



## UvA-DARE (Digital Academic Repository)

### Unsupervised, Efficient and Semantic Expertise Retrieval

Van Gysel, C.; de Rijke, M.; Worring, M.

**DOI**

[10.1145/2872427.2882974](https://doi.org/10.1145/2872427.2882974)

**Publication date**

2016

**Document Version**

Author accepted manuscript

**Published in**

WWW'16

[Link to publication](#)

**Citation for published version (APA):**

Van Gysel, C., de Rijke, M., & Worring, M. (2016). Unsupervised, Efficient and Semantic Expertise Retrieval. In *WWW'16: proceedings of the 25th International Conference on World Wide Web : May 11-15, 2016, Montreal, Canada* (pp. 1069-1079). Association for Computing Machinery. <https://doi.org/10.1145/2872427.2882974>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

# Unsupervised, Efficient and Semantic Expertise Retrieval

Christophe Van Gysel  
cvangysel@uva.nl

Maarten de Rijke  
derijke@uva.nl

Marcel Worring  
m.worring@uva.nl

University of Amsterdam, Amsterdam, The Netherlands

## ABSTRACT

We introduce an unsupervised discriminative model for the task of retrieving experts in online document collections. We exclusively employ textual evidence and avoid explicit feature engineering by learning distributed word representations in an unsupervised way. We compare our model to state-of-the-art unsupervised statistical vector space and probabilistic generative approaches. Our proposed log-linear model achieves the retrieval performance levels of state-of-the-art document-centric methods with the low inference cost of so-called profile-centric approaches. It yields a statistically significant improved ranking over vector space and generative models in most cases, matching the performance of supervised methods on various benchmarks. That is, by using solely text we can do as well as methods that work with external evidence and/or relevance feedback. A contrastive analysis of rankings produced by discriminative and generative approaches shows that they have complementary strengths due to the ability of the unsupervised discriminative model to perform semantic matching.

## Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing

## Keywords

Expertise retrieval; Language models; Semantic matching

## 1. INTRODUCTION

The transition to the knowledge and information economy [1] introduces a great reliance on cognitive capabilities [50]. It is crucial for employers to facilitate information exchange and to stimulate collaboration [17]. In the past, organizations would set-up special-purpose database systems for their members to maintain a profile [7]. However, these systems required employees to be proactive. In addition, self-assessments are known to diverge from reality [10, 32] and document collections can quickly become practically infeasible to manage manually. Therefore, there has been an active interest in automated approaches for constructing expertise

profiles [7, 61] and retrieving experts from an organization's heterogeneous document repository [15]. *Expert finding* (also known as expertise retrieval or expert search) addresses the task of finding the right person with the appropriate skills and knowledge [6]. It attempts to provide an answer to the question:

Given a topic  $X$ , who are the candidates with the most expertise w.r.t.  $X$ ?

The expertise retrieval task gained popularity in the research community during the TREC Enterprise Track [61] and has remained relevant ever since, while broadening to social media and to tracking the dynamics of expertise [5, 6, 10, 19, 21–23, 47, 49, 62]. Existing methods fail to address key challenges: (1) Queries and expert documents use different representations to describe the same concepts [26, 34]. Term mismatches between queries and experts [34] occur due to the inability of widely used maximum-likelihood language models to make use of *semantic* similarities between words [53]. (2) As the amount of available data increases, the need for more powerful approaches with greater learning capabilities than smoothed maximum-likelihood language models is obvious [63]. (3) Supervised methods for expertise retrieval [23, 47] were introduced at the turn of the last decade. However, the acceleration of data availability has the major disadvantage that, in the case of supervised methods, manual annotation efforts need to sustain a similar order of growth. This calls for the further development of *unsupervised* methods. (4) In some expertise retrieval methods, a language model is constructed for every document in the collection. These methods lack *efficient* query capabilities for large document collections, as each query term needs to be matched against every document [6]. Our proposed solution has a strong emphasis on *unsupervised model construction*, *efficient query capabilities* and *semantic matching* between query terms and candidate experts.

Specifically, we propose an unsupervised log-linear model with efficient inference capabilities for the expertise retrieval task. We show that our approach improves retrieval performance compared to vector space-based and generative language models, mainly due to its ability to perform semantic matching [34]. Our method does not require supervised relevance judgments and is able to learn from raw textual evidence and document-candidate associations alone. The purpose of this work is to provide insight in how discriminative language models can improve performance of core retrieval tasks compared to maximum-likelihood language models. Therefore, we avoid explicit feature engineering and the incorporation of external evidence in this paper. In terms of performance, the current best-performing formal language model [4] exhibits a worst-case time complexity linear in the size of the document collection. In contrast, the inference time complexity of our approach is asymptotically bounded by the number of candidate experts.

Our research questions are as follows: (1) How does our discriminative log-linear model compare to vector space-based methods and generative language models for the expert retrieval task in terms of retrieval performance? (2) What can we learn regarding the different types of errors made by generative and discriminative language models? (3) How does the complexity of inference in our log-linear model compare to vector-space based and generative models? (4) How does the log-linear model handle incremental indexing and what are its limitations?

We contribute: (i) An unsupervised log-linear model with efficient inference capabilities for the expertise retrieval task, together with an open-source implementation.<sup>1</sup> (ii) Comparison of the retrieval performance of the log-linear model with traditional vector space-based models and language model methods on well-known benchmarks. (iii) Insights in how certainty in predictions by the log-linear model correlates with performance. (iv) In-depth analysis of the inference complexity of the log-linear model. (v) Comparative error analysis between the semantic log-linear model and traditional generative language models that perform exact matching. (vi) Insight in the relative strengths of semantic matching and exact matching for the expert retrieval task.

The remainder of this paper is organized as follows. In Section 2, we briefly discuss related work. Section 3 introduces the log-linear model for the expert retrieval problem. In Section 4 we state our research questions and detail our experimental set-up and implementations. We provide an overview of our experimental results, followed by a discussion and analysis in Section 5. Section 6 concludes the paper and discusses ideas for future work.

## 2. RELATED WORK

We first discuss prior work on expert retrieval and its relation to document retrieval. Then, we review semantic search methods. The ideas we present in this paper are inspired by neural language models used in automated speech recognition (ASR) and natural language processing (NLP). Consequently, we also review work from those fields.

### 2.1 Expert retrieval

Early expert retrieval systems were often referred to as expert locator and expertise management systems [38]. These database systems often relied on people to self-assess their expertise against a predefined set of topics [39], which is known to generate unreliable results [7].

With the introduction of the P@NOPTIC system [15], and later the TREC Enterprise track [61], there has been an active research interest in automated expertise profiling methods. It is useful to distinguish between *profile-based* methods, which create a textual representation of a candidate’s knowledge, and *document-based* methods, which represent candidates as a weighted combination of documents. The latter generally performs better at ranking, while the former is more efficient as it avoids retrieving all documents relevant to a query [6, p. 221].

There has been much research on generative probabilistic models for expert retrieval [21, 49]. Such models have been categorized in candidate generation models [13], topic generation models [4, 5] and proximity-based variants [5, 54]. Of special relevance to us are the unsupervised profile-centric (Model 1) and document-centric (Model 2) models of Balog et al. [4], which focus on raw textual evidence without incorporating collection-specific information (e.g., query modeling, document importance or document structure). Supervised discriminative models [23, 47, 60] are preferred when

query-candidate relevance pairs are available for training. Unlike their generative counterparts these models have no issue combining complex and heterogeneous features (e.g., link-based features, document importance features, etc.); they resemble Learning to Rank (L2R) methods for document retrieval [6, 35]. However, a lack of training data may greatly hinder their applicability [6, p. 179]. Beyond unsupervised generative and supervised discriminative approaches, there are graph-based approaches based on random walks [55] and voting-based approaches based on data fusion [36]. Demartini et al. [19] propose a vector space-based method for the entity ranking task; their framework extends vector spaces operating on documents to entities. See [6] for a survey on the topic.

### 2.2 Latent semantic models for document retrieval

Li and Xu [34] note that the query document mismatch poses the most critical challenge in search. Semantic matching is an important attempt to remedy this problem. There has been much work on bridging the semantic gap for various different tasks [19, 26, 28, 30, 31, 41, 48, 53]. Expertise and document retrieval [6, p. 224] are closely related as performance of the latter can greatly impact that of the former [37].

Latent Semantic Models (LSM) first became popular through the introduction of Latent Semantic Indexing (LSI) [18], followed by probabilistic LSI (pLSI) [27]. LSMs based on neural networks [28, 53, 58] emerged in the last decade. Salakhutdinov and Hinton [53] employ unsupervised deep auto-encoders to map documents to bit patterns using semantic addressing. Huang et al. [28] perform semantic matching of documents and queries by leveraging click-through data optimizing for Web document ranking. Further, deep models [12, 20] have been proposed to learn to rank [35].

