

Facet Discovery for Structured Web Search: A Query-log Mining Approach

Jeffrey Pound^{*}
University of Waterloo
Waterloo, Canada
jpound@cs.uwaterloo.ca

Stelios Pappas
Microsoft Research
Mountain View, CA, USA
steliosp@microsoft.com

Panayiotis Tsaparas
Microsoft Research
Mountain View, CA, USA
panats@microsoft.com

ABSTRACT

In recent years, there has been a strong trend of incorporating results from structured data sources into keyword-based web search systems such as Google or Amazon. When presenting structured data, facets are a powerful tool for navigating, refining, and grouping the results. For a given structured data source, a fundamental problem in supporting faceted search is finding an ordered selection of attributes and values that will populate the facets. This creates two sets of challenges: First, because of the limited screen real-estate, it is important that the top few facets best match the anticipated user intent. Second, the huge scale of available data to engines like Google or Amazon demands an automated unsupervised solution.

In this paper, we model the user faceted-search behavior using the intersection of web query logs with existing structured data. Since web queries are formulated as free-text queries, a challenge in our approach is the inherent ambiguity in mapping keywords to the different possible attributes of a given entity type. We present an automated solution that elicits user preferences on attributes and values, employing different disambiguation techniques ranging from simple keyword matching, to more sophisticated probabilistic models. We demonstrate experimentally the scalability of our solution by running it on hundreds of categories of diverse entity types, and a large query log. We also present an evaluation of the quality of our results based on a user-study with Mechanical Turk.

Categories and Subject Descriptors

H.2.8 [Database Management]: Database applications—*Data Mining*

General Terms

Algorithms, Experimentation, Human Factors

1. INTRODUCTION

In recent years, search engines like Google or Bing have evolved to include in their results information from structured data sources

^{*}Work done while at Microsoft Research.

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.

SIGMOD'11, June 12–16, 2011, Athens, Greece.

Copyright 2011 ACM 978-1-4503-0661-4/11/06 ...\$10.00.

along with text documents. Furthermore, there are specialized search verticals like Amazon that rely on an even greater variety of structured data sources for their results. Such engines provide traditional keyword search capability to their users. However, users are often interested in exploring a structured collection of data rather than query for a specific item. For example, a user who is looking to buy a new digital camera is likely interested in browsing through a *consideration set* of products that meets her requirements, rather than directly target a specific model. Creating such a consideration set through a keyword search interface is a process that requires sufficient domain knowledge from the user, and several iterations over large amounts of data. Furthermore, this process does not make use of the rich meta-data associated with structured data, which can be used to facilitate the exploration of the user.

One common interaction paradigm that utilizes structured information is *faceted search*. Faceted search provides a more user-friendly visual alternative to keywords for the user to explore the structured results. Facets can be thought of as summarized data axes that allow the users to pivot on the information shown, and create a different or more restrictive consideration set. If we abstract the structured data as typed entities, facets correspond to attributes and values of such entities. A facet system consists of an ordered list of facet attributes, and an ordered list of values for each attribute. For example, for a user looking for digital cameras, exposing popular attributes such as brand, or megapixel resolution as facets, enables them to quickly zoom into the appropriate subset of products by selecting the values of interest. Similar examples exist for other types of structured data, such as auto search, or house search.

It is clear that the user experience for structured web search can benefit from facets. However for an engine to support facets effectively it needs to address two groups of challenges. (a) Given the limited screen real estate and the large number of possible facets to consider, we need select the top- k most important facets, where k is usually a small number. Facet importance in this case can be measured by the *utility* of a facet towards a user's anticipated action like a pivot or refinement. Since an entity-type can have over one hundred attributes, the challenge becomes finding the few most important attributes, and attribute values with the maximum anticipated utility. (b) There is a huge number of structured data sources currently available to engines like Google or Amazon. If we abstract the data as de-normalized entity-type tables, there are thousands of such tables to consider. So any solution that finds important attributes from these tables must be fully automated otherwise it will not scale appropriately.

The best method to learn the utility of attributes for users would be to consume user-produced signals. Such an ideal signal would be produced if all possible attributes for all possible entity-type ta-

bles are shown and then the system observes and records which ones are being selected by an extensive user study. Unfortunately, this is not feasible due to the scale of data and the number of required users for statistical significance. Instead, common practice [10, 12] has been to show attributes selected manually by experts with values populated from sampling the results. It is debatable how well a domain expert can match the average user intent, but regardless, the scale of structured data available to web engines makes it impossible to go this route with consistently good quality across all tables. As a result, existing approaches [4, 24] rely on automatically learning facets from the structured data. However, such approaches do not capture the true user intent since they focus on the structured data without incorporating user signals.

In this work we use web query logs intersected with the structured data to model the facet utility. Our intuition is that query logs act as a proxy for the exploration process of the users when navigating through a structured domain. We assume that users issue queries to refine or alter their results into an appropriate consideration set in a similar manner they would have used facets if there were available. The attributes that appear frequently in the query logs are the ones that users would use to narrow down their search, and thus they are appropriate for facets. The values of these attributes that appear frequently in the queries and are good values to populate the facets. This formulation has the advantage that it can lead to a scalable automated solution that captures the real user needs better than the opinion of an expert or the just the statistics of the structured data. Additionally, it adapts to the needs and preferences of the users as they change over time. It is worth mentioning that trying to match the user intent for attribute importance has been studied in the past, see [9, 23], with the application domain being either attribute ranking or product visibility. Although the high-level principle is similar, the underlying models and corresponding solutions are different from those for discovering the facet utility, as we present in this paper.

Using the web query logs to find frequently queried values and attributes would be trivial if the attribute-value pairs were clearly specified in the queries as in SQL statements. However, although the data is structured, web queries are free-text keywords. Finding the frequent attributes poses two hard challenges due to the inherent ambiguity of translating the keywords to attribute-value pairs: (a) Mapping the query to the appropriate structured source – e.g., map the keyword *apple* to either mp3 players, computers, or fruit; (b) identifying and disambiguating the occurrence of attributes in the query amongst attributes of the same entity-type – e.g., for the table *televisions*, disambiguate for the keyword *30inch* between the diagonal screen size, the television height, and the television width; or for the table *watches*, map the keyword *gold* to the color or the material of a watch.

In this paper, we assume that the first problem, understanding if a query targets a source of structured data and which source in particular, has been addressed from existing literature, such as [15, 21, 27, 29]. We use the existing work to find the appropriate target data source. We focus on the second problem, disambiguating between attributes within an entity-type domain. Using appropriate disambiguation, we learn the most suitable facet attributes and values. Furthermore, we also discuss how data statistics can be used to help with the disambiguation process by being tolerant to data noise and imperfections, and also which data statistics indicate potentially good facets.

Our contributions are summarized as follows: (i) We formulate the problem of query-log based facet attribute and value selection, presented in Section 2. (ii) We introduce different attribute-value identification and disambiguation techniques in Section 3, which

we extend for computing facet values, including query-dependent ones, in Section 4. (iii) We explore the effect of appropriate data statistics in facet selection, presented in Section 5. And (iv), we perform a detailed experimental study on a large scale data set with a thorough evaluation on the quality of the results from real users, presented in Section 6, where we discuss the qualitative difference between our techniques as well as show our best approach outperforming state-of-the-industry systems. We conclude with a discussion of related work in Section 7 and some final thoughts in Section 8.

2. PRELIMINARIES

In this section we formulate our problem, and provide the necessary definitions for the remainder of the paper.

2.1 Problem Formulation

We assume that structured data are organized as a collection of tables $\mathcal{T} = \{T_1, T_2, \dots, T_\tau\}$ ¹. A table T is a set of related entities defined over a set of attributes, $T.A = \{T.A_1, T.A_2, \dots, T.A_\alpha\}$. We use $T.A.V$ to denote the domain of the attribute $T.A$. We omit the prefix T and use A , whenever the table is unambiguous.

Given a table T we define a facet system $\mathcal{F} = \langle F_1, \dots, F_k \rangle$ for table T as an ordered list of k facets, where k is the facet budget. A facet $F_i = \{A_i, \langle A_i.v_1, \dots, A_i.v_{m_i} \rangle\}$ consists of an attribute A_i and an ordered list of m_i values, where m_i is the value budget for attribute A_i . For the following we use $\mathcal{A}_{\mathcal{F}}$ (or $\mathcal{F}.A$) to denote the set of attributes of \mathcal{F} , and $A.V_{\mathcal{F}}$ to denote the set of values for attribute $A \in \mathcal{A}_{\mathcal{F}}$ in facet system \mathcal{F} . We say that a facet system supports an attribute A if $A \in \mathcal{A}_{\mathcal{F}}$, and we define the attribute-support function as follows:

$$\text{ASUP}(\mathcal{F}, A) = \begin{cases} 1, & \text{if } \mathcal{F} \text{ supports } A \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

We also say that a facet attribute $A_f \in \mathcal{A}_{\mathcal{F}}$ supports value v if $v \in A_f.V_{\mathcal{F}}$, and we define the attribute-value support function as follows:

$$\text{VSUP}(A_f, v) = \begin{cases} 1, & \text{if } A_f \text{ supports } v \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

A facet system is exposed through the user interface, and it enables the user to navigate through the data in the table T by selecting values for the different attributes. The design and analysis of faceted user interfaces is an active research area in the HCI community [12]. In any faceted search system, the facet selection is paramount to the success of the system. The order of facets is also important since it determines how the facets will be displayed, with more important facets being displayed more prominently. The goal is to maximize the user engagement and satisfaction, and the i -th facet corresponds to the i -th best attribute for this task. A similar reasoning applies also to the ordering of the facet values within a facet. However, for a selected facet attribute there is usually a greater flexibility in the presentation of values (drop-down menus, slider bars, text boxes) making the value selection and ordering less critical. We thus consider the *facet attribute ordering* problem and the *facet attribute value ordering* problem separately, with the former being more important than the latter.

