning has been introduced in many information retrieval tasks because of its several advantages, such as the ability to learn word embeddings automatically and to train end-to-end [24, 25]. Ad-hoc rankings have

¹<https://dmoztools.net/>

benefited enormously from deep learning methods and great results have been achieved. DSSM [18] and CDSSM [30] embedded the query-document pair into a semantic space and ranked the results by the similarity between the embeddings of documents and queries. Severyn and Moschitti [29] also used a convolutional deep neural network to represent query-document pair and computed their semantic similarity. The model developed by Palangi et al. [26] shared a similar idea but addressed the embeddings of queries and documents using Long Short-Term Memory (LSTM) [17]. More recently, some elaborate models like DRMM [14] and K-NRM [47] were proposed to model word-level similarities and achieved even better results. But little has been studied on applying deep learning to search results personalization. Song et al. [32] adapted a global ranking model with continue-train for each individual user to personalize the results. Li et al. [22] improved in-session contextual results by involving semantic features generated by deep learning models. Different from these studies which incorporate deep learning only as a component, we intend to train a complete deep learning framework to personalize the search results.

The Applications of Hierarchical Recurrent Neural Network. The hierarchical recurrent neural network (HRNN) was firstly proposed by El Hidi and Bengio [12] aiming at modeling longer dependence in the sequential data. Since then people gradually realized that many real tasks contain hierarchically sequential dependencies, and HRNN attracts increasing attention from various communities. Quadrona et al. [27] employed HRNN to deal with cross- and in-session commodity information in a session-based recommendation. Besides, hierarchical recurrent encoder-decoder was designed to handle text generation, such as building dialogue system [28] and query suggestions [34]. In addition to basic hierarchical recurrent neural network, in this paper we also employ an attention model to generate dynamic user profile based on current query.

3 PERSONALIZATION FRAMEWORK

The one-size-fits-all retrieval models are known to be suboptimal and can be potentially improved with search results personalization. Through re-ranking the unpersonalized results for different users according to their interests, search engines can provide better search experiences to individual users. In this paper, we focus on mining sequential information and building dynamic user profiles using users' past interactions. Specifically, we employ a framework based on hierarchical recurrent neural network with query-aware attention to learn short- and long-term user profiling automatically. The short-term interest is collected from current session, and the long-term interest describes the preferences inferred by other sessions before the current session.

To start with, we formulate our problem as follows. Suppose that for each user u , there is a query log \mathcal{L}_u which includes past sessions $\mathcal{L}_u = \{S_1, \dots, S_{M-2}, S_{M-1}\}$ and current session S_M , where M is the index of the current session. Each session is defined as a sequence of queries and a list of documents retrieved for each query, denoted as $S_m = \{\{q_{m,1}, d_{m,1,1}, \dots\}, \dots, \{q_{m,n_m}, d_{m,n_m,1}, \dots\}\}$, where n_m represents the total number of queries in the session S_m . In the current session S_M , the user is issuing a query q_{M,n_M} and a document ranking list $\mathcal{D} = \{d_1, d_2, \dots\}$ is returned by the search engine, where d_i is short for $d_{M,n_M,i}$. For each document,

we want to compute a relevance score separately based on the issued query q_{M,n_M} , and past interactions \mathcal{L}_u and S_M :

$$\begin{aligned} \text{score}(d_i) &= \text{score}(d_i | q_{M,n_M}, S_M, S_{M-1}, \dots, S_1) \\ &= \text{score}(d_i | q_{M,n_M}) + \text{score}(d_i | \mathcal{L}_u) + \text{score}(d_i | S_M), \end{aligned}$$

where $\text{score}(d_i | q_{M,n_M})$ represents the relevance between document and query, $\text{score}(d_i | \mathcal{L}_u)$ and $\text{score}(d_i | S_M)$ correspond to the relevance with regard to long- and short-term user interests respectively. Finally, we re-rank the document list \mathcal{D} according to the overall score $\text{score}(d_i)$, to produce a personalized ranking list.

As shown in Figure 1, we devise a deep learning framework based on HRNN to calculate the above personalized scores using query logs and to predict the optimal document ranking list. We elaborate our model based on three main components as follows: (1) Exploiting sequence-enhanced user interests, (2) Modeling dynamic user profiles and (3) Re-ranking.

3.1 Exploiting Sequence-enhanced User Interests

As we stated in Section 1, traditional methods for building topical user profiles have several potential problems. Most of them ignored the sequential information contained in historical user behaviors. Inspired by the ability of RNN to model sequential data and the hierarchical structure of query logs, we propose to use a Hierarchical RNN to model user interests with click-through data.

Recurrent neural network has been widely used in natural language processing because of its ability to model the sequences of words and sentences. After inputs at each time point, the RNN cell is able to determine what information will be kept and passed to next state, and what to discard. More sophisticated RNN cells such as Gated Recurrent Unit (GRU) [8] and Long Short-Term Memory (LSTM) [17] were proposed to solve vanishing gradient problems and learn a long-term dependency. Specifically, we employ GRU as a basic RNN cell in our model, since it has a relatively simpler structure which is easier to train. The reasons we devise a two-level structure include: (1) The click-through data is intrinsically hierarchical with intra-session interactions conveying short-term interest and inter-session interactions reflecting long-term preference; (2) To enhance the ability of our framework to model longer dependency, specifically to learn longer user interests.

3.1.1 Modeling short-term user interests. Within a session S_m , a user issues a series of queries. Sometimes this user issues a query at first and reformulates it in the following queries. During this procedure, the user probably clicks on certain documents and may be satisfied by their content. Intuitively, both the satisfactory documents and the order of these queries are highly informative to inferring the real intent and should be utilized while tailoring the ranking of the result lists.

