

Integrating Lexical and Temporal Signals in Neural Ranking Models for Searching Social Media Streams

Jinfeng Rao,^{1,3} Hua He,¹ Haotian Zhang,² Ferhan Ture,³
Royal Sequiera,² Salman Mohammed,² and Jimmy Lin²

¹ Department of Computer Science, University of Maryland

² David R. Cheriton School of Computer Science, University of Waterloo

³ Comcast Applied AI Research Labs

ABSTRACT

Time is an important relevance signal when searching streams of social media posts. The distribution of document timestamps from the results of an initial query can be leveraged to infer the distribution of relevant documents, which can then be used to rerank the initial results. Previous experiments have shown that kernel density estimation is a simple yet effective implementation of this idea. This paper explores an alternative approach to mining temporal signals with recurrent neural networks. Our intuition is that neural networks provide a more expressive framework to capture the temporal coherence of neighboring documents in time. To our knowledge, we are the first to integrate lexical and temporal signals in an end-to-end neural network architecture, in which existing neural ranking models are used to generate query-document similarity vectors that feed into a bidirectional LSTM layer for temporal modeling. Our results are mixed: existing neural models for document ranking alone yield limited improvements over simple baselines, but the integration of lexical and temporal signals yield significant improvements over competitive temporal baselines.

1 INTRODUCTION

There is a large body of literature in information retrieval that has established the importance of modeling the temporal characteristics of documents as well as queries for various information seeking tasks [4–9, 18, 24, 25]. Such techniques are particularly important for searching real-time social media streams such as Twitter, which rapidly evolves in reaction to real-world events. In this paper, we tackle the problem of retrospective *ad hoc* retrieval over a collection of short, temporally-ordered social media posts (tweets). Given information needs expressed as queries, we aim to build systems that return high-quality ranked lists of relevant tweets.

We are motivated by Efron et al.’s temporal cluster hypothesis [8], which stipulates that in search tasks where time plays an important role (such as ours), relevant documents tend to cluster together in time, and that this property can be exploited to improve search effectiveness. Efron et al. take advantage of kernel density estimation (KDE) to infer the temporal distribution of relevant documents from an initial search; the inferred distribution is then used to rerank the original documents. Experiments show that this approach is simple yet effective [8, 29].

In this paper, we take the KDE technique as a baseline and explore an alternative approach for temporal modeling using recurrent neural networks. Such models have been successfully applied to many sequence learning tasks in natural language processing where the modeling units are temporally dependent (e.g., tagging and parsing). We draw a connection between the temporal clustering of documents, where the relevance of one document may affect its neighbors, to a sequence learning task, and explore the hypothesis that recurrent neural networks provide a rich, expressive modeling framework to capture such temporal signals. To this end, we build a unified neural network model to integrate lexical and temporal relevance signals, and we examine the effectiveness of several existing neural rankings models that consider only query-document textual similarity. We wondered how they would fare in the context of noisy social media posts.

Contributions. We view this work as having two contributions:

- (1) We examined the effectiveness of several existing neural ranking models on standard tweet test collections. Results show that, in considering only lexical signals, they yield limited improvements over simple baselines, suggesting that social media search presents a different set of challenges compared to traditional *ad hoc* retrieval (e.g., over newswire documents).
- (2) We present, to our knowledge, the first end-to-end neural network architecture that integrates lexical and temporal signals. Using the best lexical modeling component (from above), we are able to obtain significant improvements over competitive temporal baselines on standard tweet test collections.

2 RELATED WORK

2.1 Temporal Information Retrieval

There is a long thread of research exploring the role of temporal signals in search [4–9, 18, 30, 32], and it is well established that for certain tasks, better modeling of the temporal characteristics of queries and documents can lead to higher retrieval effectiveness.

