

Entity-Relationship Queries over Wikipedia

Xiaonan Li
Department of Computer
Science and Engineering
University of Texas at Arlington
xiaonan.li@mavs.uta.edu

Chengkai Li
Department of Computer
Science and Engineering
University of Texas at Arlington
cli@uta.edu

Cong Yu
Yahoo! Research New York
congyu@yahoo-inc.com

ABSTRACT

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 for accurate ranking of query answers. Experiments on INEX benchmark queries and our own crafted queries show the effectiveness and accuracy of our ranking method.

Categories and Subject Descriptors

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

General Terms

Design, Languages, Performance, Experimentation

Keywords

entity search, entity 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. 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.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SMUC'10, October 30, 2010, Toronto, Ontario, Canada.

Copyright 2010 ACM 978-1-4503-0386-6/10/10 ...\$10.00.

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 [8]. For example, a page showing a list of Silicon Valley companies may exist. However, it is unrealistic to expect such pages always 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*¹, 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, our analyst can write the following SQL-like query for Example 1.

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

```
SELECT x, y
FROM PERSON x, COMPANY y
WHERE x:["Stanford", "graduate"] // Predicate p1
      AND y:["Silicon Valley"] // Predicate p2
      AND x,y:["found"] // Predicate p3
```

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-time², instead of explicitly extracting and reasoning about complex semantic information in text before query-time [7, 5, 12, 15, 9, 16, 13]. 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 [7, 5] and entity search and ranking [11].

We implemented a prototype entity-relationship query system in Apache Lucene³. 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. Interested readers are referred to [19] for details on index and query processing.

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 [10, 11] and may not leverage redundancy [10]. 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 conduct comprehensive experiments to verify the effectiveness of the ranking method on both benchmark queries and our own crafted queries.

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). Empirical results are reported in Section 4. Section 5 discusses limitations and future work. Section 6 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 [7, 5, 12, 15, 9, 16, 13]. 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 [25, 17, 14, 6] 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 [2] for querying entities. Some of them [25, 17, 6] 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 [24] unified from WordNet and Wikipedia. Other systems apply IE techniques over Web pages to bootstrap RDF extraction [14], thus bearing the same limitation of the aforementioned DB-based approach.

¹It should not be confused with the well-known ER model.

²The effectiveness of such entity-relationship queries partially relies on the user’s capability in forming proper keyword constraints, like in IR queries.

³<http://lucene.apache.org/>

Entity ranking has gained significant interest recently [23, 3, 4, 27, 26]. However, they focus on different problem settings from ours. Take, for example, the INEX Entity Ranking track [3]. 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.

The studies most related to ours are [10, 11, 28]. Chakrabarti et al. [10] 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 [11] aggregates scores of locally evaluated co-occurrence contexts into global scores to improve ranking. The Content Query Language [28] 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 phrases⁴. 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 $\langle doc, sent, \hat{V}_p, \hat{C}_p \rangle$. doc and $sent$ 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

⁴A single keyword is treated as a phrase of length 1.

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 $\langle 9, 8, \{\langle 3, 4 \rangle\}, \{0, 2\} \rangle$.

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, ϕ , 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, $|token(\text{IBM}, s)|$ 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. Con-

text 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 $\sum_{e \in t_{p_1}} |token(e, s_1)|=2$. The scope covering the entity and the two phrases spans 5 tokens, from “Stanford” to “Yang”, thus $|scope_{p_1}(t_{p_1}, s_1)|=5$. Therefore, the proximity of Jerry Yang in s_1 is $prox_{p_1}(t_{p_1}, s_1)=\frac{2+2}{5}=0.8$. 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 $p_1=\{x, \{\text{“Stanford”, “graduate”}\}\}$ 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 $p'_1=\{x:\text{NOVEL}, \{\text{“by”, “Jane Austen”}\}\}$, 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 ϕ_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 ϕ_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 ϕ_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 Allan 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_5 : 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 ϕ_p . The credit of o in sentence 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 . $t^2=T_{p_1}^*(o_2, s)=\text{Paul Allen}$ 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, $credit_{p_1}(o_1, s)=\frac{4}{4+2}=0.67$ and $credit_{p_1}(o_2, s)=0.33$.