In our model, the inputs to the low-level RNN is extracted from a subset of the session $S'_m = \{\{q_{m,1}, \mathcal{D}_{m,1}\}, \dots, \{q_{m,n_m}, \mathcal{D}_{m,n_m}\}\}$, where $\mathcal{D}_{m,i}$ represents the set of SAT-clicked documents under query $q_{m,i}$. And we use the concatenation of query vector $q_{m,i}$ and the average vector $d_{m,i}$ of document vectors in $\mathcal{D}_{m,i}$ to feed the RNN. Note that the average document vector $d_{m,i}$ is assigned to a zero vector if the user is not satisfied by any documents. To be specific, for each query $q \in \mathbb{R}^{d_e}$ and document $d \in \mathbb{R}^{d_e}$, their

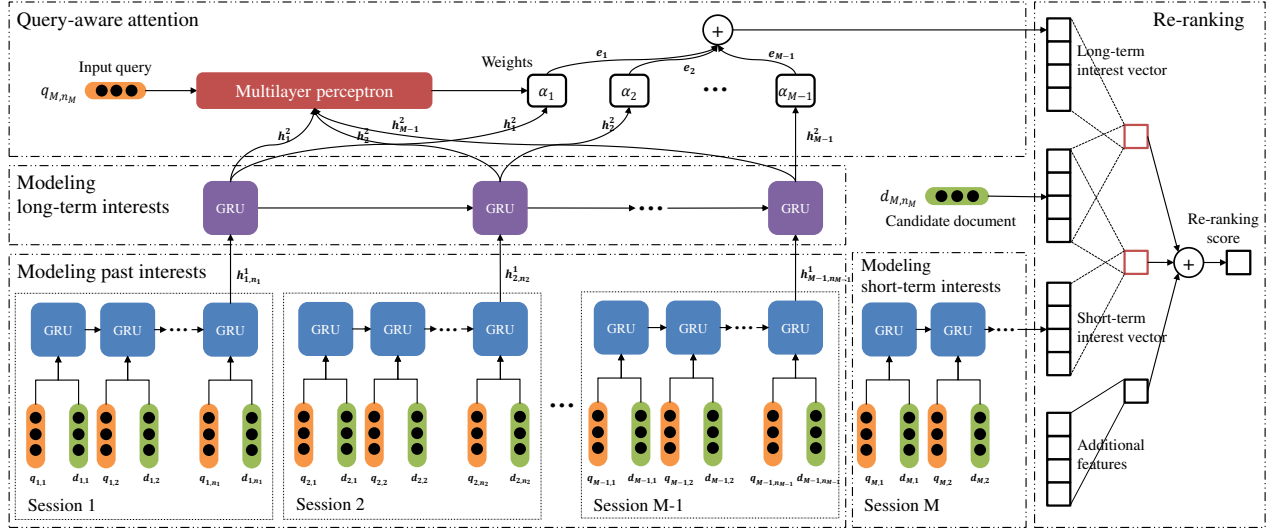


Figure 1: The architecture of our personalization framework. Given query and document representations in each past session, the interest vector of the session is calculated as the last latent state vector of a low-level RNN, which is then proceeded into a high-level RNN. Next, weights are applied on each latent state vector of the high-level RNN according to the input query and generate the long-term interest vector. By matching the interest vectors with candidate documents, and involving additional features, we compute the final personalized scores.

representations are calculated as the weighted average of the word representations $w \in \mathbb{R}^{d_e}$ multiplied by TF-IDF weights, while the words are represented by distributional vectors looked up in an embeddings matrix $W \in \mathbb{R}^{|\mathcal{V}| \times d_e}$. This part could be replaced by a more complex model such as Doc2Vec, CNN, or RNN model, which also involves more parameters and increase the difficulty of training. Then the session-level RNN calculates a series of latent state vectors according to the inputs and previous state vectors, formally:

$$h_{m,n}^1 = f(h_{m,n-1}^1, q_{m,n}, d_{m,n})$$

where $h_{m,n}^1 \in \mathbb{R}^{d_{h^1}}$ and $h_{m,0}^1$ is initialized by a zero vector. Note that we use superscript 1 and 2 to distinguish the notations in low- and high-level RNN respectively such as hidden state vectors. The RNN cell $f(\cdot)$ can be vanilla, LSTM, and GRU cell, which is implemented as GRU in our model:

$$\begin{aligned} r_{m,n}^1 &= \sigma(W_r^1[q_{m,n}; d_{m,n}] + V_r^1 h_{m,n-1}^1), \\ z_{m,n}^1 &= \sigma(W_z^1[q_{m,n}; d_{m,n}] + V_z^1 h_{m,n-1}^1), \\ c_{m,n}^1 &= \tanh(W^1[q_{m,n}; d_{m,n}] + V^1(r_{m,n}^1 \cdot h_{m,n-1}^1)), \\ h_{m,n}^1 &= (1 - z_{m,n}^1) \cdot h_{m,n-1}^1 + z_{m,n}^1 \cdot c_{m,n}^1, \end{aligned}$$

where $h_{m,n}^1 \in \mathbb{R}^{d_{h^1}}$, the reset gate $r_{m,n}^1$ and update gate $z_{m,n}^1$ control the trade-off between previous states and present inputs, $\sigma(\cdot)$ is the sigmoid function, parameters $W_r^1, V_r^1, W_z^1, V_z^1, W^1$ and V^1 are shared across all sessions and updated during training. Each latent hidden state $h_{m,n}^1$ models the short-term interest in session S_m after issuing query $q_{m,n}$. The dense vector of the last latent hidden state, denoted by h_{m,n_m}^1 , is viewed as a representation of

the whole session, i.e. the short-term interest vector, which is used as h_m^1 for short in the following parts.

3.1.2 Modeling long-term user interests. Looking at the whole search history, some interests do not change frequently alongside time. For example, if a user kept searching "JAVA runtime environment", "JAVA data type", and "JAVA string to int", the user was probably a programmer. Even though after a very long time this user issued another query "JAVA Book", there should still be a great chance that the user is looking for a book of JAVA programming instead of JAVA Island. Therefore, the parameters of such cell states of interests should change less frequently than the parameters of short-term model, and this fact motivates us to build a hierarchical structure in our framework. Therefore, we add another RNN to encode a higher level information flow.