For example, Jones and Diaz [16] studied the temporal profiles of queries, classifying queries as atemporal, temporally ambiguous, or temporally unambiguous. They showed that the temporal distribution of retrieved documents can provide an additional source of evidence to improve rankings. Building on this, Li and Croft [18] introduced recency priors that favor more-recent documents. Dakka et al. [4] proposed an approach to temporal modeling based on moving windows to integrate query-specific temporal evidence with lexical evidence. Efron et al. [7] presented several language modeling variants that incorporate query-specific temporal evidence. The most direct point of comparison to our work (as discussed in

the introduction) is the use of non-parametric density estimation to infer the temporal distribution of relevant documents from an initial list of retrieved documents [8, 29]. Most recently, Rao et al. [32] proposed alternative models that attempt to make such predictions directly from query term statistics, obviating the need for an initial retrieval stage.

There have been several other studies of time-based pseudo relevance feedback. Keikha et al. [17] represented queries and documents with their normalized term frequencies in the time dimension and used a time-based similarity metric to measure relevance. Craveiro et al. [3] exploited the temporal relationship between words for query expansion. Choi and Croft [2] presented a method to select time periods for expansion based on users’ behaviors (i.e., retweets). Rao et al. [28] proposed a continuous hidden Markov model to identify temporal burst states in order to select better query expansion terms.

In addition to ranking, modeling temporal signals has also been shown to benefit related tasks such as behavior prediction [24, 25], time-sensitive query auto-completion [35], and real-time query suggestion [19]. For example, Radinsky et al. [24, 25] built predictive models to learn query dynamics from historical user data.

2.2 Neural Information Retrieval

Following great successes in computer vision, speech recognition, and natural language processing, we have recently seen a new wave of research applying neural networks to information retrieval. Huang et al. [15] proposed a technique called Deep Structured Semantic Modeling (DSSM), which has led to follow-on work [34, 37]. The basic idea is to use a feedforward function to learn low-dimensional vector representations of queries and documents, aiming to capture latent semantic information in texts. Recently, Guo et al. [10] proposed a deep relevance matching model for *ad hoc* retrieval, pointing out differences between search and many NLP problems. Mitra et al. [20] presented a neural matching model to combine local and global interactions between queries and documents. There are many other applications of neural networks to information retrieval, for example, relevance-based word embeddings [39], voice search with hierarchical recurrent neural networks [31], reinforcement learning for query reformulation [21], and generative adversarial training for retrieval models [38].

On a slightly different thread, there has been work on modeling textual similarity between short text pairs. Severyn and Moschitti [33] proposed a convolutional neural network (CNN) for exactly this, which was further expanded and analyzed by Rao et al. [27]. He et al. [11] proposed an ensemble approach of CNNs that take advantage of different types of convolutional feature maps, pooling methods, and window sizes to capture sentence pair similarity from multiple perspectives. Rao et al. [26] extended this line of work by studying different negative sampling strategies in a pairwise ranking framework, which obtains state-of-the-art accuracy on a standard question answering benchmark dataset.

3 APPROACH

We present a neural network architecture that integrates lexical and temporal signals, shown in Figure 1. The overall architecture consists of distinct components for lexical modeling, to capture

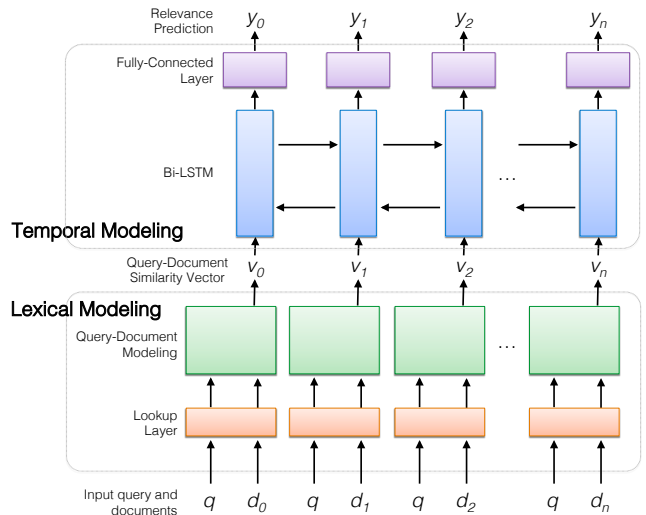


Figure 1: Our neural network architecture that integrates lexical and temporal signals. The lexical modeling component can be viewed as a black box for producing query–document similarity vectors. A temporally-ordered sequence of these vectors feed into our bidirectional LSTM for temporal modeling.

query–document similarity, and temporal modeling, to capture relevance signals contained in the temporal sequencing of documents. The two components are independent and in particular we can view the lexical modeling component as a black box, allowing us to explore different architectures. However, as we explain later, the entire model is trained end-to-end in a two-stage process.