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 ϕ_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 4, below we list three important customizations.

$$\text{COUNT } F_p(t) = \sum_{o \in O_p} (1 \sum_{s \in \phi_p(t, o)} 1) = \sum_{o \in O_p} |\phi_p(t, o)| = |\phi_p(t)|$$

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

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

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

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

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.

Table 1: Example Answers

	x	y	p_1	p_2	p_3	Π	Σ
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

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,

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

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 1. 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. EXPERIMENTS

In this section, we provide the empirical evaluation results of our prototype system implemented in Apache Lucene.

4.1 Data and Query Sets

We used the 2008-07-24 snapshot of Wikipedia⁵. After we removed all the irrelevant pages (such as category and administrative pages), there were about 2.4 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 2) and assigned about 0.75 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 [1]. 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 100 million occurrences of the 0.75 million typed entities.

Named entity recognition (NER) [22] and entity disambiguation [14] are intensively studied problems. Our hyperlink-based occurrence detection can be viewed as a rudimentary

⁵<http://download.wikimedia.org>

Table 2: Ten Types from Wikipedia

Type	(E)ntities	(O)ccurrences	O/E
AWARD	1,045	626,340	600
CITY	70,893	28,261,278	389
CLUB	15,688	5,263,865	335
COMPANY	24,191	9,911,372	409
FILM	41,344	3,047,576	74
NOVEL	16,729	1,036,596	63
PERSON	427,974	38,228,272	89
PLAYER	95,347	2,398,959	25
SONG	29,934	732,175	24
UNIVERSITY	19,717	6,141,840	311
TOTAL	742,862	95,648,273	129

entity disambiguation method. Recently we have seen advanced entity recognition and disambiguation methods using Wikipedia as entity catalog [20, 21, 18] to automatically link entities mentioned in plain text to their corresponding Wikipedia articles. One of our ongoing efforts is to use Wikify⁶ (the system based on [21]) to detect entity occurrences 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.

Two query sets were used for experiments, INEX17 and OWN28. INEX17 is adapted from the topics in the Entity Ranking track of INEX 2009 [3]. From the 60 topics in the track, we adapted the ones that are on our predefined 10 entity types. In this way, we obtained 17 queries, including 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 manually collecting complete true answers for all test queries is prohibitively costly, we alternatively adopted the depth- N pooling approach used by INEX. Basically for each query, we pooled the top N answers returned by all compared methods and manually checked their correctness. This approach allows us to evaluate precision up to rank N . Though recall cannot be measured directly with the pooled ground truth, the Mean Average Precision we used for evaluation reflects recall to certain degree. N is set to 100 in our case.

4.2 Answers to Sample Queries

In this section we use several sample queries to demonstrate the accuracy of our query answering system. 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"> • ⟨Michigan State University⟩ • ⟨Indiana University (Bloomington)⟩ • ⟨Ohio State University⟩ • ⟨University of Iowa⟩ • ⟨University of Illinois at Urbana-Champaign⟩ • ⟨Purdue University⟩ • ⟨University of Wisconsin-Madison⟩ • ⟨Pennsylvania State University⟩ • ⟨University of Michigan⟩ • ⟨Northwestern University⟩

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

SELECT x, y FROM COMPANY x, PERSON y WHERE x:["Silicon Valley"] AND y:["Stanford Graduate"]
<ul style="list-style-type: none"> • ⟨Jerry Yang, Yahoo!⟩ • ⟨Scott McNealy, Sun Microsystems⟩ • ⟨David Packard, Hewlett-Packard Company⟩ • ⟨Vinod Khosla, Sun Microsystems⟩ • ⟨William Hewlett, Hewlett-Packard Company⟩ • ⟨Bill Gates, Microsoft⟩ • ⟨Larry Page, Google⟩ • ⟨Andy Bechtolsheim, Sun Microsystems⟩ • ⟨Vinod Khosla, Kleiner Perkins Caufield & Byers⟩ • ⟨Andy Bechtolsheim, Cisco Systems⟩

