

# Entity-Relationship Queries over Wikipedia

XIAONAN LI, CHENGKAI LI

Department of Computer Science and Engineering, University of Texas at Arlington

CONG YU

Google Research

---

Wikipedia is the largest user-generated knowledge base. We propose a structured query mechanism, *entity-relationship query*, for searching entities in Wikipedia corpus by their properties and inter-relationships. An entity-relationship query consists of arbitrary number of predicates on desired entities. The semantics of each predicate is specified with keywords. Entity-relationship query searches entities directly over text rather than pre-extracted structured data stores. This characteristic brings two benefits: (1) Query semantics can be intuitively expressed by keywords; (2) It avoids information loss that happens during extraction. We present a ranking framework for general entity-relationship queries and a position-based Bounded Cumulative Model (BCM) for accurate ranking of query answers. We also explore various weighting schemes for further improving the accuracy of BCM. We test our ideas on a 2008 version of Wikipedia using a collection of 45 queries pooled from INEX entity ranking track and our own crafted queries. Experiments show that our ranking techniques and weighting schemes are both effective, particularly on multi-predicate entity-relationship queries.

Categories and Subject Descriptors: H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Retrieval models*

General Terms: Design, Languages, Performance, Experimentation

Additional Key Words and Phrases: entity search and ranking, structured entity query, Wikipedia

---

## 1. INTRODUCTION

Since its inception in January 2001, Wikipedia has risen to be the largest encyclopedia ever created, containing more than 3 million articles in English alone as of 2010. In the meantime, Wikipedia articles have amazingly evolved, from mostly plain texts at earlier stage to current ones with substantial structural annotations. It is now the primary knowledge source for many users on a wide variety of *entities*, including people, institutions, geographical locations, events, etc. For discovering and exploring the entities that fascinate them, users are in need of structured querying facilities, coupled with text retrieval capabilities, that explicitly deal with the entities, their properties, and relationships.

The prevalent manner in which users access Wikipedia is still keyword-based document search. Although keyword search has been quite effective in finding specific pages matching the keywords, there clearly exists a mismatch between its document-centric view and the aforementioned entity-centric user information needs. Users' tasks often cannot be clearly expressed with simple keyword queries and processing the query results may require substantial user efforts.

**Example 1 (Motivating Example):** Consider a business analyst investigating the development of Silicon Valley. Particularly, she is interested in this task: Find the list of *companies* and their *founders*, where the companies are in Silicon Valley and the founders are Stanford graduates.

There are two major mismatches that make keyword search unsuitable for resolving this task. First, the task focuses on *typed* entities, PERSON and COMPANY, and, in database terminology, their “join” relationships. Second, the task involves synthesizing information scattered across different places, therefore a simple list of pages is not sufficient. For instance, one page may tell the analyst that Jerry Yang is a founder of Yahoo!, but whether Yahoo! is a Silicon Valley company and whether Jerry Yang is a Stanford graduate may have to be found in other pages.

While conceptually simple, with only keyword search, tasks like the above one require substantial user efforts to perform multiple searches and assemble information from a potentially large number of articles. Our analyst may start with a search on “Silicon Valley company” and scan through the potentially long list of result articles to, hopefully, fetch a list of companies that are likely to be in Silicon Valley. She then similarly issues another search on “Stanford graduate” to find a list of people graduated from Stanford University. She then manually combine entities in these two lists and, by multiple additional searches, check if a company was founded by a person, for each pair of person and company. Alternatively, she can also go through the list of companies and, for each company, find its founders and check if Stanford is their alma mater by multiple search queries. Both are painful options and require the user to break down the task into a time-consuming, error-prone iterative procedure of searching, reading, and re-searching.

Wikipedia (and the Web) contains various tables and lists, which can be extracted into databases for powerful queries [Cafarella et al. 2008]. For example, a page listing all Silicon Valley companies may exist. However, it is unrealistic to expect such pages exist for arbitrary user tasks. Moreover, it is less common to find such tables/lists for relationships between entities, e.g., who founded which company.

We propose *entity-relationship query*, a declarative query mechanism for the aforementioned task. The results of such queries are tuples of entities that are likely to meet the semantic requirements of the query, instead of articles containing such entities. For Example 1, our analyst can write the following SQL-like query.

**Query 1 (Entity-Relationship Query for Example 1):**

```
SELECT x, y
FROM PERSON x, COMPANY y
WHERE x:["Stanford", "graduate"] // Predicate  $p_1$ 
      AND y:["Silicon Valley"] // Predicate  $p_2$ 
      AND x,y:["found"] // Predicate  $p_3$ 
```

We take a DB-IR integration approach in proposing this direction. On the one hand, entity-relationship queries have explicit structured components: *typed entity variables* (e.g.,  $x$ , bound to entities of type PERSON, and  $y$ , for entities of type COMPANY), *selection predicates* for selecting entities by their properties (e.g., predicate  $p_1$ ), and *relation predicates* for specifying relations between entities (e.g., predicate  $p_3$  for the requirement that  $y$  was founded by  $x$ ). On the other hand, the individual predicates are specified by keyword-based constraints. The query semantics dictates that entities satisfy a predicate by a simple and intuitive requirement: the entities co-occur with the keywords in some contexts. Such contexts can be sentences, windows of texts, etc. For simplicity, we use sentence as context in this paper. For example, predicate  $p_1$  requires every PERSON in the query answer to co-occur with “Stanford” and “graduate” in at least one sentence. In short, we aim

to capture entity properties and relationships through shallow syntax requirements implied by users at query-time, instead of explicitly extracting and reasoning about explicit semantic information before query-time [Brin 1998; Agichtein and Gravano 2000; Chu et al. ; Etzioni et al. 2008; Cafarella et al. 2007; Kandogan et al. 2006; DeRose et al. 2007]. Although such syntax clue is by no means rigorous or error-proof, it becomes robust when we take into account the redundancy in a corpus: true facts are more likely to be repetitively stated in multiple places. This intuition has been widely used in Web search and mining, e.g., information extraction [Brin 1998; Agichtein and Gravano 2000] and entity search and ranking [Cheng et al. 2007].

We implemented a prototype entity-relationship query system in Apache Lucene. The system consists of several components. The *Indexer* creates an Entity-Centric Index. It associates each term  $w$  with a list of entities (ordered by entity IDs) that co-occur with  $w$  somewhere in the corpus. For each entity  $e$  in the list, it further records where  $w$  and  $e$  co-occur. Leveraging this design, the *Retriever* efficiently retrieves entities and co-occurrence contexts. The results are fed to *Ranker* for ranking. This paper reports our study on ranking the related weighting problem.

**Challenges:** Entity-relationship queries can yield many false answers due to abundant accidental co-occurrences (e.g., false evidence such as “X’s partner is a Stanford graduate” for predicate  $p_1$ ). Therefore, how to rank query answers presents a critical challenge. *First*, the presence of multiple predicates in a query requires us to aggregate the rankings of entities for multiple predicates. *Second*, many true answers have small numbers of co-occurrence contexts (i.e., low redundancy) in Wikipedia. Using redundancy solely is not sufficient to tell such entities apart from false answers.

**Contributions:** To the best of our knowledge, this paper is the first attempt to study multi-predicate entity-relationship query and its ranking problem. We present a ranking framework for general entity-relationship queries. It first evaluates how well an answer satisfies individual predicates and then aggregates multiple predicate scores into an answer score. A Bounded Cumulative Model (BCM) is proposed for scoring predicates. BCM relies on redundant co-occurrence contexts for robust evaluation. To improve the ranking accuracy for answers with small numbers of supporting contexts, BCM performs refined assessment on each co-occurrence context based on three positional features—*proximity*, *ordering pattern*, and *mutual exclusion*. Contexts that are likely to be true evidence are given higher importance. Existing systems only exploit proximity feature [Chakrabarti et al. 2006; Cheng et al. 2007] and may not leverage redundancy [Chakrabarti et al. 2006]. In certain sensitive cases, BCM tends to yield sketically high scores. To address this issue, we further study two weighting schemes for BCM. The basic idea is to detect and penalize susceptible answers based on the number of supporting contexts.