Lexical Modeling. The architecture for the lexical modeling component is shown in the lower half of Figure 1, where each “slice” of the network is identical (i.e., with shared parameters). Each instance of the model takes as input a query and a document to generate a query–document similarity vector v . This is accomplished by translating an input sequence of tokens (either the query or the document) into a sequence of distributional vectors $[w_1, w_2, \dots, w_{|S|}]$, where $|S|$ is the length of the token sequence, from a word embedding lookup layer. The resulting matrix then feeds into a neural network. At a high level, this similarity model can be viewed as a black box, but we describe several instantiations below.

Temporal Modeling. The architecture of the temporal modeling component is shown in the upper half of Figure 1. We use a bidirectional LSTM where the inputs are the query–document similarity vectors from the lexical modeling component, sorted in time order. That is, documents from the training set are temporally ordered, and the lexical modeling component is applied to the query paired with each individual document to yield a collection of query–document similarity vectors $\{v_0, v_1, \dots, v_n\}$. The output of the bidirectional LSTM feeds into a fully-connected layer plus softmax to yield a prediction of document relevance y . Note that each instance of the fully-connected layer and softmax share parameters.

In what follows, we describe each of the components in detail.

3.1 Lexical Modeling Component

In this work, we considered three existing approaches to generating query–document similarity vectors. All three adopt what is commonly known as a “Siamese” structure [1], with two subnetworks processing the query and document in parallel, yielding a “joined” representation that feeds into a relevance modeling component:

DSSM [15]: The Deep Structured Semantic Model (DSSM) is an early application of neural networks to web search. One of its key features is a word hashing layer that converts all tokens into trigrams, which greatly reduces the size of the vocabulary space to help handle misspellings and other noisy text input. In parallel, the dense hashed features from either the query or the document feed into a multi-layer perceptron with a softmax on top to make the final relevance prediction. We take the intermediate semantic representation of the query and document, just before the softmax, as our query–document similarity vector.

SM [33]: The convolutional neural network (CNN) proposed by Severyn and Moschitti [33] has been previously applied to question answering as well as tweet reranking. In both the query and document subnetworks, convolutional feature maps are applied to the input embedding matrix, followed by ReLU activation and simple max-pooling, to arrive at a representation vector x_q for the query and x_d for the document. Intermediate representations are concatenated into a single vector at the join layer:

$$x_{\text{join}} = [x_q^T; x_{\text{sim}}; x_d^T; x_{\text{feat}}^T] \quad (1)$$

where x_{sim} defines the *bilinear* similarity between x_q and x_d . The final component consists of “extra features” x_{feat} derived from four word overlap measures between the query and the document.

In the original SM model, the join vector feeds into a fully-connected layer and softmax for final relevance prediction, but in our approach we use the join vector x_{join} as the query–document similarity vector.

Multi-Perspective CNN [11]: This approach was developed at roughly the same time as the SM model and can be described as an ensemble of convolutional neural networks. The “multi-perspective” idea refers to different types of convolutional feature maps, pooling methods, and window sizes to capture semantic similarity between textual inputs. Another key feature is a similarity measurement layer to explore the interactions between the learned convolutional feature maps at different levels of granularity. At the time the work was published, it achieved state-of-the-art effectiveness on several semantic modeling tasks such as paraphrase detection and question answering (although other models have improved upon it since).

As with the SM model, we take the joined representation just before the fully-connected layer and softmax as the query–document similarity vector.

3.2 Temporal Modeling Component

On top of a sequence of temporally ordered query–document similarity vectors (the output of the lexical modeling component), we layer a recurrent neural network to capture the temporal clustering of relevant documents (see Figure 1). Compared to kernel density estimation, we hypothesized that recurrent neural networks provide a richer, more expressive modeling framework to capture temporal signals that can yield more effective results.

