

Entity Retrieval

Sisay Fissaha Adafre
School of Computing, DCU
Dublin 9, Ireland
sadafre@computing.dcu.ie

Maarten de Rijke and Erik Tjong Kim Sang
ISLA, University of Amsterdam
Kruislaan 403, 1098SJ Amsterdam, The Netherlands
mdr,erikt@science.uva.nl

Abstract

Generalizing recent attention to retrieving entities and not just documents, we introduce two entity retrieval tasks: list completion and entity ranking. For each task, we propose and evaluate several algorithms. One of the core challenges is to overcome the very limited amount of information that serves as input—to address this challenge we explore different representations of list descriptions and/or example entities, where entities are represented not just by a textual description but also by the description of related entities. For evaluation purposes we make use of the lists and categories available in Wikipedia. Experimental results show that cluster-based contexts improve retrieval results for both tasks.

Keywords

Entity retrieval, Wikipedia, language modeling

1 Introduction

Both commercial systems and the information retrieval community are displaying an increasing interest in not just returning web pages or other documents in response to a user’s query but “objects,” “entities” or their properties. E.g., various web search engines recognize specific types of entity (such as books, CDs, restaurants), and list these separately from the standard document-oriented hit list. Enterprise search provides another example [5], as has also been recognized within the TREC Enterprise track. In its 2005 and 2006 editions, the track featured an expert finding task [6] where systems return a list of entities (people’s names) who are knowledgeable about a certain topic (e.g., “web standards”).

This emerging area of *entity retrieval* differs from traditional document retrieval in a number of ways. Entities are not represented directly (as retrievable units such as documents), and we need to identify them “indirectly” through occurrences in documents. Entity retrieval systems may initially retrieve documents (pertaining to a given topic or entity) but they must then extract and process these documents in order to return a ranked list of entities [20]. In order to understand the issues at hand, we propose two entity retrieval tasks (building on a proposal launched in the run-up to INEX 2006 [7] and scheduled to be implemented at INEX 2007): *list completion* and *entity ranking*.

The *list completion* task is defined as: given a topic text and a number of examples, the system has to produce further examples. I.e., given a topic description, a set of entities S and a number of example entities e_1, \dots, e_n in S that fit the description, return “more examples like e_1, \dots, e_n ” from S that fit the description. E.g., given the short description *tennis players* and two example entities such as *Kim Clijsters* and *Martina Hingis*, entities such as tennis tournaments or coaches are not relevant. Instead, the expected set should include only individuals who are or have been professional tennis players. In the *entity ranking* task, a system has to return entities that satisfy a topic described in natural language text. I.e., given a set of entities S and a topic statement t , return elements of S that satisfy t . For example, let S denote a set of Dutch people; then “Dutch actors,” “Dutch politicians,” “Dutch artists,” etc., are some of the typical topic statements t that we envisage for this task.

The main research questions we address concern the ways in which we represent entities and in which we match topics and entities. As we will see, providing a sufficiently rich description of both topics and entities to be able to rank entities in an effective manner, is one of the main challenges. We address this challenge by using several contextual models.

For evaluation purposes we make use of Wikipedia, the online encyclopedia. The decision for using Wikipedia for this task is based on practical and theoretical considerations. Wikipedia contains a large set of lists that can be used for generating the necessary test data, and also assessing the outputs of our methods. Also, with its rich structure Wikipedia offers an interesting experimental setting where we can experiment with different features, both content-based and structural. Finally, by using Wikipedia’s lists, we can avoid the information extraction task of *identifying entities* in documents and focus on the retrieval task itself, instead. Below, we will only consider entities available in Wikipedia, and we will identify each entity with its Wikipedia article.¹

The remainder of the paper is organized as follows. First, we provide background material and related work on working with Wikipedia, list questions, and contextual models. After that we turn to the list completion task, proposing and evaluating a number of algorithms. We then do the same for the entity ranking task before concluding the paper.

¹ We used the XML version of the English Wikipedia corpus made available by Denoyer and Gallinari [8]. It contains 659,388 articles, and has annotations for common structural elements such as article title, sections, paragraphs, sentences, and hyperlinks.

