

# Generalized Syntactic and Semantic Models of Query Reformulation

Amaç Herdağdelen\*  
University of Trento  
Rovereto, Italy  
amac@herdagdelen.com

Maria Holmqvist\*  
Linköpings University  
Linköpings, Sweden  
marho@ida.liu.se

Massimiliano Ciaramita  
Google  
Zürich, Switzerland  
massi@google.com

Keith Hall  
Google  
Zürich, Switzerland  
kbhall@google.com

Daniel Mahler  
Google  
Zürich, Switzerland  
mahler@google.com

Stefan Riezler  
Google  
Zürich, Switzerland  
riezler@google.com

## ABSTRACT

We present a novel approach to query reformulation which combines syntactic and semantic information by means of generalized Levenshtein distance algorithms where the substitution operation costs are based on probabilistic term rewrite functions. We investigate unsupervised, compact and efficient models, and provide empirical evidence of their effectiveness. We further explore a generative model of query reformulation and supervised combination methods providing improved performance at variable computational costs.

Among other desirable properties, our similarity measures incorporate information-theoretic interpretations of taxonomic relations such as specification and generalization.

## Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Query formulation, Search process, Retrieval models.

## General Terms

Algorithms, Experimentation.

## Keywords

Query reformulation, query rewriting, generalized edit distance, similarity metrics.

## 1. INTRODUCTION

Query reformulation is the process of iteratively modifying a query to improve the quality of a search engine results, in order to satisfy one's information need. Search engines

\*Work carried out during internships at Google.

support users in this task explicitly; e.g., by suggesting related queries or query completions, and implicitly; e.g., by expanding the query to improve quality and recall of organic and sponsored results. The close interaction between users and algorithms makes this a central topic in search technology and research [10, Ch. 6].

Successful refinements are closely related to the original query [22]. This is not surprising as reformulations involve spelling corrections, morphological variations and tend to reuse parts of the previous query. More precisely, reformulations are close to the previous query both *syntactically*, as sequences of characters or terms,<sup>1</sup> and *semantically*, often involving transparent taxonomic relations. As an example, for the query “becoming a dentist”, the reformulation “becoming an oral surgeon” might have a higher chance of producing relevant results than “becoming a doctor”. In this paper we address the following question: how can we model query reformulation as a process involving syntactic and semantic operations within a unified and principled framework?

String distance metrics model the similarity between two queries as a function of the edit operations (insertion, deletion, substitution) that are necessary to generate one string from the other. Jones *et al.* [14] noticed that edit distance can be an accurate query similarity measure as it approximates well the users' conservative disposition in query refinement. Semantic approaches are based on the linguistic notion that similar words (queries) occur in similar contexts; an intuition that can be captured by statistical association measures extracted from simple document counts [5], or involving deeper analyses; e.g., of search results snippets [24]. Here we investigate a class of models for query reformulation which combines the syntactic and semantic aspects. We call these models *generalized* in the sense that they aim at capturing both syntactic and semantic properties of a reformulation. These models build upon the generalized edit distance framework. In our formulation, the cost of an edit operation, rather than being fixed, is weighted by probabilistic interpretations of the semantic relation between two terms. Our approach, while conceptually simple, unsupervised, and efficient, outperforms several competitors and baselines. We

<sup>1</sup>Through the paper with the term *syntactic* we refer purely to the surface properties of queries as sequences of symbols, without any reference to their constituent structure.

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.

SIGIR'10, July 19–23, 2010, Geneva, Switzerland.

Copyright 2010 ACM 978-1-60558-896-4/10/07 ...\$10.00.

provide empirical evidence from extensive evaluations on two datasets in Section 6.

Pushing the framework further, we investigate a generative model previously applied to biological sequence alignment problems [20]. We show that in this direction improved performance can be expected, although at increased computational cost and additional complexities in parameter estimation, leaving room for further research. While most of the focus is on single unsupervised signals for query reformulation, we show that our measures provide mutually complementary information: weighted combinations further improve performance.

The paper also touches upon a related topic. Recently, Boldi *et al.* [4] proposed the idea of capturing explicitly the relation between two queries with respect to a taxonomic representation (e.g., specification, generalization, etc.) to improve query reformulations. With respect to this issue, we show how asymmetric and symmetric probabilistic similarity measures, and their combinations, can be loosely interpreted as information-theoretic approximations of categorical notions such as "generalization" or "specification".

