

Learning to Combine Collection-centric and Document-centric Models for Resource Selection

Krisztian Balog
University of Stavanger
krisztian.balog@uis.no

Abstract: This paper describes our participation in the Federated Web Search track at TREC 2014. Our main focus is on the resource selection task, where we employ a learning-to-rank approach to combine various (instantiations of) resource ranking models. Further, we show that vertical selection can be run on the output from resource selection, and that it directly benefits from the improvements of thereof.

1 Introduction

We describe our participation in the Federated Web Search track at TREC 2014. Specifically, we took part in the *resource selection* and *vertical selection* tasks. For resource selection, our focus was on finding a way to effectively combine two principal strategies, Collection-centric (CC) and Document-centric (DC), we developed in prior work (Balog, 2014). We employ a learning-to-rank approach, where various instantiations of the CC and DC models, using different representations and relevance cutoff values, are used as features. We present our approach and results in Section 2. We base our *vertical selection* runs on the outcomes of resource selection step. Specifically, we use the estimated collection relevance scores as binary judgments, thereby essentially delegating the “selection” problem to the resource ranking component. The method and the results are described in Section 3.

2 Resource selection

In prior work, we presented two approaches to the resource selection task based on generative language modeling techniques (Balog, 2014). According to the Collection-centric (CC) model, each collection is represented as a term distribution, which is estimated from all sampled documents. The second model, Document-centric (DC), first scores individual sampled documents, then considers the top-K ranked ones to determine collection relevance. Despite its relative simplicity, the DC model delivers solid performance; at TREC 2013 it came very close to the top performing runs on

all metrics (Demeester et al., 2014). We also experimented with the combination of the CC and DC strategies in our participation last year, using a linear mixture model, but it did not improve over the DC model. This year our aim is to find a way to effectively combine the CC and DC models. To this end, we employ learning-to-rank techniques.

2.1 Approach

We use the scores estimated by the CC and DC models as features. Specifically, we consider a number of different configurations, based on the type of document representation (title, snippet, page) and the cutoff value (K , only for the DC model). In the following subsections, we briefly present the CC and DC models; for a more detailed description we refer to Balog (2014). Additionally, we take collection size to be a feature as well (previously, it was incorporated as a prior collection probability). Table 1 lists our features (36 in total).

2.1.1 Collection-centric Model

Drawing on Callan et al. (1995) and Si et al. (2002), this approach treats each collection as a single, large document. Under the language modeling framework, the probability of the collection generating the query is expressed as follows:

$$P(q|c) = \prod_{t \in q} \left\{ (1 - \lambda) \left(\sum_{d \in c} P(t|d)P(d|c) \right) + \lambda P(t) \right\}^{n(t,q)}, \quad (1)$$

where $n(t, q)$ is the number of times term t is present in the query q , $P(t|d)$ and $P(t)$ are maximum-likelihood estimates of the probability of observing term t given the document and background language models, respectively, and λ is a smoothing parameter. The background language model is estimated from all sampled documents. Here, all documents are assumed to be equally important within a given collection, therefore, $P(d|c)$ is set to $1/|c|$, where $|c|$ is the number of (sampled) documents in collection c .

2.1.2 Document-centric Model

Instead of creating a direct term-based representation of collections, we model and query individual (sampled) docu-

Table 1: List of features used for resource selection.

Feature	Description
$CC_r(q, c)$	$P(q c)$ estimated using the CC model (Eq. 1) representations: $r = \{\text{title, snippet}\}$
$DC_{r,K}(q, c)$	$P(q c)$ estimated using the DC model (Eq. 2) representations: $r = \{\text{title, snippet, document}\}$ cutoff values: $K = \{10, 20, 50, 75, 100, 150, 200, 250, 300, 500, 1000\}$
$\text{snippets}(c)$	Number of snippets in the sample of c

ments, then aggregate their relevance estimates. This approach closely resembles the ReDDE collection selection algorithm (Si and Callan, 2003). Formally:

$$P(q|c) = \sum_{d \in c} P(d|c) \prod_{t \in q} ((1 - \lambda)P(t|d) + \lambda P(t))^{n(t,q)}, \quad (2)$$

where, as before, $P(t|d)$ and $P(t)$ and the document and background term probabilities, λ is the smoothing parameter, and $P(d|c)$ is the importance of the document given the collection. Additionally, we apply a rank-based cut-off and consider only the top K most relevant documents in the sample index for the computation of Eq. 2.

2.1.3 Combining Models

We employ a listwise learning-to-rank approach, LambdaMART (Wu et al., 2010). For training the machine learning model we use data from prior editions of the TREC Fed-Web track. Our results in §2.2 indicate that the choice of the training material has a major impact on performance.