2 Background

Mining/Retrieval against Wikipedia Wikipedia has attracted interest from researchers in disciplines ranging from collaborative content development to language technology, addressing aspects such as information quality, users motivation, collaboration pattern, network structures, e.g., [25]. Several publications describe the use of Wikipedia as a resource for question answering and other types of IR systems; see e.g., [1, 10, 17]. Wikipedia has been used for computing word semantic relatedness, named-entity disambiguation, text classification, and as a document collection in various retrieval and knowledge representation tasks, e.g., [11].

Entity Retrieval List queries are a common types of web queries [22]. The TREC Question Answering track has recognized the importance of list questions [23]; there, systems have to return two or more instances of the class of entities that match the description in the list question. List questions are often treated as (repeated) factoids, but special strategies are called for as answers may need to be collected from multiple documents [4].

Recognizing the importance of list queries, Google Sets allows users to enter some instances of a concept and retrieve others that closely match the examples provided [13]. Ghahramani and Heller [12] developed an algorithm for completing a list based on examples using machine learning techniques. A proposed INEX entity retrieval task, with several tasks will likely be run during 2007 [7].

Our entity retrieval tasks are related to ontological relation extraction [14], where a combination of large corpora with simple manually created patterns are often used. Wikipedia, as a corpus, is relatively small, with much of the information being presented in a concise and non-redundant manner. Therefore, pattern-based methods may have limited coverage for the entity retrieval tasks that we consider.

Document expansion and contextual IR Enriching the document representation forms an integral part of the approach we propose in this paper. Though, in the past, application of document expansion techniques, particularly document clustering, has shown mixed results in document retrieval settings, recent studies within the language modelling framework provide new supporting evidence of the advantages of using document clusters [19]. Due to the nature of the tasks defined in this paper, the cluster hypothesis which states that “closely associated documents tend to be relevant to the same request” [16] provides for an intuitive starting point in designing our methods. Specifically, for each entity (or article) a precomputed cluster will be used to supply it with contextual information, much in the spirit of the work done by Azopardi [2] and Liu and Croft [19].

3 Task 1: List Completion

The main challenge of the list completion task is that the topic statement, example entity descriptions, and, more generally, entity descriptions in Wikipedia, tend

to be very short. Therefore, a straightforward retrieval baseline may suffer from poor recall. Hence, in our modeling we will address several ways of representing the topic statement and example entities.

We model the list completion task as follows: *what is the probability of a candidate e belonging to the list defined by the topic statement t and example entities e_1, \dots, e_n ?* We determine $p(e|t, e_1, \dots, e_n)$ and rank entities according to this probability. To estimate $p(e|t, e_1, \dots, e_n)$, we proceed in two steps: (1) select candidate entities, and (2) rank candidate entities. More formally,

$$p(e|t, e_1, \dots, e_n) \propto \chi_C \cdot \text{rank}(e; t, e_1, \dots, e_n),$$

where χ_C is a characteristic function for a set of selected candidate entities C and $\text{rank}(\cdot)$ is a ranking function. Below, we consider alternative definitions of the function χ_C and we describe two ranking functions. First, though, we define so-called entity neighborhoods that will be used in the candidate selection phase: to each individual entity e they associate additional entities based on e 's context, both in terms of link structure and contents.

3.1 Entity Neighborhoods

In the context of a hypertext documents, identification of a cluster typically involves searching for graph structures, where co-citations and bibliographic couplings provide importance features. Fissaha Adafre and de Rijke [9] describe a Wikipedia specific clustering method called *LTRank*. Their clustering method primarily uses the co-citation counts. We provide a slight extension that exploits the link structure (both incoming and outgoing links), article structure, and content. In Wikipedia, the leading few paragraphs contain essential information about the entity being described in the articles serving as summary of the content of the article; we use the first five sentences of the Wikipedia article as a representation of the content of the article. Our extension of the *LTRank* method for finding the neighborhood $\text{neighborhood}(e)$ of an entity e is summarized in Figure 1. With this definition we can return to the first phase in our approach: *candidate entity selection*.

3.2 Candidate Entity Selection

To perform the candidate entity selection step, we use a two part representation of entities (Wikipedia articles). Each entity e is represented using (1) the textual content of the corresponding article a_e , and (2) the list of all entities in the set of $\text{neighborhood}(e)$ defined above. We propose four candidate entity selection methods, that exploit this representation in different ways.

B-1. Baseline: Retrieval Here we rank entities by the similarity of their content part to a query consisting of the topic statement t and the titles t_{e_1}, \dots, t_{e_n} of the example entities. We used a simple vector space retrieval model for computing the similarity. The top n retrieved documents constitute the baseline candidate set C_1 .