## 2. RELATED WORK

Query reformulation is an important topic in web search as a large fraction of the queries issued to search engines are modified after examination of the results [12]. Query modification is supported in several ways to improve search experience; e.g., via automatic spelling correction [6]. Query reformulation also requires editing or expanding the query. Several techniques have been proposed based on relevance feedback, query log analysis and distributional similarity [2, 18, 23, 24, 28, 29]. A related task to query reformulation is session segmentation [9, 13].

As relevance and pseudo relevance feedback impose additional cognitive load on the user, and can lead to query drift or costly computations, Jones *et al.* [14] proposed to pre-compute reformulations by ranking candidate queries extracted from query logs, using several types of features and learning methods. Interestingly, they notice how simple linear combinations of just a few edit distance features provide powerful ranking functions, comparable to more complex methods. Previously, Wen *et al.* [26] clustered queries combining several sources of information such as click and traditional IR document similarity, including string distance. They also suggest using a smaller fixed cost for pairs of terms occurring in Wordnet in the edit distance computation, but did not carry out a systematic evaluation. Generalizations of string matching metrics are the focus of our study. Generalized Levenshtein distance algorithms [17] have been intensively investigated in bioinformatics for solving sequence alignment problems. Oommen and Kashyap proposed a model which generates the probability of a string being rewritten into another accounting for all possible combinations of edit operations [20] that has been used successfully in peptide classification [3] and optical character recognition [15].

Previous studies have teased apart the semantic aspects of query reformulations. Rieh and Xie [22] (see also [11, 16]) analyzed query transitions in terms of syntactic and semantic operations and found that when reformulating previous queries users adopt several tactics including generalization, replacement with synonyms, parallel movement (approximately 50% of the time) and specification (approximately

30% of the time). Boldi *et al.* [4] proposed a query reformulation approach based on classifying reformulation types (QRTs) as belonging to a small taxonomy. They represent query transitions in a feature space including properties extracted from sessions and similarity features including edit distance, Jaccard and term vector cosine. Hence, they build a supervised QRT decision tree classifier which achieves 92% accuracy in a four-class task (specialization, generalization, correction, parallel move). Their methodology includes an unspecified feature selection process, thus we don't know the contribution of each feature. However, the high accuracy suggests that a few features, at least partially based on simple string matching metrics, can go a long way in capturing taxonomic aspects of query reformulations (see also Huang and Efthimiadis [11] for a related unsupervised approach). In evaluation they find that recommendations limited to specializations provide the best accuracy while introducing other types of QRTs decreases the quality of the recommendations.

## 3. PRELIMINARIES

Let  $(q_s, q_t)$  be an ordered pair where  $q_t$  is a candidate reformulation of a query  $q_s$ . We call  $q_s$  the *source* and  $q_t$  the *target*. A similarity measure between two queries is a function  $f : (q_s, q_t) \rightarrow \mathbb{R}$  which takes  $(q_s, q_t)$  as input and returns a score. In particular, we are interested in functions which correlate well with human judgments of how good a reformulation  $q_t$  is for  $q_s$ .

### 3.1 Semantic similarity

For several of the similarity measures described below, we employ *pointwise mutual information* (PMI) as a measure of the association between two terms or queries. PMI has been applied extensively to model semantic similarity – e.g., Turney [25] uses it to discover synonyms on web data – and correlates well with human judgments [21]. Let  $x$  and  $y$  be two strings that we want to measure the amount of association between. Let  $p(x)$  and  $p(y)$  be the probability of observing  $x$  and  $y$  in a given model; e.g., relative frequencies estimated from occurrence counts in a corpus. We also define  $p(x, y)$  as the joint probability of  $x$  and  $y$ ; i.e., the probability of the two strings occurring together. An abstract definition of PMI for our purposes is as follows:

$$\text{PMI}(x, y) = \log \left( \frac{p(x, y)}{p(x)p(y)} \right). \quad (1)$$

PMI can yield negative values, if  $p(x, y) < p(x)p(y)$ . For the purposes of normalization, below in this section, we discard negative PMI values and assign zero to such cases. PMI is also used as a basis for the substitutability score of two terms (see Section 5.2). Limiting PMI to positive values is further motivated by the assumption that substituting two terms occurring together less frequently than random should not be penalized more than two unrelated terms.<sup>2</sup>