In summary, we make the following contributions in this paper:

- We propose the concept of entity-relationship queries, for structured querying of entities directly over Wikipedia text with complex constraints (i.e., multiple predicates).
- We propose a ranking framework and a position-based Bounded Cumulative Model for ranking the answers to entity-relationship queries.

- We present two weighting schemes, maximal-support weighting and document-frequency weighting, for improving the accuracy of BCM.
- We conduct comprehensive experiments, and demonstrate the effectiveness of the ranking and weighting methods.

The rest of the paper is organized as follows. Section 2 reviews related work. Section 3.1 formally defines the concept of entity-relationship query and its ranking problem. In Section 3.2, we discuss three position-based features, which are used in our Cumulative Model (Section 3.3) and Bounded Cumulative Model (Section 3.4). Section 4 explores the weighting schemes aimed for improving BCM’s performance. Empirical results are reported in Section 5. Section 6 discusses limitations and future work. Section 7 concludes the paper.

## 2. RELATED WORK

Previous studies on structured querying of the Web focus on DB-based approach that explicitly extracts structured information into databases [Brin 1998; Agichtein and Gravano 2000; Chu et al. ; Etzioni et al. 2008; Cafarella et al. 2007; Kandogan et al. 2006; DeRose et al. 2007]. This approach lends itself to the rich and mature techniques of database querying. The DB-based approach is constrained by the capability of the information extraction (IE) and natural language processing (NLP) techniques. Particularly, it requires explicit identification of the “names” of entity relationships. For example, if a “found” relation between Jerry Yang and Yahoo! was not detected during the extraction phase, such information is lost and could not be queried.

Some systems [Suchanek et al. 2007; Kasneci et al. 2008; Dill et al. 2003; Auer et al. 2007] explicitly encode entities and their relations (and general knowledge) in RDF, the W3C recommendation of data model for Semantic Web. They can thus leverage the rich expressiveness of query languages like SPARQL<sup>1</sup> for querying entities. Some of them [Suchanek et al. 2007; Kasneci et al. 2008; Auer et al. 2007] only capture structured and semi-structured information, e.g., infoboxes in Wikipedia, leaving out the implicit information in unstructured text. For example, YAGO only supports around 100 relations [Suchanek 2009] unified from WordNet and Wikipedia. Other systems apply IE techniques over Web pages to bootstrap RDF extraction [Dill et al. 2003], thus bearing the same limitation of the aforementioned DB-based approach.

Entity ranking has gained significant interest recently [Petkova and Croft 2007; Zaragoza et al. 2007; Vercoustre et al. 2008]. However, they focus on different problem settings from ours. For example, in the entity ranking tracks of both INEX<sup>2</sup> and TREC<sup>3</sup>, the participant systems should find named entities relevant to some narrative descriptions. The focus is to accurately understand the descriptions and rank entities accordingly. There are often type constraints on the entities as well. However, they do not have the concept of “predicates” and do not deal with multiple predicates.

<sup>1</sup><http://www.w3.org/TR/rdf-sparql-query>

<sup>2</sup><http://www.inex.otago.ac.nz/tracks/entity-ranking/entity-ranking.asp>

<sup>3</sup><http://ilps.science.uva.nl/trec-entity/guidelines/>

The studies most related to ours are [Chakrabarti et al. 2006; Cheng et al. 2007; Zhou et al. 2010]. Chakrabarti et al. [Chakrabarti et al. 2006] learns an optimal scoring function on proximity feature, but it only scores entities by single context. It makes no attempt to integrate information found in multiple documents. Leveraging the redundancy on the Web, EntityRank [Cheng et al. 2007] aggregates scores of locally evaluated co-occurrence contexts into global scores to improve ranking. The Content Query Language [Zhou et al. 2010] extends EntityRank with more context matching patterns. All three systems only focus on queries comparable to our single-predicate queries and thus do not study multi-predicate queries.

### 3. POSITION-BASED RANKING

#### 3.1 Problem Statement

We formalize an **entity-relationship query** as  $q = \langle V, P \rangle$ .  $V$  is a set of entity variables. Each  $v \in V$  is bound to entities of certain type, e.g., PERSON.  $P$  is a set of predicates. Each  $p \in P$  is a pair  $\langle V_p, C_p \rangle$ , where  $V_p \subseteq V$  and  $C_p$  is a set of phrases<sup>4</sup>. Query 1 is thus formalized as  $q_1 = \langle V, P \rangle$ , where  $V = \langle x:\text{PERSON}, y:\text{COMPANY} \rangle$  and  $P = \{p_1, p_2, p_3\}$ . Among the predicates,  $p_1 = \langle \{x\}, \{\text{"Stanford"}, \text{"graduate"}\} \rangle$  and  $p_2 = \langle \{y\}, \{\text{"Silicon Valley"}\} \rangle$  are selection predicates, and  $p_3 = \langle \{x, y\}, \{\text{"found"}\} \rangle$  is a relation predicate.

An **answer** to a query  $q$  is a tuple of entities, denoted by  $t$ . For each  $v \in V$ , there is a corresponding entity  $e \in t$  instantiated from  $v$ , e.g.,  $t = \langle \text{Jerry Yang}, \text{Yahoo!} \rangle$  for Query 1. Given a predicate  $p = \langle V_p, C_p \rangle$ , we use  $t_p$  to represent the sub-tuple of  $t$  such that each entity  $e \in t_p$  is instantiated from a corresponding  $v \in V_p$ . Take  $p_1$  in Query 1 for example.  $t_{p_1} = \langle \text{Jerry Yang} \rangle$  because  $V_{p_1}$  has only one variable  $x$  and Jerry Yang is instantiated from  $x$ . Similarly,  $t_{p_3} = t$ .

Given a predicate  $p$ , if a sentence contains all the phrases in  $C_p$  and one entity for each variable in  $V_p$ , it is a (co-occurrence) context for  $p$ . These entities in whole are said to satisfy  $p$ . Suppose three sentences are found in the corpus:

- $s_1$ : Stanford University graduates Jerry Yang and ...
- $s_2$ : ...a senior manager at Yahoo! in Silicon Valley.
- $s_3$ : Jerry Yang co-founded Yahoo!.

Jerry Yang satisfies  $p_1$  by sentence  $s_1$ ; Yahoo! satisfies  $p_2$  by sentence  $s_2$ ; and they together satisfy  $p_3$  by  $s_3$ . Assembling the information together, the entity tuple  $\langle \text{Jerry Yang}, \text{Yahoo!} \rangle$  is composed as an answer to the query since it satisfies all the query predicates. Note that in  $s_1$ , Stanford University is treated as plain text since it is neither a PERSON nor a COMPANY.

A co-occurrence context of answer  $t$  for predicate  $p = \langle V_p, C_p \rangle$  is a quadruple  $\langle \text{doc}, \text{sent}, \hat{V}_p, \hat{C}_p \rangle$ .  $\text{doc}$  and  $\text{sent}$  refer to the document ID and the sentence number that together identify a unique sentence in the corpus.  $\hat{V}_p$  are the positions of entities in the aforementioned sub-tuple  $t_p$  and  $\hat{C}_p$  are the positions of phrases in  $C_p$ . Suppose the aforementioned  $s_1$  is the 8th sentence of document 9. In this context, Jerry Yang spans from position 3 to 4 and the two phrases ("Stanford" and "graduate") are at positions 0 and 2. Hence, it is represented as  $\langle 9, 8, \{\langle 3, 4 \rangle\}, \{0, 2\} \rangle$ .

<sup>4</sup>A single keyword is treated as a phrase of length 1.

Note that there can be multiple contexts of  $t_{p_1}$ , each being a sentence containing Jerry Yang, “Stanford”, and “graduate”. We denote all contexts of  $t_p$  by  $\phi_p(t)$ . Without loss of generality, we use sentence and context interchangeably unless distinction is needed.

**Problem Statement:** Denote all answers to query  $q=\langle V, P \rangle$  by  $A$ . Our goal is to rank the answers in  $A$  according to information provided by  $\phi=\{\phi_p|p \in P\}$ , where  $\phi_p=\bigcup_{t \in A} \phi_p(t)$ .

Since the information that is used for ranking,  $\phi$ , is primarily position information (i.e., documents IDs, sentence numbers, entity spans and phrase positions), the problem is called *position-based* ranking problem.