- Given a Wikipedia article a_e of an entity e , collect the titles of pages with links to or from a_e , as well as the words in the first five sentences of a_e . Let $long(a_e)$ be the resulting bag of terms; this is the *long* representation of a_e .
- Given a Wikipedia article a_e , rank all articles w.r.t. their content similarity to $long(a_e)$; we use a simple vector space model for the ranking. This produces a ranked list $L_{a_e} = a_{e_1}, \dots, a_{e_n}, \dots$.
- Given a Wikipedia article a_e , consider the titles t_1, \dots, t_k of the top k articles in the list L_{a_e} . Represent a_e as the bag of terms $short(a_e) = \{t_1, \dots, t_k\}$; we call this the *short* representation of a_e .
- For each Wikipedia article a_e , rank the short representations of other Wikipedia articles w.r.t. their content similarity to $short(a_e)$; again, we use a simple vector space model for the ranking. This produces a ranked list L'_{a_e} . The *neighborhood*(e) is defined to be the set of top l articles in L'_{a_e} whose similarity score is above some threshold α .

Fig. 1: An extension of LTRank [9]. Our extension is in the first step, where we add outgoing links and the first 5 sentences of a_e . For the experiments in this paper, we took $k = 10$, $l = 100$, and $\alpha = 0.3$.

B-2. Neighborhood search Our second candidate selection method matches the titles of the example entities against the neighborhoods of Wikipedia articles.

$$C_2 = \{e | \bigvee_i (e_i \in neighborhood(e))\}$$

B-3. Neighborhood and Topic statement search Here we take the union of the entities retrieved using the topic statement, and method B-2 described above. First, we rank entities by the similarity of their content part to a query which corresponds to the topic statement t . Here again, we used a simple vector space similarity measure to compute the similarity. We take the top k entities ($k = 200$ in this paper) which constitute the first set, $C_{3.1}$. We then take all entities that contain at least one example entity in their neighborhood as with B-2, i.e.,

$$C_{3.2} = \{e | \bigvee_i (e_i \in neighborhood(e))\}.$$

The final candidate set is simply the union of these two sets, i.e., $C_3 = C_{3.1} \cup C_{3.2}$.

B-4. Neighborhood and Definition search This method is similar to the method B-3. But instead of taking the topic statement t as a query for ranking entities (in the set $C_{3.1}$ above), we take the definitions of the example entities e_1, \dots, e_n , where the first sentence of the Wikipedia article a_e of an entity e to be its definition; stopwords are removed.

3.3 Candidate Entity Ranking

We compare two methods that make use of the content of articles for ranking the entities generated by the previous step. Particularly, we apply the following two methods: Bayesian inference [12] and relevance-based language models [18]. Both methods provide a

mechanism for building a model of the concept represented by the example set. These two algorithms are developed for a task which closely resembles our task definitions, i.e., given a limited set of examples, find other instances of the concept represented by the examples. In the next paragraphs, we briefly discuss these methods.

C-1. Bayesian Inference Ghahramani and Heller [12] addressed the entity ranking task in the framework of Bayesian Inference. Given n example entities, e_1, \dots, e_n , and candidate entity e , the ranking algorithm is given by

$$score(e) = \frac{P(e, e_1, \dots, e_n)}{P(e)P(e_1, \dots, e_n)}. \quad (1)$$

To compute Eq. 1, a parameterized density function is posited. We list all terms $t_{e_{1.1}}, \dots, t_{e_{1.k_1}}, \dots, t_{e_{n.k_n}}$ occurring in the example entities. Then, each candidate entity e is represented as a binary vector where vector element $e_{i,j}$ corresponds to the j -th term from article a_{e_i} of the i -th example instance and assumes 1 if $t_{e_{i,j}}$ appears in the article for the entity e and 0 otherwise. It is assumed that the terms $e_{i,j}$ are independent and have a Bernoulli distribution θ_j with parameters α_j and β_j ; see [12]. In sum, Eq. 1 is rewritten to:

$$score(e) = c + \sum_{j=1}^N q_j e_{.,j},$$

where the summation ranges over the binary vector representation of e , and

$$c = \sum_j (\log(\alpha_j + \beta_j) - \log(\alpha_j + \beta_j + n) + \log(\beta_j + n - \sum_{i=1}^n e_{i,j}) - \log(\beta_j)),$$

while

$$q_j = \log(\alpha_j + \sum_{i=1}^n e_{i,j}) - \log(\alpha_j) + \log(\beta_j) - \log(\beta_j + n - \sum_{i=1}^n e_{i,j})$$

For given values of α_j and β_j , the quantity q_j assigns more weights to terms that occur in most of the example entities. Therefore, a candidate instance e_i will be ranked high if it contains many terms from the example instances and the $e_{i,j}$ receive high weights from the q_j s.