### 2.3 Neural language modeling

Large-vocabulary neural probabilistic language models for modeling word sequence distributions have become very popular recently [8, 43, 44]. These models learn continuous-valued distributed representations for words, also known as embeddings [42, 45, 48], in order to fight the curse of dimensionality and increase generalization by introducing the expectation that similar word vectors signify semantically or syntactically similar words. Recurrent neural language models have shown to perform well in ASR [40]. Collobert et al. [14] propose a unified neural network architecture for various NLP tasks. Even more recently, there has been a surge in multimodal neural language models [31], which lend themselves to the task of automated image captioning.

What we add on top of the related work described above is the following. In this work we model the conditional probability of the expertise of a candidate given a single query term (contrary to binary relevance given a character-based n-gram [28]). In the process we learn a distributed vector representation (similar to LSI, pLSI and semantic hashing) for both words and candidates such that nearby representations indicate semantically similar concepts.

We propose a log-linear model that is similar to neural language models. The important difference is that we predict a candidate expert instead of the next word. To the best of our knowledge, we are the first to propose such a solution. We employ an embedding layer in our shallow model for the same reasons as mentioned above: we learn continuous word representations that incorporate semantic and syntactic similarity tailored to an expert’s domain.

<sup>1</sup><https://github.com/cvangysel/SERT>

### 3. A LOG-LINEAR MODEL FOR EXPERT SEARCH

In the setting of this paper we have a document collection  $D$  and a predefined set of candidate experts  $C$  (entities to be retrieved). Documents  $d \in D$  are represented as a sequence of words  $w_1, \dots, w_{|d|}$  originating from a vocabulary  $V$ , where  $w_i \in V$  and the operator  $|\cdot|$  denotes the document length in tokens. For every document  $d \in D$  we write  $C_d$  to denote the set of candidates  $c \in C$  associated with it (i.e.,  $C = \bigcup_{d \in D} C_d$ ). These document-candidate associations can be obtained explicitly from document meta-data (e.g., the author of an e-mail) or implicitly by mining references to candidates from the document text. Notice that some documents might not be associated with any candidate. When presented with a query  $q$  of constituent terms  $t_1, \dots, t_{|q|}$ , the expert retrieval task is to return a list of candidates  $\rho(C)$  ordered according to topical expertise. We generate this ranking using a relatively shallow neural network which directly models  $P(c | q)$ .

We employ vector-based distributed representations [26], for both words (i.e., word embeddings) and candidate experts, in a way that motivates the unsupervised construction of features that express regularities of the expertise finding domain. These representations can capture the similarity between concepts (e.g., words and candidate experts) by the closeness of their representations in vector space. That is, concepts with similar feature activations are interpreted by the model as being similar, or even interchangeable.

#### 3.1 The model

To address the expert search task, we model  $P(c_j | q)$  and rank candidates  $c_j$  accordingly for a given  $q$ . We propose an unsupervised, discriminative approach to obtain these probabilities. We construct our model solely from textual evidence: we do not require query-candidate relevance assessments for training and do not consider external evidence about the corpus (e.g., different weightings for different sub-collections), the document (e.g., considering certain parts of the document more useful) nor link-based features.

Let  $e$  denote the size of the vector-based distributed representations of both words in  $V$  and candidate experts in  $C$ . These representations will be learned by the model using gradient descent [42] (Section 3.2). For notational convenience, we write  $P(c | \cdot)$  for the (conditional) probability distribution over candidates, which is the result of vector arithmetic. We define the probability of a candidate  $c_j$  given a single word  $w_i \in V$  as the log-linear model

$$P(c | w_i) = \frac{1}{Z_1} \exp(W_c \cdot (W_p \cdot v_i) + b_c), \quad (1)$$

where  $W_p$  is the  $e \times |V|$  projection matrix that maps the one-hot representation (i.e., 1-of- $|V|$ ) of word  $w_i$ ,  $v_i$ , to its  $e$ -dimensional distributed representation,  $b_c$  is a  $|C|$ -dimensional bias vector and  $W_c$  is the  $|C| \times e$  matrix that maps the word embedding to an unnormalized distribution over candidates  $C$ , which is then normalized by  $Z_1 = \sum_{j=1}^{|C|} [\exp(W_c \cdot (W_p \cdot v_i) + b_c)]_j$ . If we consider Bayes' theorem, the transformation matrix  $W_c$  and bias vector  $b_c$  can be interpreted as the term log-likelihood  $\log P(w_i | c)$  and candidate log-prior  $\log P(c)$ , respectively. The projection matrix  $W_p$  attempts to soften the curse of dimensionality introduced by large vocabularies  $V$  and maps words to word feature vectors [8]. Support for large vocabularies is crucial for retrieval tasks [28, 53].

We then assume conditional independence of a candidate's expertise given an observation of data (i.e., a word). Given a sequence

of words  $w_1, \dots, w_k$  we have:

$$\begin{aligned} P(c | w_1, \dots, w_k) &= \frac{1}{Z_2} \tilde{P}(c | w_1, \dots, w_k) = \frac{1}{Z_2} \prod_{i=1}^k P(c | w_i) \\ &= \frac{1}{Z_2} \exp\left(\sum_{i=1}^k \log(P(c | w_i))\right) \end{aligned} \quad (2)$$

where  $\tilde{P}(c | w_1, \dots, w_k)$  denotes the unnormalized score and  $Z_2 = \sum_{j=1}^{|C|} \exp\left(\sum_{i=1}^k \log(P(c_j | w_i))\right)$  is a normalizing term. The transformation to log-space in (2) is a well-known trick to prevent floating point underflow [46, p. 445]. Given (2), inference is straight-forward. That is, given query  $q = t_1, \dots, t_k$ , we compute  $P(c | t_1, \dots, t_k)$  and rank the candidate experts in descending order of probability.

Eq. 1 defines a neural network with a single hidden layer. We can add additional layers. Preliminary experiments, however, show that the shallow log-linear model (1) performs well-enough in most cases. Only for larger data sets did we notice a marginal gain from adding an additional layer between projection matrix  $W_p$  and the softmax layer over  $C$  ( $W_c$  and bias  $b_c$ ), at the expense of longer training times and loss of transparency.

#### 3.2 Parameter estimation

The matrices  $W_p$ ,  $W_c$  and the vector  $b_c$  in (1) constitute the parameters of our model. We estimate them using error back propagation [51] as follows. For every document  $d_j \in D$  we construct an ideal distribution over candidates  $\mathbf{p} = P(c | d_j)$  based on the document-candidate associations  $C_{d_j}$  such that

$$P(c | d_j) = \begin{cases} \frac{1}{|C_{d_j}|}, & c \in C_{d_j} \\ 0, & c \notin C_{d_j} \end{cases}$$

We continue by extracting n-grams where  $n$  remains fixed during training. For every n-gram  $w_1, \dots, w_n$  generated from document  $d$  we compute  $\tilde{\mathbf{p}} = P(c | w_1, \dots, w_n)$  using (2). During model constructing we then optimize the cross-entropy  $H(\mathbf{p}, \tilde{\mathbf{p}})$  (i.e., the joint probability of the training data if  $|C_{d_j}| = 1$  for all  $j$ ) using batched gradient descent. The loss function for a single batch of  $m$  instances with associated targets  $(\mathbf{p}^{(i)}, \tilde{\mathbf{p}}^{(i)})$  is as follows:

$$\begin{aligned} L(W_p, W_c, b_c) &= \frac{1}{m} \sum_{i=1}^m \frac{|d_{\max}|}{|d^{(i)}|} H(\mathbf{p}^{(i)}, \tilde{\mathbf{p}}^{(i)}) \\ &\quad + \frac{\lambda}{2m} \left( \sum_{i,j} W_{p,i,j}^2 + \sum_{i,j} W_{c,i,j}^2 \right) \\ &= -\frac{1}{m} \sum_{i=1}^m \frac{|d_{\max}|}{|d^{(i)}|} \sum_{j=1}^{|C|} P(c_j | d^{(i)}) \log(P(c_j | w_1^{(i)}, \dots, w_n^{(i)})) \\ &\quad + \frac{\lambda}{2m} \left( \sum_{i,j} W_{p,i,j}^2 + \sum_{i,j} W_{c,i,j}^2 \right), \end{aligned} \quad (3)$$

where  $d^{(i)}$  refers to the document from which n-gram  $w_1^{(i)}, \dots, w_n^{(i)}$  was extracted,  $d_{\max} = \arg \max_{d \in D} |d|$  indicates the longest document in the collection, and  $\lambda$  is a weight regularization parameter. The update rule for a particular parameter  $\theta$  ( $W_p$ ,  $W_c$  or  $b_c$ ) given a single batch of size  $m$  is:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha^{(t)} \odot \frac{\partial L(W_p^{(t)}, W_c^{(t)}, b_c^{(t)})}{\partial \theta}, \quad (4)$$

where  $\alpha^{(t)}$  and  $\theta^{(t)}$  denote the per-parameter learning rate and parameter  $\theta$  at time  $t$ , respectively. The learning rate  $\alpha$  consists of the same number of elements as there are parameters; in the case of a global learning rate, all elements of  $\alpha$  are equal to each other. The derivatives of the loss function (3) are given in the Appendix.