As pointed out in [16, 22] query transitions tend to correlate with taxonomic relations such as generalization and specialization. Boldi *et al.* [4] show how knowledge of transition types can positively impact query reformulation. We would like to exploit this information as well. However, rather than

<sup>2</sup>As a matter of fact, such occurrences are extremely rare in our data and within noise levels.

building a dedicated supervised classifier for this task we try to capture it directly at the source. We propose that by manipulating PMI we can directly model taxonomic relations to some extent. In the following definitions we interpret  $(x, y)$  as a transition from  $x$  (i.e. source) to  $y$  (i.e. target) to break the symmetry without loss of generalization.

### 3.1.1 Joint normalization

The first type of normalization, called *joint normalization*, uses the negative log joint probability and is defined as:

$$\text{PMI}(\text{J})(x, y) = \text{PMI}(x, y) / -\log(p(x, y)). \quad (2)$$

As we limit PMI to positive values the normalization bounds the range between 0 and 1. The jointly normalized PMI(J) is a symmetric measure between  $x$  and  $y$  in the sense that  $\text{PMI}(\text{J})(x, y) = \text{PMI}(\text{J})(y, x)$ . Intuitively it is a measure of the amount of shared information between the two strings relative to the sum of individual strings information.

### 3.1.2 Specialization normalization

To capture asymmetries in the relation between two strings, we apply two non-symmetric normalizations also bounding the measure between 0 and 1. The first asymmetric normalization is called *specialization* and is defined as:

$$\text{PMI}(\text{S})(x, y) = \text{PMI}(x, y) / -\log(p(x)). \quad (3)$$

The reason we call it specialization is that PMI(S) favors pairs where the second one is a specialization of the first one. For instance, PMI(S) is at its maximum when  $p(x, y) = p(y)$  and that means the conditional probability  $p(x|y)$  is 1 which is an indication of a specialization relation.

### 3.1.3 Generalization normalization

The second asymmetric normalization is called *generalization* and is defined in the reverse direction as:

$$\text{PMI}(\text{G})(x, y) = \text{PMI}(x, y) / -\log(p(y)). \quad (4)$$

PMI(G) is at maximum when  $p(y|x)$  is 1.

The three normalizations provide a richer representation of the association between two strings and approximate the generalization-specialization dimension from an information-theoretic perspective. As an example, for the query transition “apple” to “mac os”  $\text{PMI}(\text{G})=0.2917$  and  $\text{PMI}(\text{S})=0.3686$ ; i.e., there is more evidence for a specialization. Conversely for the query transition “ferrari models” to “ferrari” the measures yield  $\text{PMI}(\text{G})=1$  and  $\text{PMI}(\text{S})=0.5558$ ; i.e., the target is a “perfect” generalization of the source<sup>3</sup>.

## 3.2 Syntactic similarity

Let  $V$  be a finite vocabulary and  $\xi$  be the null symbol. An edit operation: insertion, deletion or substitution, is a pair  $(a, b) \in \{V \cup \{\xi\} \times V \cup \{\xi\}\} - \{(\xi, \xi)\}$ . An alignment between two sequences  $x$  and  $y$  is a sequence of edit operations  $\omega = (a_1, b_1), \dots, (a_n, b_n)$ . Given a non-negative cost function  $c$  the cost of an alignment is  $c(\omega) = \sum_{i=1}^n c(\omega_i)$ . The Levenshtein distance, or edit distance, defined over  $V$ ,  $d_V(x, y)$  between  $x$  and  $y$  is the cost of the least expensive sequence of edit operations which transforms  $x$  into  $y$  [17]. The distance computation can be performed via dynamic programming in time  $O(|x||y|)$ . Edit distance captures the amount of overlap between the queries as sequences of symbols and have been previously used in information retrieval [4, 14, 27].

<sup>3</sup>The values are computed from Web counts.