C-2. Relevance Models Lavrenko and Croft [18] proposed so-called relevance-based language models for information retrieval. Given n example entities, e_1, \dots, e_n , and the candidate e from the candidate set C , the ranking function is given by the KL-divergence between two relevance models:

$$score(e) = KL(P_{e_1, \dots, e_n} || P_e),$$

where P_{e_1, \dots, e_n} is the relevance model of the example entities, and P_e is the language model induced from the Wikipedia article for entity e . The relevance models are given by

$$\begin{aligned} P(w|e_1, \dots, e_n) &= \sum_{e \in W} P(w|e) \cdot P(e|e_1, \dots, e_n) \\ P(e|e_1, \dots, e_n) &= \begin{cases} 1/n & \text{if } e \in \{e_1, \dots, e_n\} \\ 0 & \text{otherwise} \end{cases} \\ P(w|e) &= \frac{\#(w, e)}{|e|}, \end{aligned}$$

where W is the collection (Wikipedia), and w represents the terms in the Wikipedia article for entity e . The KL divergence will be small for entities that more closely resemble the example entities in terms of their descriptions.

Summary Both of the ranking methods outlined above return a ranked list of candidate entities. We normalize the scores using

$$\text{score}_{\text{norm}} = \frac{\text{score}_{\text{MAX}} - \text{score}}{\text{score}_{\text{MAX}} - \text{score}_{\text{MIN}}},$$

and take those candidate entities for which the normalized score lie above empirically determined threshold ($\text{score}_{\text{norm}} > 0.5$). The resulting set will be assessed.

3.4 Experimental Set-up

The performance of our approach to the list completion task depends on the performance of the two sub-components: candidate selection and candidate ranking. We conduct two sets of experiments, one to determine the effectiveness of the candidate selection methods, and a second to determine the effectiveness of the overall approach. We are especially interested in the contribution of using the neighborhoods of entities.

The Wikipedia lists serve as our gold standard. We selected a random sample of 30 lists (the topics) from Wikipedia. We chose relatively homogeneous and complete lists, and excluded those that represent a mixture of several concepts. We take 10 example sets for each topic. Each example set consists of a random sample of entities from the Wikipedia list for the topic. We run our system using each of these 10 example sets as a separate input. The final score for each topic is then the average score over the ten separate runs. In the experiments in this section, we assume that each example set contains two example instances. This choice is mainly motivated by our assumption that users are unlikely to supply many examples.

The results are assessed based on the following scores: $P@20$ (number of correct entities that are among the top 20 in the ranked list), *precision* (P; number of correct entities that are in the ranked list, divided by size of the ranked list), *recall* (R; number of correct entities that are in the ranked list, divided by the number of entities in the Wikipedia list) and *F-scores* (F; harmonic mean of the recall and precision values).

In order to test if the differences among the methods measured in terms of F-scores is statistically significant, we applied the two-tailed Wilcoxon matched pair signed-ranks test (for $\alpha = 0.05$ and $\alpha = 0.005$).

3.5 Results

First, we assess the methods we used for candidate selection. Following this, we present the evaluation results of the overall system.

Candidate selection Table 1 shows results of the evaluation of the candidate selection module. The figures are averages over all topics and all sets of example entities. The values are relatively low. Retrieving additional candidates using terms derived either from the

Selection method	P	R	Result set size
B-1 (Top $k = 500$)	0.042	0.235	500
B-2	0.142	0.236	206
B-3	0.089	0.311	386
B-4	0.093	0.280	367

Table 1: Performance on the candidate selection sub-task.

Candidate selection	Candidate ranking	P	R	F	P@20
B-1	C-1	0.100	0.068	0.058	0.128
	C-2	0.203	0.046	0.060	0.144
B-2	C-1	0.172	0.163	0.136	0.205
	C-2	0.227	0.142	0.137	0.231
B-3	C-1	0.121	0.236	0.136	0.196
	C-2	0.188	0.210	0.151	0.249
B-4	C-1	0.140	0.202	0.142	0.201
	C-2	0.204	0.209	0.158	0.248

Table 2: Performance on the entire list completion task. Best scores per metric in boldface.

definition of the entities or topic statement improves recall to some extent. The recall values for method B-3 are the best. This suggests that the terms in the topic are more accurate than the terms automatically derived from the definitions.