Given a query  $q=\langle V, P \rangle$ , our **ranking framework** consists of three scoring functions  $F^S$ ,  $F^R$  and  $F^A$ , such that for each answer  $t$ : (1) its score on a selection predicate  $p \in P$  is given by  $F_p^S(t)$ ; (2) its score on a relation predicate  $p \in P$  is given by  $F_p^R(t)$ ; and (3) its final score  $F^A(t)$  (the answer score) aggregates all predicate scores obtained via  $F^S$  and  $F^R$ . In this framework, the scores of different predicates are computed independently from each other. The intuition can be explained as follows. In Query 1, whether a PERSON is a Stanford graduate ( $p_1$ ) is independent from whether she founded any COMPANY ( $p_3$ ) and certainly irrelevant to whether a COMPANY is in Silicon Valley ( $p_2$ ).

The rest of this section proposes our position-based ranking method following this framework.

### 3.2 Position-Based Features

This section studies three position-based features that are derivable from a co-occurrence context. These features are the key components in our Cumulative Model (CM) and Bounded Cumulative Model (BCM) that are introduced later.

**3.2.1 Proximity.** Intuitively, if the entities in  $t_p$  and the keywords in  $C_p$  are close to each other in a context  $s \in \phi_p(t)$ , they likely belong to the same grammatical unit of the corresponding sentence (e.g., a phrase like *Stanford University graduate Jerry Yang*) and thus form a piece of true evidence. Given predicate  $p$ , we define the proximity of  $t_p$  in  $s$  as

$$prox_p(t, s) = prox_p(t_p, s) = \frac{\sum_{e \in t_p} |token(e, s)| + \sum_{c \in C_p} |c|}{|scope_p(t_p, s)|}$$

where  $|token(e, s)|$  is the number of tokens in  $s$  representing entity  $e$ ;  $|c|$  is the number of tokens in phrase  $c$ ;  $scope_p(t_p, s)$  is the smallest scope in  $s$  covering all the entities in  $t_p$  and all the phrases in  $C_p$  (a scope is a consecutive sequence of tokens in  $s$ ); and consequently  $|scope_p(t_p, s)|$  is the total number of tokens in the scope. Note that the proximity value is in the range of  $[0,1]$  by this definition.

Different representations may be used in various places to refer to the same entity and may have different numbers of tokens. For example, the entity IBM may be represented by “IBM”, “Big Blue”, or “International Business Machine”. Hence,  $|token(\underline{IBM}, s)|$  may be 1, 2, or 3 in different  $s$ .

**Example 2:** The following two sentences are both contexts of the underlined entities for predicate  $p_1$  in Query 1. Context  $s_1$  is true evidence, supporting a true positive, while  $s_4$  is false, supporting a false positive.

$s_1$ : *Stanford University graduates Jerry Yang and ...*

$s_4$ : *A professor at Stanford University, Colin Marlow had a relationship with Cristina Yang before she graduated ...*

Predicate  $p_1$  has two phrases, “Stanford” and “graduate”, each with one token, hence  $\sum_{c \in C_{p_1}} |c| = 2$ . In  $s_1$ , the PERSON Jerry Yang is represented by two tokens, “Jerry” and “Yang”, hence  $\sum_{e \in t_{p_1}} |\text{token}(e, s_1)| = 2$ . The scope covering the entity and the two phrases spans 5 tokens, from “Stanford” to “Yang”, thus  $|\text{scope}_{p_1}(t_{p_1}, s_1)| = 5$ . Therefore, the proximity of Jerry Yang in sentence  $s_1$  is  $\text{prox}_{p_1}(t_{p_1}, s_1) = \frac{2+2}{5} = 0.8$ . Similarly, the proximity of Colin Marlow in  $s_4$  is  $\frac{2+2}{13} = 0.31$ . Based on proximity alone, we say that  $s_1$  is more likely to be true evidence and therefore, Jerry Yang is more likely to satisfy  $p_1$  than Colin Marlow, given no other context.

**3.2.2 Ordering Pattern.** An ordering pattern refers to the order of entities and phrases in a co-occurrence context. Consider again predicate  $p_1 = \langle \{x\}, \{\text{“Stanford”}, \text{“graduate”}\} \rangle$  in Query 1. Let  $c_1$  be the first phrase (“Stanford”) and  $c_2$  the second (“graduate”). This predicate has six different ordering patterns ( $xc_1c_2$ ,  $xc_2c_1$ ,  $c_1xc_2$ ,  $c_2xc_1$ ,  $c_1c_2x$  and  $c_2c_1x$ ). Generally, if we denote all possible patterns of a predicate  $p$  by  $O_p$ , we have  $|O_p| = (|V_p| + |C_p|)!$ . Note that, extra tokens and punctuations between entities and phrases are irrelevant to the patterns. Hence, “Stanford University graduate, Jerry Yang” and “Stanford graduate Jerry Yang” follow the same pattern,  $c_1c_2x$ .

We observe that some ordering patterns are better indicators of true evidence than others. For example, to express that somebody is a graduate of Stanford University, true evidence usually follows the pattern  $c_1c_2x$  (e.g.,  $s_1$ ). Context following another pattern,  $c_1xc_2$ , is likely to be false evidence (e.g.,  $s_4$ ). To distinguish good patterns (those that tend to indicate true evidence) from others, we may assign a different weight to each pattern, so that entities supported by contexts following good patterns are scored higher. However, it is impossible to pre-determine the weights since the goodness of ordering patterns are predicate-dependent. To illustrate,  $c_1c_2x$  is a good pattern for predicate  $p_1$  in Query 1, but may not be equally good for another predicate  $p'_1 = \langle \{x:\text{NOVEL}\}, \{\text{“by”}, \text{“Jane Austen”}\} \rangle$ , because it is less common to see true evidence such as

... *written by Jane Austen, Pride and Prejudice ...*

In our approach, the weights of ordering patterns for a predicate  $p$  are dynamically derived from  $\phi_p$ , the set of all co-occurrence contexts for  $p$ . Denoting  $\phi_p(o)$  as the subset of contexts following pattern  $o$ , we define the weight of  $o$  for predicate  $p$  as its frequency in  $\phi_p$ ,

$$f_p(o) = |\phi_p(o)| / |\phi_p|$$

This definition assumes that good patterns appear more often than bad ones. Although in theory there might be a pattern frequently appearing in false evidence, making a bad pattern more common, we do not observe such case in our experiments.

Another possible direction is leveraging Machine Learning techniques to predict which patterns lead to better results. While we are also exploring this direction as future work, we note here that one significant challenge of the Machine Learning

approach is the need to obtain training data, which can be costly in terms of human effort.

**3.2.3 Mutual Exclusion.** Given a predicate  $p$ , multiple contexts in  $\phi_p$  may have the same  $\langle doc, sent \rangle$  value (i.e., come from the same sentence). They are contexts of different entities and may follow different ordering patterns in that sentence. The co-existence of different patterns in one sentence is called *collision* and the patterns are referred to as *colliding patterns*. The mutual exclusion rule dictates that, when collision happens, at most one colliding pattern is effective and the sentence is only considered evidence following that pattern.

**Example 3:** The following sentence illustrates mutual exclusion rule for  $p_1$  in Query 1. The sentence appears as three contexts, one for each underlined entity. Ric Weiland follows the pattern  $o_1 = xc_2c_1$ . Paul Allen and Bill Gates follow  $o_2 = c_2c_1x$ . Semantically, the former pattern is the effective pattern and the sentence is only evidence of Ric Weiland.

*s<sub>5</sub>: After Ric Weiland graduated from Stanford University, Paul Allen and Bill Gates hired him in 1975 ...*

Without understanding the semantics, it is difficult to decide which colliding pattern is absolutely effective. Therefore, we relax the rule with a *credit* mechanism, where every colliding pattern is considered partially effective, and patterns with higher credits are more likely to be effective than those with lower credits. We assume each sentence  $s$  (that is a context of at least one sub-tuple  $t_p$  for predicate  $p$ ) has a total credit of 1, meaning that there is only one effective pattern. Given a predicate  $p$ , denote the colliding patterns in  $s$  by  $O_p(s) \subseteq O_p$ . Each  $o \in O_p(s)$  gets a credit  $credit_p(o, s)$ , and  $\sum_{o \in O_p(s)} credit_p(o, s) = 1$ .