In this work we are interested in constructing a facet system using web search queries. We assume that we have a query log $\mathcal{Q} = \{q_1, \dots, q_n\}$, consisting of queries posed to a generic web

¹The organization of data into tables is purely conceptual and orthogonal to the underlying storage layer: the data can be physically stored in XML files, relational tables, retrieved from remote web services, etc. Our assumption is that a mapping between the storage layer and the “schema” of table collection \mathcal{T} has been defined.

search engine. For each query q_i we also have a weight w_i denoting the importance of the query, e.g., the number of times that the query appears in the query log.

Although queries are posed to the search engine, it is common that users make queries that are targeted to structured data. For example, the query “*canon powershot*” can be interpreted as targeting a table with structured data about digital cameras. For a given table T let $\mathcal{Q}_T \subseteq \mathcal{Q}$ denote the subset of queries that target table T . We will drop the subscript T whenever the table is unambiguous. We assume that each query $q \in \mathcal{Q}_T$ specifies a set of attribute-value pairs $\mathcal{AV}_q = \{AV_i\}$ of the form $AV_i = (A_i, A_i.v_j)$. For example, the query “*canon powershot*” specifies the value *canon* for attribute “brand”, and the value *powershot* for attribute “product line”. We can thus map the query $q = \text{“canon powershot”}$ to the set $\mathcal{A}_q = \{(\text{brand}, \text{canon}), (\text{product line}, \text{powershot})\}$. We use \mathcal{A}_q and \mathcal{V}_q to denote the set of attributes and values specified in the query. Figure 1(a) depicts graphically an example. Note that we only assume that such a mapping exists; obtaining it from the web log is a challenging problem as we elaborate below.

We now define the *attribute utility* of a facet system \mathcal{F} . For the definition we make use of the attribute support function defined in Equation 1.

DEFINITION 1 (ATTRIBUTE UTILITY). *Given a facet system \mathcal{F} , and a query q , we define the attribute utility of \mathcal{F} for q as*

$$U_a(\mathcal{F}, q) = \sum_{A \in \mathcal{A}_q} \text{ASUP}(\mathcal{F}, A)$$

the number of attributes in \mathcal{A}_q supported by \mathcal{F} . The attribute utility of \mathcal{F} for query log \mathcal{Q} is defined as

$$U_a(\mathcal{F}, \mathcal{Q}) = \sum_{q \in \mathcal{Q}} w_q U_a(\mathcal{F}, q)$$

We can similarly define the *attribute-value utility* of a facet attribute A_f , using the attribute-value support function defined in Equation 2.

DEFINITION 2 (ATTRIBUTE-VALUE UTILITY). *Given a facet attribute A_f and a query q , we define the attribute-value utility of A_f for q as*

$$U_v(A_f, q) = \sum_{(A, v) \in \mathcal{AV}_q: A=A_f} \text{VSUP}(A_f, v)$$

the number of values in \mathcal{AV}_q supported by A_f . The attribute-value utility of A_f for query log \mathcal{Q} is defined as

$$U_v(A_f, \mathcal{Q}) = \sum_{q \in \mathcal{Q}} w_q U_v(A_f, q)$$

As we have already argued, queries act as a proxy for the exploration process of the users over the structured data. Users include in their queries attributes over which they want to dissect the table. A good faceted system should anticipate the information need of the user and select attributes (and values) that maximize the utility to the users. We thus have the following definition of the *query-based facet selection problems*.

PROBLEM 1 (QUERY-BASED FACET ATTRIBUTE SELECTION). *Given a table T , and a query log \mathcal{Q}_T over the table T , create an ordering $\mathcal{L}(T.\mathcal{A})$ of the attributes $T.\mathcal{A}$ such that for any $k \geq 1$, a facet system \mathcal{F} with attribute set $\mathcal{F}.\mathcal{A} = \mathcal{L}_k(T.\mathcal{A})$, the top- k attributes according to ordering \mathcal{L} , has maximum attribute utility $U_a(\mathcal{F}, \mathcal{Q}_T)$ over all facet systems with attribute set of size k .*

PROBLEM 2 (QUERY-BASED FACET VALUE SELECTION). *Given a table T , an attribute $A \in T.\mathcal{A}$, and a query log \mathcal{Q}_T over the table T , create an ordering $\mathcal{L}(A.\mathcal{V})$ of the values in $A.\mathcal{V}$ such that for any $m \geq 1$, a facet system \mathcal{F} that contains A with attribute value set $A.\mathcal{V}_{\mathcal{F}} = \mathcal{L}_m(A.\mathcal{V})$, the top- m attributes according to ordering \mathcal{L} , has maximum attribute-value utility $U_v(A, \mathcal{Q}_T)$ over all facet systems that contain A with attribute-value set of size m .*

If we are given the mapping \mathcal{AV}_q from queries to attribute value pairs, as in Figure 1(a) then the facet selection problems we defined above are easy to solve. This is the case for example with structured query logs (e.g. SQL), where the user specifies explicitly the attributes and values they are interested in. Then, it is straightforward to see that the optimal solution for Problem 1 is to order the attributes according to their *popularity* in the query log. The popularity $F_{\mathcal{Q}}(A_i)$ of attribute A_i over query log \mathcal{Q} is computed as

$$F_{\mathcal{Q}}(A_i) = \sum_{q \in \mathcal{Q}} \sum_{(A_i, v) \in \mathcal{AV}_q} w_q \quad (3)$$

The popularity $F_{\mathcal{Q}}(A_i)$ is thus the number of attribute-value pairs in \mathcal{Q} that contain attribute A_i weighted by the weight of the query in which they appear.

Similarly, the optimal solution for Problem 2 is to order the values within attribute A_i according to their popularity in the query log. The popularity of a value v_j of attribute A_i is computed as

$$F_{\mathcal{Q}}(A_i.v_j) = \sum_{q \in \mathcal{Q}} \sum_{(A_i, v_j) \in \mathcal{AV}_q} w_q$$

the number of attribute value pairs in the query log, where A_i is the attribute and v_j is the value.

Unfortunately, in the case of web query logs this mapping is not given to us and we need to discover it from the data. We elaborate on this challenge in the following section.

2.2 Mapping Web Queries to Structured Data

In order to extract attribute value pairs from the queries in the unstructured web query log, we need to address two issues: (a) for every table T , identify the queries in \mathcal{Q} that target table T , i.e., determine the set \mathcal{Q}_T , and (b) identify and disambiguate the attribute occurrences in the queries. As previously discussed, in this paper we focus on the latter problem.

Given a query $q \in \mathcal{Q}_T$ we now need to identify the occurrences of the attributes in $T.\mathcal{A}$ in the query q . We assume a *tagging process* TAGGER that performs this task. Define a *token* to be a string, and let L_q denote the set of all tokens in q consisting of one or more contiguous words. The TAGGER produces all tokens $t \in L_q$, and matches them with attributes in $T.\mathcal{A}$, producing attribute-token pairs, $AT = (A, t)$. For a categorical attribute $A_c \in T.\mathcal{A}$, the TAGGER outputs an attribute-token pair $AT = (A_c, t)$, if the token t appears in $A_c.\mathcal{V}$. Approximate, or synonymous matching is also possible as an extension of the same process, although discussing it in detail is outside the scope of this paper. For a numerical attribute A_n with unit u (e.g., inches, Kg, or x zoom), the TAGGER produces an attribute-token pair $AT = (A_n, t)$ for every token t that consists of a number followed by a unit u . For example, the tagger output for the query “*5x nikon*” over the Digital Cameras table is $\{(\text{digital zoom}, “5x”), (\text{optical zoom}, “5x”), (\text{brand}, “nikon”)\}$.

For a query q we use $\mathcal{AT}_q = \{(A, t)\} \subseteq T.\mathcal{A} \times L_q$ to denote the set of attribute-token pairs output by the TAGGER. Note that attribute-token pairs are different from attribute-value pairs, since a token may not be a value in the table T . This is the case with synonyms and approximate matches, as well as with numeric attributes: in the query “*32 inch lcd TV*” the token “*32 inch*” will

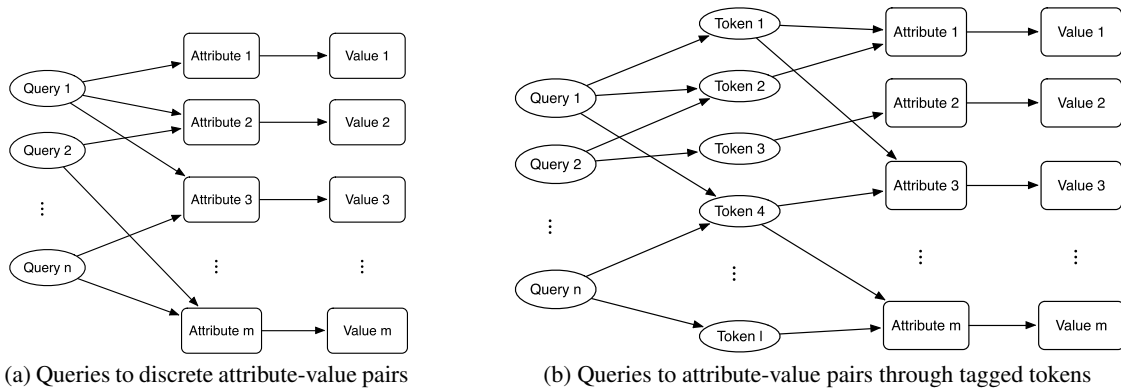


Figure 1: Query to Attribute-Value Pairs Mappings

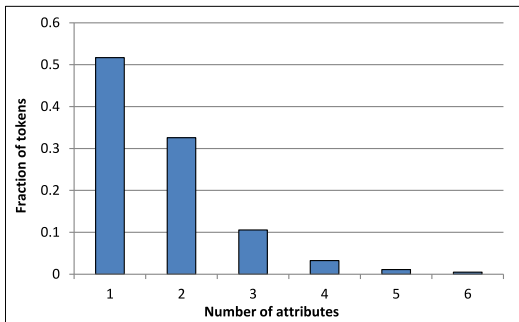


Figure 2: Fraction of tokens mapped to varying numbers of attributes.

map to both height and diagonal screen size in the Television table, even if there are no TVs with an exact value of 32 inches for height or diagonal.