The neighborhood-based methods achieve better recall values while returning fewer number of candidates (cf. the last column of Table 1).

Overall results Table 2 shows the scores resulting from applying the two ranking methods C-1 and C-2 on the output of different candidate selection methods. The first column of Table 2 shows the different candidate selection methods; the second column shows the ranking methods.

The neighborhood-based combinations outperform the baselines at the $\alpha = 0.005$ significance level (when considering F-scores). The combination of C-2 (*Relevance model*) with B-4 (*Neighborhood plus Definition Terms*) input outperforms both the B-2 + C-1 and B-2 + C-2 combinations at the $\alpha = 0.05$ significance level. Generally, the C-2 ranking method has a slight edge over the C-1 method on most inputs. Furthermore, retrieving additional candidates using either the topic statement or the definition terms improves results, especially when used in combination with the C-2 ranking method.

3.6 Error Analysis

A closer look at the results for the individual topics reveals a broad range of recall values. The recall values for the topics *North European Jews*, *Chinese Americans*, *French people*, and *Miami University alumni* are very low. On the other hand, the topics *Indian Test cricketers*, *Revision control software*, *Places in Norfolk*, and *Cities in Kentucky* receive high recall scores. For the neighborhood-based methods, there is some correlation between the composition of the neighborhoods corresponding to the example entities and the results obtained. For example, the neighborhoods corresponding to the example entities for the topic *Indian Test cricketers* contain Indian cricket players. On the

other hand, the neighborhoods corresponding to the example entities for the topic *Chinese Americans* contain individuals from the USA, most of whom are not Chinese Americans, and have very little in common except for the features identified by the topic titles, which are too specific.

4 Task 2: Entity Ranking

The goal of the entity ranking task is to retrieve a subset of a given set of entities that satisfy a topic statement. More formally, let E , a set of entities, be given. We rank entities according to the probability $p(t|e)$, where e ranges over elements of E and t is a topic statement. We present different methods of estimating $p(t|e)$. These methods are organized along two dimensions; along one we consider richer representations of the topic statement t , along the other we consider different ways of representing entities.

4.1 Topic Representations

We compare two types of topic representation which we describe below.

F-1. Baseline As our baseline, we only remove stopwords from the topic statements. No further processing is done on the topic statement.

F-2. Topic expansion In addition to removing stopwords, we enrich the topic by incorporating additional terms based on the method proposed in [21]. We assume the top n ($n = 5$) articles returned using the *Collection smoothing method* (see below) with $\lambda = 0.9$ as being relevant. Extra terms are added based on the log ratio of their likelihood in terms of the model for relevant articles to their likelihood in terms of the model for whole entity set.

4.2 Entity Representations

We now introduce several ways of representing entities, all in terms of two or three part mixture models. We start with our baseline approach.

G-1. Baseline As explained in the introduction, the entities we consider are titles of Wikipedia articles. Hence, the simplest representation of an entity e is its associated Wikipedia article a_e . As usual, the topic t is represented by a set of terms: $t = \{t_1, \dots, t_k\}$; we write $c(t_i, a_e)$ to indicate the number of times t_i occurs in a_e . Each topic term is assumed to be generated independently, and so the topic likelihood is obtained by taking the product across all the terms in the topic:

$$p(t|e) = \prod_{t_i \in t} p(t_i|a_e)^{c(t_i, t)}.$$

In our baseline approach, we estimate $p(t_i|a_e)$ by taking the maximum likelihood estimate of t_i in a_e :

$$p_{baseline}(t_i|a_e) = p_{MLE}(t_i|a_e) = \frac{c(t_i, e)}{|a_e|},$$

where $|a_e|$ the total number of term occurrences in a_e .

G-2. Collection smoothing Since $p_{MLE}(t_i|a_e)$ may contain zero probabilities it is standard to employ smoothing [24]. Therefore, we smooth the maximum likelihood estimate, i.e., $p_{MLE}(t_i|e)$, against a general model estimated from the whole Wikipedia collection as follows:

$$p(t_i|a_e) = \lambda \cdot p_{MLE}(t_i|a_e) + (1 - \lambda) \cdot p_{MLE}(t_i|W), \quad (2)$$

where the latter is the maximum likelihood estimate of t_i in W , the entire Wikipedia corpus.