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:["Academy Award"] AND y:["Australian" "actor"] AND x,y:["star"]
<ul style="list-style-type: none"> • ⟨Brokeback Mountain, Heath Ledger⟩ • ⟨Gladiator, Russell Crowe⟩ • ⟨A Beautiful Mind, Russell Crowe⟩ • ⟨Braveheart, Mel Gibson⟩ • ⟨American Gangster, Russell Crowe⟩ • ⟨Munich, Eric Bana⟩ • ⟨The Adventures of Priscilla, Hugo Weaving⟩ • ⟨The Adventures of Priscilla, Guy Pearce⟩ • ⟨The Insider, Russell Crowe⟩ • ⟨L.A. Confidential, Russell Crowe⟩

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"> • ⟨Golden State Warriors, Rick Barry⟩ • ⟨Chicago Bulls, Michael Jordan⟩ • ⟨Los Angeles Lakers, Magic Johnson⟩ • ⟨Los Angeles Lakers, Kareem Abdul-Jabbar⟩ • ⟨Boston Celtics, Larry Bird⟩ • ⟨San Antonio Spurs, Tim Duncan⟩ • ⟨Detroit Pistons, Isiah Thomas⟩ • ⟨Los Angeles Lakers, Shaquille O’Neal⟩ • ⟨Los Angeles Lakers, James Worthy⟩ • ⟨Houston Rockets, Moses Malone⟩

4.3 Comparing Ranking Methods

In this section, we compare and analyze the multiple ranking methods discussed earlier, namely COUNT, PROX, MEX, CM and BCM. All the methods differ in how they compute predicate scores, i.e., $F_p(t)$. For multi-predicate queries, the same aggregate function, product, is used to compute answer scores, $F^A(t)$. We compare these ranking methods using three popular measures: $nDCG$, MAP , and $Precision-at-k$.

nDCG (Normalized Discounted Cumulative Gain): The first block in Table 3 shows the average nDCG on single-predicate queries (Single-11), multi-predicate queries (Multi-6), and all queries (All-17) from INEX17. Both MEX and PROX improve over COUNT, by 0.02-0.05 across all three cases. PROX appears to be more effective than MEX. CM and BCM are comparable to PROX on Single-11, but further improve by more than 0.02 on Multi-6. We only observe minor difference between CM and BCM.

MAP (Mean Average Precision): The second block of Table 3 shows the MAP on INEX17. The observations are

⁶<http://wdm.cs.waikato.ac.nz:8080/>

Table 3: MAP and nDCG on INEX17/OWN28

Query	COUNT	MEX	PROX	CM	BCM	ER
nDCG on INEX17						
Single-11	0.889	0.911	0.920	0.920	0.920	0.904
Multi-6	0.880	0.918	0.932	0.954	0.958	0.927
All-17	0.886	0.913	0.924	0.932	0.933	0.912
MAP on INEX17						
Single-11	0.756	0.812	0.843	0.844	0.842	0.779
Multi-6	0.772	0.820	0.852	0.885	0.894	0.809
All-17	0.762	0.815	0.846	0.859	0.860	0.790
nDCG on OWN28						
Single-16	0.917	0.943	0.947	0.953	0.954	0.923
Multi-12	0.800	0.812	0.836	0.844	0.878	0.781
ALL-28	0.867	0.887	0.899	0.906	0.922	0.862
MAP on OWN28						
Single-16	0.758	0.825	0.838	0.858	0.853	0.760
Multi-12	0.579	0.620	0.660	0.684	0.748	0.521
ALL-28	0.681	0.738	0.762	0.783	0.808	0.658

mostly similar to those from the nDCG analysis. Note that a larger distinction between CM and BCM is observed on Multi-6, with BCM about 0.01 better than CM.

For further investigation, we repeat the above experiments on OWN28 and provide the results in the bottom half of Table 3. Most results are consistent with INEX17. However, on multi-predicate queries in OWN28 (Multi-12), BCM shows clear advantage over CM in terms of both nDCG (by 0.034) and MAP (by 0.064). The different observations on INEX17 and OWN28 is because, we believe, OWN28 has more multi-predicate queries than INEX17 and the advantage of BCM is more stably observed on OWN28.