To allocate credits to the colliding patterns  $O_p(s)$ , we adopt the intuition that patterns followed by more prominent entities are more likely to be effective. Specifically, let  $T_p(o, s)$  be all sub-tuples on  $p$  following pattern  $o$  in  $s$ . For each  $o \in O_p(s)$ , we choose a representative from  $T_p(o, s)$ , denoted by  $T_p^*(o, s)$ , which is the one with the highest proximity value, i.e.,  $T_p^*(o, s) = \arg \max_{t_p \in T_p(o, s)} prox_p(t_p, s)$ . We compare the representatives (and thus the patterns that they follow) by how prominent they are, i.e., by their overall numbers of contexts in  $\phi_p$ . The credit of  $o$  in  $s$  is

$$credit_p(o, s) = \frac{|\phi_p(T_p^*(o, s))|}{\sum_{o' \in O_p(s)} |\phi_p(T_p^*(o', s))|}$$

where  $\phi_p(T_p^*(o, s))$  is the set of contexts of  $T_p^*(o, s)$  for predicate  $p$ . Note that we choose the most proximate sub-tuple as the representative of a colliding pattern and allocate credits based on representatives only. The intuition is that the most proximate sub-tuple is most likely to form a grammatical unit with phrases in  $C_p$ , and hence the most reliable one for allocating credits.

In Example 3,  $t^1 = T_{p_1}^*(o_1, s) = \text{Ric Weiland}$  (i.e., the representative of pattern  $o_1$  is Ric Weiland) since he is the only PERSON in  $s$  following pattern  $o_1$ ; and  $t^2 = T_{p_1}^*(o_2, s) = \text{Paul Allen}$  because he has higher proximity (0.67) than Bill Gates (0.44), though both follow  $o_2$ . Suppose Ric Weiland is found in 4 contexts ( $|\phi_{p_1}(t^1)| = 4$ ) and Paul Allen in 2 ( $|\phi_{p_1}(t^2)| = 2$ ). Then, the credits of their corresponding patterns in  $s$  would be  $\frac{4}{4+2} = 0.67$  (for  $o_1$ ) and 0.33 (for  $o_2$ ).



Note that the pattern credit here is different from the weight of pattern in Section 3.2.2. The weight of pattern  $o$  is a global measure (aggregated over  $\phi_p$ ) of how frequent, and thus how reliable, pattern  $o$  is. The credit of  $o$ , on the contrary, is a local measure particular to each sentence  $s$ , indicating how likely  $o$  is the effective pattern in  $s$ .

### 3.3 Single-Predicate Scoring

So far, we have introduced all the position-based features for assessing individual contexts. Integrating these features together, this section presents Cumulative Model (CM) for scoring an answer on a single predicate. We assume that  $F^S$  is the same as  $F^R$  (i.e., the same function is used for scoring all predicates), hence for brevity, we use  $F_p(t)$  instead of  $F_p^S(t)$  and  $F_p^R(t)$ .

Let  $\phi_p(t, o) \subseteq \phi_p(t)$  be all contexts of  $t$  for predicate  $p$  that follow pattern  $o \in O_p$ . Our Cumulative Model (CM) is

$$F_p(t) = \sum_{o \in O_p} (f_p(o) \sum_{s \in \phi_p(t, o)} \text{prox}_p(t, s) \text{credit}_p(o, s))$$

where  $f_p(o)$  is the weight of pattern  $o$ ;  $\text{prox}_p(t, s)$  is  $t_p$ 's proximity in context  $s$ ;  $\text{credit}_p(o, s)$  is the credit of  $o$  in  $s$ .

The model divides  $\phi_p(t)$ ,  $t$ 's contexts for  $p$ , into  $|O_p|$  groups,  $\{\phi_p(t, o) | o \in O_p\}$ , so that contexts in each group follow the same pattern. For each group  $\phi_p(t, o)$ , the model computes a *group score* (the inner summation). The group scores are linearly combined using weights  $f_p(o)$  (the outer summation), such that the group scores of better patterns account more in  $F_p(t)$ . The kernel of the function,  $\text{prox}_p(t, s) \text{credit}_p(o, s)$ , assesses how likely  $s$  is true evidence of  $t$  for predicate  $p$ . It is monotonic to both the proximity of  $t_p$  and the credit of  $t_p$ 's pattern  $o$ . Answers supported by contexts having higher proximities and pattern credits will accumulate higher scores and thus ranked higher.

It is interesting to note that CM can be customized easily by switching on and off its component features, so that we can evaluate the effectiveness of individual features. While detailed evaluations are presented in Section 5, below we list three important customizations.

$$\begin{aligned} \text{COUNT } F_p(t) &= \sum_{o \in O_p} (1 \sum_{s \in \phi_p(t, o)} 1) = \sum_{o \in O_p} |\phi_p(t, o)| = |\phi_p(t)| \\ \text{PROX } F_p(t) &= \sum_{o \in O_p} \sum_{s \in \phi_p(t, o)} \text{prox}_p(t, s) = \sum_{s \in \phi_p(t)} \text{prox}_p(t, s) \\ \text{MEX } F_p(t) &= \sum_{o \in O_p} \sum_{s \in \phi_p(t, o)} \text{credit}_p(o, s) = \sum_{s \in \phi_p(t)} \text{credit}_p(o, s) \end{aligned}$$

### 3.4 Multi-Predicate Scoring

We extend our single-predicate scoring model to handle multi-predicate queries. Given a query answer, CM computes a score on each predicate. However, it remains unclear how to derive the final score,  $F^A(t)$ , from predicate scores.

With CM, predicate scores are unbounded, i.e., the more contexts the higher scores. When multiple predicate scores are aggregated, some could be so high

Table I. Example Answers

	$x$	$y$	$p_1$	$p_2$	$p_3$	$\Pi$	$\Sigma$
$t_1$	Jerry Yang	Yahoo!	0.8	0.7	0.8	0.448	2.3
$t_2$	Larry Page	Google	0.6	0.5	0.6	0.18	1.7
$t_3$	Scott McNealy	Cisco	0.9	0.8	0.2	0.144	1.9
$t_4$	Bill Gates	IKEA	0.3	0.1	0.2	0.006	0.6

that they dominate the aggregate score, which is called *predicate dominance*. To alleviate this problem, we propose **Bounded Cumulative Model (BCM)** for better scoring of individual predicates:

$$F_p(t) = \sum_{o \in O_p} (f_p(o) [1 - \prod_{s \in \phi_p(t,o)} (1 - \text{prox}_p(t,s) \text{credit}_p(o,s))]) \quad (1)$$

BCM uses the same three features as CM does, but differs from CM in the computation of group scores, each of which is computed from a set of contexts  $\phi_p(t,o)$ . Basically, BCM bounds all group scores in the range  $[0,1]$ , and consequently it bounds the predicate scores within  $[0,1]$ , since  $\sum_{o \in O_p} f_p(o) = 1$  according to Section 3.2.2.

Given an answer  $t$  to query  $q = \langle V, P \rangle$ ,  $t$ 's final score,  $F^A(t)$ , is computed as the product of its scores on all predicates,

$$F^A(t) = \prod_{p \in P} F_p(t) \quad (2)$$

where  $F_p(t)$  can be either BCM or CM. For our problem, product is a more reasonable aggregate function than summation, another common aggregate function, because it favors answers with balanced predicate scores over those with polarized ones. To illustrate why balanced scores should be favored, consider the case in Table I. The table shows four answers to the query of Query 1. For each answer, it lists all three predicate scores (by BCM), as well as the final scores using product and summation, respectively. The two aggregates agree on the ranking of  $t_1$  and  $t_4$ , which get unanimously (i.e., balanced) high and low predicate scores, but disagree on  $t_2$  and  $t_3$ . The true positive,  $t_2$ , gets modest and balanced scores on all the predicates. It is correctly ranked higher than  $t_3$ , a false positive, by product, but loses the comparison by summation. Answer  $t_3$  gains high scores on  $p_1$  and  $p_2$  (Both are indeed satisfied by  $t_3$ .), but low score on  $p_3$  (In reality, it does not satisfy  $p_3$ .). However, the final score of  $t_3$  by summation is dominated by the high scoring predicates and thus  $t_3$  is mistakenly ranked above  $t_2$ .