## 4. QUERY-LEVEL MEASURES

We calculate the PMI for a pair  $(q_s, q_t)$  using the number of documents retrieved by a search engine for  $q_s$ ,  $q_t$  and  $q_{s,t}$  where  $q_{s,t}$  is a shorthand for the concatenation of  $q_s$  and  $q_t$ . Formally, let  $N_s$  and  $N_t$  be the number of documents retrieved for  $q_s$  and  $q_t$  respectively. Similarly, let  $N_{s,t}$  be the number of documents retrieved for the concatenated joint query. We define the *probability* of, respectively, the two queries and the joint query as  $p(q_s) = \frac{N_s}{N}$ ,  $p(q_t) = \frac{N_t}{N}$ , and  $p(q_s, q_t) = \frac{N_{s,t}}{N}$  where  $N$  is a constant large enough to approximate the total number of documents that can be retrieved. In our implementation we use Google’s search engine. The number of results returned determine  $N_s$ ,  $N_t$ , and  $N_{s,t}$ . We denote this PMI measure between two queries by  $\text{PMI}_{Web}(q_s, q_t)$ . We renormalize the PMI values as described above thus generating three query-level similarity measures.

## 5. TERM-LEVEL MEASURES

### 5.1 Syntactic measures

We use two Levenshtein distance models as basic syntactic measures. The first, called Edit1 (E1), employs a unit cost function for each of the three operations. That is, given a finite vocabulary  $T$  of all terms occurring in queries:

$$\forall a, b \in T, c_{E1}(a, b) = 1 \text{ if } a \neq b, 0 \text{ otherwise.} \quad (5)$$

The second, called Edit2 (E2), uses unit costs for insertion and deletion, but computes the character-based edit distance between two terms to determine the substitution cost. If two terms are similar at the character level, the cost of substitution is lower. Given the vocabulary  $T$  of terms and a finite vocabulary  $A$  of characters the cost function is defined as:

$$\forall a, b \in T, c_{E2}(a, b) = d_A(a, b) \text{ if } a \wedge b \neq \xi, 1 \text{ otherwise} \quad (6)$$

where  $0 \leq d_A(a, b) \leq 1$ , normalizing by  $\max(|a|, |b|)$ .

We also investigate a variant in which the input sequences are alphabetically sorted before the edit distance computation. The motivation is the observation that queries may be often formulated as sets of terms in which the order of the terms is irrelevant. Thus, “Brooklyn pizza” and “pizza Brooklyn” may denote same user intent but the edit distance is unable to capture the similarity. By presorting the terms in the queries we compute an order-free version of edit distance. We prefix the names of these models with “Sorted”.

### 5.2 Generalized edit distance

Extending the Levenshtein distance framework to take into account semantic similarities between terms is conceptually simple. As in the Edit2 model above we use a modified cost function. For our purposes, the cost function should have the following properties: whenever there is evidence of semantic association between two terms, it should be “cheaper” to substitute these terms instead of deleting one and inserting the other. For an unrelated term pair, a combination of insertion and deletion should always be less costly than a substitution. We also assume that the cost of the substitution of a term with itself (i.e. identity substitution) is always 0. Considering these requirements, we define the cost function as a *cost matrix*  $S$  based on the normalized PMI measures defined above. Given a normalized similarity measure  $f$ , an entry in a cost matrix  $S$  for a term pair

$(w_i, w_j)$  is defined as:

$$s(w_i, w_j) = 2 - 2f(w_i, w_j) + \epsilon \quad (7)$$

The  $\epsilon$  correction, coupled with unit insertion and deletion costs, guarantees that these requirements are fulfilled. We call these models GenEdit (GE). Given a finite term vocabulary  $T$  and cost matrix  $S$  the cost function is defined as:

$$\forall a, b \in T, c_{GE}(a, b) = s(a, b) \text{ if } a \wedge b \neq \xi, 1 \text{ otherwise.} \quad (8)$$

### 5.2.1 Cost matrix estimation

To estimate a cost matrix we used session logs consisting of actual transitions of consecutive queries. The data consists of approximately 1.3 billion English queries generated from the U.S. A session is defined as a sequence of queries from the same user within a controlled time interval. Let  $q_s$  and  $q_t$  be a query pair observed in the session data where  $q_t$  is issued immediately after  $q_s$  in the same session. Let  $q'_s = q_s \setminus q_t$  and  $q'_t = q_t \setminus q_s$ , where  $\setminus$  is the set difference operator. The co-occurrence count of two terms  $w_i$  and  $w_j$  from a query pair  $q_s, q_t$  is denoted by  $n_{i,j}(q_s, q_t)$  and is defined as:

$$n_{i,j}(q_s, q_t) = \begin{cases} 1 & \text{if } w_i = w_j \wedge w_i \in q_s \wedge w_j \in q_t \\ 1/(|q'_s| |q'_t|) & \text{if } w_i \in q'_s \wedge w_j \in q'_t \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