Given the query log \mathcal{Q} the TAGGER outputs a set of attribute-token pairs $\mathcal{AT}_{\mathcal{Q}}$. It would appear that our job is now done, since we can compute the popularity of attributes and values as before, simply by interchanging tokens with values, and attribute-value pairs with attribute-token pairs. However, this is not the case due to *ambiguity* in the token interpretation. The domains of attributes in a table often overlap, and as a result a token t can map to more than one attribute. For example, in the query "32 inch lcd TV", the token "32 inch" is potentially mapped to diagonal, width, height, and depth attributes, since all these attributes are numeric with the same unit; similarly in the query "camera with 3x zoom", "3x" can be either optical or digital zoom. Such confusion is not limited to just numeric attributes. For example in the query "gold watch", the token "gold" can specify either the material, or the color attribute; similarly for "leather dress shoes" the token "leather" can map to either the top shoe material or the shoe sole. There are many such examples that appear in the web queries. We computed the number of attributes of the same table that a token matches on average, over all queries. As we can see from Figure 2, approximately half the tokens are confused over at least two or more attributes. So ambiguity is a real problem that needs to be addressed.

Estimating attribute popularity in the presence of ambiguity becomes problematic. Our data no longer looks like Figure 1(a), but instead like Figure 1(b). Computing the attribute popularity over attribute-token pairs using Equation 3 will lead to misleading results: TV height will be deemed equally important to TV diagonal, and camera digital zoom the same as optical zoom. In order to

obtain correct estimates for the true popularity of an attribute, we need to *disambiguate* the ambiguous tokens.

Let \tilde{t} denote a token in a query q . We say that token \tilde{t} is ambiguous for a table T , if there are at least two attribute-token pairs (A_i, \tilde{t}) involving \tilde{t} . We use $\mathcal{A}_{\tilde{t}}$ to denote the set of candidate attributes to which the token \tilde{t} is mapped. Our goal is to disambiguate between the attributes in $\mathcal{A}_{\tilde{t}}$. We do this by estimating the probability of an attribute given the ambiguous token. We thus have the following problem definition.

PROBLEM 3 (ATTRIBUTE DISAMBIGUATION). *Given a table T , a query log \mathcal{Q} over T , and an ambiguous token $\tilde{t} \in L_{\mathcal{Q}}$ that maps to the set of attributes $\mathcal{A}_{\tilde{t}} \subseteq T.A$, compute the disambiguation probability distribution $P(A|\tilde{t})$ over the set of candidate attributes $A \in \mathcal{A}_{\tilde{t}}$.*

Given $P(A|\tilde{t})$ we now have a measure of likelihood for the mapping of token \tilde{t} to attribute A . We can use this probability to perform either hard or soft disambiguation. In hard disambiguation, we map token \tilde{t} to the most likely attribute

$$A_{\tilde{t}} = \arg \max_{A \in \mathcal{A}_{\tilde{t}}} P(A|\tilde{t})$$

In hard disambiguation our data will take the form of Figure 1(a), and we can apply directly Equation 3 to estimate the attribute popularity. In soft disambiguation, we modify Equation 3 such that each occurrence of attribute A with the ambiguous token \tilde{t} is weighted by the probability $P(A|\tilde{t})$. We thus have

$$\tilde{F}_{\mathcal{Q}}(A_i) = \sum_{q \in \mathcal{Q}} \sum_{(A_i, t) \in \mathcal{AT}_q} w_q P(A_i|t) \quad (4)$$

We have that $P(A|t) = 1$ in the case token t maps unambiguously to attribute A .

3. ATTRIBUTE DISAMBIGUATION

In this section we describe four different approaches for attribute disambiguation. The first three rely on the principle that the best solution is the one that better *explains* the attribute-token pairs we observe. If we assume that the attributes in the table are generating queries and tokens, we identify the attribute that is most likely to have generated the tokens we are observing. The three different algorithms differ in the granularity at which they try to explain the data, becoming progressively more complex. The last approach relies on user feedback to disambiguate between different interpretations of a token.

3.1 Token-level Disambiguation

The first algorithm we consider treats each token independently and tries to find the attribute that is most likely to have generated this specific token. Let \tilde{t} denote an ambiguous token and let $\mathcal{A}_{\tilde{t}}$ the set of candidate attributes for \tilde{t} . We want to estimate $P(A|\tilde{t})$ for all $A \in \mathcal{A}_{\tilde{t}}$. Using Bayes rule we have

$$P(A|\tilde{t}) = \frac{P(\tilde{t}|A)P(A)}{\sum_{A_i \in \mathcal{A}_{\tilde{t}}} P(\tilde{t}|A_i)P(A_i)}$$

We will estimate the right-hand side using the data in T . Assuming that all attributes are equally likely in the table, then $P(A|\tilde{t})$ is proportional to $P_T(\tilde{t}|A)$, the probability that the token \tilde{t} is generated from the distribution of A in the table. We have that

$$P_T(\tilde{t}|A_i) = \frac{|T(A, \tilde{t})|}{|A|}$$

where $T(A, \tilde{t})$ denotes the set of entries in the table where the attribute A takes the value \tilde{t} , and $|A|$ is the number of table entries with some value for attribute A . Numeric attributes are similarly described using histogram estimations that are better suited than discrete counts to the continuous nature of numbers.

Therefore, according to this algorithm, the most likely attribute is the one that is most likely to have generated token \tilde{t} . Note that this is also the attribute whose value distribution can encode token \tilde{t} with the fewest number of bits, i.e., the one that reduces the uncertainty the most for \tilde{t} .

3.2 Cluster-level Disambiguation

The token-level disambiguation approach considers the probability of each token mapping to a candidate attribute independent of the rest. A natural extension to this model is to consider the ambiguity among a cluster of tokens confused with the same set of attributes. In this cluster-based approach, we aggregate over all tokens in the query-log to find clusters of ambiguous attribute-token pairs. We then resolve which attributes in the cluster are better candidates to model the distribution of confused tokens over the *full* cluster.

Consider the bi-partite graph formed by the set of attribute-token pairs in $\mathcal{AT}_{\mathcal{Q}}$ (e.g., see Figure 3 (a)). Each time a token is mapped to a set of attributes, it *supports* the ambiguity of these attributes. For example, if the token “8x” is mapped to both digital zoom and optical zoom, it supports the ambiguity of these two attributes. If the token “2x” is mapped to digital zoom, optical zoom, and model it supports the ambiguity of all three of these attributes. By aggregating over all tokens in the query log, we can find the groups of attributes most commonly confused. In particular, if we consider the set of all tokens that form a bi-clique with a set of attributes, then these attributes are all pairwise ambiguous over the same set of tokens, and the support for the ambiguity of this cluster is proportional to the number of tokens in the bi-clique. It is these strongly supported bi-cliques that we aim to find.

To find clusters of ambiguous attributes, we first construct an *attribute ambiguity graph* on a per table basis. For a table T , consider a graph $G = (V, E, w)$ where there exists a vertex $v_a \in V$ for every attribute $a \in \mathcal{T.A}$. We create an edge $e = \langle v_{A_1}, v_{A_2} \rangle$ between two vertices v_{A_1} and v_{A_2} if some token in the query log is mapped to both A_1 and A_2 over all interpretations of queries. The weight function $w(e)$ assigns a weight to the edge e proportional to the number of tokens confused between (the attributes denoted by) v_{A_1} and v_{A_2} . We then enumerate all cliques in the attribute ambiguity graph, with the support being the minimum edge weight in the clique. For example, in Figure 3 (b), the support of the 2-clique containing digital zoom and optical zoom is ten, while the

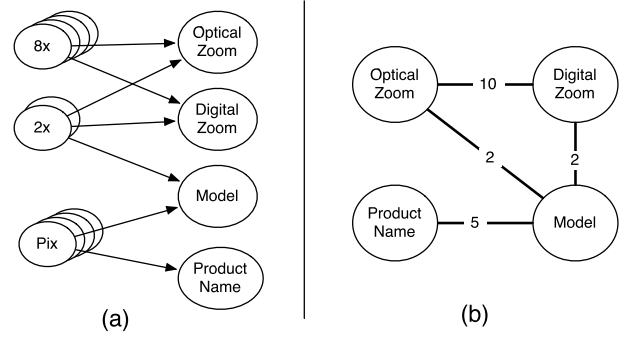


Figure 3: (a) a bi-partite token-attribute graph and (b) the resulting attribute ambiguity graph.

support of the 3-clique containing digital zoom, optical zoom, and model is two. While clique enumeration is in general NP-hard, the attribute ambiguity graph for a given table is relatively small and approximation algorithms can be used to improve performance if necessary. Note however that this is an offline process and ambiguous clusters need only be computed once.