G-3. Context models 1: A generic approach In this paragraph and the next, we introduce two context models, both give rise to three part mixture models, involving the entity, the context, and the collection. The intuition behind these models is that a more focused context should be more accurate in capturing the topic of the entity, thus producing a more meaningful representation of the entity than the entire collection. The first context model we consider is generic, and does not exploit special features of the Wikipedia corpus. Specifically, we use probabilistic latent semantic analysis (PLSA, [15]) to induce a context for every entity e . Given an entity e , a latent class z is selected with probability $p(z|e)$, and given the class z , terms t_i are generated with probability $p(t_i|z)$. Then the following context model is obtained:

$$p_{PLSA}(t_i|e) = \sum_{z \in Z} p(t_i|z) \cdot p(z|e), \quad (3)$$

where Z is the set of latent variables considered (in our experimental evaluation we fix $|Z| = 20$). The probabilities $p(t_i|z)$ and $p(z|e)$ are estimated using the EM algorithm as described in [15]. Putting Eq. 3 together with the smoothed model (Eq. 2), we obtain the following:

$$p_{TOPIC}(t_i|e) = \lambda_1 \cdot p_{MLE}(t_i|e) + \lambda_2 \cdot p_{PLSA}(t_i|e) + (1 - \lambda_1 - \lambda_2) \cdot p_{MLE}(t_i|W), \quad (4)$$

where $\lambda_1, \lambda_2 \in [0, 1]$ and $\lambda_1 + \lambda_2 \leq 1$.

G-4. Context models 2: A Wikipedia-specific approach The second context model we consider in this paper exploits specific features of the Wikipedia corpus. We use the method summarized in Figure 1 for estimating the Wikipedia specific context model. Specifically, given an entity e , consider the neighborhood of e as produced by the algorithm in Figure 1. Assume $neighborhood(e) = d_1(e), \dots, d_k(e)$ for e . Then,

$$p_{WIKI}(t_i|e) = \lambda_1 \cdot p_{MLE}(t_i|e) + \lambda_2 \cdot p_{LTS}(t_i|d(e)_1, \dots, d(e)_k) + (1 - \lambda_1 - \lambda_2) \cdot p_{MLE}(t_i|W), \quad (5)$$

where, as before, $\lambda_1, \lambda_2 \in [0, 1]$ and $\lambda_1 + \lambda_2 \leq 1$. $p_{LTS}(t_i|d_1, \dots, d_k)$ is the context model, which gives the likelihood of the term t_i in the cluster consisting of the context documents, d_1, \dots, d_k .

4.3 Experimental Set-up

The experiments in this section are aimed at gaining insight into the contributions (for the *Entity Ranking* task) of the different topic and document representation methods introduced previously. We used

Document representation		Topic representation			
		F-1		F-2	
↓	Parameters	P@10	R-Prec	P@10	R-Prec
G-1	–	0.587	0.399	0.217	0.211
G-2	$\lambda = 0.9$	0.567	0.428	0.567	0.413
G-3	$\lambda_1 = 0.7, \lambda_2 = 0.2$	0.583	0.448	0.570	0.426
G-4	$\lambda_1 = 0.7, \lambda_2 = 0.2$	0.623	0.476	0.580	0.464

Table 3: *Entity ranking results: average values over all topics.*

Wikipedia’s hierarchical categories for generating the data for evaluating the methods. We selected a random sample of 30 Wikipedia lists, i.e., *main entity sets*. For each *main entity set*, we selected a subset of entities and the associated topic. Each of the alternative approaches presented in this section rank entities in the *main entity set*. The ranked list is assessed based on the following precision scores: R-Precision (the fraction of the number of correct entities for each topic that are among the top n entities returned, where n is the size of the sublist we are seeking), and $p@10$ (number of correct entities for each topic that are among the top 10 entities returned).

We applied the two-tailed Wilcoxon matched pair signed-ranks test to determine whether the differences among the methods as measured in terms of R-Precision scores are statistically significant ($\alpha = 0.05$).