For our task, we used a variant of recurrent neural networks, bidirectional LSTM (Bi-LSTM) [14], which have been successfully applied to text similarity tasks [12, 13]. One key feature of LSTMs is their ability to capture long-range dependencies, and a bidirectional LSTM consists of two LSTMs that run in parallel in opposite directions: one (forward LSTM^f) on the input sequence and the other (backward LSTM^b) on the reverse of the sequence. At time step t , the Bi-LSTM hidden state h_t^{bi} is a concatenation of the hidden state h_t^{for} of LSTM^f and the hidden state h_t^{back} of LSTM^b, representing the neighboring contexts of input v_t in the temporal sequence.

Given Bi-LSTM output h_t^{bi} , the prediction output y_t of our temporal ranking model at time step t is obtained by passing the Bi-LSTM output through a fully-connected layer and softmax as follows:

$$g_t = \sigma(W^m \cdot h_t^{\text{bi}} + b^m) \quad (2)$$

$$y_t = \text{softmax}(W^p \cdot g_t + b^p) \quad (3)$$

where the output y_t indicates the relevance of the document at time step t . W^* and b^* are learned weight matrices and biases.

3.3 Model Training

Although our neural network architecture breaks down into two distinct components, we train the entire model end-to-end in a two-stage manner, with stochastic gradient descent to minimize negative log-likelihood loss of the entire model. In each epoch, we first train the lexical modeling component independently, and then use the results to generate inputs to the temporal modeling layer. The losses from all documents are summed together to train the Bi-LSTM and the top layers, while the underlying lexical component is held constant. The reason for this two-stage approach is to restrict the search space during model optimization, since we have limited labeled data for training.

At inference time, we first retrieve candidate documents from the collection using a standard ranking function. These documents are then ordered chronologically and fed into the model. The classification scores outputted by each step of the Bi-LSTM (corresponding to the processing of that document) are used to resort the ranked list, which we take as final output for evaluation.

4 EXPERIMENTS

4.1 Experimental Setup

We evaluated our proposed models on Twitter test collections from the TREC 2011 and 2012 Microblog Tracks (49 topics and 59 topics, respectively). Both use the Tweets2011 collection, which consists of an approximately 1% sample (after some spam removal) of tweets from January 23, 2011 to February 7, 2011 (inclusive), totaling approximately 16 million tweets. Relevance judgments were made on a 3-point scale (“not relevant”, “relevant”, “highly relevant”), but in this work we treated both higher grades as relevant. We removed all the retweets in our experiments since they are by definition not relevant according to the assessment guidelines.

To rule out the effects of different preprocessing strategies during collection preparation (i.e., stemming, stopword removal, etc.), we used the open-source implementations of tweet search provided by the TREC Microblog API¹ to retrieve up to 1000 tweets per topic

¹<https://github.com/lintool/twitter-tools>

ID	Method	P15	P30	P100	AP
1	Query Likelihood (QL) [23]	0.381	0.329	0.234	0.200
Temporal Baselines					
2	uniform	0.366	0.326	0.243	0.203
3	KDE [8] score-based	0.383	0.334	0.244	0.203
4	rank-based	0.387	0.337	0.244	0.204
5	oracle	0.411 ^{1,4}	0.389 ^{1,4}	0.260 ^{1,4}	0.229 ^{1,4}
Neural Ranking Approaches					
6	DSSM [15]	0.187	0.168	0.153	0.102
7	SM [33]	0.203	0.188	0.170	0.116
8	Multi-Perspective CNN [11]	0.401 ¹	0.356 ¹	0.252 ¹	0.197
Neural Ranking + Temporal Modeling					
9	SM [33] + Temporal	0.222	0.196	0.169	0.116
10	Multi-Perspective CNN [11] + Temporal	0.418 ^{1,4,8}	0.366 ^{1,4}	0.257 ^{1,4,8}	0.203 ^{1,8}

Table 1: Results from the TREC 2011/2012 Microblog Track test collections, using TREC 2011 data for training and TREC 2012 data for evaluation. Superscripts indicate the row indexes from which the metric difference is statistically significant ($p < 0.05$) using Fisher’s two-sided, paired randomization test.

using query likelihood (QL) for scoring. These initial results were then reranked using our proposed models. For effectiveness, we measured average precision (AP) and precision at 15, 30, and 100 (P15, P30, and P100); note that P30 was the official metric used in the TREC Microblog Tracks. Since all the models required training, we used the TREC 2011 topics for training and the TREC 2012 topics for evaluation. Additionally, we randomly selected 5% of query-document pairs from the training set as the development set; those selected samples were removed from the training set.