Precision-at- k : To further analyze how various methods perform at different ranks, we plot precision-at- k curves. Figure 1(a)(b) shows the results for $k=10$. 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 nDCG and MAP analysis. BCM is consistently the best among all, while CM has inconsistent performance on INEX17 and OWN28. Figure 1(c)(d) show the results for $k=50$. Each curve shows the average precision of the corresponding method for queries that returned 50 or more answers, including 7 queries in INEX17 and 18 in OWN28. In Figure 1(c), CM and BCM excel before $k=10$ and BCM is slightly better. PROX is the best after $k=10$ but is significantly worse than BCM at top ranks. In Figure 1(d), BCM is clearly the best among all, although a little worse than CM between 10 and 25.

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

4.4 BCM vs. Other Entity Ranking Methods

As a first study on multi-predicate entity-relationship query, we did not find directly comparable systems. Instead, we chose three state-of-the-art systems proposed for related problems: EntityRank, INEX and INRIA. All three systems used Wikipedia as corpus and entity catalog, though INEX and INRIA used different snapshots from ours.

EntityRank (ER) [11] outperforms another closely related method [10] by a large margin, in term of MRR. We re-implemented ER as a plugin 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, the predicate scores computed by ER are aggregated with the same function, product, for computing answer scores (F^A).

We tested ER on our data set. In Table 3, both CM and BCM outperform ER by large margins. The advantage of CM/BCM over ER is more clear on multi-predicate queries than on single-predicate ones. The peak margin (0.22) in terms of MAP is observed on Multi-12 from OWN28, between BCM and ER. In Figure 1, ER rivals PROX, CM, and BCM at top-2, verifying the high MRR reported in [11]. However, it deteriorates very fast when $k > 2$, dropping below 0.7 around $k=5$, while BCM remains above 0.7 even 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.

INEX Entity Ranking track [3] focuses 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 [26] works on the same problem as INEX. Unlike INEX participants, it is not based on co-occurrence of entities and query inputs. Rather, it ranks entities by 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 on INEX17 is 0.860. We acknowledge that this comparison is not strictly fair. First, the results are based on different query sets (INEX17 is a subset of INEX Entity Ranking topics) 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. 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 demonstrate 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

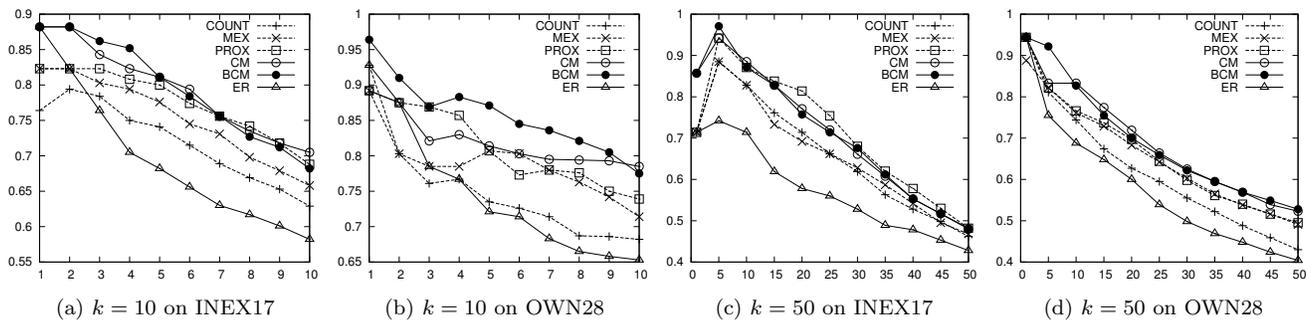


Figure 1: Precision-at- k on INEX17/OWN28

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.

6. 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 on both benchmark queries and our own crafted queries.