The high-level RNN takes the short-term interest vectors (for past sessions) $\{h_{1,n}^1, \dots, h_{M-1,n}^1\}$ outputted by the low-level RNN as inputs. It computes the sequence of user representations $\{h_{1,n}^2, \dots, h_{M-1,n}^2\}$ at the end of every session. Formally:

$$h_m^2 = f(h_{m-1}^2, h_m^1),$$

where $h_m^2 \in \mathbb{R}^{d_{h^2}}$, and h_0^2 is initialized by a zero vector. The RNN cell $f(\cdot)$ is similarly defined in the low-level RNN:

$$\begin{aligned} r_m^2 &= \sigma(W_r^2 h_m^1 + V_r^2 h_{m-1}^2), \\ z_m^2 &= \sigma(W_z^2 h_m^1 + V_z^2 h_{m-1}^2), \\ c_m^2 &= \tanh(W^2 h_m^1 + V^2(r_m^2 \cdot h_{m-1}^2)), \\ h_m^2 &= (1 - z_m^2) \cdot h_{m-1}^2 + z_m^2 \cdot c_m^2, \end{aligned}$$

where parameters $W_r^2, V_r^2, W_z^2, V_z^2, W^2$, and V^2 are shared across all users and updated during training. Analogously, each state h_m^2 of this layer models long-term interest of user u after the session S_m ends. A simple way to obtain the long-term user interest vector is to adopt the last latent state vector h_{M-1}^2 directly. However, this could introduce unnecessary noise as we stated in Section 1. To avoid this issue, we implement an attention model to assign query-aware weights to long-term interest vectors $\{h_1^2, \dots, h_{M-1}^2\}$ in the end of different sessions, so as to highlight different parts of historical search behaviors dynamically. This intuitive idea is proved effective and better than plain RNN model in our experiments.

3.2 Building Dynamic User Profiles

As discussed in Section 1, the same past interactions could contribute differently to personalized rankings under different search circumstances. We believe that discriminating previous behaviors is necessary in building an effective user profile. Vu et al. built user profiles using dynamic group formation based on input query [40]. In our model, we deploy an attention model to apply different weights for each previous session based on current query. Note that in an extreme case, the weights for some interactions should be very high when a user issues a historical query to re-find the information they have searched before. Such re-finding problem has been well studied by Teevan et al. as personal navigation [37] and is a special case naturally incorporated in our model.

Attention mechanism was first proposed in machine translation to deal with the limitation of fixed-length representation of input sentences [1]. Similarly we hypothesize that encoding the click through data into a fixed-length representation is not reasonable when we are personalizing search results under different situations. Therefore, with respect to the current query, we assign query-aware weights to past interactions and compute a more precise long-term user profile. In other words, we suppose that the long-term interest vector is influenced by the current information need dynamically. At the end of session S_{M-1} , we calculate weights $\{\alpha_1, \dots, \alpha_{M-1}\}$ for each past interest vector in $\{h_1^2, \dots, h_{M-1}^2\}$ as follow.

$$e_i = \phi(q_{M, n_M}, h_i^2),$$

$$\alpha_i = \frac{\exp(e_i)}{\sum_{j=1}^{M-1} \exp(e_j)},$$

where $\phi(\cdot)$ is a multilayer perceptron (MLP) with $\tanh(\cdot)$ as activation function in our model, which is updated during training and could be replaced by more complex functions in the future. Then the query-aware long-term interest vector $h_{M-1}^{2,q}$ is computed by a weighted linear combination of $\{h_1^2, \dots, h_{M-1}^2\}$:

$$h_{M-1}^{2,q} = \sum_{i=1}^{M-1} \alpha_i h_i^2.$$

In sum, we denote the final short- and long-term interest vectors as h_{M, n_m}^1 and $h_{M-1}^{2,q}$ respectively and use them to re-rank the search results in the following.

3.3 Re-ranking

Finally we re-rank the original results using the personalized information we collect. Given the short- and long-term interest vectors, h_{M, n_m}^1 and $h_{M-1}^{2,q}$, we calculate the personalized ranking scores of a document by measuring its similarity to these two interest vectors:

$$\text{score}(d_i | \mathcal{L}_u) = \text{sim}((h_{M-1}^{2,q})^T W_L, d_i),$$

$$\text{score}(d_i | S_M) = \text{sim}(h_{M, n_m}^1 W_S, d_i),$$

where $W_S \in \mathbb{R}^{d_{h^1} \times d_e}$, $W_L \in \mathbb{R}^{d_{h^2} \times d_e}$ are two similarity matrices whose functions are to project the interest vectors into the same semantic space as the documents, which are optimized during training. The similarity function $\text{sim}(\cdot)$ is defined as:

$$\text{sim}(X, Y) = \frac{X^T Y}{\|X\| \cdot \|Y\|}.$$

In addition to personalization scores calculated by our model, we also incorporate query-document relevance feature and click-based features as additional features. Since the original query-document features are inaccessible in our dataset, here we use the original position of the document as a feature. Also, the click features include the total number of historical clicks on the candidate document by the user, the number of clicks on candidate document under the input query by the user and the click entropy of input query. The reason we incorporate click entropy is because the value of personalization varies a lot across different queries, and indiscriminately applying personalization on all queries could produce an adverse effect on the overall quality [11, 36]. Note that as for fair comparison, the baseline model is a more complex model that incorporates more than these three features. And our idea is focused on building topic-based user profile instead of investigating on click-based features as discussed in Section 1. These additional features are fed into a multilayer perceptron (MLP) with $\tanh(\cdot)$ as the activation function. Finally, we sum up the scores calculated by different parts as the final relevance score.

We choose a basic ranking algorithm, LambdaRank [5], to train the whole framework. We generate training pairs from query logs by treating the SAT-clicked documents as relevant samples and the others as irrelevant ones. Then we use the representations of document pairs to calculate the loss. Take a pair of relevant document d_i and irrelevant document d_j as an example. The loss function is the product of cross entropy between desired probabilities and predicted probabilities and the change of metrics, Δ , while swapping the positions of the two documents, defined as:

$$\text{loss} = (-\bar{p}_{ij} \log(p_{ij}) - \bar{p}_{ji} \log(p_{ji})) |\Delta|,$$

where p_{ij} represents the predicted probability that d_i is more relevant than d_j , and \bar{p}_{ij} represents the real probability. Specifically, the predicted probabilities are computed by a logistic function,

$$p_{ij} = \frac{1}{1 + \exp(-(\text{score}(d_i) - \text{score}(d_j)))}.$$

4 EXPERIMENT SETUP

4.1 Dataset and Evaluation

The dataset in our experiment is sampled randomly by users from the logs of a commercial search engine, comprising click-through

Table 1: Basic statistics of the dataset.