After identifying the clusters of ambiguous attributes, we want to compute how likely each attribute is to model the set of confused tokens. We do this by computing the KL-divergence (Kullback-Leibler divergence) between the value distribution of the attribute and the set of confused tokens. Let $C = (\mathcal{A}_C, L_C)$ denote a attribute-token cluster discovered by our algorithm. We compute the disambiguation probability of an attribute $A \in \mathcal{A}_C$ given the cluster C as follows.

$$P(A|C) = \frac{KL(A||C)}{\sum_{A_i \in C} KL(A_i||C)}$$

where the KL-divergence between two distributions P and Q is defined as follows.

$$KL(P||Q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

For every ambiguous attribute-token pair $(A, \tilde{t}) \in L_C$, we set the disambiguation probability $P(A|\tilde{t}) = P(A|C)$.

The intuition in this kind of disambiguation is that the attribute which forms a better model for the query token distribution is more likely to be the correct attribute users had in mind when formulating these queries. For example, if a collection of tokens around measurements are confused among television height and diagonal, we would expect to find the distribution of values in the diagonal attribute to be a better model for the set of confused query tokens.

As an example, consider Figure 3. There are eight tokens with the same connectivity of the token labeled “8x”, two tokens with the same connectivity as “2x”, and five tokens with the same connectivity as “Pix.” The ambiguous attribute clusters correspond to the 2-clique {“Digital Zoom”, “Optical Zoom”} with support 10, the 2-clique {“Product Name”, “Model”} with support five, and lastly the 3-clique {“Digital Zoom”, “Optical Zoom”, “Model”} with a support of two. The 3-clique will be used to disambiguate among tokens confused with all three attributes, while the sub-clique will be used to disambiguate between tokens confused among only those two attributes. We may also find the remaining two 2-cliques {“Digital Zoom”, “Model”} and {“Optical Zoom”, “Model”} however in our implementation we are not concerned with these as there are no token that confuse “Model” with one of the other attributes and not the third (i.e., all tokens confused with model and

optical zoom are also confused with digital zoom, and all tokens confused with model and digital zoom are also confused with optical zoom). This means that there are no tokens in our query-log confused among only two of the three attributes.

3.3 Query-log-level Disambiguation

We now present an approach that considers the tokens in the full query log, and estimates the probability $P(A|\tilde{t})$ for ambiguous tokens \tilde{t} , such that the likelihood of the full observed query log is maximized.

Let \mathcal{Q} be a query log over table T . Let $L_{\mathcal{Q}}$ denote the set of all tokens generated by the TAGGER that appear in the attribute-token pairs $\mathcal{AT}_{\mathcal{Q}}$. Each token t is assigned a weight

$$w_t = \sum_{q \in \mathcal{Q}} \sum_{(A_i, t) \in \mathcal{AT}_q} w_q,$$

the total weight of queries in which token t appears. For simplicity we can consider w_t to be the frequency of the token t in $\mathcal{AT}_{\mathcal{Q}}$.

We assume that the tokens in $L_{\mathcal{Q}}$ were generated according to a generative model, which generates tokens as follows: an attribute A_i is selected with probability $P(A_i)$, and then a token t is drawn from the data distribution $P_T(t|A_i)$. Thus token t is generated with probability

$$P(t) = \sum_{A_i \in T.A} P_T(t|A_i)P(A_i)$$

The probability of observing the full set of tokens $L_{\mathcal{Q}}$ is

$$P(L_{\mathcal{Q}}) = \prod_{t \in \mathcal{Q}} P(t)^{w_t}$$

The generative model we described has α parameters $\pi_i = P(A_i)$, one for each attribute A_i in table T . Parameter π_i is the probability that attribute A_i is activated for the generation of a token, and we have that $\sum_{i=1}^{\alpha} \pi_i = 1$. Let π denote the vector of parameters $\pi_1, \dots, \pi_{\alpha}$. We perform Maximum Likelihood Estimation of π , and we look for the value of π that maximizes the probability that we observe the tokens in $L_{\mathcal{Q}}$. Using log-likelihood notation we want to minimize

$$\mathcal{L}(L_{\mathcal{Q}}) = -\log P(L_{\mathcal{Q}}) = -\sum_{t \in \mathcal{Q}} w_t \log P(t)$$

We use an iterative EM algorithm for finding the Maximum Likelihood Estimate for π . The algorithm initializes the values of the parameters π_i to an arbitrary distribution, in our implementation uniform. Then it alternates between the following two steps.

In the E-step given an estimation for the parameters π , for each attribute-token pair (A_i, t_j) in $L_{\mathcal{Q}}$ we compute the posterior probability of the attribute A_i given token t_j

$$P(A_i|t_j) = \frac{P(t_j|A_i)\pi_i}{\sum_{A_{\ell} \in \mathcal{AT}.A} P(t_j|A_{\ell})\pi_{\ell}} \quad (5)$$

In the M-step, we can now estimate new values for the parameters in π as follows

$$\pi_i = \sum_{(A_i, t) \in \mathcal{AT}_{\mathcal{Q}}} P(A_i|t)P(t) = \frac{w_t}{W} \sum_{(A_i, t) \in \mathcal{AT}_{\mathcal{Q}}} P(A_i|t) \quad (6)$$

where $W = \sum_{t \in L_{\mathcal{Q}}} w_t$ is the total weight of all tokens.

We repeat these two steps until the algorithm converges. Given that the optimization function is convex we know that it will reach a global maximum. After the algorithm has converged, for an ambiguous token \tilde{t} we can now compute $P(A|\tilde{t})$ using Equation 5.

The Maximum Likelihood Estimation of the probabilities $P(A_i)$ finds a set of values that best explains the full query log we observe. Intuitively what this means is that when observing a token t , even though attribute A_1 may have the highest probability $P(t|A_1)$, the token t may be attributed to another attribute A_2 since this is overall more likely to occur in the log. For example, consider a query log for queries the television table, and let L_{inch} denote the set of all tokens of the form number followed by the unit ‘‘inch’’. If the overall distribution of these tokens agrees better with the attribute ‘‘diagonal’’ than ‘‘height’’ then for a token like ‘‘32 inch’’ in L_{inch} the attribute ‘‘diagonal’’ will have higher probability $P(A|t)$ even though the value is more likely in the ‘‘height’’ attribute. This agrees well with our intuition that ‘‘32 inch’’ is more likely to refer to the diagonal rather than the height, since it is typical for users to query for diagonal when querying for televisions. The iterative algorithm captures nicely this intuition of a typical query.

3.4 Clicks-based Disambiguation

The ambiguity problem is essentially an intent issue. Given an ambiguous token, and multiple interpretations of the token in the data table we have no indication which interpretation was intended by the user. A strong indication of intent in an online environment is clicks. A query-click log $\mathcal{QC} = \{(q_1, p_1), \dots, (q_n, p_n)\}$ over a table T is a set of query-entity pairs, where q is a query posed in a commercial search engine, and p is an entity in table T that was shown to the user as a result of the query and the user chose to click on it to obtain further information about it. Such query-click logs are judiciously collected by commerce search engines and reveal a great amount of information for the behavior of the users.

An entity p is an entry in table T . We represent it as a set of attribute-value pairs $p = \{(A_i, A_i.v_j)\}$ for every attribute A_i in table T . Given a query-click pair (q, p) , we make the assumption that the user who posed the query intended to access an entity with the attributes of clicked p . Since we now have an indication for the intent of the user, the disambiguation problem becomes significantly easier to tackle. Given an ambiguous token \tilde{t} that maps to attributes $\mathcal{A}_{\tilde{t}} = \{A_{i_1}, \dots, A_{i_k}\}$ we disambiguate by selecting from entity p the attribute $A^* \in \mathcal{A}_{\tilde{t}}$ that has value $A^*.v$ that best *matches* token \tilde{t} .

The notion of a match depends on the type of token that we consider. In the case of categorical tokens we require that $A^*.v = \tilde{t}$. In such cases we have that $P(A^*|\tilde{t}) = 1$, since we assume complete certainty in our mapping. If there is more than one attribute that match token \tilde{t} then we assume that \tilde{t} maps to all of them, and we assign equal probability to all. It is worth noting that, although possible, this event is not very probable. In our experiments, we did not observe any case where multiple categorical attributes match on the same token for the same clicked entity.

In the case of numerical tokens, we map token \tilde{t} to the attribute A^* that has the closest value, i.e., it minimizes the difference $|\tilde{t} - A^*.v|$. we compute the disambiguation probability as

$$P(A^*|\tilde{t}) = \frac{\exp(-|\tilde{t} - A^*.v|)}{\sum_{A \in \mathcal{A}_{\tilde{t}}} \exp(-|\tilde{t} - A^*.v|)}$$

favoring the attributes with smallest distance.

Clicks has been used extensively in web search for eliciting the preferences of users [2, 6, 16, 21, 30]. It is well known that they contain a very strong signal for the intent of the user, but they also come with different caveats, such as presentation bias (users click only on what they are shown), ranking bias (users tend to click on the top positions), and noise. We discuss some of these issues in the experimental section.

4. FACET VALUE SELECTION

Equally important to finding important facets, is the problem of finding important values to populate the facets. We apply the same philosophy for values as we do for facets; popular values in a query log will be strong facet values that are important to users. In the following section, we discuss two types of facet value selection: category-dependent facet value selection, and dynamic query-dependent facet value selection.

4.1 Category-dependent Value Selection

The core problem of identifying to which facet a value belongs deals with the same ambiguity challenge discussed in Section 3. We apply the token-disambiguation approach to compute probabilistic frequency counts of values per facet attribute. Thus, for each table, we have a category-dependent set of facet values which can be displayed without a query or other context beyond the table category. Figure 4 shows an example for the golf clubs table. These values are popular golf club brands.