In the next section we will discuss our experimental setup, followed by an overview of our experimental results and further analysis in Section 5.

## 4. EXPERIMENTAL SETUP

### 4.1 Research questions

As indicated in the introduction, we seek to answer the following research questions:

**RQ1** How does our discriminative log-linear model compare to vector space-based methods and generative language models for the expert retrieval task in terms of retrieval performance?

In particular, how does the model perform when compared against vector space-based (LSI and TF-IDF) and generative approaches (profile-centric Model 1 and document-centric Model 2)?

**RQ2** What can we learn regarding the different types of errors made by generative and discriminative language models?

Does the best-performing generative model simply perform slightly better on the topics for which the other models perform decent as well, or do they make very different errors? If the latter holds, an ensemble of the rankings produced by both model types might exceed performance of the individual rankings.

**RQ3** How does the complexity of inference in our log-linear model compare to vector-space based and generative models?

The worst-case inference cost of document-centric models makes them unattractive in online settings where the set of topics is not defined beforehand and the document collection is large. Profile-centric methods are preferred in such settings as they infer from one language model per candidate expert for every topic (i.e., a pseudo-document consisting of a concatenation of all documents associated with an expert) [6]. Vector space-based methods [19] have similar problems due to the curse of dimensionality [29] and consequently their inferential time complexity is likewise asymptotically bounded by the number of experts.

**RQ4** How does the log-linear model handle incremental indexing and what are its limitations?

### 4.2 Benchmarks

The proposed method is applicable in the setting of the Expert Search task of the TREC Enterprise track from 2005 to 2008 [61]. We therefore evaluate on the W3C and CERC benchmarks released by the track. The W3C dataset [16] is a crawl of the W3C’s sites in June 2004 (mailing lists, web pages, etc.). The CSIRO Enterprise Research Collection (CERC) [2] is a dump of the intranet of Australia’s national science agency. Additionally, we evaluate our method on a smaller, more recent benchmark based on the employee database of Tilburg University (TU) [10], which consists of bi-lingual, heterogeneous documents. See Table 1.

Mining document-candidate associations and how they influence performance has been extensively covered in previous work [4, 6] and is beyond the scope of this work. For TU, the associations are part of the benchmark. For W3C, a list of possible candidates is given and we extract the associations ourselves by performing a case-insensitive match of full name or e-mail address [4]. For CERC, we make use of publicly released associations [3].

As evaluation measures we use Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), Normalized Discounted Cumula-

tive Gain at rank 100 (NDCG@100) and Precision at rank 5 (P@5) and rank 10 (P@10).

### 4.3 Baselines

We compare our approach to existing unsupervised methods for expert retrieval that solely rely on textual evidence and static document-candidate associations. (1) Demartini et al. [19] propose a generic framework to adapt vector spaces operating on documents to entities. We compare our method to TF-IDF (raw frequency and inverse document frequency) and LSI (300 latent topics) variants of their vector space model for entity ranking (using cosine similarity). (2) In terms of language modeling, Balog et al. [4] propose two models for expert finding based on generative language models. The first takes a profile-centric approach comparing the language model of every expert to the query, while the second is document-centric. We consider both models with different smoothing configurations: Jelinek-Mercer (jm) smoothing with  $\lambda = 0.5$  [4] and Dirichlet (d) smoothing with  $\beta$  equal to the average document length [5] (see Table 1). Significance of results produced by the baselines (compared to our method) is determined using a two-tailed paired randomization test [59].

### 4.4 Implementation details

The vocabulary  $V$  is constructed from each corpus by ignoring punctuation, stop words and case; numbers are replaced by a numerical placeholder token. During our experiments we prune  $V$  by only retaining the  $2^{16}$  most-frequent words so that each word can be encoded by a 16-bit unsigned integer. Incomplete n-gram instances are padded by a special-purpose token.

In terms of parameter initialization, we sample the initial matrices  $W_c$  and  $W_p$  (1) uniformly in the range

$$\left[ -\sqrt{\frac{6.0}{m+n}}, \sqrt{\frac{6.0}{m+n}} \right]$$

for an  $m \times n$  matrix, as this initialization scheme improves model training convergence [25], and take the bias vector  $b_c$  to be null. The projection layer  $W_p$  is initialized with pre-trained word representations trained on Google News data [41]; the number of word features is set to  $e = 300$ , similar to pre-trained representations.

We used adadelta ( $\rho = 0.95$ ,  $\epsilon = 10^{-6}$ ) [64] with batched gradient descent ( $m = 1024$ ) and weight decay  $\lambda = 0.01$  during training on NVidia GTX480 and NVidia Tesla K20 GPUs. We only iterate once over the entire training set for each experiment.

## 5. RESULTS AND DISCUSSION

We start by giving a high-level overview of our experimental results and then address issues of scalability, provide an error analysis and discuss the issue of incremental indexing.

### 5.1 Overview of experimental results

We evaluate the log-linear model on the W3C, CERC and TU benchmarks (Section 4.2). During training we extract non-overlapping n-grams for the W3C and CERC benchmarks and overlapping n-grams for the TU benchmark. As the TU benchmark is considerably smaller, we opted to use overlapping n-grams to counter data sparsity. The architecture of each benchmark model (e.g., number of candidate experts) is inherently specified by the benchmarks themselves (see Table 1). However, the choice of n-gram size during training remains open. Errors for input  $w_1, \dots, w_n$  are propagated back through  $W_c$  until the projection matrix  $W_p$  is reached; if a single word  $w_i$  causes a large prediction error, then this will influence its neighboring words  $w_1, \dots, w_{i-1}, w_{i+1}, \dots, w_n$  as well.

	W3C		CERC		TU	
Number of documents	331,037		370,715		31,209	
Average document length <sup>a</sup>	1,237.23		460.48		2,454.93	
Number of candidates <sup>b</sup>	715		3,479		977	
Number of document-candidate associations	200,939		236,958		36,566	
Number of documents (with $C_d > 0$ )	93,826		123,934		27,834	
Number of associations per document <sup>c</sup>	$2.14 \pm 3.29$		$1.91 \pm 3.70$		$1.13 \pm 0.39$	
Number of associations per candidate	$281.03 \pm 666.63$		$68.11 \pm 1,120.74$		$37.43 \pm 61.00$	
Queries	49	(2005)	50	(2007)	1,662	(GT1)
	50	(2006)	77	(2008)	1,266	(GT5)

<sup>a</sup> Measured in number of tokens.

<sup>b</sup> Only candidates with at least a single document association are considered.

<sup>c</sup> Only documents with at least one association are considered.

**Table 1: An overview of the three datasets (W3C, CERC and TU) used for evaluation and analysis.**

This allows the model to learn continuous word representations tailored to the expert retrieval task and the benchmark domain.

A larger window size has a negative impact on batch throughput during training. We are thus presented with the classic trade-off between model performance and construction efficiency. Notice, however, that the number of n-grams decreases as the window size increases if we extract non-overlapping instances. Therefore, larger values of  $n$  lead to faster wall-clock time model construction for the W3C and CERC benchmarks in our experiments.

We sweep over the window width  $n = 2^i$  ( $0 \leq i < 6$ ) for all three benchmarks and their corresponding relevance assessments. We report MAP and MRR for every configuration (see Figure 1). We observe a significant performance increase between  $n = 1$  and  $n = 2$  on all benchmarks, which underlines the importance of the window size parameter. The increase in MAP implies that the performance achieved is not solely due to initialization with pre-trained representations (Section 4.4), but that the model efficiently learns word representations tailored to the problem domain. The highest MAP scores are attained for relatively low  $n$ . As  $n$  increases beyond  $n = 8$  a gradual decrease in MAP is observed on all benchmarks. In our remaining experiments we choose  $n = 8$  regardless of the benchmark.

Words  $w_i$  that mainly occur in documents associated with a particular expert are learned to produce distributions  $P(c | w_i)$  with less uncertainty than words associated with many experts in (1). The product of  $P(c | w_i)$  in (2) aggregates this expertise evidence generated by query terms. Hence, queries with strong evidence for a particular expert should be more predictable than very generic queries. To quantify uncertainty we measure the normalized entropy [56] of  $P(c | q)$ :

$$\eta(c | q) = -\frac{1}{\log(|C|)} \sum_{j=1}^{|C|} p(c_j | q) \log(p(c_j | q)). \quad (5)$$

Equation 5 can be interpreted as a similarity measure between the given distribution and the uniform distribution. Importantly, Figure 2 shows that there is a statistically significant negative correlation between query-wise normalized entropy and average precision for all benchmarks.