Item	Statistic	Item	Statistic
#days	58	#distinct queries	1,624,496
#users	33,204	#sessions	654,776
#queries	2,665,625	#SAT-clicks	1,228,028

data of the users between 1st January 2013 and 28th February 2013. Each piece of data in the log contains an anonymous user identifier, a query, a session identifier, query issued time, the top 20 URLs retrieved by the search engine, clicks and dwelling time. The logs are collected when personalization support was not applied, so that our results are guaranteed not to be biased toward other personalization signals. The basic statistics is shown in Table 1.

Following [4, 41], we similarly regard the click that has a dwelling time of more than 30 seconds or is the last one in the session as a satisfied click (SAT-click). We use the first six weeks of data to generate basic user profiles and use the remaining two weeks to train and test the models. Note that we use not only sessions from the first six weeks but all sessions before current session to calculate long-term interest vector. We split the last two-week data into training and test sets according to the sessions instead of the dates. The motivation behind this is that two weeks are a relatively short time span and the distribution of queries sometimes is uneven across dates. Besides, sessions can be viewed as search activities with independent intents of users, thus can be reasonably divided into different sets. For each user, we divide the sessions into 5:1 as training and test set respectively in the time order of sessions and treat the last one fifth sessions of training set as validation set. For each URL in the logs we retrieve its content and we remove the URLs that cannot be found anymore. Then we preprocess the documents by removing stopwords and punctuations. To ensure an effective segmentation of training and test dataset, we also remove the users who had less than 4 sessions.

We measure the quality of each ranking results using mean average precision (MAP), mean reciprocal rank (MRR), precision@1 (P@1) and average click position (Avg. Click). In addition to these mainstream metrics, we further evaluate the rankings by measuring the actual improvements on inverse document pairs [20], i.e. a clicked document and the skipped documents ranked before it. Our reason behind this is that clicks are not only based on the relevance of documents but also the positions [10]. This kind of position bias may make mainstream metrics somewhat problematic. For example, the low-ranked documents have lower probabilities to be examined. Even though some of them may be relevant, they are not clicked because they are not examined by users. These results are labeled as irrelevant by traditional metrics. Consequently, if we boost the positions of the unexamined documents, the traditional evaluation metrics cannot reflect the real changes. Therefore, we use the strategy, which considers clicked documents are better than skipped documents, proposed by Joachims et al. [20] to collect inverse document pairs. We compute the number (#Better) and percentage of improved pairs (P-Improve) in each personalized ranking to evaluate the reliable improvements made by personalization methods.

4.2 Baselines

In addition to the original ranking generated by the search engine, which is usually of very high quality, we further reproduce several state-of-the-arts personalization models as follows.

P-Click: Users often issue an identical query to re-find the previously viewed information. This observation was confirmed by Teevan et al [37] as a safe opportunity to improve the quality of search engine. Dou et al. proposed P-Click [11] as a basic personalization strategy. As for a user u , P-Click calculates the personalized score of a document d as the percentage of clicks on this document from past queries that are identical to the input query q , formally:

$$\text{score}(d|q, u) = \frac{|\text{clicks}(q, d, u)|}{|\text{clicks}(q, \bullet, u)| + \beta},$$

where $|\text{clicks}(q, d, u)|$ represents the click counts on d from query q by user u in the past, $|\text{clicks}(q, \bullet, u)|$ is the total click counts on query q by user u , and β is a smoothing factor set as 0.5. P-Click re-ranks the results based on this personalized score and combines the personalized ranking with original ranking using Borda' ranking fusion method. This model is also used as a baseline in [23]. Note that this method is purely based on the click information.

Besides, we also care about the quality of user profile built by our model and past studies. Many past studies built user profiles based on clicked documents over a topic space. They assume that the topic that users are interested in is reflected by the topics of documents they clicked [33]. These past studies tailored the original ranking either directly using similarity between user profile and documents [40] or training a supervised model with diverse features [4]. However, the intrinsic idea of these methods is the same - using the aggregation of topics of past documents as user profiles.

SLTB: Bennett et al [4] implemented a personalization method by extracting diverse features from short- and long-term behavior (SLTB). SLTB generates personalized features through the combination of four options: (1) feature types (click-based or topic-based), (2) document coverage (across all queries; under queries identical/generalizations/specializations to the current query.), (3) temporal angle (historic, session and aggregate) and (4) temporal decay (able or disable). The decay is defined by the function $0.95^{p(q_r)-1}$. Here $p(q_r)$ refers to the number of queries preceding the current query and 0.95 is a chosen decay factor. As for the topical representation of the documents, it uses the top two levels of the Open Directory Project hierarchy as the 207 labels of documents and trains a classifier to predict the categories of documents. Besides, SLTB also implements click entropy [11], topic entropy [33], and other features. All these features are fed into a learning-to-rank model, LambdaMART [46], to generate a personalized ranking. As for the parameters used in our experiments, especially those in the Learning-to-rank models, we initialize them with parameters same as in [4] and then tune them using the validation set.

PTM: Learning document representations from manual labels is problematic according to Carmen [7], since the categories of some documents are missing in the ontology. Instead, they proposed to use an unsupervised approach to learn a multinomial distribution for documents on the latent topics, and this method is also used by Vu et al. [40] later. Here we reproduce the PTM model proposed by [16] as our baseline. PTM calculated personalized scores based on

Table 2: Overall performances of models. Bold indicates the main model proposed by this paper, which is also the best among all compared models. The improvements achieved by HRNN and HRNN+QA on MAP, MRR, P@1 and Avg.Click are significantly larger than improvements made by any baseline model with paired t-test at p-value<0.01.