7. REFERENCES

- [1] <http://tartarus.org/martin/porterstemmer/>.
- [2] <http://www.w3.org/tr/rdf-sparql-query>.
- [3] INEX 2009 entity-ranking track. <http://www.inex.otago.ac.nz/tracks/entity-ranking/entity-ranking.asp>.
- [4] TREC 2009 entity track: Searching for entities and properties of entities. <http://ilps.science.uva.nl/trec-entity/guidelines/>.
- [5] E. Agichtein and L. Gravano. Snowball: Extracting relations from large plain-text collections. In *DL*, 2000.
- [6] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. DBpedia: A nucleus for a Web of open data. In *Int'l Semantic Web Conf.*, 2007.
- [7] S. Brin. Extracting patterns and relations from the world wide web. In *WebDB*, 1998.
- [8] M. J. Cafarella, A. Halevy, D. Z. Wang, E. Wu, and Y. Zhang. Webtables: exploring the power of tables on the web. *Proc. VLDB Endow.*, 1(1):538–549, 2008.
- [9] M. J. Cafarella, C. Ré, D. Suciu, O. Etzioni, and M. Banko. Structured querying of Web text. In *CIDR*, pages 225–234, 2007.
- [10] S. Chakrabarti, K. Punyani, and S. Das. Optimizing

scoring functions and indexes for proximity search in type-annotated corpora. In *WWW*, 2006.

- [11] T. Cheng, X. Yan, and K. C.-C. Chang. EntityRank: searching entities directly and holistically. In *VLDB*, pages 387–398, 2007.
- [12] E. Chu, A. Baid, T. Chen, A. Doan, and J. Naughton. A relational approach to incrementally extracting and querying structure in unstructured data. In *VLDB*, pages 1045–1056, 2007.
- [13] P. DeRose, W. Shen, F. Chen, A. Doan, and R. Ramakrishnan. Building structured Web community portals: a top-down, compositional, and incremental approach. In *VLDB*, 2007.
- [14] S. Dill, N. Eiron, D. Gibson, D. Gruhl, R. Guha, A. Jhingran, T. Kanungo, S. Rajagopalan, A. Tomkins, J. A. Tomlin, and J. Y. Zien. SemTag and seeker: bootstrapping the semantic Web via automated semantic annotation. In *WWW*, 2003.
- [15] O. Etzioni, M. Banko, S. Soderland, and D. S. Weld. Open information extraction from the Web. *Commun. ACM*, 51(12):68–74, 2008.
- [16] E. Kandogan, R. Krishnamurthy, S. Raghavan, S. Vaithyanathan, and H. Zhu. Avatar semantic search: a database approach to information retrieval. In *SIGMOD*, pages 790–792, 2006.
- [17] G. Kasneci, F. Suchanek, G. Ifrim, M. Ramanath, and G. Weikum. NAGA: Searching and ranking knowledge. In *ICDE*, pages 953–962, 2008.
- [18] S. Kulkarni, A. Singh, G. Ramakrishnan, and S. Chakrabarti. Collective annotation of Wikipedia entities in Web text. In *KDD*, pages 457–466, 2009.
- [19] X. Li, C. Li, and C. Yu. Structured querying of annotation-rich web text with shallow semantics. Technical report, Univ. of Texas at Arlington, 2010.
- [20] R. Mihalcea and A. Csomai. Wikify!: linking documents to encyclopedic knowledge. In *CIKM*, 2007.
- [21] D. Milne and I. H. Witten. Learning to link with wikipedia. In *CIKM*, 2008.
- [22] Nadeau, David, Sekine, and Satoshi.
- [23] D. Petkova and W. B. Croft. Proximity-based document representation for named entity retrieval. In *CIKM*, 2007.
- [24] F. Suchanek. *Automated Construction and Growth of a Large Ontology*. PhD thesis, Saarland University, 2009.
- [25] F. M. Suchanek, G. Kasneci, and G. Weikum. YAGO: a core of semantic knowledge unifying WordNet and Wikipedia. In *WWW*, 2007.
- [26] A.-M. Vercoustre, J. A. Thom, and J. Pehcevski. Entity ranking in wikipedia. In *SAC*, 2008.
- [27] H. Zaragoza, H. Rode, P. Mika, J. Atserias, M. Ciaramita, and G. Attardi. Ranking very many typed entities on Wikipedia. In *CIKM*, 2007.
- [28] M. Zhou, T. Cheng, and K. C.-C. Chang. Data-oriented content query system: searching for data into text on the web. In *WSDM*, 2010.