Table 2 presents a comparison between the log-linear model and the various baselines (Section 4.3). Our unsupervised method significantly ( $p < 0.01$ ) outperforms the LSI-based method consistently. In the case of the TF-IDF method and the profile-centric

generative language models (Model 1), we always perform better and statistical significance is achieved in the majority of cases. The document-centric language models (Model 2) perform slightly better than our method on two (out of six) benchmarks in terms of MAP and NDCG@100: (1) For the CERC 2007 assessment we match performance of the document-centric generative model with Jelinek-Mercer smoothing. (2) For TU GT1 the generative counterpart seems to outperform our method.

Notice that over all assessments, the log-linear model consistently outperforms all profile-centric approaches and is only challenged by the smoothed document-centric approach. In addition, for the precision-based measures ( $P@k$  and MRR), the log-linear model consistently outperforms all other methods we compare to.

Next, we turn to a topic-wise comparative analysis of discriminative and generative models. After that, we analyze the scalability and efficiency of the log-linear model and compare it to that of the generative counterparts, and address incremental indexing.

## 5.2 Error analysis

How does our log-linear model achieve its superior performance over established generative models? Figure 3 depicts the per-topic differences in average precision between the log-linear model and Model 2 (with Jelinek-Mercer smoothing) on all benchmarks. For each plot, the vertical bars with a positive AP difference correspond to test topics for which the log-linear model outperforms Model 2 and vice versa for bars with a negative AP difference.

The benefit gained from the projection matrix  $W_p$  is two-fold. First, it avoids the curse of dimensionality introduced by large vocabularies. Second, term similarity with respect to the expertise domain is encoded in latent word features. When examining words nearby query terms in the embedding space, we found words to be related to the query term. For example, word vector representations of *xml* and *nonterminal* are very similar for the W3C benchmark ( $l_2$  norm). This can be further observed in Figure 1: log-linear models trained on single words perform significantly worse compared to those that are able to learn from neighboring words.

We now take a closer look at the topics for which the log-linear model outperforms Model 2 and vice versa. More specifically, we investigate textual evidence related to a topic and whether it is considered relevant by the benchmark. For the log-linear model, we examine terms nearby topic terms in  $W_p$  ( $l_2$ -norm), as these terms are considered semantically similar by the model and provide a means for semantic matching. For every benchmark, we first consider top-

W3C	2005					2006				
	MAP	NDCG@100	MRR	P@5	P@10	MAP	NDCG@100	MRR	P@5	P@10
LSI	0.135	0.266	0.306	0.192	0.196	0.245	0.371	0.482	0.287	0.338
TF-IDF	0.243	0.426	0.541	0.384	0.350	0.343	0.531	0.650	0.492	0.498
Model 1 (d)	0.192	0.358	0.433	0.276	0.266	0.321	0.491	0.635	0.478	0.449
Model 1 (jm)	0.190	0.352	0.390	0.272	0.276	0.311	0.483	0.596	0.502	0.437
Model 2 (d)	0.198	0.369	0.429	0.288	0.272	0.261	0.419	0.551	0.441	0.404
Model 2 (jm)	0.211	0.380	0.451	0.332	0.296	0.260	0.423	0.599	0.449	0.429
Log-linear (ours)	<b>0.248</b>	<b>0.444</b>	<b>0.618*</b>	<b>0.412</b>	<b>0.361</b>	<b>0.484***</b>	<b>0.667***</b>	<b>0.833***</b>	<b>0.713***</b>	<b>0.644***</b>

CERC	2007					2008				
	MAP	NDCG@100	MRR	P@5	P@10	MAP	NDCG@100	MRR	P@5	P@10
LSI	0.031	0.107	0.060	0.016	0.014	0.038	0.099	0.106	0.042	0.055
TF-IDF	0.332	0.486	0.463	0.196	0.141	0.269	0.465	0.525	0.332	0.277
Model 1 (d)	0.287	0.427	0.384	0.156	0.096	0.181	0.355	0.388	0.200	0.172
Model 1 (jm)	0.278	0.420	0.384	0.156	0.084	0.170	0.347	0.339	0.181	0.159
Model 2 (d)	0.352	0.495	0.454	0.180	0.138	0.264	0.461	0.510	0.281	0.244
Model 2 (jm)	<b>0.361</b>	<b>0.500</b>	0.467	0.192	0.138	0.274	0.463	0.517	0.278	0.239
Log-linear (ours)	0.344	0.493	<b>0.513</b>	<b>0.215</b>	<b>0.150</b>	<b>0.342***</b>	<b>0.519**</b>	<b>0.656**</b>	<b>0.381*</b>	<b>0.299</b>

TU	GT1					GT5				
	MAP	NDCG@100	MRR	P@5	P@10	MAP	NDCG@100	MRR	P@5	P@10
LSI	0.095	0.205	0.153	0.060	0.051	0.097	0.208	0.129	0.043	0.036
TF-IDF	0.216	0.356	0.324	0.131	0.097	0.233	0.378	0.288	0.108	0.079
Model 1 (d)	0.171	0.308	0.258	0.103	0.082	0.241	0.385	0.292	0.109	0.081
Model 1 (jm)	0.189	0.325	0.277	0.112	0.085	0.231	0.373	0.271	0.100	0.075
Model 2 (d)	0.154	0.284	0.228	0.087	0.070	0.191	0.334	0.233	0.084	0.065
Model 2 (jm)	<b>0.234</b>	<b>0.370</b>	0.342	0.136	0.101	0.253	0.394	0.302	0.108	0.081
Log-linear (ours)	0.219	0.356	<b>0.351</b>	<b>0.145*</b>	<b>0.105</b>	<b>0.287***</b>	<b>0.425***</b>	<b>0.363***</b>	<b>0.134***</b>	<b>0.092***</b>

**Table 2: Evaluation results for models trained on the W3C, CERC and TU benchmarks. Suffixes (d) and (jm) denote Dirichlet and Jelinek-Mercer smoothing, respectively (Section 4.3). Significance of results is determined using a two-tailed paired randomization test [59] (\*\* $p < 0.01$ ; \*\* $p < 0.05$ ; \* $p < 0.1$ ) with respect to the log-linear model (adjusted using the Benjamini-Hochberg procedure for multiple testing [9]).**

ics where exact matches (Model 2) perform best, followed by examples which benefit from semantic matching (log-linear model). Topic identifiers are between parentheses.

**W3C** Topics *P3P specification* and *CSS3* (EX8 and EX69, respectively) should return candidates associated with the definition of these standards. The log-linear model, however, considers these close to related technologies such as CSS2 for CSS3 and UTF-8 for P3P. Semantic matching works for topics *Semantic Web Coordination* and *Annotea server protocol* (EX1 and EX103), where the former is associated with RDF libraries, RDF-related jargon and the names of researchers in the field, while the latter is associated with implementations of the protocol and the maintainer of the project.

**CERC** For CSIRO, topic *nanohouse* (CE-035) is mentioned in many irrelevant contexts (i.e., spam) and therefore semantic matching fails. The term *fish oil* (CE-126) is quickly associated with different kinds of fish, oils and organizations related to marines and fisheries. On the other hand, we observe *sensor networks* (CE-018) to be associated with sensor/networking jargon and sensor platforms. Topic *forensic science workshop* (CE-103) expands to syntactically-similar terms (e.g., plural), the names of science laboratories and references to support/law-protection organizations.

**TU** The TU benchmark contains both English and Dutch textual evidence. Topics *sustainable tourism* and *interpolation* (1411 and 4882) do not benefit from semantic matching due to a semantic gap: *interpolation* is associated with the polynomial kind while the relevance assessments focus on stochastic methods. Interestingly, for the topic *law and informatization/computerization* (1719) we see that the Dutch translation of *law* is very closely related. Similar terms to *informatization* are, according to the log-linear model, Dutch words related to cryptography. Similar dynamics are at work for *legal-political space* (12603), where translated terms and semantic-syntactic relations aid performance.