Model	MAP	MRR	P@1	Avg. Click	#Better	P-Improve
Original Ranking	.7226	.7334	.5931	2.292	-	-
P-Click	.7348	.7467	.6015	2.138	4,834	.1419
PTM	.6679	.6801	.5244	2.578	9,684	.2845
SLTB	.7776	.7881	.6698	2.054	16,257	.4777
SLTB+PTM	.7830	.7929	.6716	1.998	16,602	.4878
HRNN	.7989	.8107	.7039	1.930	18,166	.5338
HRNN+QA	.8017	.8135	.7067	1.904	18,609	.5468

the likelihood of documents given both query and the user as:

$$\text{score}(d|q, u) \propto P(d) \prod_{w \in q} \sum_z P(w|z) P(u|z)^\lambda P(z|d),$$

where λ balances the weight of user's topical interest influencing on the overall ranking, and $P(d)$ is estimated with Dirichlet smoothing based on the relative frequency of clicks on d in the whole log:

$$\hat{P}(d) = \frac{\#clicks(d) + \sigma \frac{1}{|\mathcal{D}|}}{\sum_{d_i} \#clicks(d_i) + \sigma}.$$

SLTB+PTM: We also replace the topical features in SLTB by the features generated by the topic model from PTM, and keep the other features the same as in SLTB. We use this method as the fourth baseline.

4.3 Our Models

To verify the effects of our personalization framework, we train two personalization models formed with different compositions but under a same learning setting. Specifically:

Hierarchical RNN+Query-aware Attention (HRNN+QA): The complete model stated in Section 3;

Hierarchical RNN (HRNN): The Query-aware Attention is disabled and the long-term interest vector is represented by the last latent vector of the high-level RNN layer;

To determine a group of appropriate parameters, we apply a grid search according to the performance of the model on the validation set. We empirically use a two-layer network for the attention and additional MLP. We use a range of word embedding sizes $d_e \in \{300, 1000\}$, sizes of short-term interest vector $d_{s1} \in \{100, 200, 300, 500\}$, sizes of long-term interest vector $d_{s2} \in \{200, 400, 600, 1000\}$, number of hidden units in attention MLP $d_a \in \{512, 1024\}$, number in additional MLP $d_f \in \{32, 64, 128\}$, and learning rates $\lambda \in \{1e^{-4}, 1e^{-3}, 1e^{-2}\}$. No noticeable differences are observed while changing learning rates. When embedding size is larger, a relatively smaller size of RNN state will yield worse results. Considering the balance of efficiency and result quality, we finally chose a combination of parameters as $d_e = 300$, $d_{s1} = 300$, $d_{s2} = 600$, $d_a = 1024$, $d_f = 64$ and $\lambda = 1e^{-3}$.

Similar to [29], we initialize the word embedding matrix W with a pre-trained unsupervised model and keep it fixed during the training. In this experiment, we train an embedding matrix on the documents from training dataset using the Google word2vec tool

². Finally, as we state at Section 3, we use Mean Average Precision as our metric to calculate Δ in the learning-to-rank model, since MAP encourages to rank the most relevant documents at top positions, but we still evaluate the final rankings with four different metrics. Besides, we adopt an early stop strategy to end the training when the average loss on validation set stops decreasing in three continuous epochs.

5 EXPERIMENTAL RESULTS AND ANALYSIS

As discussed in Section 1, the sequences of previous queries and sessions are valuable to personalization and a static user profile is not enough. Therefore, we are very interested in two research questions: (1) Is hierarchical recurrent neural network able to mine sequential information from query logs and learn a better topical user profile? (2) Does attention mechanism help to highlight different parts of search histories dynamically? To answer the questions, we firstly evaluate the overall performances of models stated in the previous section. Then we will discuss the effectiveness of attention mechanism by visualizing the provided weights. In addition, we will analyze the performances of baseline models and our models under different situations, including on queries with different click entropies, on repeated and non-repeated queries, and on queries at different positions in a session.

5.1 Overall Performance

We evaluate the results generated by the search engine, baseline models and our proposed models using MAP, MRR, P@1, Avg. Click, #Better, and P-Improve. In addition to original ranking, the baseline models include P-Click, PTM, SLTB and SLTB+PTM, and our frameworks include HRNN and HRNN+QA. All the scores are computed over all test queries. Results are shown in Table 2. We have the following observations:

(1) All personalization methods except PTM successfully improve the original ranking. Directly using user profiles to re-rank the results might cause many problems. For example, applying personalization on queries with low click entropies may increase the risk of personalization. It is observed that P-Click only fixes about half of inverse document pairs of those improved by PTM. In contrast, though very simple, P-Click improves the original ranking in a relatively safe way. SLTB and SLTB-PTM significantly improve

²<https://code.google.com/archive/p/word2vec/>

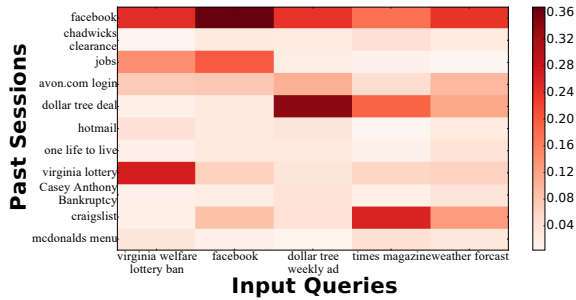


Figure 2: The weights of past sessions when different queries are issued. A darker area indicates a larger weight.

original ranking with paired t-test at $p < 0.01$. This proves that a well-trained learning-to-rank model with rich features can effectively yield personalization under different circumstances.

(2) As shown in Table 2, our models generate great improvements over all baseline models and the differences are statistically significant with paired t-test at $p < 0.01$ level. More specifically, HRNN+QA obtains rankings that are 0.0187 higher in MAP than the rankings generated by SLTB+PTM. As for reliable improvements, HRNN+QA produces a 12% increase in the number of improved inverse document pairs than SLTB+PTM. This outcome demonstrates that our models are able to learn more precise user profiles, and consequently yield better personalization.

(3) We find that attention mechanism works on HRNN and improves its results with 443 more inverse document pairs fixed. However, such improvement is not statistically significant. One possible reason is that the attention model is difficult to apply on a relatively long sequence (around 80 queries were issued per user). Since the attention layer calculates weights that despite the time span, some old behaviors might be assigned undeserved weights in the final profile and destroy the sequential information. One potential replacement of attention model is to apply Reinforcement Learning to discretely select the search sessions and queries, which successfully worked on sentence classification task with noise data [13].