More formally, we score a value for a given attribute by the probability that the value actually represents the attribute. Let A denote an attribute, v denote a value and \mathcal{A}_v the set of mapped attributes for v . We want to estimate $P(A|v)$ for all values v found in the query-log. Using Bayes rule we have

$$P(A|v) = \frac{P(v|A)P(A)}{\sum_{A_i \in \mathcal{A}_v} P(v|A_i)P(A_i)}$$

Following the token-disambiguation scheme, we estimate the right-hand side using the data in the table T . Since we assume that all attributes are equally likely in the table, $P(A|v)$ is proportional to $P_T(v|A)$, the probability that the value v is generated from the distribution of A in the table. It follows that

$$P_T(v|A_i) = \frac{|T(A, v)|}{|A|}$$

where $T(A, v)$ denotes the set of entries in the table where the attribute A takes the value v , and $|A|$ is the number of table entries with some value for attribute A . The popularity of the value for a query log \mathcal{Q} is then estimated by the following.

$$F_{\mathcal{Q}}(A, v) = \sum_{q \in \mathcal{Q}} \sum_{(A, v) \in \mathcal{AV}_q} w_q P(A|v)$$

4.2 Query-dependent Value Selection

Given a query, in the form of an annotated keyword query or a facet selection, it is desirable to update the values of facets to reflect the context given by the query. For example, consider a faceted search engine which includes “brand”, “model,” and “golf club type” among its facets for the golf club table. The category-dependent values would include popular brands, popular models (from different brands), and the most popular club types. If a user has submitted the query “odyssey golf clubs”, then any model attribute values occurring in the model facet that are not from the brand “Odyssey” are irrelevant. This could be particularly confusing to a user if they had chosen “Odyssey” as a facet selection from the “brand” facet. A subsequent selection of the model facet for a model corresponding to a different brand would produce an empty result. This is a clear motivation for conditional, query-dependent facet values.

Pre-computing conditional values however, is an expensive task for large data sets, since we need to consider all pairwise values in the table. Furthermore, even after this expensive computation, supporting the retrieval of such information for real-time search is a non-negligible investment. This is a significant effort to support

updating facet values in response to searches. The alternative is to support the expensive run-time computation of conditional facet values, which can have an overall impact on response time.

For our approach to query-dependent facets, we again turn to the query log. We compute co-occurrence frequencies for all pairs of values that co-occur in a query, over the entire log. We then validate the co-occurring values against the database to ensure a non-zero entity count in the result. This approach has a number of benefits. First, it allows us to find the most popular values co-occurring with a given selection condition. This benefits users by putting the most relevant facet values at the top of the list for a given selection.

We define the conditional score of a value for a given attribute by the probability that the value represents the attribute in a query that co-occurs with a given attribute value pair. Let $(A_c, A_c.v_c)$ denote a conditional attribute-value pair, A denote an attribute, v denote a value and \mathcal{A}_v the set of mapped attributes for v . We want to estimate $P(A|v, A_c = A_c.v_c)$ for all values v found in the query-log. Since tokens are tagged independent of one another, this simplifies to $P(A|v)$ which we can compute as shown Section 4.1. The difference is in which queries are used from the log. Given a query log \mathcal{Q} and attribute-value pair $AV = (A_c, A_c.v_c)$, let \mathcal{Q}_{AV} denote the subset of queries $\{q \mid q \in \mathcal{Q} \wedge (A_c, A_c.v_c) \in \mathcal{AV}_q\}$. We can then estimate the popularity of a value over all queries satisfying the given condition.

$$F_{\mathcal{Q}}(A, v) = \sum_{q \in \mathcal{Q}_{(A_c, A_c.v_c)}} \sum_{(A, v) \in \mathcal{AV}_q} w_q P(A|v)$$

Returning to our golf example, our category-dependent list has the values “drivers” as the most popular value for “golf club type”. But given the example query “odyssey golf clubs”, the most popular “golf club type” is “putter”.² The other main advantage to query log based discovery of conditional values is that it restricts the space of pairwise values (to count and store) to only those values appearing in the log. This can significantly prune the space of conditional counts, making pairwise co-occurrence statistics feasible with limited resources. An example of our conditional values can be seen in Table 4 for the query “golf putters”. The values in the query-dependant list are brands that are popular for their putters.

The last consideration for query-dependent facet values is the case where there are multiple selection conditions specified in the query. This may occur either as multiple recognized attribute values in a keyword query, or as multiple facet selections. Supporting n -way conditionals with exact information is impractical. The data-driven approach suffers the combinatorial explosion of an already large data set, and the query log driven approach suffers from the sparseness of queries. After filtering by multiple selection conditions, the number of queries from which to compute popularity frequencies becomes too small. For this scenario, we adopt a simple approach of intersecting value lists and aggregating popularity counts by addition. This will find values relevant to multiple selections, and score them proportionally to how relevant they are to their respective conditional selections. Handling the case of an empty intersection is more of a user interface design decision. One could eliminate the facet entirely, or default to the category-dependent values.

5. FACETS AND DATA STATISTICS

An additional challenge to the problem of attribute selection comes from the web nature of our structured data. Real-world data is often noisy and incomplete, and we need to account for that when

²This is consistent with the Odyssey golf brand, which is well known for putters.

Category: Golf Clubs	Query: "golf putters"
Nike	Odyssey
Ping	Ping
TaylorMade	Scotty Cameron

Figure 4: Example top-3 facet values for golf club brand.

selecting attributes to be exposed as facets. Furthermore, even in the presence of clean data, some attributes are not appropriate for being selected as facets since they are not helpful in the exploration process of the user, because they have very little information content. For example, if almost all cameras have optical image stabilization, exposing this feature gives little information to the user in narrowing down their search. We thus need to eliminate noisy or uninformative attributes from being considered as candidate facets.

The query logs can reveal which attributes in the table are popular among the users. However, just because an attribute is popular does not necessarily mean that it is good to be used as a facet. In this section we discuss how to use some signals from the data to discover which attributes should be considered as candidates for the facets.

Information Content: Recall that the goal of facets is to enable the users to quickly zoom into a smaller subset of products that are of interest to them. Therefore, in order for an attribute A to be a good candidate for being a facet, it should have high information content, that is, the knowledge of the value selection v for the attribute A should give information about the tuples of interest to the user. This property is naturally captured in the *entropy* of the attribute A . For an attribute A defined over the domain $A.V$, the entropy of A is defined as

$$H(A) = - \sum_{v \in V} P_A(v) \log P_A(v)$$

where $P_A(v)$ is the fraction of entries in the table where attribute A takes value v , over the number of entries in the table where the attribute A takes any value (i.e., is non-null).

Attributes with low entropy do not make good facets. For such attributes the distribution of values is highly skewed, and the knowledge of the value gives very little information for the subset of entities of interest to the user. This is clear in the extreme example when the entropy is zero, meaning that all entities take the same value. In this case, the attribute is useless for navigation, since it gives no information about which entities the user is interested in. For example, the attribute "color" in the televisions table conveys no useful information as a facet if all televisions are black.

An alternative interpretation of $H(A)$ is that it captures the decrease in the uncertainty for the preference of the user for the entities (rows) in table T . When no facet attribute has been selected, any entity $e \in E$ is equally likely to be of interest to the user, therefore, if N is the number of entities in the table, each entity has probability $P(e) = 1/N$. The uncertainty is captured in the entropy of the random variable E , which in this case is maximum, $H(E) = \log N$. Now assume that the user is allowed to use attribute A to zoom in on a smaller subset of entries. The uncertainty for E decreased given the knowledge of the value in A is measured as $H(E) - H(E|A)$, where $H(E|A)$ is the conditional entropy of E given A . It is well known that $H(E|A) = H(E, A) - H(A)$. Since each entity is distinct from the rest, and assuming that A takes a value for all entities, we have that $H(E, A) = H(E) = \log N$. Therefore, $H(E|A) = H(E) - H(A)$. Thus the entropy $H(A)$ captures how much our uncertainty for the entities that the user is interested will decrease when we have the knowledge for attribute

A . If $H(A) = 0$, then the knowledge of A gives no extra information for E and thus it is redundant.

It follows naturally, that attributes with low entropy should not be used as facets. We thus compute the entropy for all attributes and we apply a threshold for eliminating attributes with low information content, that is, $H(A) < \theta_H$.

Sparsity: Real-world data is often noisy with missing or incomplete data. As a result there are often cases where some attributes are only sparsely populated. If such an attribute is selected for faceted navigation, then selecting any value will immediately discard most of the entities since they do not have a values for the sparse attribute. This could be ok if this corresponded to a rare feature, in which case missing information corresponds to negative information for the existence of the feature. However, it is often the case that sparsity is due to noise in the data collection, in which case, the entities that are eliminated are valid for the selected facet, yet will never be visible to users. Furthermore, noisy attributes are often likely to contain incorrect values, confusing the tagging process of the queries and the corresponding probabilities.

Therefore, we exclude from the candidate attributes the ones that are very sparse. The sparsity of attribute A is defined as

$$R(A) = \frac{|T(A)|}{|T|}$$

where $T(A)$ is the set of entries in which the attribute has a non-null value, and $|T|$ is the total number of entries in the table. We use as candidate facets only the attributes that have $R(A) \geq \theta_R$.

6. EXPERIMENTAL EVALUATION

We abstract the structured data as sparse tables, each containing entities of the same type. For our experiments we have 1164 such tables that we crawled using the MSN shopping XML API³. The tables correspond to products used to answer shopping queries and are similar to the data used by sites like Amazon or Google Product Search. We consider each category of products to be a table of entities of the same type. The available range covers entities from electronics like digital cameras or televisions to golf clubs and soft goods like shoes and shirts. In total, there were around 10 million structured distinct product entities occupying approximately 300GB on disk when stored in our database.