In order to further quantify the effect of the embedding matrix  $W_p$ , we artificially expand benchmark topic terms by  $k$  nearby terms. We then examine how the performance of a profile-centric generative language model [4, Model 1] evolves for different values of  $k$  (Figure 4). The purpose of this analysis is to provide further insight in the differences between maximum-likelihood language models and the log-linear model. Figure 4 shows that, for most benchmarks, MAP increases as  $k$  goes up. Interestingly enough, the two benchmarks that exhibit a decrease in MAP for larger  $k$  (CERC 2007 and TU GT1) are likewise those for which generative language models outperform the log-linear model in Table 2. This suggests that the CERC 2007 and TU GT1 benchmarks require ex-

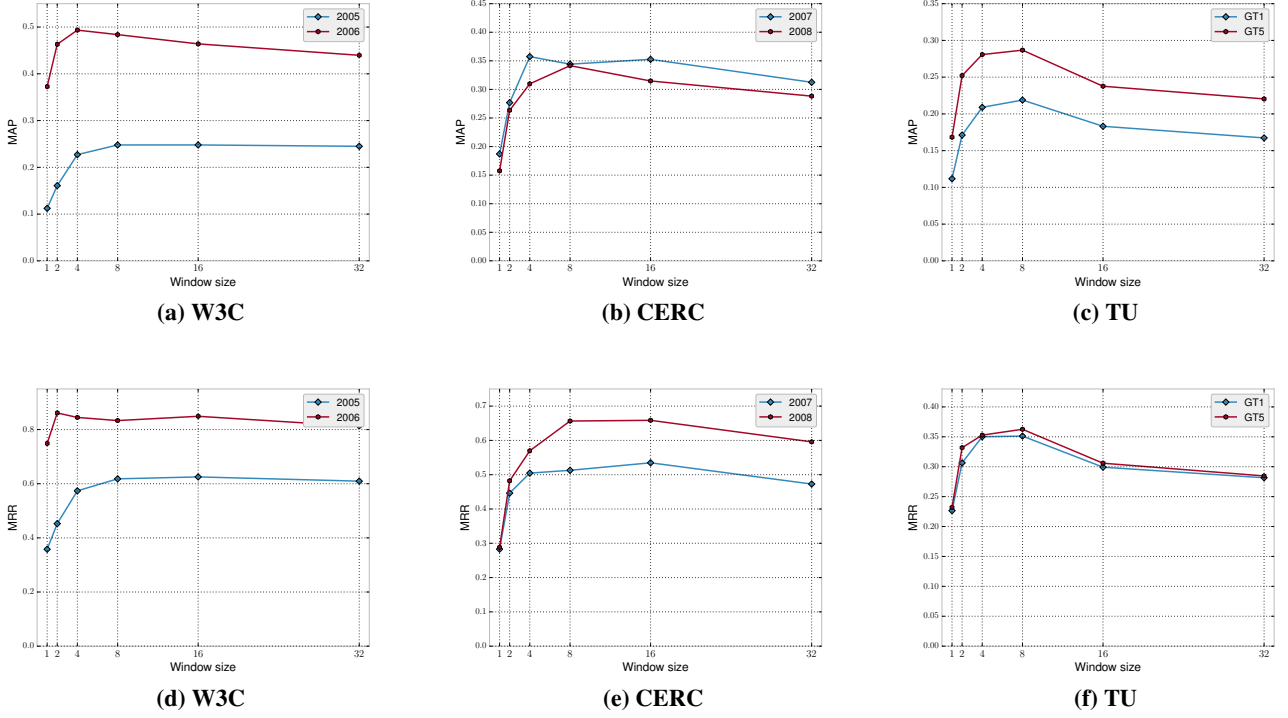


Figure 1: Sensitivity analysis for window size (n-gram) during parameter estimation (3) for W3C, CERC and TU benchmarks.

act term matching, while the remaining four benchmarks benefit greatly from the semantic matching provided by our model.

The per-topic differences suggest that Model 2 and the log-linear model make very different errors: Model 2 excels at retrieving exact query matches, while the log-linear model performs semantic matching. Based on these observations we hypothesize that a combination of the two approaches will raise retrieval performance even further. To test this hypothesis, we propose a simple ensemble of rankings generated by Model 2 and the log-linear model by re-ranking candidates according to the multiplicatively-combined reciprocal rank:

$$\text{rank}_{\text{ensemble}}(c_j, q_i) \propto \frac{1}{\text{rank}_{\text{model 2}}(c_j, q_i)} \cdot \frac{1}{\text{rank}_{\text{log-linear}}(c_j, q_i)}, \quad (6)$$

where  $\text{rank}_M(c_j, q_i)$  denotes the position of candidate  $c_j$  in a ranking generated by model  $M$  for answering query  $q_i$ . Equation (6) is equivalent to performing data fusion using CombSUM [57] where the scores are given by the logarithm of the reciprocal ranks of the experts. Table 3 compares the result of this ensemble to that of its constituents. Compared to the supervised methods of Fang et al. [23], we conclude that our fully unsupervised ensemble matches the performance of their method on the CERC 2007 benchmark and outperforms their method on the W3C 2005 benchmark. The superior performance of the ensemble suggests the viability of hybrid methods that combine semantic and exact matching.

### 5.3 Scalability and efficiency

Inference in the log-linear model is expressed in linear algebra operations (Section 3). These operations can be efficiently performed by highly optimized software libraries and special-purpose hardware (i.e., GPUs). But the baseline methods against which we compare do not benefit from these speed-ups. Furthermore, many implementation-specific details and choice of parameter values can

influence runtime considerably (e.g. size of the latent representations). Therefore, we opt for a theoretical comparison of the inference complexity of the log-linear model and compare these to the baselines (Section 4.3).

The log-linear model generates a ranking of candidate experts by straight-forward matrix operations. The look-up operation in the projection matrix  $W_p$  occurs in constant time complexity, as the multiplication with the one-hot vector  $v_i$  comes down to selecting the  $i$ -th column from  $W_p$ . Multiplication of the  $|C| \times e$  matrix  $W_e$  with the  $e$ -dimensional word feature vector exhibits  $O(|C| \cdot e)$  runtime complexity. If we consider addition of the bias term and division by the normalizing function  $Z_1$ , the time complexity of (1) becomes

$$O(\underbrace{|C| \cdot (e + (e - 1))}_{\text{matrix-vector multiplication}} + \underbrace{|C|}_{\text{bias term}} + \underbrace{2 \cdot |C| - 1}_{Z_1}).$$

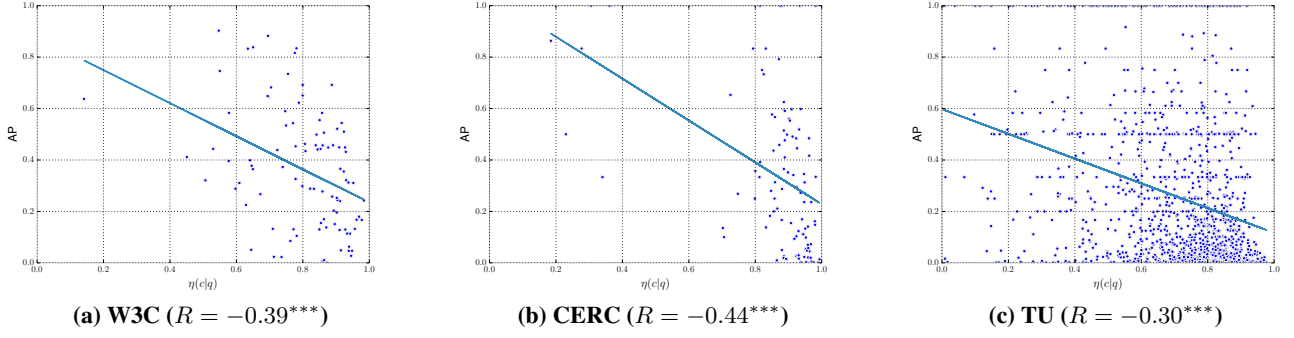
Notice, however, that the above analysis considers sequential execution. Modern computing hardware has the ability to parallelize common matrix operations [24, 33]. The number of candidate experts  $|C|$  is the term that impacts performance most in the log-linear model (under the assumption that  $|C| \gg e$ ).

If we consider  $n$ , where  $n$  is the query length during inference or the window size during training, then the complexity of (2) becomes

$$O(\underbrace{n \cdot |C| \cdot (2 \cdot e - 1) + n \cdot (3 \cdot |C| - 1)}_{n \text{ forward-passes}} + \underbrace{(n - 1) \cdot |C|}_{\text{factor product}} + \underbrace{2 \cdot |C| - 1}_{Z_2})$$

Notice that  $Z_2$  does not need to be computed during inference as it does not affect the candidate expert ranking.