The overall performances provide us a general evaluation on these methods. **Observing the significantly improvement generated by our model, it is safe to draw conclusion that hierarchical recurrent neural network is able to model the sequential information and build better user profiles than traditional methods.** To further analyze the function of attention model introduced in our framework, it is necessary to find out on what queries the model applies a higher weight. In the next section, we give a real example and visualize the weights.

5.2 Visualization of Weights Assigned by Attention Model

While personalizing search results in a specific situation, it is unnecessary for search engine to use the whole search history. Depending on the input query, some previous information is useless. Blindly using all data could bring noise into user profiles. Intuitively, sessions that have more similar intents should contribute more in building the current user profile. In this paper, we implement a query-aware attention mechanism to highlight different past interactions on the basis of input query and generate user profile dynamically. Table 2 shows that using this kind of dynamic user

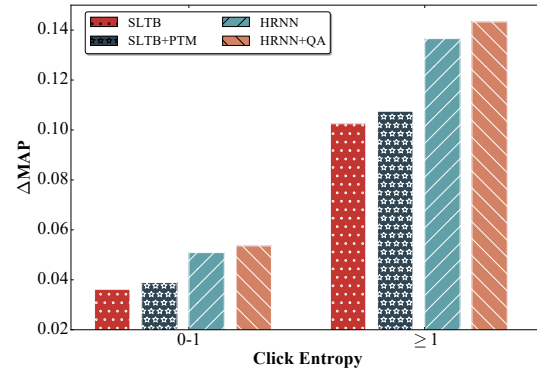


Figure 3: The improvements over original ranking on queries with different click entropies.

profiles improves personalization. To further analyze the influence of attention mechanism, we sample a user from query logs who has relatively rich search history, and visualize the weights applied on different sessions. To make it clearer, we represent the intent of each past session with a typical query from it and sum up the weights of different sessions with the same intent. We select five input queries in the test data and remove the past sessions that have weights less than 0.01 for any of the queries.

As shown in Figure 2, we find that the long-term interest vector attends to past sessions that are more relevant to input queries. In general, the attention can filter the irrelevant information in the query logs. For example, when a user issues a query "virginia welfare lottery ban", the past session containing "virginia lottery" gains the largest weight and the other sessions obtain lower weights. We also find that the past sessions for "facebook" are weighted highly in all five input queries, which is contrary to our intuition. One possible reason is that queries on "facebook" have been issued frequently in the past, so the aggregation of these different past sessions leads to a higher weight for "facebook" sessions. **Input queries such as "times magazine" and "weather forest" whose topics did not occur in the search history obtain weights that do not vary too much. The visualization of weights on past sessions shows that the attention model is able to distinguish the relatively important sessions in the user query logs.**

In order to further verify the effects of our model, it is worthwhile to analyze on what kinds of queries that our methods could achieve larger improvements. In the following sections, we focus on comparing our frameworks with two SLTB baseline models. We use ΔMAP , the change of MAP over original ranking, as the main metric to describe results.

5.3 Performance on Navigational and Non-navigational Queries

A larger click entropy often indicates a higher potential for personalization, as it represents a larger probability of diverse user intents [11, 36]. To study this problem, we firstly group the queries with cutoff of click entropy at 1.0. Teevan et al. [37] used it as an indicator of distinguishing general navigational queries. Then we compare the improvements made by different models on two groups of queries.

As seen in Figure 3, all personalization models produce improved rankings on both groups of queries, and the average improvements

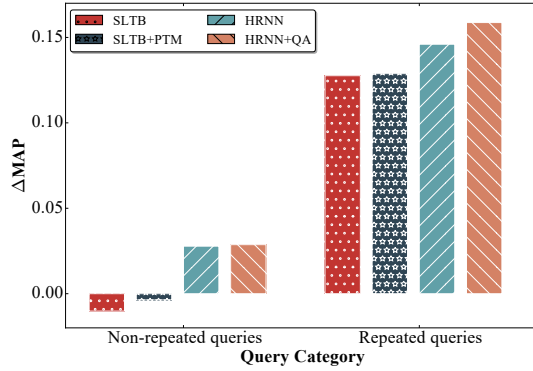


Figure 4: The improvements over original ranking on repeated queries and non-repeated queries.

on non-navigational queries (click entropy ≥ 1) are much larger than on navigational queries (click entropy < 1). In general, our model based on HRNN+QA outperforms any other methods and all of our models outperform the baseline methods ($p < 0.01$). Specifically, as for navigational queries, the Δ MAP of our framework based on HRNN+QA is 0.0124 higher than the best baseline model's (SLTB+PTM). Such improvement further increases to 0.0337 when click entropy is larger than 1.0. Also, all of our two models have improvements of about 0.03 on the non-navigational queries. These results confirm that our framework works better on non-navigational queries than navigational queries.

5.4 Performance on Repeated and Non-repeated Queries

In this experiment, we categorize test queries into two groups: repeated queries and non-repeated queries. A repeated query is a query that the user has issued in the past, and this indicates that the user probably issues the current query to re-find the same information. The click-based personalization methods we discuss in Section 2 highly depend on the information extracted from these repeated queries. If a model never saw the query in the past, most click-based features will be disabled. Therefore, it is worthwhile to evaluate our framework on these non-repeated queries when click based features do not work.

As shown in Figure 4, all methods significantly improve the original results on repeated queries, while our models also outperform the baseline models. The MAP of ranking generated by HRNN+QA on these queries is around 0.03 larger than by SLTB+PTM. With the topical information, the model can better promote the clicked results. Note that personalization on the non-repeated queries is very challenging when baseline models fail improving the search results. On the contrary, our framework successfully generalizes to the queries that are never seen before. Due to the lack of click information for non-repeated queries, only topical features can play a role in tailoring the search results. This observation shows that it is insufficient to simply use an aggregation of clicked documents, and it demonstrates the effectiveness of our model in learning a better user profile. In addition, it shows that applying attention model or not does not make a big difference on non-repeated queries. The majority of the improvements achieved by attention model occurs on repeated queries, which means attention model further highlight the very similar previous queries and reduce the noise.