We considered several lexical and temporal baselines to evaluate our models. The standard query likelihood (QL) approach [23] was used as the lexical baseline. We used the kernel density estimation techniques of Efron et al. [8] as our temporal baseline (with the implementation from Rao et al. [29]). They proposed four different weighting schemes to estimate feedback parameters: uniform, score-based, rank-based, and oracle. The first three take advantage of document timestamp distributions from an initial retrieval, while the oracle method requires actual human relevance judgments. The oracle is useful to illustrate upper bound effectiveness for KDE-based techniques.

The neural ranking approaches were implemented using the Torch deep learning toolkit (in Lua). For the SM model² [33] and the multi-perspective CNN³ [11], we took advantage of existing open-source implementations; DSSM is our own re-implementation. We used existing 300-dimensional GloVe [22] word embeddings to encode each word, which was trained on 840 billion tokens and freely available. The vocabulary size of our dataset is 90.3K, with around 37% words not found in the GloVe word embeddings. Unknown words were randomly initialized with values uniformly sampled from $[-0.05, 0.05]$. During training, we used stochastic gradient descent together with RMS-PROP to iteratively update the model. The output size of the Bi-LSTM layer is 400 and the hidden layer size is 150. The learning rate was initially set to 0.001, and then decreased by a factor of three when the development set

loss stopped decreasing for three epochs. The maximum number of training epochs was 25.

4.2 Experimental Results

Table 1 shows our experimental results, with each row representing an experimental condition (numbered for convenience). For each method, we performed significance testing against the lexical baseline (QL) and the best-performing temporal KDE model (rank-based). In addition, we tested the significance of differences between each pair of lexical-only model vs. lexical + temporal model. In all cases, we used Fisher’s two-sided, paired randomization test [36]. Superscripts indicate the row indexes for which the metric difference is statistically significant ($p < 0.05$).

From the block in Table 1 labeled “Temporal Baselines”, we see that the KDE approaches (with the exception of the oracle condition) yield limited improvements over the QL baseline.⁴ Looking at the block of Table 1 labeled “Neural Ranking Approaches”, we find that the SM model and DSSM do not appear to be as effective as the multi-perspective CNN; in particular, the first two models actually perform worse than the simple QL baseline.

In Table 1, under “Neural Ranking + Temporal Modeling”, we report results from combining the SM model and the multi-perspective CNN with our Bi-LSTM temporal model. In the first case, the improvement is minor over the SM model alone, but with the multi-perspective CNN, the addition of a temporal layer yields significant improvements over the multi-perspective CNN alone (condition 8) and also rank-based KDE (condition 4). We also note that the multi-perspective CNN + Bi-LSTM model approaches the effectiveness of the oracle KDE condition (and in the case of P15, exceeds it, albeit not significantly). This suggests that neural networks offer an expressive framework for integrating lexical and temporal signals, potentially beyond what is available to non-parametric density estimation techniques alone, even with oracle input.

²<https://github.com/castorini/SM-CNN-Torch>

³<https://github.com/castorini/MP-CNN-Torch>

⁴These results are consistent with those of Rao et al. [29]; although those experiments affirmed the overall effectiveness of the KDE techniques, results from individual configurations (such as a particular train/test split) may not yield significant improvements.