**Figure 2: Scatter plot of the normalized entropy of distribution  $P(c | q)$  (2) returned by the log-linear model and per-query average precision for W3C, CERC and TU benchmarks. Pearson’s  $R$  and associated  $p$ -value (two-tailed paired permutation test: \*\*\*  $p < 0.01$ ; \*\*  $p < 0.05$ ; \*  $p < 0.1$ ) are between parentheses. The depicted linear fit was obtained using an ordinary least squares regression.**

W3C			2005					2006		
	MAP	NDCG@100	MRR	P@5	P@10	MAP	NDCG@100	MRR	P@5	P@10
Model 2 (jm)	0.211	0.380	0.451	0.332	0.296	0.260	0.423	0.599	0.449	0.429
Log-linear (ours)	0.248	0.444	0.618	0.412	0.361	<b>0.484**</b>	<b>0.667**</b>	<b>0.833</b>	<b>0.713**</b>	<b>0.644**</b>
Ensemble	<b>0.291***</b>	<b>0.479**</b>	<b>0.668</b>	<b>0.440</b>	<b>0.378</b>	0.433	0.634	0.825	0.657	0.586
CERC			2007					2008		
	MAP	NDCG@100	MRR	P@5	P@10	MAP	NDCG@100	MRR	P@5	P@10
Model 2 (jm)	0.361	0.500	0.467	0.192	0.138	0.274	0.463	0.517	0.278	0.239
Log-linear (ours)	0.344	0.493	0.513	0.215	0.150	0.342	0.519	0.656	0.381	0.299
Ensemble	<b>0.452**</b>	<b>0.589***</b>	<b>0.627***</b>	<b>0.248*</b>	<b>0.160</b>	<b>0.395***</b>	<b>0.593***</b>	<b>0.716</b>	<b>0.459**</b>	<b>0.357***</b>
TU			GT1					GT5		
	MAP	NDCG@100	MRR	P@5	P@10	MAP	NDCG@100	MRR	P@5	P@10
Model 2 (jm)	0.234	0.370	0.342	0.136	0.101	0.253	0.394	0.302	0.108	0.081
Log-linear (ours)	0.219	0.356	0.351	0.145	0.105	0.287	0.425	0.363	0.134	0.092
Ensemble	<b>0.271***</b>	<b>0.417***</b>	<b>0.403***</b>	<b>0.165***</b>	<b>0.121***</b>	<b>0.331***</b>	<b>0.477***</b>	<b>0.402***</b>	<b>0.156***</b>	<b>0.105***</b>

**Table 3: Comparison of Model 2, the log-linear model and an ensemble of the former on W3C, CERC and TU benchmarks. Significance of results is determined using a two-tailed paired randomization test [59] (\*\*\*  $p < 0.01$ ; \*\*  $p < 0.05$ ; \*  $p < 0.1$ ) with respect to the ensemble ranking (adjusted using the Benjamini-Hochberg procedure for multiple testing [9]).**

In terms of space complexity, parameters  $W_p$ ,  $W_c$  and  $b_c$ , in addition to the intermediate results, all require memory space proportional to their size. Considering (2) for a sequence of  $k$  words and batches of  $m$  instances, we require  $O(m \cdot k \cdot |C|)$  floating point numbers for every forward-pass to fit in-memory. While such an upper bound seems reasonable by modern computing standards, it is a severely limiting factor when considering large-scale communities and while utilizing limited-memory GPUs for fast computation.

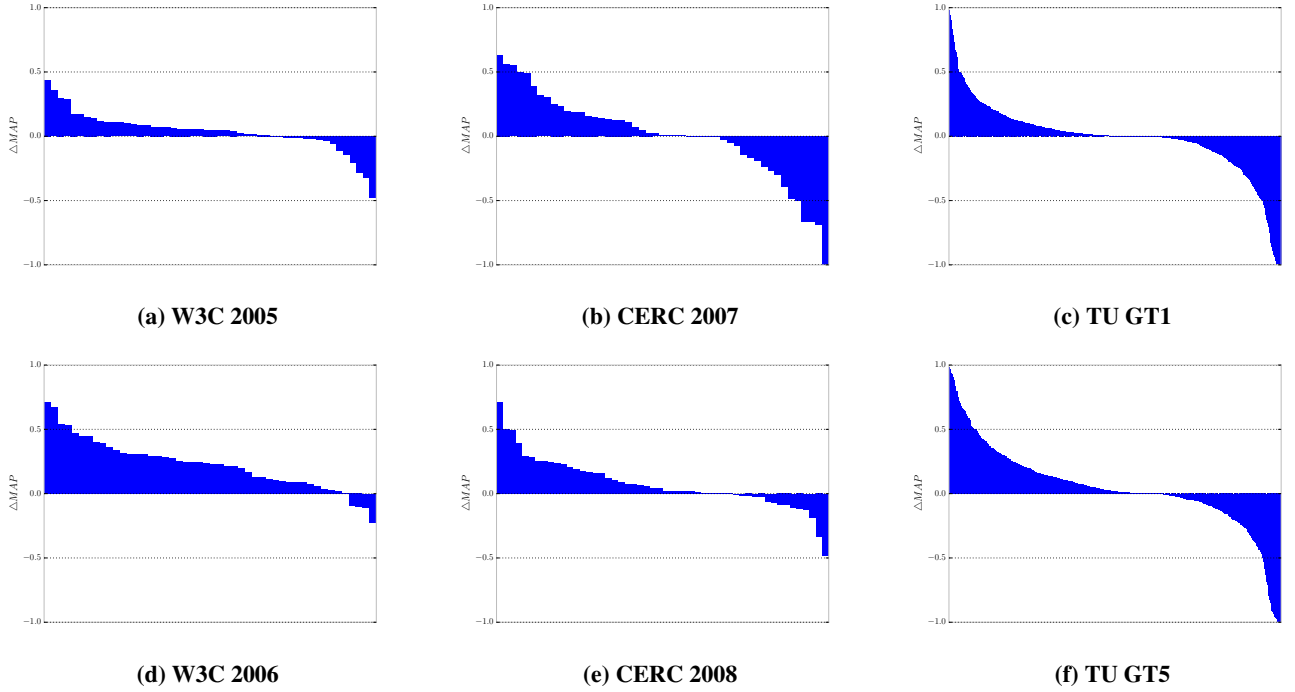
The inferential complexity of the vector space-based models for entity retrieval [19] depends mainly on the dimensionality of the vectors and the number of candidate experts. The dimensionality of the latent entity representations is too high for efficient nearest neighbor retrieval [29] due to the curse of dimensionality. Therefore, the time complexity for the LSI- and TF-IDF-based vector space models are respectively  $O(\gamma \cdot |C|)$  and  $O(|V| \cdot |C|)$ , where  $\gamma$  denotes the number of latent topics in the LSI-based method. As hyperparameters  $e$  and  $\gamma$  both indicate the dimensionality of latent entity representations, the time complexity of the LSI-based method is comparable to that of the log-linear model. We note that  $|V| \gg |C|$  for all benchmarks ( $|V|$  is between 18 to 91 times larger than  $|C|$ ) we consider in this paper and therefore conclude that the TF-IDF method loses to the log-linear model in terms of efficiency.

Compared to the unsupervised generative models of Balog et al., we have the profile-centric Model 1 and the document-centric Model 2 with inference time complexity  $O(n \cdot |C|)$  and  $O(n \cdot |D|)$ , respectively, with  $|D| \gg |C|$ . In the previous section we showed that the log-linear model always performs better than Model 1 and nearly always outperforms Model 2. Hence, our log-linear model generally achieves the expertise retrieval performance of Model 2 (or higher) at the complexity cost of Model 1 during inference.

## 5.4 Incremental indexing

Existing unsupervised methods use well-understood maximum-likelihood language models that support incremental indexing. We now briefly discuss the incremental indexing capabilities of our proposed method. Extending the set of candidate experts  $C$  requires the log-linear model to be re-trained from scratch as it changes the topology of the network. Moreover, every document associated with a candidate expert is considered as a negative example for all other candidates. While it is possible to reiterate over all past documents and only learn an additional row in matrix  $W_c$ , the final outcome is unpredictable.

If we consider a stream of documents instead of a predefined set  $D$ , the log-linear model can be learned in an online fashion. However, stochastic gradient descent requires that training examples are



**Figure 3: Difference of average precision between log-linear model and Model 2 [4] with Jelinek-Mercer smoothing per topic for W3C, CERC and TU benchmarks.**

picked at random such that the batched update rule (4) behaves like the empirical expectation over the full training set [11]. While we might be able to justify the assumption that documents arrive randomly, the  $n$ -grams extracted from those documents clearly violate this requirement.

Considering a stream of documents leads to the model forgetting expertise evidence as an (artificial) shift in the underlying distribution of the training data occurs. While such behavior is undesirable for the task considered in this paper, it might be well-suited for temporal expert finding [22, 52], where expertise drift over time is considered. However, temporal expert finding is beyond the scope for this paper and left for future work.