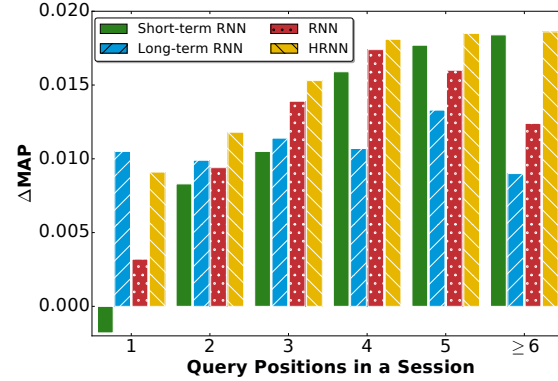


Figure 5: The improvements of models vs. positions of query in sessions.

5.5 Performances on Queries at Different Positions in Sessions

In Table 2, we find that hierarchical structure achieved the best compared with baseline models, but the roles of long- and short-term interest vectors in our framework are still unclear. In this experiment we analyze the performances of each part on queries at different positions in a session. To avoid the influence of additional features, we focus our experiment on non-repeated queries, where click-based features are no more useful. Specifically, we train a plain RNN model whose hierarchical RNN layer is replaced by a one-layer RNN. The model unpacks the session-segmented inputs of HRNN and feeds the whole historical data to the RNN to calculate evolving latent state vectors. In addition, we set a short-term RNN model by disabling the high-level RNN and using the short-term interest vector only, and also a long-term RNN model in the similar way. Note that the long-term RNN model still depends on the low-level RNN to modeling past interests. We analyze the performances of above frameworks and compare them with the complete HRNN.

From Figure 5, we find that short-term RNN steadily improves the results more with the increasing of query positions, because more information is available at a higher position during the session. This observation is consistent with Bennett et al. [4]. In contrast, long-term RNN in our model generates relatively stable benefits on queries at different positions. On the queries at first position in each session, short-term RNN fails to improve the search results, since no session information is available. In that case, the zero short-term vector is used, which turns out to be ineffective. **By analyzing the queries at first positions, we find that hierarchical structure is better in modeling long-term dependency, and it produces improvements two times larger than plain structure.**

6 CONCLUSION

Large amount of data is generated every single day alongside the interactions between users and search engines. Previous studies have demonstrated that by extracting diverse features from large-scale query logs, search engines are able to tailor the original ranking to satisfy individual users. However, none of these studies successfully exploit the sequential information contained among queries and sessions. In this paper, we propose to apply a deep learning framework to solve this problem. Specifically, we deploy a hierarchical recurrent neural network with query-aware attention mechanism

to dynamically generate the topical user profiles. The evaluation on the query logs from a commercial search engine shows that our framework significantly outperforms the existing non-deep learning methods. In contrast to previous studies which assigned a timespan-based decay on historical data or ignored the sequential information, our model is able to automatically decide what to keep in the user profiles and consequently generate better user profiles. Our experiments also confirm that the query-aware attention model is able to highlight the important parts from previous queries and sessions. To further improve our framework, we can replace the attention model with a Reinforcement Learning model, leveraging its ability to discretely select past data. Besides, how to leverage documents skipped by users is also worth further studying.

ACKNOWLEDGMENTS

Zhicheng Dou is the corresponding author. This work was supported by National Key R&D Program of China No. 2018YFC0830703, National Natural Science Foundation of China No. 61872370 and No. 61502501, and the Beijing Natural Science Foundation No. 4162032.