Besides the structured data collection, we also use a web query log and a vertical click log. The *web query log* contains queries posed on a major web search engine on a period of five months from July 2009 to November 2009. The web queries are approximately 31 million distinct queries all in the English language. As query weight we use the aggregated impressions count – each time a query is asked it increments its impression count. We limit ourselves to queries with at least 10 aggregated impressions to reduce noise in the queries. The total aggregate impression weight of the log is approximately 4.2 billion impressions. The average query length is 3.42 words and 22.04 characters. The *click log* is available via toolbar usage logs collected for users that have installed the browser add-on Bing toolbar. It contains queries and clicks on Amazon products over a period of one year from March 2009 to February 2010. Since our structured data set is in the shopping domain, we were able to map the Amazon click log to our product entities using unique identifier references like UPC and EIN codes. The format of the log is query, entity-id and number of clicks as the query weight. The total number of distinct queries is in the order of 2.2 million

³See <http://shopping.msn.com/xml/v1/getresults.aspx?text=digital+cameras> for a table of digital cameras and <http://shopping.msn.com/xml/v1/getspecs.aspx?itemid=25236859> for an example of a camera entity with its attributes.

with average length of 3.67 and 24.1 characters. It is worth noting that the click log is smaller although the period used is longer. This is attributed to two reasons: there are more queries on a search engine than on Amazon, and the toolbar, as an add-on, means that it can only capture a fraction of the click activity since it is opt-in only and many users do not install such add-ons. However both logs provide very valuable information.

For our experiments, we use the work described in [27] to map queries to tables in our collection. We run the tool on our web query log and kept only the interpretation to tables that are above the threshold of $\theta = 1$. Since a query can map to more than one table, we took all the plausible tables and normalized their probabilities to fractions summing to 1. Then we used each fraction to map the query to the table while appropriately adjusting the query weight. Thus we produced a subset of the query log for each table, with each query having a weight that is the corresponding fraction of its impressions multiplied by the normalized table probability.

We tested the various techniques we proposed in the paper and we present the results below. The naming scheme is as follows: *DataProb* is the probabilistic disambiguation from Section 3.1, *ClusterProb* the cluster-based approach from Section 3.2, *Iterative* is the approach described in Section 3.3, and *Clicks* utilizes the click log as described in Section 3.4.

In addition to these techniques, we implemented a simpler one, shown as *NoProb*, that does not use any token based probabilities for disambiguation. It computes the score using, for each query mapped to a table T, all the plausible token-attribute pairs from the table T while assigning to all of them equally the query weight multiplied by the normalized table weight. *NoProb* is not blind counting because it does consider the query weight and probability from the interpretation algorithm [27] that we used in assigning the queries to tables, thus we consider it a fair baseline. The main difference with our proposed disambiguation techniques is that it does not use an explicit probabilistic disambiguation on the possible token-attribute mappings within each table.

6.1 Relevance Judgments

To measure the effectiveness of our results we conducted an extensive study using Amazon Mechanical Turk⁴. Assessing the entire dataset and all possible attributes would have been prohibitively expensive. Instead we chose a diverse subset of our tables covering a variety of different areas representing as much as possible the entire data spectrum. The tables used in our evaluation were televisions, microwave ovens, cell phones, golf clubs, mp3 players, shoes, laptop computers, printers, watches, video games, engagement rings, digital cameras, skin cleansers, and camera bags and cases.

We ran all our approaches on the full set of tables and produced a set of facet attributes for each table and each algorithm. Then we took the results for the evaluation tables and created a single pool of attributes for judging by merging the results of all algorithms. We posted Mechanical Turk HITs (Human Intelligence Tasks) to obtain user judged relevance for the quality of attributes for search. Users were asked to judge how important a given attribute was for searching within a given category (with the example of faceted search explained) on a three-point (one to three) scale. The scale items were labeled “highly important”, “somewhat important”, and “not important.”

Producing high quality judgments in a crowd-sourced environment like Mechanical Turk is a challenging problem on its own. We employed a number of methods to ensure quality judgments.

First, we created HITs of ten judgments, plus an additional judgment used for spam detection. The additional judgment was drawn from a pool of judgments we performed manually, and for which we felt there was a clear “highly important” or “not important” answer. The deviation of workers from our gold standard tasks was often a clear indication of spam. However, to ensure fair judgment, we manually inspected the results of any candidate spammer for consistency. The average completion time of a worker was also used as an aid in detecting unreliable workers. Lastly, each task was repeated by nine unique workers. This allowed us to find majority decisions for roughly two-thirds of all attribute judgments. For the remaining third without a majority, we took the average score. The result is a single valuation of the importance of each attribute.

Using these graded judgments we computed the normalized discounted cumulative gain (NDCG) for the ranked output of each approach. Given a list of results L , DCG was computed using linear gain and a logarithmic rank discount as follows.

$$DCG(L) = \sum_{i=1}^k \frac{rel(L_i)}{\log_2(i+1)}$$

Where $rel(L_i)$ denotes the judged relevance of result L_i . Let I denote the ideal result ordering (the list of all possible results ordered by judged relevance), then NDCG is defined as follows.

$$NDCG(L) = \frac{DCG(L)}{DCG(I)}$$

We report on various values of k , the number of attributes returned. If a system returned less than k attributes it was penalized with a relevance score of 0 for each missing value.

Validating Mturk Judgements: While we were confident in our quality assurance techniques, we further validated the collection for significant outliers, by manually judging five categories for which we have personal knowledge of the domain through shopping experience. The table in Figure 5 shows the confusion matrix for the judgments. Across the top are the Mechanical Turk judgments, and

	judgements			mturk		
	1	2	3	1	2	3
manual judgments	1	6	15	3		
	2	6	34	4		
	3	2	17	50		

Figure 5: Confusion matrix between mturk and manual judgments.

on the left are the manual judgments. Entry (i, j) is the number of times the manual judgment was i , and the Mechanical Turk was j . We round averaged Mechanical Turk valuations to the nearest integer for this computation. We can see from this that the Mechanical Turk workers often favor the “safe” middle valuation. The effect of averaging judgements that do not reach a majority may also push valuations to the middle. The manual judgments tend to distribute more judgments to the “highly important” and “not important” valuations. Overall, we see exact agreement on 66% of judgments, with opposing 1 vs 3 valuations on only 4%. The Mechanical Turk valuations are therefore somewhat more coarse grained than our careful manual evaluations, but are still similar in overall trend. We have also run all of our experiments over the manual judgments and found the results to be equivalent in terms of trend and rank order of the systems. We omit these graphs for brevity.

6.2 Attribute Pre-selection with Data Filters

We start with an experiment that measures the effect of data

⁴<http://www.mturk.com>

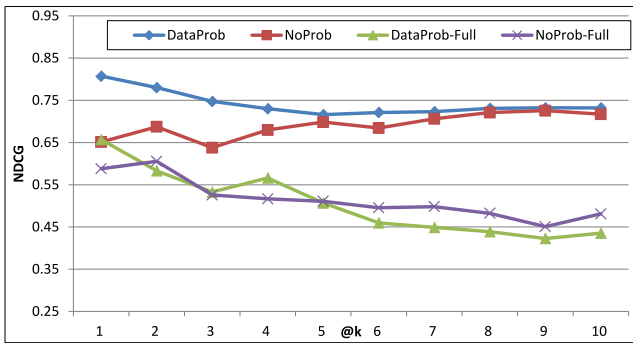


Figure 6: Attribute pre-selection using data filters.

statistics on the produced facets such as selectivity and entropy as described in Section 5. The results are summarized in Figure 6.

The first step was to take all attributes for each table⁵ and perform a run on our two simpler algorithms, shown as *NoProb-Full* and *DataProb-Full* in Figure 6. Surprisingly, *DataProb-Full* does not improve much over *NoProb-Full*. In fact it actually performs worse at higher values of k . Upon further investigation, this is understandable due to how certain attributes affect the computed probabilities. For example, there is an attribute called ‘video input type’ on digital cameras where the value is almost always ‘digital camera’. The probability $P(t|A)$ for the token ‘digital camera’ was very high, affecting the disambiguation process whenever it was recognized in a query. As a result the *DataProb-Full* got skewed in counting the proper attribute importance incorrectly. There were other such attributes with similar characteristics that produced a skewed behavior. We observed that such attributes had in common certain data statistics such as low value entropy and low sparsity across entities in a table.

We did a second run where we did not load attributes that appear in less than 10% ($\theta_R = 0.1$) of the entities and have entropy less than 0.1 ($\theta_H = 0.1$). Using this data set we run again the same two algorithms, shown as *NoProb* and *DataProb* in Figure 6. The pre-selection step reduced noise in the data significantly and also removed attributes that have little information content (low entropy) and are not appropriate for facets. As a result the same algorithm with the pre-selected data set perform much better. Furthermore, now the actual data probabilities are more meaningful and *DataProb* performs better than *NoProb*.

Entropy and sparsity can be seen as continuous discount values. We saw in practice that using them as parameterized filters is easier and produces good results. The specific values we used capture the problems with our particular data set and might not be optimal for all data sets. But the point we wish to make is that one has to consider such data filters to pre-select good attribute candidates in the mining algorithms. The full attribute set triggers bad results even on the more complex algorithms. It is not shown here for presentation simplicity. The remainder of the experimental section uses the same pre-selected set of attributes for all our techniques.

6.3 Explicit Disambiguation

Using the above mentioned preselected attribute set we run all our disambiguation algorithms. The results are shown in Figure 7.