2.2 Runs and results

We submitted the following runs:

NTNUIsrs1 Document-centric model using the entire document text ($r = \text{document}$) and a cutoff value of $K = 500$. This particular setting was chosen based on a (non-extensive) set of experiments performed on the FedWeb’13 collection.

NTNUIsrs2 Learning-to-rank approach trained on the FedWeb’13 data set.

NTNUIsrs3 Learning-to-rank approach trained on the FedWeb’12 and ’13 data sets.

Table 2 presents the results. We find that the learning-to-rank approach trained on FedWeb’13 outperforms the DC model by over 13% in terms of the official metric, nDCG@20 (NTNUIsrs2 vs. NTNUIsrs1). Interestingly, when training was done on both FedWeb’12 and ’13 performance dropped

Table 2: Results for our official resource selection runs. Best scores for each metric are in boldface.

Run	nDCG@20	nDCG@10	P@1	P@5
NTNUIsrs1	0.306	0.225	0.148	0.195
NTNUIsrs2	0.348	0.281	0.206	0.257
NTNUIsrs3	0.248	0.205	0.202	0.189

substantially (NTNUIsrs3 vs. NTNUIsrs1). Discriminative learning is indeed a promising direction for this task, but further research is needed to understand how the training material should be composed. It is also left to future work to experiment with different learning-to-rank algorithms, specifically pointwise and pairwise approaches.

3 Vertical selection

3.1 Approach

Our choice of method for the vertical selection task is closely tied to our resource selection approach. We assume that resource selection produces a relevance score $s(q, c)$ for each collection such that

$$s(q, c) = \begin{cases} > 0 & c \text{ is relevant} \\ \leq 0 & c \text{ is nonrelevant} \end{cases} \quad (3)$$

Then, we simply select all collections that have a positive relevance score:

$$V(q) = \{c | s(q, c) > 0\}, \quad (4)$$

where $V(q)$ denotes the set of selected verticals for query q . In a way, we delegate the “selection” problem to the resource ranking component.

3.2 Runs and results

We submitted the following runs:

NTNUIsrs2 Based on resource selection run NTNUIsrs2.

NTNUIsrs3 Based on resource selection run NTNUIsrs3.

Table 3 displays precision (P), recall (R), and F1-measure (F1) for our submitted runs. Based on these results, we make the not surprising observation that better resource selection indeed leads to better vertical selection. The scores, however, are quite low in absolute terms, which suggests that the scores produced by the resource selection approach may not satisfy the criteria that we have specified regarding the signs of collection scores (cf. Eq. 3). We hypothesize that using a simple score-based thresholding (i.e., changing the value 0 to a parameter in Eq. 3) might alleviate this issue. It might also be the case that the underlying resource selection step needs to be casted as a classification task as opposed to a ranking problem.

Table 3: Results for our official resource selection runs. Best scores for each metric are in boldface.

Run	P	R	F1
NTNUIsVs2	0.157	0.406	0.205
NTNUIsVs3	0.145	0.281	0.177

4 Conclusions

We described our participation in the TREC 2014 Federated Web Search track. For resource selection we have experimented with a discriminative learning approach for combining numerous instantiations of resource selection models. We have shown that it can outperform a competitive baseline model, but is sensitive to the choice of the underlying training material. We have used the estimated collection relevance scores, as binary judgments, to make a selection of verticals. We have found that improvements in resource selection indeed translate to better vertical selection performance. At the same time, making a binary judgement about the relevance of a collection remains to be challenging, given that resource selection is approached as a ranking problem, and not as a classification task.

5 References

- Balog, K. (2014). Collection and document language models for resource selection. In *Proceedings of the Twenty-Second Text REtrieval Conference (TREC 2013)*.
- Callan, J. P., Lu, Z., and Croft, W. B. (1995). Searching distributed collections with inference networks. In *Proceedings of SIGIR'95*, pages 21–28.
- Demeester, T., Trieschnigg, D., Nguyen, D., and Hiemstra, D. (2014). Overview of the TREC 2013 federated web search track. In *Proceedings of the Twenty-Second Text REtrieval Conference (TREC 2013)*.
- Si, L. and Callan, J. (2003). Relevant document distribution estimation method for resource selection. In *Proceedings of SIGIR'03*, pages 298–305.
- Si, L., Jin, R., Callan, J., and Ogilvie, P. (2002). A language modeling framework for resource selection and results merging. In *Proceedings of CIKM'02*, pages 391–397.
- Wu, Q., Burges, C. J., Svore, K. M., and Gao, J. (2010). Adapting boosting for information retrieval measures. *Inf. Retr.*, 13(3).