## 6. CONCLUSIONS

We have introduced an unsupervised discriminative, log-linear model for the expert retrieval task. Our approach exclusively employs raw textual evidence. Future work can focus on improving performance by feature engineering and incorporation of external evidence. Furthermore, no relevance feedback is required during training. This renders the model suitable for a broad range of applications and domains.

We evaluated our model on the W3C, CERC and TU benchmarks and compared it to state-of-the-art vector space-based entity ranking (based on LSI and TF-IDF) and language modeling (profile-centric and document-centric) approaches. The log-linear model combines the ranking performance of the best maximum-likelihood language modeling approach (document-centric) with inference time complexity linear in the number of candidate experts. We observed a notable increase in precision over existing methods. Analysis of our model’s output reveals a negative correlation between the per-query performance and ranking uncertainty: higher confidence (i.e., lower entropy) in the rankings produced by our approach often occurs together with higher rank quality.

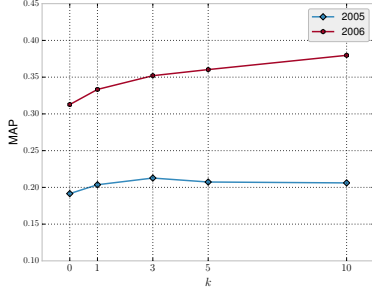
An error analysis of the log-linear model and traditional language models shows that the two make very different errors. These errors are mainly due to the semantic gap between query intent and the raw textual evidence. Some benchmarks expect exact query matches, others are helped by our semantic matching. An ensemble of methods employing exact and semantic matching generally outperforms the individual methods. This observation calls for further research in the area of combining exact and semantic matching.

One current limitation of our work is its scalability with respect to the number of candidate experts. We have started investigating trade-offs between model performance and time/space complexity. In the future we hope to apply scalable variants of this method on large-scale social media communities, for the purpose of determining topic ownership. While in this work we focus on expertise retrieval, the ideas we proposed can easily be transferred to the more general entity retrieval task. Moreover, our approach is likely to be applicable to authorship attribution and various other entity retrieval and prediction tasks.

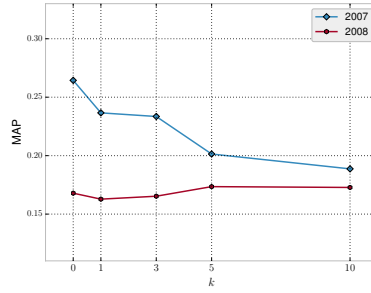
**Acknowledgments.** We thank Isaac Sijarmanual, Manos Tsagkias, Tom Kenter, Zhaochun Ren and Ke Tran for their useful comments and insights.

This research was supported by Amsterdam Data Science, the Dutch national program COMMIT, Elsevier, the European Community’s Seventh Framework Programme (FP7/2007-2013) under grant agreement nr 312827 (VOX-Pol), the ESF Research Network Program ELIAS, the Royal Dutch Academy of Sciences (KNAW) under the Elite Network Shifts project, the Microsoft Research Ph.D. program, the Netherlands eScience Center under project number 027.012.105, the Netherlands Institute for Sound and Vision, the Netherlands Organisation for Scientific Research (NWO) under project nrs 727.011.005, 612.001.116, HOR-11-10, 640.006.013, 612.066.-930, CI-14-25, SH-322-15, 652.002.001, 612.001.551, the Yahoo Faculty Research and Engagement Program, and Yandex. All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

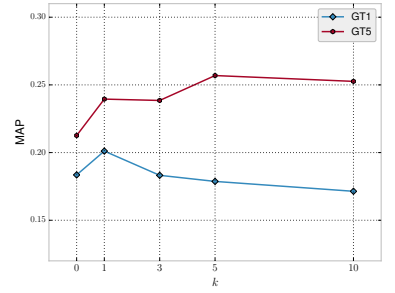
Computing resources were provided by the Netherlands Organisation for Scientific Research (NWO) through allocation SH-322-15 of the Cartesius system and by the Advanced School for Computing and Imaging (ASCI) by allocation of the Distributed ASCII Supercomputer 4 (DAS-4) system.



(a) W3C



(b) CERC



(c) TU

**Figure 4: Effect of query expansion by adding nearby terms in  $W_p$  (1) in traditional language models (Model 1 [4] with Jelinek-Mercer smoothing) for W3C, CERC and TU benchmarks.**

## APPENDIX

The derivative of (3) w.r.t. bias term  $b_c$  equals

$$\frac{\partial L(W_p, W_c, b_c)}{\partial b_c} = -\frac{1}{m} \left( \sum_{i=1}^m \frac{|d_{\max}|}{|d^{(i)}|} \sum_{j=1}^{|C|} P(c_j | d^{(i)}) \frac{\partial \log(P(c_j | w_1^{(i)}, \dots, w_n^{(i)}))}{\partial b_c} \right)$$

and w.r.t. an arbitrary matrix parameter  $\theta$  ( $W_p$  or  $W_c$ ):

$$\begin{aligned} \frac{\partial L(W_p, W_c, b_c)}{\partial \theta} &= -\frac{1}{m} \left( \sum_{i=1}^m \frac{|d_{\max}|}{|d^{(i)}|} \sum_{j=1}^{|C|} P(c_j | d^{(i)}) \frac{\partial \log(P(c_j | w_1^{(i)}, \dots, w_n^{(i)}))}{\partial \theta} \right) \\ &\quad + \frac{\lambda}{m} \sum_{i,j} \theta_{i,j}. \end{aligned}$$

Further differentiation for parameter  $\theta$  ( $W_p$ ,  $W_c$  or  $b_c$ ):

$$\begin{aligned} \frac{\partial \log(P(c_j | w_1, \dots, w_n))}{\partial \theta} &= \frac{1}{P(c_j | w_1, \dots, w_n)} \frac{\partial P(c_j | w_1, \dots, w_n)}{\partial \theta} \\ &= \frac{\frac{\partial \tilde{P}(c_j | w_1, \dots, w_n)}{\partial \theta} Z_2 - \tilde{P}(c_j | w_1, \dots, w_n) \frac{\partial Z_2}{\partial \theta}}{Z_2^2} \\ \frac{\partial Z_2}{\partial \theta} &= \sum_k \frac{\partial \tilde{P}(c_k | w_1, \dots, w_n)}{\partial \theta} \\ \frac{\partial \tilde{P}(c_j | w_1, \dots, w_n)}{\partial \theta} &= \sum_k \frac{\partial P(c_j | w_k)}{\partial \theta} \prod_{i \neq k} P(c_j | w_i) \end{aligned}$$

For a given candidate  $c_j$  and word  $w_i$ , following (1) we have

$$\begin{aligned} P(c_j | w_i) &= \frac{\tilde{P}(c_j | w_i)}{Z_1} \\ &= \frac{\exp((\sum_{k=1}^e W_{c_j,k} W_{p_k,i}) + b_{c_j})}{\sum_{l=1}^{|C|} \exp((\sum_{k=1}^e W_{c_l,k} W_{p_k,i}) + b_{c_l})} \end{aligned}$$

and consequently, with  $\mathbf{W}_{p_i}^\top$  denoting the  $i$ -th column of matrix  $W_p$ ,

$$\begin{aligned} \frac{\partial P(c_j | w_i)}{\partial \mathbf{W}_{c_j}} &= \frac{(Z_1 - \tilde{P}(c_j | w_i)) \tilde{P}(c_j | w_i) \mathbf{W}_{p_i}^\top}{Z_1^2} \\ \frac{\partial P(c_j | w_i)}{\partial b_{c_j}} &= \frac{(Z_1 - \tilde{P}(c_j | w_i)) \tilde{P}(c_j | w_i)}{Z_1^2} \quad (7) \\ \frac{\partial P(c_j | w_i)}{\partial \mathbf{W}_{p_i}^\top} &= \frac{(\mathbf{W}_{c_j} - \sum_{l=1}^{|C|} \mathbf{W}_{c_l}) \tilde{P}(c_j | w_i)}{Z_1} \end{aligned}$$

As can be seen in (7), the distributed representations of candidates  $c_j$  at time  $t+1$  are updated using the representation of words  $w_i$  at time  $t$  and vice versa.