The first observation is that *ClusterProb* and *DataProb* have similar curves. Although *ClusterProb* seemed like a very good idea for disambiguation by looking at token-attribute clusters across many queries, we see that it does not improve things in practice. There are a number of reasons why this happens. First, attributes encod-

⁵We exclude metadata and boolean attributes with values like ‘yes/no’, ‘true/false’

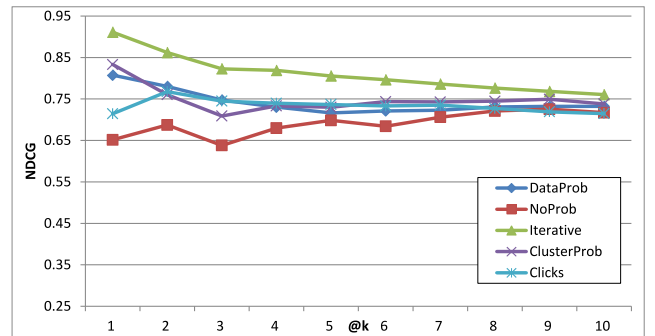


Figure 7: Token-attribute disambiguation comparison

ing exact numeric measurements do not always reflect the product class to which they belong. Going back to our television example, a ‘32 inch’ television may actually have a diagonal closer to 31 or 33 inches. In this situation, our *ClusterProb* approach may incorrectly reinforce the height attribute rather than the diagonal given that many modern 50 inch televisions in the data set have a height closer to this value. Also, this problem is exacerbated by the discretization of continuous domains for the KL-divergence computation. Lastly, KL-divergence is only well defined for distributions that have non-zero probability over the entire domain of values, meaning even attributes like depth are assigned a small probability for values like 50 inches.

Iterative tries to detect correlations between data and queries and it exploits the overall bias of users towards the correct attributes. Intuitively, the reason it disambiguates the best is due to how users enter queries. Although the system can be confused with multiple interpretations, users actually know which attributes they are looking for when they enter their values. When examined across the full query log, the user behavior tends to match the intended attribute. For example, although token ‘32 inch’ can be considered closer to a TV height than diagonal screen size for many modern televisions, users query much more frequently for diagonal size including other values like ‘50inch’ or ‘55inch’. This creates a bias for *inch* tokens towards the attribute diagonal, which *Iterative* applies to the ‘32 inch’ token. So diagonal size gets preferred over height even for that particular value.

One surprise in the disambiguation methods is that *Clicks* do more or less similarly to *DataProb*. At first, one would think that clicks offer the perfect disambiguation method. A user click on a particular entity id allows us to select only the attributes that closely match that entity id disambiguating attributes almost perfectly. However, clicks have problems in the way they were produced that make them less than optimal for learning good facets. First, a user can only click on results shown to them, so clicks incorporate the engine bias and do not represent the true popularity of how many times all possible queries are asked. Hence some attributes are under-represented in a click log due to the engine bias of what results were shown. Second, there are far fewer queries with clicked results than generally asked, meaning many categories have very few queries. This means the mined facets do not have much support and are somewhat erratic. Finally, the queries that tend to trigger clicks are very specific queries looking for one or few products, e.g. ‘canon eos 50d digital camera’. More general category based queries are more likely to be followed by a refinement since they return a large result set that users cannot easily consider, e.g. ‘12 megapixel digital camera’. This creates a skew of the importance learned towards attributes that act as unique keys. However these attributes are not well suited for facets. If a user

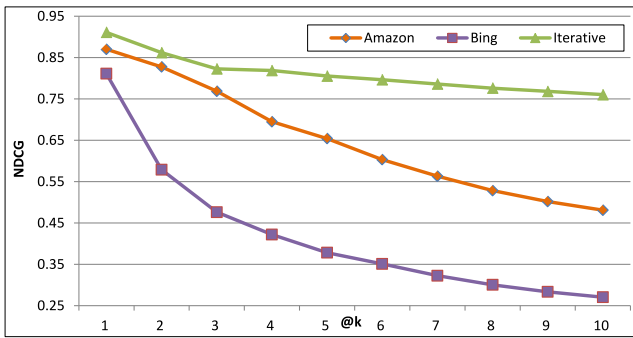


Figure 8: Comparing with Amazon and Bing.

already knows the unique key they are looking for (like camera model), they can simply select it from the results and do not need a facet for it. Instead such attributes are better suited for entity snippets, that is summarized views that help users differentiate amongst entities.

6.4 Commercial Faceted Search

Commercial web engines are already supporting faceted search. Although we do not know the technical details of industrial implementations, we felt the best way to test the real world effectiveness of our solution is by comparing against state of the industry engines. Since our data is shopping based, we considered Amazon as the world’s most popular shopping engine. Furthermore, since our data comes from the public MSN Shopping API, we considered Bing Shopping (redirects from MSN Shopping) as a shopping engine that is using the same data set as our experiments. Both Amazon and Bing Shopping show facets. To test them, we crawled the attributes shown on Amazon and Bing for our test categories and submitted them in the same pool of attributes that we labeled with mturk judges. To remove any bias, we dropped any affiliation information on the web site. Furthermore, we did not consider the handful of generic shopping attributes shown on both Amazon and Bing, like price or shipping options. Although very valuable facets there is no need for an automated technique to discover them, as such attributes can easily be added to all categories. Instead we limited our comparison to only *category specific* attributes from each solution. The results are summarized in Figure 8.

Our best solution performs better than Amazon and significantly better than Bing Shopping. For many categories Amazon shows very few attributes and Bing shows even less. As a result, both Amazon and Bing drop dramatically as the size k increases in our experiments. It is important to note that Amazon does show up to 12 category-dependent facets in some cases (e.g., watches) so screen real-estate is not a limiting factor. For small values of k ($k \leq 3$) we still do better but the difference, particularly with Amazon, is smaller. We are not familiar with the details of their facet selection approach and so their techniques cannot be contrasted to ours. However, since they are a popular site with lots of query traffic and domain knowledge, we expect their facet selection to be of the highest caliber. Having our best solution perform even better than Amazon, demonstrates the true value of a pure query-based automated solution and how well it approximates the desired facet utility. Furthermore, it shows that our best method deals more than adequately with ambiguity amongst attributes within each table.

6.5 Values

Our approach populates facets with the most popular values found in the query-log. This is done in a query-independent way, by look-

Judgment	Query-depdnt	Category-depdnt	Equivalent
Preferred	46.3%	23.96%	29.74%

Figure 9: Comparison of facet values for query-dependent vs. category-dependent approach.

ing at all values for all attributes as described in Section 4. So for every category, for the top facets we have computed the top values in a category-dependent list. It would be very expensive to perform an NDCG experiment on values, as the potential judgements needed to capture the value relevancy, in a statistically significant way, is very large. Furthermore, given how we find the popular values, we believe that the quality of those values will correlate with the results of our disambiguation techniques and such NDCG experiments would be moot. However, we wanted to validate the benefit of showing the query-dependent conditional values, which were computed off-line as described in Section 4.

We created a synthetic online experiment by generating queries using the top values for the top facets for our test categories. For example, using the brand for golf clubs, we generated ‘odyssey golf clubs’ or ‘callaway golf clubs’ as test queries. For each query, we judged the quality of the query-dependent conditional value list vs. the category-dependent value list, for the top facets. The judgement favored one or the other list for each facet, based on the relevance of the values shown given the query. If the values seemed equally relevant, they were judged to be equal. In many cases, the query itself may dictate whether the values are correct or not based on what is present in the data. For example, if the condition specifies the brand “taylormade”, then a value specific to a “nike” product is simply incorrect. Other times, results require specialized knowledge of the domain, for example, that the brand “odyssey” is well known for putters or that “callaway” is famous for drivers. The results are summarized in Figure 9.

As expected, query-dependent values are preferred more than category-dependent. However, at first look, it is surprising to see the category-dependent gaining so much preference. This happens because category-dependent attribute are quite often a super set of the query-dependent ones. The nature of conditional computation has smaller support thus eliminating some values. For example, the conditional list may show values {mens, womens} while the category-dependent list shows {mens, womens, juniors} which is perceived as better or more complete. A solution to this to take the cases where the query-dependent is a subset of the category-dependent list, show the superset but force the conditional ones to the top. We performed the same experiment after we applied this simple heuristic and found out that only 1.98% of the queries went to category-dependent whereas the query-dependent jumped to 65.58%. There was no significant change in the equivalent cases, as they are generally similar value sets with marginal order changes.

7. RELATED WORK

Our problem of discovering important facets in a query log is related to a number of active research areas, including faceted search, attribute ordering, query understanding, and search over structured data. In the following section we survey a sample of the relevant work in these areas.

Faceted Search The problem of automatically discovering facets (or facet hierarchies) from text has been studied in the past in [7, 8, 20, 28]. The principle challenges in this domain are in finding the terms which can be used to describe (sub)categories, and in building the hierarchy of categories. This is in contrast to our setting, where the actual facets have a direct correspondence to the struc-

tured data attributes. The process of extracting facets from text often makes use of structured resources like WordNet⁶ [7, 28] or semi-structured sources like Wikipedia⁷ [7], but only for categorization and hierarchy extraction, rather than the extraction of attributes. In these works, the facet selection is achieved by ranking based on a navigation cost [8, 20], or by frequency-based metrics [7].

Recently, there is also work on facet selection for structured data [4, 18, 19, 24, 26], which looks for structured data attributes as potential facets. Facet ranking is generally done by estimating the user cost to identify a result when searching with the facets [4], or by ordering the facets according to some property in the data, such as the partitioning balance, cardinality, or attribute frequency [24]. This does not have the benefit of approximating the true user intent as in our query-log based approach. The use of data properties however, is similar to our data filtering. In [19] facets are ranked based on a history of user preference for documents. However the facet-value pairs are given explicitly for each document so there are no challenges in resolving ambiguity. In [18], a cost model driven by user actions is approximated using user relevance judgements over attributes. This is in contrast to our extraction of attribute relevance from query-logs.