5 FUTURE WORK AND CONCLUSIONS

While our results are certainly encouraging, there are a number of unresolved issues and open questions; these are avenues for future work. First, we have only experimented on a single collection of tweets, and thus there are questions about the robustness of our results. Second, we have yet to perform detailed error analysis to uncover the differences between the three neural ranking models we examined, and thus have not answered the *why* questions: For example, what characteristics of the multi-perspective CNN allow it to serve as an effective ranker while the SM model and DSSM do not appear to work? As a start, a topic-by-topic analysis might uncover some insights. Third, it is interesting to note that our approaches improve early precision more than average precision: it is unclear if this is due to inherent properties of our model, our reranking setup, or some other reason.

To conclude, we believe that this work is most valuable in providing a general architecture for integrating lexical and temporal signals for information seeking on time-ordered documents. We have already shown that different lexical modeling components can be “plugged in”—our experiments examined three neural network models, but more can be straightforwardly explored. Similarly, we can imagine different temporal models beyond the Bi-LSTM approach proposed here. In addition, we have shown that our combined lexical and temporal models can be trained end to end, which yields an integrated, flexible, and expressive ranking framework.

6 ACKNOWLEDGMENTS

This work was supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada, with additional contributions from the U.S. National Science Foundation under CNS-1405688. Any findings, conclusions, or recommendations expressed do not necessarily reflect the views of the sponsors.