4.4 Results

Table 3 shows the result of the different runs. In the tables, the columns *Parameters*, $p@10$ and *R-Prec* correspond to the parameter settings, precision at 10, and R-Precision. The parameter settings are the optimal mixing values for the given model. As the results show, the baseline method, which uses the maximum likelihood estimation without term expansion (F-1 + G-1), performs relatively well. However, term expansion hurts performance of the baseline method due to the MLE estimation (the extended topic tends to be assigned zero probability). All methods outperform the F-2 + G-1 combination. The ranking method that uses the *Wikipedia Specific Context model* (G-4) outperforms the *Collection-based context* and the MLE method at a significance level of $\alpha = 0.05$. G-4 performs better than G-3 at the significance level of $\alpha = 0.1$. Term expansion tends to hurt performance as can be seen from the general pattern in Table 3.

5 Discussion

Entity retrieval vs information extraction

The tasks considered in this paper, i.e., *list completion* and *entity ranking*, share a common overall goal. They both aim at identifying entities that share certain characteristics. In this respect, they resemble tasks commonly addressed in Information Extraction (IE), such as *named entity recognition* and *relation extraction*. However, there are important distinctions between traditional IE and the entity retrieval tasks we consider. First, in typical IE scenarios, the entities are embedded in a text, and the aim is to extract or recognise

occurrences of these entities in the text. Systems commonly use surrounding contextual information, and redundancy information to recognise the entities in the text. The inputs to these systems are documents that may contain one or more occurrences of the target entities. In contrast, in the entity retrieval tasks that we consider, the entities are represented by documents which provide descriptive information about them—typically, there is a one-to-one relation between the entities and the documents. In our setting, then, we abstract away from the recognition phase so that we are able to zoom in on the retrieval task only—unlike, e.g., the expert finding scenarios currently being explored at TREC, that do require participating systems to create effective combinations of extraction and retrieval [3].

One or two tasks? Although the list completion and entity ranking tasks are similar at an abstract level, a closer look at the specific details reveals important differences which necessitated task-specific approaches. One aspect concerns the size of the input; for the list completion task, the inputs are example entities with/without topic statements, and the candidates are all Wikipedia entries. On the other hand, the inputs for the entity ranking task consist of the topic statements only, and the candidates are entities in a particular Wikipedia list, such as, e.g., the List of Countries, which is obviously much smaller and more homogeneous than the entire Wikipedia collection.

The result of the list completion task shows that traditional information retrieval methods significantly underperform for selecting initial candidates from all of Wikipedia. This affects the overall score of the method as subsequent processing makes use of the output of this step. On the other hand, preclustering of Wikipedia articles led to much better performance. The re-ranking methods showed comparable performance results, with the relevance feedback method having a slight edge over the Bayesian method.

In the entity ranking task, we compared different ways of enriching the topic statements and document representations. As to the former, we added more terms to the topic description, and in the latter, we applied document modeling techniques that capture natural groupings that may exist in the target list. The results show that automatic addition of terms using relevance feedback methods seems to hurt performance. Here again, our notion of neighbourhood seems to capture the natural groupings in the target list better than the topic modeling method we considered in this paper.

By comparing the absolute scores of the two tasks, it seems safe to conclude that the richer input used for the entity ranking task (working with a specific list rather than all of Wikipedia) leads to higher scores.

6 Conclusion

We described, and proposed solutions for, two types of entity retrieval tasks, *list completion* and *entity ranking*. We conducted two sets of experiments in order to assess the proposed methods, which focused on enriching the two key elements of the retrieval tasks, i.e.,

Topic statements and *Example entities*.

For the list completion task, the methods that used the titles of the example entities and the topic statements or definition terms performed better. All methods that used a context set consisting of related articles significantly outperformed a simple document-based retrieval baseline that does not use the related articles field.

For the entity ranking task, the method that used a context set of related articles also performed better than most of the alternatives we considered. Here, we used the related articles to provide contextual information for the entity description when computing the similarity between the topic statement and entity description. Our notion of related articles improves results when used both as a means of retrieving initial candidates and for providing contextual information during similarity computations.

Our results are limited in a number of ways. For example, entities are represented primarily by the combination of the content of their Wikipedia articles (as a bag of words) and a precomputed set of related articles. We need to explore other—rich—representations of the content, e.g., phrases or anchor text, and also other concepts of relatedness, e.g., the Wikipedia categories.