Attribute Ordering The problem of ordering or ranking attributes of structured data repositories is a related problem to facet discovery. The goals of attribute ordering problems however, often differ to that of ranking facets. In [9, 17], the aim was to choose a subset attributes that are most “useful” to display to a user when rendering structured search results with limited screen space. Usefulness was formalized as the attributes most influential in the search systems ranking function. A related problem is product visibility, in which attributes are chosen that will distinguish one’s product from others in the database. In [23] the product visibility problem was approached by finding the k attributes that cause a product to appear in the maximum number of results over a given query log.

The idea of using a query-log as a source of important attributes has also been studied in previous work. Miah et al. proposed to compute the frequency of the attribute names in a query log as a signal of importance for the product visibility problem [23]. However, they assume binary attributes, such that each attribute is essentially an attribute/value pair in our data model (e.g., the attribute “four door” describes the boolean property of a car having four doors or not, as opposed to an attribute like “doors” having values two, three, four, or five). This is different from our approach in that we mine attribute instances from the query-log, and thus have to deal with the ambiguity in mapping instances to their correct attributes. Chaudhuri et al. have also used query-log mining to rank attributes based on their probability of being correlated with a query, with the intention of ranking structured results [5]. However, their work considered structured queries, and thus a structured query-log which does not suffer from ambiguity in query tokens.

Query Understanding A number of recent works have looked at query understanding using structured data [11, 22, 25, 27]. Because our approach rests on having query interpretations, this field is particularly relevant. As part of our experiments we implemented the work in [27] and benefit from it in finding plausible query interpretations and assigning queries to tables. However, the facet discovery problem poses more strenuous requirements than previous query understanding work in disambiguating amongst multiple query to attribute mappings within the same table. In other words, it is not enough to just find the plausible interpretations of a query

on a table. For facet discovery we have to actually be able to distinguish amongst them with high confidence. An approximation of this effect is the results of *NoProb* vs. all the other disambiguation techniques we discussed – *NoProb* can be viewed as just using the interpretation score given by our implementation of [27]. As such, our work is different and can be considered an extension to previous work.

A related problem to query understanding is search over structured data, a field that has received a lot of attention in recent years. Search over structured data is also the core motivation for our work of supporting faceted search over structured collections. In this context, each unique mapping of a keyword query into a structured database represents an understanding of the query. Early work focused on how to assemble search results, with ranking based on candidate network size [3, 14], graph-based weight propagation [1], or integrating IR style ranking functions [13]. However, no particular effort has been made in this field to disambiguate among possible mappings of query tokens to attributes.

8. CONCLUSIONS

Structured web search in the form of a unified ecosystem like Google or Bing, or in the form of a specialized vertical like Amazon can take advantage of rich meta-data to produce a better user experience with the appropriate facets. An important problem in this setting is that of selecting the attributes and values for the facet system. Discovering facets must also deal with the scale of data and the restricted screen real estate that only allows for the ones best matching the anticipated user intent.

In this paper we proposed a model that exploits the intersection of web query logs and structured data to find the facet utility. We showed a variety of techniques that dealt with disambiguating amongst different overlapping attribute-value pairs per query. We extended them to capture values and conditional values, that give a dynamic query-dependent facet experience, and we also discussed how data statistics play a role in producing better facets. Finally, we showed experimentally that our approach is scalable, it works well in practice, and even outperforms popular commercial systems that offer facet functionality.

9. REFERENCES

- [1] B. Aditya, G. Bhalotia, S. Chakrabarti, A. Hulgeri, C. Nakhe, P. Parag, and S. Sudarshan. Banks: browsing and keyword searching in relational databases. In *Proc. VLDB Conf.*, pages 1083–1086. VLDB Endowment, 2002.
- [2] R. Agrawal, A. Halverson, K. Kenthapadi, N. Mishra, and P. Tsaparas. Generating labels from clicks. In *WSDM*, pages 172–181, 2009.
- [3] S. Agrawal, S. Chaudhuri, and G. Das. Dbxplorer: A system for keyword-based search over relational databases. In *Proc. ICDE Conf.*, page 5, Washington, DC, USA, 2002. IEEE Computer Society.
- [4] S. Basu Roy, H. Wang, G. Das, U. Nambiar, and M. Mohania. Minimum-effort driven dynamic faceted search in structured databases. In *Proc. CIKM*, pages 13–22, New York, NY, USA, 2008. ACM.
- [5] S. Chaudhuri, G. Das, V. Hristidis, and G. Weikum. Probabilistic information retrieval approach for ranking of database query results. *ACM Trans. Database Syst.*, 31(3):1134–1168, 2006.
- [6] N. Craswell, O. Zoeter, M. J. Taylor, and B. Ramsey. An experimental comparison of click position-bias models. In *WSDM*, pages 87–94, 2008.

⁶<http://wordnet.princeton.edu/>

⁷<http://en.wikipedia.org/>

- [7] W. Dakka and P. G. Ipeirotis. Automatic extraction of useful facet hierarchies from text databases. In *Proc. ICDE Conf.*, pages 466–475, Washington, DC, USA, 2008. IEEE Computer Society.
- [8] W. Dakka, P. G. Ipeirotis, and K. R. Wood. Automatic construction of multifaceted browsing interfaces. In *Proc. CIKM*, pages 768–775, New York, NY, USA, 2005. ACM.
- [9] G. Das, V. Hristidis, N. Kapoor, and S. Sudarshan. Ordering the attributes of query results. In *Proc. SIGMOD Conf.*, pages 395–406, New York, NY, USA, 2006. ACM.
- [10] P. Dmitriev, P. Serdyukov, and S. Chernov. Enterprise and desktop search. In *Proc. WWW Conf.*, pages 1345–1346, April 2010.
- [11] R. Fagin, B. Kimelfeld, Y. Li, S. Raghavan, and S. Vaithyanathan. Understanding queries in a search database system. In *Proc. PODS Symposium*, pages 273–284, New York, NY, USA, 2010. ACM.
- [12] M. A. Hearst. *Search User Interfaces*. Cambridge University Press, 2009.
- [13] V. Hristidis, L. Gravano, and Y. Papakonstantinou. Efficient ir-style keyword search over relational databases. In *Proc. VLDB Conf.*, pages 850–861. VLDB Endowment, 2003.
- [14] V. Hristidis and Y. Papakonstantinou. Discover: keyword search in relational databases. In *Proc. VLDB Conf.*, pages 670–681. VLDB Endowment, 2002.
- [15] P. Ipeirotis and L. Gravano. Distributed search over the hidden web: Hierarchical database sampling and selection. In *Proc. VLDB Conf.*, August 2002.
- [16] T. Joachims and F. Radlinski. Search engines that learn from implicit feedback. *IEEE Computer*, 40(8):34–40, 2007.
- [17] N. Kapoor, G. Das, V. Hristidis, S. Sudarshan, and G. Weikum. Star: A system for tuple and attribute ranking of query answers. *Proc. ICDE Conf.*, 0:1483–1484, 2007.
- [18] A. Kashyap, V. Hristidis, and M. Petropoulos. Facetor: cost-driven exploration of faceted query results. In *Proc. Information and knowledge management, CIKM '10*, pages 719–728, New York, NY, USA, 2010. ACM.
- [19] J. Koren, Y. Zhang, and X. Liu. Personalized interactive faceted search. In *Proc. WWW Conf.*, pages 477–486, New York, NY, USA, 2008. ACM.
- [20] C. Li, N. Yan, S. B. Roy, L. Lisham, and G. Das. Facetedpedia: dynamic generation of query-dependent faceted interfaces for wikipedia. In *Proc. WWW Conf.*, pages 651–660, New York, NY, USA, 2010. ACM.
- [21] X. Li, Y.-Y. Wang, and A. Acero. Learning query intent from regularized click graphs. In *Proc. SIGIR Conf.*, pages 339–346, 2008.
- [22] X. Li, Y.-Y. Wang, and A. Acero. Extracting Structured Information from User Queries with Semi-supervised Conditional Random Fields. In *Proc. SIGIR Conf.*, pages 572–579, 2009.
- [23] M. Miah, G. Das, V. Hristidis, and H. Mannila. Standing out in a crowd: Selecting attributes for maximum visibility. In *Proc. ICDE Conf.*, pages 356–365, Washington, DC, USA, 2008. IEEE Computer Society.
- [24] E. Oren, R. Delbru, and S. Decker. Extending faceted navigation for RDF data. *The Semantic Web-ISWC 2006*, pages 559–572, 2006.
- [25] S. Papanizos, A. Ntoulas, J. Shafer, and R. Agrawal. Answering web queries using structured data sources. In *Proceedings of the 35th SIGMOD international conference on Management of data, SIGMOD '09*, pages 1127–1130, New York, NY, USA, 2009. ACM.
- [26] S. Roy, H. Wang, U. Nambiar, G. Das, and M. Mohania. Dynacet: Building dynamic faceted search systems over databases. In *Proc. ICDE Conf.*, pages 1463–1466, 2009.
- [27] N. Sarkas, S. Papanizos, and P. Tsaparas. Structured annotations of web queries. In *Proc. SIGMOD Conf.*, June 2010.
- [28] E. Stoica, M. A. Hearst, and M. Richardson. Automating creation of hierarchical faceted metadata structures. In *Proc. Human Language Technology Conf. (NAACL HLT)*, 2007.
- [29] A. Sugiura and O. Etzioni. Query routing for web search engines: Architecture and experiments. In *Proc. WWW Conf.*, April 2000.
- [30] K. Wang, T. Walker, and Z. Zheng. Pskip: estimating relevance ranking quality from web search clickthrough data. In *KDD*, pages 1355–1364, 2009.