#### 4. WEIGHTING

We have introduced our multi-predicate ranking method based on Bounded Cumulative Model. In this section, we will look further into BCM and explore several weighting schemes to improve its accuracy.

As discussed earlier, BCM is proposed to alleviate the problem of predicate dominance. But its multiplicative nature (the product in formula 1) also makes predicate scoring sensitive to a few good contexts. As an extreme example, if only one ordering pattern is involved (thus,  $f_p(o) = 1$ ), the predicate score will become 1 (the maximum) as long as any context is scored 1, regardless of other supporting contexts. In multi-predicate queries, such sensitivity of predicate scores will be further

magnified by formula 2.

Our idea for solving the sensitivity problem is to penalize susceptible answers by weighting. The new weighted scoring model for general entity-relationship queries is shown by formula 3. Compared to formula 2, an exponential weight  $W_p(t)$  is applied over each  $F_p(t)$ , the BCM score of answer  $t$  for predicate  $p$ . In the ideal case, the weight should adjust itself according to the degree of susceptibility of each answer. Intuitively, if  $F_p(t)$  is computed according to a small number of contexts, it is susceptible (even though the score may be high) and should receive a penalty for that. Briefly speaking, the weight  $W_p(t)$  needs to be larger for more susceptible answers (i.e., answers supported by less contexts). Note that  $F_p(t)$  is bounded between  $[0, 1]$  by BCM, hence larger weight will lower the predicate score.

$$F^A(t) = \prod_{p \in P} F_p(t)^{W_p(t)} \quad (3)$$

The rest of this section discusses two weighting schemes that match our intention—higher weights  $W_p(t)$  for answers with less contexts. With regard to a predicate, let us define the *support* of a candidate answer  $t$  to be the number of supporting contexts for  $t$ , which is denoted by  $|\phi_p(t)|$ . We measure the significance of  $t$  by comparing its support with some “best-possible” support. In the first weighting scheme, maximal-support weighting, we compare the support of  $t$  with the largest support among all the answers. The larger the difference is, the higher the weight will be. In the second scheme, corpus-frequency weighting, the support of  $t$  is compared with its “corpus frequency”, which refers to the number of contexts in the corpus that contain  $t$ .

#### 4.1 Maximal-Support Weighting

The maximal-support weight is defined as

$$W_p(t) = \alpha_p(t) = \frac{\max_{t'} \log(|\phi_p(t')| + 1)}{\log(|\phi_p(t)| + 1)}$$

The denominator is the logarithm of answer  $t$ ’s support (plus 1) for predicate  $p$  and the numerator is the maximum of such logarithm values among all the answers. The smoothing constant 1 is used to avoid zero-denominator when  $|\phi_p(t)| = 1$ . By this weighting scheme, answers with less support (i.e., smaller denominators) will receive higher penalties. Suppose there are two answers for predicate  $p$ , where  $t_1$  has 3 contexts ( $\log(3+1) = 2$ ) and  $t_2$  has much more contexts, say 127 ( $\log(127+1) = 7$ ). Their maximal-support weights are  $\alpha_p(t_1) = \frac{\max\{2,7\}}{2} = 3.5$  and  $\alpha_p(t_2) = \frac{7}{7} = 1$ , respectively. Thus,  $t_1$  is penalized more with weight 3.5. Note that, although we use base 2 logarithm in the calculation for illustration purpose, the weight is actually independent of the choice of base.

For an individual predicate, if two answers have the same support (hence, the same denominator), their maximal-support weights will be the same, and their relative ranking is preserved in spite of the weighting, according to the monotonicity that  $x < y \Leftrightarrow x^w < y^w$  for  $w > 0$ .

The numerator of  $\alpha_p(t)$  corresponds to the maximum of observed support for a predicate, it is a hint on how much support a true answer could get from the corpus. Here, we implicitly assume that the maximally supported answer is true, which is

usually the case. We utilize this hint to evaluate the significance of the supports of candidate answers and weights them accordingly. To illustrate, let us consider an answer  $t$  with 3 contexts. We observed two situations in our corpus: (1) the maximal support is much higher than 3 (e.g., 127), and (2) the maximal support is also low (e.g., 7). The former indicates that true answers tend to have many contexts, while the latter indicates that it is unlikely to find many contexts for this predicate. Intuitively, answer  $t$  should be penalized more in the former situation since it does not look like a true answer in terms of support. In this regard, the definition of  $\alpha_p(t)$  tends to penalize poorly-supported answers more severely.

It is important to note that the maximal support (i.e., the numerator), has no effect on ranking single-predicate queries because it is constant for all answers. But for multi-predicate queries, each predicate has a different maximal support. Predicates with higher maximal support tend to yield higher weights, and penalize the answer more severely given the same support value, and thus impact the ranking of candidate answers. In sum, maximal-support weighting applies distinguished treatment of predicates according to their maximal supports, relying on predicates with high maximal supports to penalize their poorly-supported answers.

The maximal-support weighting has its own concern as the maximal support itself is sensitive to outliers. In an undesirable scenario, it is possible that most true answers of a predicate has less than 10 contexts, except for one outlier with 100 contexts. In this case, the maximum support is 100, resulting in high penalty to other true answers. A more robust solution could use the average of the largest  $k$  supports as the numerator in  $\alpha_p(t)$ . We leave this issue for future study.

#### 4.2 Corpus-Frequency Weighting

The corpus-frequency weighting takes into consideration the corpus background of an answer for evaluating the significance of its support. Given the same support, the more an answer appears in the corpus, the less significant the support is. For example, for predicate  $p_2 = \langle \{y\}, \text{“Silicon Valley”} \rangle$  on COMPANY, NASDAQ is supported by 3 contexts out of its 845 occurrences in the corpus while Mayfield Fund has 3 out of 11. Given its high corpus frequency, NASDAQ is quite susceptible because its co-occurrence with “Silicon Valley” likely happens just by chance. In this sense, we consider the support of Mayfield Fund more significant than that of NASDAQ. This intuition is reflected in our corpus-frequency weighting as

$$W_p(t) = \beta_p(t) = \frac{\log(N(t_p) + 1)}{\log(|\phi_p(t)| + 1)}$$

where  $t_p$  is the projection of  $t$  onto the variables in  $p$ .  $N(t_p)$  is  $t_p$ 's corpus frequency. Note that  $|\phi_p(t)| \leq N(t_p)$ .

As mentioned above, maximal-support weighting will weight two answers equally if they have the same support. This is not desired. By looking at the corpus frequency,  $\beta_p(t)$  is able to capture the difference in case of tied support.

Corpus-frequency weighting shares the same intuition as IDF (inverse document frequency) weighting in document retrieval. Both weighting schemes favor specificity over generality. The difference is that IDF weight depends on query terms only and is the same for all hit answers, while our corpus-frequency weight is dependent on both query predicates and the returned answers.

Table II. Ten Types from Wikipedia

Type	(E)ntities	(O)ccurrences	O/E
AWARD	979	83,001	85
CITY	54,296	2,269,771	42
CLUB	13,365	701,042	52
COMPANY	20,966	842,849	40
FILM	37,455	595,774	16
NOVEL	15,385	181,156	12
PERSON	368,833	5,985,231	16
PLAYER	74,171	750,708	10
SONG	27,805	251,186	9
UNIVERSITY	17,491	677,140	39
TOTAL	630,746	11,742,084	19

A concern is how to weight between two answers of the same specificity, e.g., (1)  $|\phi_p(t_1)|=1$  with  $N(t_{1p})=5$ , versus (2)  $|\phi_p(t_2)|=20$  and  $N(t_{2p})=100$ . Here, the specificity is measured by the ratio between support and corpus frequency. While the conclusion is unclear in general, for our problem and corpus, the latter case should be favored. The intuition is that, since  $t_2$  already has large support, its BCM score is not as susceptible as  $t_1$  and hence the corpus frequency becomes less important in determining its significance. Our definition of corpus-frequency weight is capable of capturing this intuition.