## REFERENCES

- [1] The knowledge-based economy. Techn. report, Organisation for Economic Co-operation and Development, 1996.
- [2] P. Bailey, A. P. De Vries, N. Craswell, and I. Soboroff. Overview of the TREC 2007 enterprise track. In *TREC*, 2007.
- [3] K. Balog. *People Search in the Enterprise*. PhD thesis, University of Amsterdam, 2008.
- [4] K. Balog, L. Azzopardi, and M. de Rijke. Formal models for expert finding in enterprise corpora. In *SIGIR*, pages 43–50, 2006.
- [5] K. Balog, L. Azzopardi, and M. de Rijke. A language modeling framework for expert finding. *IPM*, 45:1–19, 2009.
- [6] K. Balog, Y. Fang, M. de Rijke, P. Serdyukov, and L. Si. Expertise retrieval. *Found. & Tr. in Information Retrieval*, 6(2-3):127–256, 2012.
- [7] I. Becerra-Fernandez. Role of artificial intelligence technologies in the implementation of People-Finder knowledge management systems. *Knowledge-Based Systems*, 13(5):315–320, 2000.
- [8] Y. Bengio, R. Ducharme, P. Vincent, and C. Janvin. A neural probabilistic language model. *JMLR*, 3:1137–1155, 2003.
- [9] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *JSTOR*, pages 289–300, 1995.
- [10] R. Berendsen, M. de Rijke, K. Balog, T. Bogers, and A. van den Bosch. On the assessment of expertise profiles. *JASIST*, 64(10):2024–2044, 2013.
- [11] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *COMPSTAT*, pages 177–186. Springer, 2010.
- [12] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender. Learning to rank using gradient descent. In *ICML*, pages 89–96, 2005.
- [13] Y. Cao, J. Liu, S. Bao, and H. Li. Research on Expert Search at Enterprise Track of TREC 2005. In *TREC*, pages 2–5, 2005.
- [14] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing (almost) from scratch. *JMLR*, 12(Aug):2493–2537, 2011.
- [15] N. Craswell, D. Hawking, A.-M. Vercoustre, and P. Wilkins. P@noptic expert: Searching for experts not just for documents. In *Ausweb Poster Proceedings*, pages 21–25, 2001.
- [16] N. Craswell, A. P. de Vries, and I. Soboroff. Overview of the TREC 2005 enterprise track. In *TREC*, 2005.
- [17] T. H. Davenport and L. Prusak. *Working Knowledge*. Harvard Business Review Press, 1998.
- [18] S. C. Deerwester, S. T. Dumais, and R. A. Harshman. Indexing by latent semantic analysis. *JASIS*, 1990.

- [19] G. Demartini, J. Gaugaz, and W. Nejdl. A vector space model for ranking entities and its application to expert search. In *ECIR*, pages 189–201. Springer, 2009.
- [20] L. Deng, X. He, and J. Gao. Deep stacking networks for information retrieval. In *ICASSP*, pages 3153–3157, 2013.
- [21] H. Fang and C. Zhai. Probabilistic models for expert finding. In *ECIR*, pages 418–430, Berlin, Heidelberg, 2007. Springer-Verlag.
- [22] Y. Fang and A. Godavarthy. Modeling the dynamics of personal expertise. In *SIGIR*, pages 1107–1110, 2014.
- [23] Y. Fang, L. Si, and A. P. Mathur. Discriminative models of integrating document evidence and document-candidate associations for expert search. In *SIGIR*, pages 683–690, 2010.
- [24] K. Fatahalian, J. Sugerman, and P. Hanrahan. Understanding the efficiency of gpu algorithms for matrix-matrix multiplication. In *SIGGRAPH HWS*, pages 133–137. ACM, 2004.
- [25] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, pages 249–256, 2010.
- [26] G. E. Hinton. Learning distributed representations of concepts. In *8th Annual Conference of the Cognitive Science Society*, volume 1, page 12, Amherst, MA, 1986.
- [27] T. Hofmann. Probabilistic latent semantic indexing. In *SIGIR*, pages 50–57. ACM, 1999.
- [28] P.-s. Huang, N. M. A. Urbana, X. He, J. Gao, L. Deng, A. Acero, and L. Heck. Learning deep structured semantic models for web search using clickthrough data. In *CIKM*, pages 2333–2338, 2013.
- [29] P. Indyk and R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *STOC*, pages 604–613. ACM, 1998.
- [30] M. Karimzadehgan and C. Zhai. Estimation of statistical translation models based on mutual information for ad hoc information retrieval. In *SIGIR*, pages 323–330. ACM, 2010.
- [31] R. Kiros, R. Salakhutdinov, and R. Zemel. Multimodal neural language models. In *ICML*, pages 595–603, 2014.
- [32] J. Kruger and D. Dunning. Unskilled and unaware of it: how difficulties in recognizing one’s own incompetence lead to inflated self-assessments. *J. Personality and Social Psych.*, 77(6):1121, 1999.
- [33] J. Krüger and R. Westermann. Linear algebra operators for gpu implementation of numerical algorithms. *ACM Transactions on Graphics*, 22(3):908–916, 2003.
- [34] H. Li and J. Xu. Semantic matching in search. *Found. & Tr. in Information Retrieval*, 7(5):343–469, June 2014.
- [35] T.-Y. Liu. *Learning to Rank for Information Retrieval*. Springer, 2011.
- [36] C. MacDonald and I. Ounis. Voting for candidates: adapting data fusion techniques for an expert search task. In *CIKM*, pages 387–396, 2006.
- [37] C. Macdonald and I. Ounis. Expert search evaluation by supporting documents. In *ECIR*, pages 555–563. Springer, 2008.
- [38] M. T. Maybury. Expert finding systems. Techn. Report MTR-06B000040, MITRE, 2006.
- [39] D. W. McDonald and M. S. Ackerman. Expertise recommender. In *CSCW*, pages 231–240, 2000.
- [40] T. Mikolov, M. Karafiat, L. Burget, J. Cernocky, and S. Khudanpur. Recurrent neural network based language model. In *Interspeech*, pages 1045–1048, 2010.
- [41] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, pages 3111–3119, 2013.
- [42] T. Mikolov, G. Corrado, K. Chen, and J. Dean. Efficient estimation of word representations in vector space. arXiv 1301.3781, 2013.
- [43] A. Mnih and G. Hinton. Three new graphical models for statistical language modelling. In *ICML*, pages 641–648, 2007.
- [44] A. Mnih and G. Hinton. A scalable hierarchical distributed language model. In *NIPS*, pages 1081–1088, 2008.
- [45] A. Mnih and K. Kavukcuoglu. Learning word embeddings efficiently with noise-contrastive estimation. In *NIPS*, pages 2265–2273, 2013.
- [46] G. Montavon, G. B. Orr, and K.-R. Müller. *Neural Networks: Tricks of the Trade*. Springer, 2012.
- [47] C. Moreira, B. Martins, and P. Calado. Using rank aggregation for expert search in academic digital libraries. In *Simpósio de Informática, INForum*, pages 1–10, 2011.
- [48] J. Pennington, R. Socher, and C. D. Manning. GloVe: Global Vectors for Word Representation. In *EMNLP*, pages 1532–1543, 2014.
- [49] D. Petkova and W. B. Croft. Hierarchical language models for expert finding in enterprise corpora. In *ICTAI ’06*, pages 599–606, 2006.
- [50] W. W. Powell and K. Snellman. The knowledge economy. *Annual review of sociology*, pages 199–220, 2004.
- [51] D. Rumelhart, G. Hinton, and R. Williams. Learning internal representations by back propagation. In *Parallel Distributed Processing*, pages 318–362. MIT Press, 1986.
- [52] J. Rybak, K. Balog, and K. Nørnvåg. Temporal expertise profiling. In *ECIR*, pages 540–546. Springer, 2014.
- [53] R. Salakhutdinov and G. Hinton. Semantic hashing. *Int. J. Approximate Reasoning*, 50(7):969–978, 2009.
- [54] P. Serdyukov and D. Hiemstra. Modeling documents as mixtures of persons for expert finding. In *ECIR*, pages 309–320. Springer, 2008.
- [55] P. Serdyukov, H. Rode, and D. Hiemstra. Modeling multi-step relevance propagation for expert finding. In *CIKM*, pages 1133–1142, 2008.
- [56] C. Shannon. A mathematical theory of communication. *Bell System Technical J.*, 27:379–423, 623–656, 1948.
- [57] J. A. Shaw, E. A. Fox, J. A. Shaw, and E. A. Fox. Combination of multiple searches. In *TREC*, pages 243–252, 1994.
- [58] Y. Shen, X. He, J. Gao, L. Deng, and G. Mesnil. A latent semantic model with convolutional-pooling structure for information retrieval. In *CIKM*, pages 101–110, 2014.
- [59] M. D. Smucker, J. Allan, and B. Carterette. A comparison of statistical significance tests for information retrieval evaluation. In *CIKM*, pages 623–632. ACM, 2007.
- [60] P. Sorg and P. Cimiano. Finding the right expert: Discriminative models for expert retrieval. In *KDIR*, pages 190–199, 2011.
- [61] TREC. Enterprise Track, 2005–2008.
- [62] D. van Dijk, M. Tsagkias, and M. de Rijke. Early detection of topical expertise in community question and answering. In *SIGIR*, 2015.
- [63] V. Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.
- [64] M. D. Zeiler. Adadelta: An adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.