REFERENCES

- [1] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. 1993. Signature Verification Using a “Siamese” Time Delay Neural Network. In *NIPS*. 737–744.
- [2] Jaeho Choi and W. Bruce Croft. 2012. Temporal Models for Microblogs. In *CIKM*. 2491–2494.
- [3] Olga Craveiro, Joaquim Macedo, and Henrique Madeira. 2014. Query Expansion with Temporal Segmented Texts. In *ECIR*. 612–617.
- [4] Wisam Dakka, Luis Gravano, and Panagiotis G. Ipeirotis. 2012. Answering General Time-Sensitive Queries. *TKDE* 24, 2 (2012), 220–235.
- [5] Anlei Dong, Yi Chang, Zhaohui Zheng, Gilad Mishne, Jing Bai, Ruiqiang Zhang, Karolina Buchner, Ciya Liao, and Fernando Diaz. 2010. Towards Recency Ranking in Web Search. In *WSDM*. 11–20.
- [6] Anlei Dong, Ruiqiang Zhang, Pranam Kolari, Jing Bai, Fernando Diaz, Yi Chang, Zhaohui Zheng, and Hongyuan Zha. 2010. Time is of the Essence: Improving Recency Ranking Using Twitter Data. In *WWW*. 331–340.
- [7] Miles Efron and Gene Golovchinsky. 2011. Estimation Methods for Ranking Recent Information. In *SIGIR*. 495–504.
- [8] Miles Efron, Jimmy Lin, Jiyin He, and Arjen de Vries. 2014. Temporal Feedback for Tweet Search with Non-Parametric Density Estimation. In *SIGIR*. 33–42.
- [9] Jonathan L. Elsas and Susan T. Dumais. 2010. Leveraging Temporal Dynamics of Document Content in Relevance Ranking. In *WSDM*. 1–10.
- [10] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-hoc Retrieval. In *CIKM*. 55–64.
- [11] Hua He, Kevin Gimpel, and Jimmy Lin. 2015. Multi-Perspective Sentence Similarity Modeling with Convolutional Neural Networks. In *EMNLP*. 1576–1586.
- [12] Hua He and Jimmy Lin. 2016. Pairwise Word Interaction Modeling with Deep Neural Networks for Semantic Similarity Measurement. In *HLT-NAACL*. 937–948.
- [13] Hua He, John Wieting, Kevin Gimpel, Jinfeng Rao, and Jimmy Lin. 2016. UMD-TTIC-UW at SemEval-2016 Task 1: Attention-Based Multi-Perspective Convolutional Neural Networks for Textual Similarity Measurement. In *SemEval*. 1103–1108.
- [14] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9, 8 (1997), 1735–1780.
- [15] 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*. 2333–2338.
- [16] Rosie Jones and Fernando Diaz. 2007. Temporal Profiles of Queries. *TOIS* 25, 3 (2007), Article 14.
- [17] Mostafa Keikha, Shima Gerani, and Fabio Crestani. 2011. TEMPER: A Temporal Relevance Feedback Method. In *ECIR*. 436–447.
- [18] Xiaoyan Li and W. Bruce Croft. 2003. Time-Based Language Models. In *CIKM*. 469–475.
- [19] Gilad Mishne, Jeff Dalton, Zhenghua Li, Aneesh Sharma, and Jimmy Lin. 2012. Fast Data in the Era of Big Data: Twitter’s Real-Time Related Query Suggestion Architecture. In *SIGMOD*. 1147–1157.
- [20] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. 2017. Learning to Match using Local and Distributed Representations of Text for Web Search. In *WWW*. 1291–1299.
- [21] Rodrigo Nogueira and Kyunghyun Cho. 2017. Task-Oriented Query Reformulation with Reinforcement Learning. *arXiv:1704.04572*.
- [22] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global Vectors for Word Representation. In *EMNLP*. 1532–1543.
- [23] Jay M. Ponte and W. Croft. 1998. A Language Modeling Approach to Information Retrieval. In *SIGIR*. 275–281.
- [24] Kira Radinsky, Krysta Svore, Susan Dumais, Jaime Teevan, Alex Bocharov, and Eric Horvitz. 2012. Modeling and Predicting Behavioral Dynamics on the Web. In *WWW*. 599–608.
- [25] Kira Radinsky, Krysta M. Svore, Susan T. Dumais, Milad Shokouhi, Jaime Teevan, Alex Bocharov, and Eric Horvitz. 2013. Behavioral Dynamics on the Web: Learning, Modeling, and Prediction. *TOIS* 31, 3 (2013), Article 16.
- [26] Jinfeng Rao, Hua He, and Jimmy Lin. 2016. Noise-Contrastive Estimation for Answer Selection with Deep Neural Networks. In *CIKM*. 1913–1916.
- [27] Jinfeng Rao, Hua He, and Jimmy Lin. 2017. Experiments with Convolutional Neural Network Models for Answer Selection. In *SIGIR*.
- [28] Jinfeng Rao and Jimmy Lin. 2016. Temporal Query Expansion Using a Continuous Hidden Markov Model. In *ICITR*. 295–298.
- [29] Jinfeng Rao, Jimmy Lin, and Miles Efron. 2015. Reproducible Experiments on Lexical and Temporal Feedback for Tweet Search. In *ECIR*. 755–767.
- [30] Jinfeng Rao, Xing Niu, and Jimmy Lin. 2016. Compressing and Decoding Term Statistics Time Series. In *ECIR*. 675–681.
- [31] Jinfeng Rao, Ferhan Ture, Hua He, Oliver Jojic, and Jimmy Lin. 2017. Talking to Your TV: Context-Aware Voice Search with Hierarchical Recurrent Neural Networks. *arXiv:1705.04892*.
- [32] Jinfeng Rao, Ferhan Ture, Xing Niu, and Jimmy Lin. 2017. Mining Temporal Statistics of Query Terms for Searching Social Media Posts. In *ICITR*.
- [33] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to Rank Short Text Pairs with Convolutional Deep Neural Networks. In *SIGIR*. 373–382.
- [34] 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*. 101–110.
- [35] Milad Shokouhi and Kira Radinsky. 2012. Time-Sensitive Query Auto-Completion. In *SIGIR*. 601–610.
- [36] Mark D. Smucker, James Allan, and Ben Carterette. 2007. A Comparison of Statistical Significance Tests for Information Retrieval Evaluation. In *CIKM*. 623–632.
- [37] Yang Song, Ali Mamdouh Elkahky, and Xiaodong He. 2016. Multi-Rate Deep Learning for Temporal Recommendation. In *SIGIR*. 909–912.
- [38] Jun Wang, Lantao Yu, Weinan Zhang, Yu Gong, Yinghui Xu, Benyou Wang, Peng Zhang, and Dell Zhang. 2017. IRGAN: A Minimax Game for Unifying Generative and Discriminative Information Retrieval Models. *arXiv:1705.10513*.
- [39] Hamed Zamani and W. Bruce Croft. 2017. Relevance-based Word Embedding. In *SIGIR*.