Acknowledgments

This research was supported by the Netherlands Organisation for Scientific Research (NWO) under project numbers 017.001.190, 220-80-001, 264-70-050, 354-20-005, 600.065.120, 612-13-001, 612.000.-106, 612.066.302, 612.069.006, 640.001.501, 640.002.-501, and by the E.U. IST programme of the 6th FP for RTD under project MultiMATCH contract IST-033104.

References

- [1] K. Ahn, J. Bos, J. R. Curran, D. Kor, M. Nissim, and B. Webber. Question answering with QED at TREC-2005. In *Proceedings of TREC 2005*, 2005.
- [2] L. Azzopardi. Incorporating context within the language modeling approach for ad hoc information retrieval. *SIGIR Forum*, 40(1):70–70, 2006.
- [3] K. Balog and M. de Rijke. Finding experts and their details in e-mail corpora. In *15th International World Wide Web Conference (WWW2006)*, 2006.
- [4] J. Chu-Carroll, K. Czuba, J. Prager, A. Ittycheriah, and S. Blair-Goldensohn. IBM's PIQUANT II in TREC 2004. In *Proceedings TREC 2004*, 2004.
- [5] N. Craswell, D. Hawking, A. M. Vercoustre, and P. Wilkins. P@noptic expert: Searching for experts not just for documents. In *Ausweb*, 2001.
- [6] N. Craswell, A. de Vries, and I. Soboroff. Overview of the TREC 2005 Enterprise Track. In *Proceedings of TREC 2005*, 2006.
- [7] A. de Vries and N. Craswell. XML entity ranking track, 2006. URL: <http://inex.is.informatik.uni-duisburg.de/2006/xmlSearch.html>.
- [8] L. Denoyer and P. Gallinari. The Wikipedia XML Corpus. *SIGIR Forum*, 2006.
- [9] S. Fissaha Adafre and M. de Rijke. Discovering missing links in Wikipedia. In *Proceedings of LinkKDD-2005 Workshop*, 2005.
- [10] S. Fissaha Adafre and M. de Rijke. Estimating importance features for fact mining (with a case study in biography mining). In *RIAO 2007*, 2007.
- [11] E. Gabrilovich and S. Markovitch. Overcoming the brittleness bottleneck using wikipedia: Enhancing text categorization with encyclopedic knowledge. In *AAAI '06*, 2006.
- [12] Z. Ghahramani and K. A. Heller. Bayesian sets. In *NIPS 2005*, 2005.
- [13] Google, 2006. GoogleSets. URL: <http://labs.google.com/sets>, accessed on 04-10-2006.
- [14] M. A. Hearst. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of COLING 1992*, pages 539–545, 1992.
- [15] T. Hofmann. Probabilistic latent semantic analysis. In *Proceedings of Uncertainty in Artificial Intelligence, UAI'99*, 1999.
- [16] N. Jardine and C. J. van Rijsbergen. The use of hierarchic clustering in information retrieval. *Information Storage Retrieval*, 7(5):217–240, 1971.
- [17] V. Jijkoun and M. de Rijke. WiQA: Evaluating Multi-lingual Focused Access to Wikipedia. In *EVIA 2007*, 2007.
- [18] V. Lavrenko and W. B. Croft. Relevance models in information retrieval. In *Language Modeling for Information Retrieval*. Kluwer Academic Publishers, 2003.
- [19] X. Liu and W. B. Croft. Cluster-based retrieval using language models. In *SIGIR '04*, pages 186–193, 2004.
- [20] D. Petkova and W. B. Croft. Hierarchical language models for expert finding in enterprise corpora. In *Proceedings ICTAI 2006*, pages 599–608, 2006.
- [21] J. M. Ponte. Language models for relevance feedback. In *Advances in Information Retrieval*, pages 73–96. 2000.
- [22] D. E. Rose and D. Levinson. Understanding user goals in web search. In *WWW '04*, pages 13–19, 2004.
- [23] E. Voorhees. Overview of the TREC 2004 question answering track. In *Proceedings of TREC 2004*, 2005. NIST Special Publication: SP 500-261.
- [24] C. Zhai and J. Lafferty. A study of smoothing methods for language models applied to ad hoc information retrieval. In *SIGIR' 01*, pages 334–342, 2001.
- [25] V. Zlatić, M. Božičević, H. Štefančić, and M. Domazet. Wikipedias: Collaborative web-based encyclopedias as complex networks. *Physical Review E (Statistical, Nonlinear, and Soft Matter Physics)*, 74(1), 2006.