As a real example from the predicate on Silicon Valley companies, Apple (a true answer) has 19 contexts out of its 2639 occurrences with specificity  $\frac{19}{2639}=0.72\%$ , while SAP (a false answer) has 2 out of 211 with specificity 0.94%. Even though Apple has very large corpus frequency and thus lower specificity, its weight is  $\frac{\log(2639+1)}{\log(19+1)}=2.6$ , compared to SAP's 4.9.

## 5. EXPERIMENTS

In this section, we provide the empirical evaluation of our prototype system implemented in Apache Lucene. We have updated the pre-processing module used for [Li et al. ], fixing some errors in parsing Wikipedia articles. This results in a slightly different dataset. Therefore, we re-evaluated all ranking methods reported earlier. Despite this change, our conclusions remain the same. We will also report our new experimental results on different weighting schemes.

### 5.1 Data and Query Sets

We used the 2008-07-24 snapshot of Wikipedia. After we removed all the irrelevant pages (such as category and administrative pages), there were about 1.8 million articles. This article set is used as the entity catalog. Each article is the description of an entity, by Wikipedia's nature of being an encyclopedia, and the article title corresponds to the entity name. We predefined 10 entity types (Table II) and assigned about 0.63 million entities to these types based on simple hand-crafted rules, mainly using their categories in Wikipedia. For example, if an article belongs to a category whose name ends with "novels" (e.g., *British novels*) it is treated as an entity of type NOVEL. This simple method turns out sufficiently accurate for our experiments.

The same article set is also used as the corpus. For each article, we removed

its section titles, tables, infoboxes, references, etc., retaining only the main textual content. The main text is segmented into sentences. We removed punctuation marks and stemmed all words using the Porter Stemmer. We consider the *internal links*, hyperlinks from one Wikipedia article to other Wikipedia articles, as occurrences of the link targets (entities). In this way, we collected nearly 12 million occurrences of the 0.63 million typed entities. There is significant reduction in total occurrences from previous report [Li et al. ]. The reduced occurrences are mainly those appearing in lists and tables that are un-removed previously due to parsing errors. Therefore, despite this significant reduction, our new experiments find out that the supporting contexts of retrieved answers have only suffered small changes.

Named entity recognition (NER) [Nadeau et al. ] and entity disambiguation [Dill et al. 2003] are intensively studied problems. Our hyperlink-based annotation can be viewed as a rudimentary entity disambiguation method. Recently we have seen advanced entity recognition and disambiguation methods using Wikipedia as entity catalog [Mihalcea and Csomai ; Milne and Witten ; Kulkarni et al. 2009] to automatically link entities mentioned in plain text to their corresponding Wikipedia articles. One of our ongoing efforts is to use Wikify [Milne and Witten ] to annotate entities in articles that are not hyperlinks. This method will give us more comprehensive entity occurrences. Furthermore, it can be applied on generic Web pages, thus enabling entity-relationship queries on Web corpus.

We have previously collected two query sets, INEX17 and OWN28. INEX17 is adapted from the topics in the Entity Ranking track of INEX 2009. From the 60 available topics, we adapted the ones that are on our predefined 10 entity types. We obtained 11 single-predicate queries and 6 multi-predicate queries. OWN28 contains our own crafted 28 queries, including 16 single-predicate queries and 12 multi-predicate queries. Since both query sets are relatively small, we merged them into one set for more robust evaluation. It contains 27 single-predicates (Single-27) and 18 multi-predicate queries (Multi-18).

**Ground Truth:** We manually checked query answers returned by our system to collect the ground truth. However, some of the queries return hundreds of answers. Exhaustive checking is prohibitively time consuming. Therefore, for such queries, we adopted the depth- $N$  pooling approach used by INEX. Basically for each query, we only check the top  $N$  answers returned by each ranking method. Since different methods overlap greatly in their top  $N$  answers, a lot of time is saved. This pooling approach allows us to evaluate ranking accuracy up to rank  $N$ . As a typical choice,  $N$  is set to 100 in our experiments.

## 5.2 Answers to Sample Queries

In this section we use several sample queries to demonstrate the accuracy of our query evaluation system, based on BCM scoring model. The top-10 answers to each query are displayed. The true answers are marked by bullets.

**Case 1 (Big Ten Universities):** A student is interested in the list of big ten universities. She can use a single-predicate query on entity type UNIVERSITY with keyword constraint “Big Ten”, as follows.

SELECT x FROM UNIVERSITY x WHERE x:["Big Ten"]
<ul style="list-style-type: none"> <li>• ⟨Indiana University (Bloomington)⟩</li> <li>• ⟨Ohio State University⟩</li> <li>• ⟨University of Michigan⟩</li> <li>• ⟨Pennsylvania State University⟩</li> <li>• ⟨University of Iowa⟩</li> <li>• ⟨Purdue University⟩</li> <li>• ⟨Michigan State University⟩</li> <li>• ⟨University of Wisconsin-Madison⟩</li> <li>• ⟨University of Illinois at Urbana-Champaign⟩</li> <li>• ⟨Northwestern University⟩</li> </ul>

**Case 2 (Business Analyst):** Our motivating Example 1– Silicon Valley companies founded by Stanford graduates. The corresponding query is Query 1.

<ul style="list-style-type: none"> <li>• ⟨Jerry Yang, Yahoo!⟩</li> <li>• ⟨David Packard, Hewlett-Packard Company⟩</li> <li>• ⟨Scott McNealy, Sun Microsystems⟩</li> <li>• ⟨Vinod Khosla, Sun Microsystems⟩</li> <li>• ⟨Bill Gates, Microsoft⟩</li> <li>• ⟨William Hewlett, Hewlett-Packard Company⟩</li> <li>• ⟨Vinod Khosla, Kleiner Perkins Caufield &amp; Byers⟩</li> <li>• ⟨Larry Page, Google⟩</li> <li>• ⟨Andy Bechtolsheim, Sun Microsystems⟩</li> <li>• ⟨Sergey Brin, Google⟩</li> </ul>
--

**Case 3 (Movie Lover):** A movie lover interested in Academy Award winning films starring Australian actors.

SELECT x, y FROM FILM x, PERSON y WHERE x:["win "Academy Award"] AND y:["Australian" "actor"] AND x,y:["star"]
<ul style="list-style-type: none"> <li>• ⟨Braveheart, Mel Gibson⟩</li> <li>• ⟨Brokeback Mountain, Heath Ledger⟩</li> <li>• ⟨The Adventures of Priscilla, Guy Pearce⟩</li> <li>• ⟨L.A. Confidential, Guy Pearce⟩</li> <li>• ⟨Mad Max, Mel Gibson⟩</li> <li>• ⟨The Year of Living Dangerously, Mel Gibson⟩</li> <li>• ⟨Elizabeth, Geoffrey Rush⟩</li> <li>• ⟨A Beautiful Mind, Russell Crowe⟩</li> <li>• ⟨The Adventures of Priscilla, Hugo Weaving⟩</li> <li>• ⟨Gladiator, Russell Crowe⟩</li> </ul>

**Case 4 (Basketball Fan):** A basketball fan looking for team leaders of NBA champions. Particularly, she is interested in those NBA Finals MVPs.

SELECT x, y FROM CLUB x, PERSON y WHERE x:["NBA" "champion"] AND y:["Finals MVP"] AND x,y:["led"]
<ul style="list-style-type: none"> <li>• ⟨Chicago Bulls, Michael Jordan⟩</li> <li>• ⟨Los Angeles Lakers, Magic Johnson⟩</li> <li>• ⟨San Antonio Spurs, Tim Duncan⟩</li> <li>• ⟨Los Angeles Lakers, Kareem Abdul-Jabbar⟩</li> <li>• ⟨Boston Celtics, Larry Bird⟩</li> <li>• ⟨Los Angeles Lakers, Shaquille O'Neal⟩</li> <li>• ⟨Houston Rockets, Moses Malone⟩</li> <li>• ⟨Los Angeles Lakers, James Worthy⟩</li> <li>• ⟨Philadelphia 76ers, Kareem Abdul Jabbar⟩</li> <li>• ⟨Detroit Pistons, Isiah Thomas⟩</li> </ul>

Table III. Comparing Different Ranking Methods

Query	COUNT	MEX	PROX	CM	BCM	ER
<b>MAP</b>						
Single-27	0.826	0.858	0.901	0.909	0.911	0.829
Multi-18	0.573	0.603	0.682	0.708	0.763	0.703
<b>nDCG</b>						
Single-27	0.888	0.908	0.932	0.937	0.940	0.901
Multi-18	0.729	0.753	0.806	0.826	0.855	0.812

### 5.3 Comparing Ranking Methods

In this section, we compare and analyze the multiple ranking methods discussed in Section 3, namely COUNT, MEX, PROX, CM and BCM. All the methods differ in how they compute predicate scores, i.e.,  $F_p(t)$ . For multi-predicate queries, the same formula 2 is used to aggregate predicate scores into answer score. We compare these ranking methods using three popular measures: *MAP*, *nDCG*, and *Precision-at-k* ( $P@k$  for short). Since it is fairly difficult to obtain the complete ground truth in Wikipedia, we are not able to evaluate recall or F-measure.