REFERENCES

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [2] Paul N Bennett, Filip Radlinski, Ryan W White, and Emine Yilmaz. 2011. Inferring and using location metadata to personalize web search. In *Proceedings of the SIGIR'2011*. ACM, 135–144.
- [3] Paul N Bennett, Krysta Svore, and Susan T Dumais. 2010. Classification-enhanced ranking. In *Proceedings of the WWW'2010*. ACM, 111–120.
- [4] Paul N Bennett, Ryan W White, Wei Chu, Susan T Dumais, Peter Bailey, Fedor Borisov, and Xiaoyuan Cui. 2012. Modeling the impact of short- and long-term behavior on search personalization. In *Proceedings of the SIGIR'2012*. ACM, 185–194.
- [5] Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. 2005. Learning to rank using gradient descent. In *Proceedings of ICML'2005*. ACM, 89–96.
- [6] Fei Cai, Shangsong Liang, and Maarten De Rijke. 2014. Personalized document re-ranking based on bayesian probabilistic matrix factorization. In *Proceedings of the SIGIR'2014*. ACM, 835–838.
- [7] Mark J. Carman, Fabio Crestani, Morgan Harvey, and Mark Baillie. 2010. Towards query log based personalization using topic models. In *Proceedings of the CIKM'2010*. 1849–1852.
- [8] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In *EMNLP'2014*. 1724–1734.
- [9] Kevyn Collins-Thompson, Paul N Bennett, Ryan W White, Sebastian De La Chica, and David Sontag. 2011. Personalizing web search results by reading level. In *Proceedings of the CIKM'2011*. ACM, 403–412.
- [10] Nick Craswell, Onno Zoeter, Michael J. Taylor, and Bill Ramsey. 2008. An experimental comparison of click position-bias models. In *WSDM'2008*.
- [11] Zhicheng Dou, Ruihua Song, and Ji-Rong Wen. 2007. A large-scale evaluation and analysis of personalized search strategies. In *WWW'2007*. ACM, 581–590.
- [12] Salah El Hahi and Yoshua Bengio. 1996. Hierarchical recurrent neural networks for long-term dependencies. In *NIPS'1996*. 493–499.
- [13] Jun Feng, Minlie Huang, Li Zhao, Yang Yang, and Xiaoyan Zhu. 2018. Reinforcement Learning for Relation Classification From Noisy Data. In *AAAI'2018*.
- [14] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. A deep relevance matching model for ad-hoc retrieval. In *CIKM'2016*. ACM, 55–64.
- [15] Aniko Hannak, Piotr Sapiezynski, Arash Molavi Kakhki, Balachander Krishnamurthy, David Lazer, Alan Mislove, and Christo Wilson. 2013. Measuring personalization of web search. In *WWW'2013*. ACM, 527–538.
- [16] Morgan Harvey, Fabio Crestani, and Mark J Carman. 2013. Building user profiles from topic models for personalised search. In *CIKM'2013*. ACM, 2309–2314.
- [17] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [18] 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'2013*. ACM, 2333–2338.
- [19] Thorsten Joachims. 2002. Optimizing search engines using clickthrough data. In *SIGKDD'2002*. ACM, 133–142.
- [20] Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, and Geri Gay. 2005. Accurately interpreting clickthrough data as implicit feedback. In *SIGIR'2005*. 154–161.
- [21] Lin Li, Zhenglu Yang, Botao Wang, and Masaru Kitsuregawa. 2007. Dynamic adaptation strategies for long-term and short-term user profile to personalize search. *Advances in Data and Web Management* (2007), 228–240.
- [22] Xiujun Li, Chenlei Guo, Wei Chu, Ye-Yi Wang, and Jude Shavlik. 2014. Deep learning powered in-session contextual ranking using clickthrough data. In *NIPS'2014*.
- [23] Nicolaas Matthijs and Filip Radlinski. 2011. Personalizing web search using long term browsing history. In *WSDM'2011*. ACM, 25–34.
- [24] Bhaskar Mitra and Nick Craswell. 2017. Neural Models for Information Retrieval. *arXiv preprint arXiv:1705.01509* (2017).
- [25] Kezban Dilek Onal, Ye Zhang, Ismail Sengor Altıngövdü, Md Mustafizur Rahman, Pinar Karagoz, Alex Braylan, Brandon Dang, Heng-Lu Chang, Henna Kim, Quinten McNamara, and others. 2018. Neural information retrieval: At the end of the early years. *Information Retrieval Journal* 21, 2-3 (2018), 111–182.
- [26] Hamid Palangi, Li Deng, Yelong Shen, Jianfeng Gao, Xiaodong He, Jianshu Chen, Xinying Song, and Rabab Ward. 2016. Deep sentence embedding using long short-term memory networks: Analysis and application to information retrieval. *TASLP* 24, 4 (2016), 694–707.
- [27] Massimo Quadrana, Alexandros Karatzoglou, Balázs Hidasi, and Paolo Cremonesi. 2017. Personalizing Session-based Recommendations with Hierarchical Recurrent Neural Networks. In *RecSys'2017*. 130–137.
- [28] Iulian Vlad Serban, Alessandro Sordoni, Yoshua Bengio, Aaron C Courville, and Joelle Pineau. 2016. Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models. In *AAAI'2016*. 3776–3784.
- [29] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to rank short text pairs with convolutional deep neural networks. In *SIGIR'2015*. ACM, 373–382.
- [30] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. A latent semantic model with convolutional-pooling structure for information retrieval. In *CIKM'2014*. ACM, 101–110.
- [31] Ahu Sieg, Bamshad Mobasher, and Robin Burke. 2007. Web search personalization with ontological user profiles. In *CIKM'2007*. ACM, 525–534.
- [32] Yang Song, Hongning Wang, and Xiaodong He. 2014. Adapting deep ranknet for personalized search. In *WSDM'2014*. ACM, 83–92.
- [33] David Sontag, Kevyn Collins-Thompson, Paul N Bennett, Ryan W White, Susan Dumais, and Bodo Billerbeck. 2012. Probabilistic models for personalizing web search. In *WSDM'2012*. ACM, 433–442.
- [34] Alessandro Sordoni, Yoshua Bengio, Hossein Vahabi, Christina Lioma, Jakob Grue Simonsen, and Jian-Yun Nie. 2015. A hierarchical recurrent encoder-decoder for generative context-aware query suggestion. In *CIKM'2015*. ACM, 553–562.
- [35] Mirco Speretta and Susan Gauch. 2005. Personalized search based on user search histories. In *Web Intelligence, 2005. Proceedings. The 2005 IEEE/WIC/ACM International Conference on*. IEEE, 622–628.
- [36] Jaime Teevan, Susan T Dumais, and Daniel J Liebling. 2008. To personalize or not to personalize: modeling queries with variation in user intent. In *SIGIR'2008*. ACM, 163–170.
- [37] Jaime Teevan, Daniel J Liebling, and Gayathri Ravichandran Geetha. 2011. Understanding and predicting personal navigation. In *WSDM'2011*. ACM, 85–94.
- [38] Maksims Volkovs. 2015. Context models for web search personalization. *arXiv preprint arXiv:1502.00527* (2015).
- [39] Thanh Vu, Dat Quoc Nguyen, Mark Johnson, Dawei Song, and Alistair Willis. 2017. Search personalization with embeddings. In *ECIR'2017*. Springer, 598–604.
- [40] Thanh Vu, Dawei Song, Alistair Willis, Son Ngoc Tran, and Jingfei Li. 2014. Improving search personalisation with dynamic group formation. In *SIGIR'2014*. 951–954.
- [41] Thanh Vu, Alistair Willis, Son N Tran, and Dawei Song. 2015. Temporal latent topic user profiles for search personalisation. In *ECIR'2015*. Springer, 605–616.
- [42] Hongning Wang, Xiaodong He, Ming Wei Chang, Yang Song, Ryan W. White, and Wei Chu. 2013. Personalized ranking model adaptation for web search. In *SIGIR'2013*. 323–332.
- [43] White, W Ryan, Bennett, N Paul, Dumais, and T Susan. 2010. Predicting short-term interests using activity-based search context. (2010), 1009–1018.
- [44] Ryan W White, Peter Bailey, and Liwei Chen. 2009. Predicting user interests from contextual information. In *SIGIR'2009*. ACM, 363–370.
- [45] Ryan W White, Wei Chu, Ahmed Hassan, Xiaodong He, Yang Song, and Hongning Wang. 2013. Enhancing personalized search by mining and modeling task behavior. In *WWW'2013*. ACM, 1411–1420.
- [46] Qiang Wu, Chris JC Burges, Krysta M Svore, and Jianfeng Gao. 2008. *Ranking, boosting, and model adaptation*. Technical Report. Technical report, Microsoft Research.
- [47] Chenyan Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan Liu, and Russell Power. 2017. End-to-End Neural Ad-hoc Ranking with Kernel Pooling. In *SIGIR'2017*.