The **MAP** (Mean Average Precision) is the arithmetic mean of average precisions for a set of queries. It is a good indicator of both precision and recall. The top part of Table III shows MAP for Single-27 and Multi-18. Both MEX and PROX improve over COUNT, by 0.032 and 0.075 on Single-27, respectively, and by 0.03 and 0.109 on Multi-18. PROX is shown to be more effective than MEX. CM and BCM have similar performance on Single-27, but BCM excels over CM by 0.055 on Multi-18. This demonstrates that BCM is a better model for ranking multi-predicate queries.

The **nDCG** (Normalized Discounted Cumulative Gain) discounts the gain of a true answer by its rank in the result list and normalizes the cumulated gain with that of a perfect ranking. The bottom part of Table III shows the average nDCG on Single-27 and Multi-18. The observation is similar to the case of MAP.

Overall, our conclusion is that, CM and BCM are comparable for ranking single-predicates and BCM has clear advantage on multi-predicate queries. We refer interested readers to [Li et al. ] for separate evaluation of INEX17 and OWN28 query sets in earlier experiments.

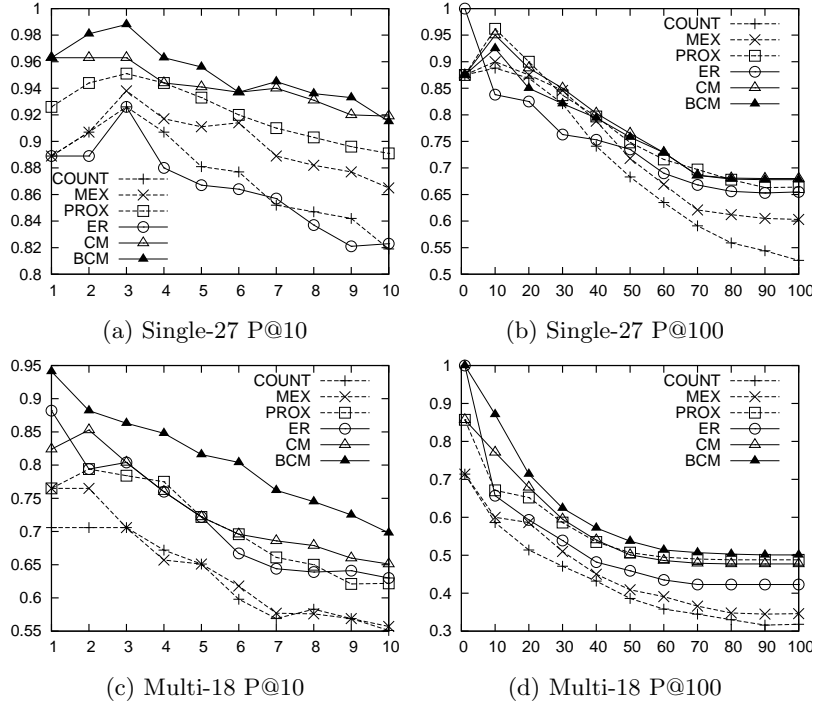
To further analyze how various methods perform at different ranks, we plot **precision-at- $k$**  curves.

Figure 1(a) shows the results for Single-27 with  $k$  up to 10 (denoted as  $P@10$ ). COUNT has the worst performance. PROX is consistently better than MEX across all ranks, but worse than CM and BCM, agreeing with the conclusion drawn from MAP and nDCG analysis. BCM is the best among all, especially for top 5. We further collected a subset of queries in Single-27, which returned more than 100 answers, and plot the  $P@100$  curve in Figure 1(b). While COUNT and MEX are still the worst, the other methods are fairly close to each. Our observation on Single-27 indicates that, for single-predicate queries, BCM is good at ranking top answers, but does not have advantage over the long range.

Figure 1(c,d) repeats the similar experiment on Multi-18. It can be easily observed that BCM has clear advantage in ranking multi-predicate queries for top answers and beyond. This reinforces our conclusion drawn from Table III.

Note that, if all true answers are ranked before position  $K$ , the precision after  $K$



Fig. 1. Precision-at- $k$  for Single-27 and Multi-18

will not change. This is reflected by the flat tails in (b) and (d).

In summary, the individual features are effective for ranking entity-relationship queries and they work best in concert when integrated and BCM. BCM rivals CM on single-predicate queries, and excels on multi-predicate queries. This is because BCM alleviates the predicate dominance problem as discussed in Section 3.4.

#### 5.4 BCM vs. Other Entity Ranking Methods

As a first study on multi-predicate entity-relationship query, we did not find directly comparable systems. To our best effort, we have chosen three state-of-the-art systems proposed for related problems: EntityRank, INEX and INRIA. All the three systems work on the entity ranking problem, though under different task settings. Besides, they are all evaluated using Wikipedia as corpus and entity catalog, though INEX and INRIA worked on different snapshots from ours.

EntityRank (ER) [Cheng et al. 2007] outperforms another closely related system [Chakrabarti et al. 2006] by a large margin, in term of MRR (Mean Reciprocal Rank). We re-implemented ER in our system for scoring individual predicates ( $F^S, F^R$ ). As ER focuses on single-predicate queries, performance on such queries can be fairly compared. For multi-predicate queries, we can also use ER to compute predicate scores and aggregate them by the same formula 2.

We tested ER with our data set. In Table III, ER is worse than BCM on Single-27, by 0.082 in MAP and 0.039 in nDCG. It is only marginally better than the baseline COUNT. For multi-predicate queries (Multi-18), ER shows relatively better performance than PROX, however, it is still worse than BCM by a large margin

Table IV. Compare Effectiveness of Different Weighting Schemes

Query	BCM	BCM- $\alpha$	BCM- $\beta$	BCM- $\gamma$
<b>MAP</b>				
Single-27	0.911	0.910	0.913	0.915
Multi-18	0.763	0.771	0.785	0.788
<b>nDCG</b>				
Single-27	0.940	0.938	0.942	0.944
Multi-18	0.855	0.863	0.866	0.871

(0.06 in MAP and 0.043 in nDCG).

From Figure 1, it can be seen that ER is good at ranking top 2 answers (with the exception of (a)), rivaling CM and BCM. However, it deteriorates very fast when  $k > 2$ . In (c), ER drops below 0.7 around  $k=6$ , while BCM remains above 0.7 at  $k=10$ . It indicates that BCM is more robust for queries with multiple true answers. This is because BCM exploits more features than ER and is thus able to promote the ranking of some *hard* true answers indistinguishable by ER.

The INEX Entity Ranking track targets on a different problem setting. INEX queries are specified as narrative descriptions on the desired entities. Participating systems can use any techniques to answer the queries, but need to understand the query descriptions, which itself is challenging, thus their MAPs may tend to be low. The MAP achieved by the best system participating in the 2009 track is 0.517. To avoid the overhead of assessing participating systems, INEX used a sampling strategy to estimate their MAPs.

INRIA [Vercoustre et al. 2008] works on the same problem as INEX. Unlike INEX participants, it is not based on co-occurrence of entities and query inputs. Rather, it finds the Wikipedia articles that best match query descriptions, through link analysis and TF-IDF weighting. It achieves MAP of 0.390 on 18 topics adapted from INEX 2006 ad hoc track.

In comparison with INEX and INRIA, the MAP achieved by BCM is 0.852 (averaged over all 45 queries). We acknowledge that this comparison is not strictly fair. First, the results are based on different query sets and snapshots of Wikipedia. Second, they focus on different query styles (structured query vs. narrative description). However, our argument is that the high MAP of BCM at least indicates that the structured entity-relationship queries can be highly effective in reality.

### 5.5 Effectiveness of Weighting

In this part, we evaluate different weighting schemes introduced in Section 4. We will compare plain BCM with BCM- $\alpha$  (applying maximal-support weight only) and BCM- $\beta$  (applying corpus-frequency weight only). We also consider the case of combining the two weights as  $\gamma_p(t) = \alpha_p(t) + \beta_p(t)$ . Table IV compares the results using both MAP and nDCG.

The benefits of the two weighting schemes are not perceived on Single-27. One major reason is that, the accuracy of BCM on single-predicate queries is already quite high, any further improvement would be very difficult. Recall that both weighting schemes are designed to penalize susceptible answers instead of promoting promising ones. Their effects tend to diminish when there are few susceptible answers ranked at top positions, which is the case of Single-27. Nevertheless, the combined weighting, BCM- $\gamma$ , does show some marginal improvement over BCM.

The result on Multi-18 is different. Both  $\text{BCM-}\alpha$  and  $\text{BCM-}\beta$  achieve better performance than BCM, by 0.008 and 0.022 respectively in term of MAP. The corpus-frequency weighting ( $\text{BCM-}\beta$ ) appears to be more effective than maximal-support weighting ( $\text{BCM-}\alpha$ ).  $\text{BCM-}\gamma$  benefits from both weighting components, achieving the highest MAP (0.788) and nDCG (0.871). The reason for the improvements is two fold. On the one hand, the sensitivity of formula 2 bring more susceptible answers into top ranks in case of multi-predicate queries. On the other hand, as discussed in Section 4.1, the maximal support (the numerator) in  $\text{BCM-}\alpha$  becomes effective for multi-predicate queries.

## 6. FUTURE WORK

This paper is the first work in the direction of multi-predicate entity-relationship queries. The distinguishing characteristic of our approach is to combine IR-style keyword constraints with SQL-style structured query constructs. As the initial exploratory study along this line, our work has mainly focused on designing an accurate ranking framework and building a prototype system, to demonstrates the effectiveness and promises of the approach. We have thus only supported basic keyword constraints in the queries. To deploy a production query system with good usability in practice, we need to support more advanced query features. Here we briefly enumerate several such features.

Our current definition of entity-relationship query does not support querying entities by name (e.g., find entities whose names contain “Michael”) or querying the relationship between entities (e.g., find the relationship between Bill Gates and Microsoft). One direction of our future work is to support these types of queries.

Although the predicates in entity-relationship query are keyword-based, it can be a burden for users to choose the right keywords. We are exploring two extensions to address this concern. First, query suggestion can directly help users find the proper keywords. For example, after the user types “Stanford” for  $p_1$  of Query 1, the system could suggest a list of keywords that commonly co-occur with “Stanford” in the corpus or in the query log, such as “graduate”, “professor”, etc. Second, the strict keyword matching can be relaxed by query expansion, e.g., allowing a context for  $p_1$  to contain “alumni”, if not “graduate”. Synonym thesaurus and paraphrases mined from the corpus may be used for this purpose. New challenges on ranking shall arise with query expansion.

In addition to position-based features, the assessment of co-occurrence contexts can be improved with more features such as syntactic information (e.g., part-of-speech tags) and lexicographic information. The BCM model may be extended with new factors for these features.

## 7. CONCLUSION

*Entity-relationship query* is a structured facility to query entities over Wikipedia. It distinguishes itself by (1) allowing multiple keyword-based predicates in a query and (2) searching directly in corpus instead of pre-extracted data stores. As a result, entity-relationship query supports semantics expressed with keywords and avoids information loss that happens when pre-extracting facts into data stores. We presented a ranking framework for entity-relationship queries and a Bounded

Cumulative Model under this framework. Our ranking method exploits three intuitive positional features, which are shown to be effective in the experiments. To address the sensitivity problem in our ranking model, we further studied how to detect and weight susceptible answers. We explored two weighting schemes, both leveraging the intuition that predicate scores supported with less contexts are more susceptible and should be penalized accordingly.

## REFERENCES

- AGICHTEIN, E. AND GRAVANO, L. 2000. Snowball: Extracting relations from large plain-text collections. In *DL*.
- AUER, S., BIZER, C., KOBILAROV, G., LEHMANN, J., CYGANIAK, R., , AND IVES, Z. 2007. DBpedia: A nucleus for a Web of open data. In *Int'l Semantic Web Conf*.
- BRIN, S. 1998. Extracting patterns and relations from the world wide web. In *WebDB*.
- CAFARELLA, M. J., HALEVY, A., WANG, D. Z., WU, E., AND ZHANG, Y. 2008. Webtables: exploring the power of tables on the web. *Proc. VLDB Endow.* 1, 1, 538–549.
- CAFARELLA, M. J., RÉ, C., SUCIU, D., ETZIONI, O., AND BANKO, M. 2007. Structured querying of Web text. In *CIDR*.
- CHAKRABARTI, S., PUNIYANI, K., AND DAS, S. 2006. Optimizing scoring functions and indexes for proximity search in type-annotated corpora. In *WWW*.
- CHENG, T., YAN, X., AND CHANG, K. C.-C. 2007. EntityRank: searching entities directly and holistically. In *VLDB*. 387–398.
- CHU, E., BAID, A., CHEN, T., DOAN, A., AND NAUGHTON, J. A relational approach to incrementally extracting and querying structure in unstructured data.
- DEROSE, P., SHEN, W., CHEN, F., DOAN, A., AND RAMAKRISHNAN, R. 2007. Building structured Web community portals: a top-down, compositional, and incremental approach. In *VLDB*.
- DILL, S., EIRON, N., GIBSON, D., GRUHL, D., GUHA, R., JHINGRAN, A., KANUNGO, T., RAJAGOPALAN, S., TOMKINS, A., TOMLIN, J. A., AND ZIEN, J. Y. 2003. SemTag and seeker: bootstrapping the semantic Web via automated semantic annotation. In *WWW*.
- ETZIONI, O., BANKO, M., SODERLAND, S., AND WELD, D. S. 2008. Open information extraction from the Web. *Commun. ACM* 51, 12, 68–74.
- KANDOGAN, E., KRISHNAMURTHY, R., RAGHAVAN, S., VAITHYANATHAN, S., AND ZHU, H. 2006. Avatar semantic search: a database approach to information retrieval. In *SIGMOD*. 790–792.
- KASNECI, G., SUCHANEK, F., IFRIM, G., RAMANATH, M., AND WEIKUM, G. 2008. NAGA: Searching and ranking knowledge. In *ICDE*. 953–962.
- KULKARNI, S., SINGH, A., RAMAKRISHNAN, G., AND CHAKRABARTI, S. 2009. Collective annotation of Wikipedia entities in Web text. In *KDD*. 457–466.
- LI, X., LI, C., AND YU, C. Entity-relationship queries over wikipedia.
- MIHALCEA, R. AND CSOMAI, A. Wikify!: linking documents to encyclopedic knowledge.
- MILNE, D. AND WITTEN, I. H. Learning to link with wikipedia.
- NADEAU, DAVID, SEKINE, AND SATOSHI. *Linguisticae Investigationes*.
- PETKOVA, D. AND CROFT, W. B. 2007. Proximity-based document representation for named entity retrieval. In *CIKM*.
- SUCHANEK, F. 2009. Automated construction and growth of a large ontology. Ph.D. thesis, Saarland University.
- SUCHANEK, F. M., KASNECI, G., AND WEIKUM, G. 2007. YAGO: a core of semantic knowledge unifying WordNet and Wikipedia. In *WWW*.
- VERCOUSTRE, A.-M., THOM, J. A., AND PEHCEVSKI, J. 2008. Entity ranking in wikipedia. In *SAC*.
- ZARAGOZA, H., RODE, H., MIKA, P., ATSERIAS, J., CIARAMITA, M., AND ATTARDI, G. 2007. Ranking very many typed entities on Wikipedia. In *CIKM*.
- ZHOU, M., CHENG, T., AND CHANG, K. C.-C. 2010. Data-oriented content query system: searching for data into text on the web. In *WSDM*.