If a term occurs in both queries, it has a co-occurrence count of 1. For all other pairs we make sure the sum of co-occurrence counts for a term  $w_i \in q_s$  is 1 for a given query pair. The normalization is an attempt to avoid the under-representation of terms occurring in both queries. The final co-occurrence count of two arbitrary terms  $w_i$  and  $w_j$  is denoted by  $N_{i,j}$  and it is defined as the sum over all query pairs in the session logs,  $N_{i,j} = \sum_{q_s, q_t} n_{i,j}(q_s, q_t)$ . Let  $N = \sum_{i,j} N_{i,j}$  be the sum of co-occurrence counts over all term pairs. Then we define a joint probability for a term pair as  $p(x, y) = \frac{N_{i,j}}{N}$ . Similarly, we define the single-occurrence counts and probabilities of the terms by computing the marginalized sums over all term pairs. Namely, the probability of a term  $w_i$  occurring in the source query is  $p(i, \cdot) = \sum_j N_{i,j}/N$  and similarly the probability of a term  $w_j$  occurring in the target query is  $p(\cdot, j) = \sum_i N_{i,j}/N$ . Plugging these values in Eq. (1), we obtain the PMI( $i, j$ ) for term pair  $w_i$  and  $w_j$ , which are further normalized as described in Section 3.1. Any term pair that is not co-occurring in the session data is considered to be unrelated and is assigned a PMI value of zero.

## 5.3 A generative model

The edit distance measures considered so far generate a score based on the least-costly alignment of two queries. This can be viewed as finding the shortest path in a query space constructed on the atomic edit operations. A natural extension would calculate the probability of producing the target from the source not only considering the least-costly alignment of two queries but computing over all possible ways the target can be obtained from the source. Oommen and Kashyap [20], proposed a syntactic transition probability model (referred to as the OK model) which shows how the probability of a string, in our case a query, rewrite can be computed with a generative model consisting of random insertion, deletion, and substitution operations. The model has been successfully applied to problems such as peptide classification and OCR correction [3, 15].

Let  $x = (x_1 x_2 \dots x_n)$  and  $y = (y_1 y_2 \dots y_m)$  be the source and target strings respectively such that  $x \in V^*$  and  $y \in V^*$  where  $V$  is the finite alphabet of symbols. We introduce two additional symbols  $\xi$  and  $\lambda$ , which are not in  $V$ , as input and output null symbols, respectively. The OK model computes the probability of obtaining  $y$  from  $x$  under a generative model which takes as input two probability distributions  $G$  and  $S$ , and works in the following steps.

1. Distribution  $G$  specifies the number of insertions to be applied to the source. In each independent string generation process  $z$  terms are inserted with probability  $G(z)$ . The intermediate output at the end of this step is  $x' = x'_1 x'_2 \dots x'_{n+z}$  where  $x'$  is the string  $x$  modified by inserting  $z$   $\xi$  symbols at random positions of  $x$ .
2. The distribution  $S$  is over  $\{T \cup \{\xi\} \times T \cup \{\lambda\}\}$ . The value  $S(y_j | x_i)$  for two symbols  $x_i$  and  $y_j$  is the probability that  $x_i$  is substituted by  $y_j$ . The output of this step is  $y'$  which is obtained by substituting all symbols in  $x'$  according to the probabilities specified by  $S$ .
3. Last step mirrors the original deletion operation: all  $\lambda$  characters remaining in  $y'$  are deleted. The remaining string is  $y$ , the output of the process.

There are two constraints on  $S$ . The first states that  $\forall x_i \sum_{y_j \in T \cup \{\lambda\}} S(y_j | x_i) = 1$ . This ensures that each symbol in  $x$  is either substituted, left intact or deleted. The second constraint  $S(\lambda | \xi) = 0$  guarantees that exactly  $z$  insertions are made and no input null symbol inserted at the first step is deleted in the second step (i.e.  $\xi$  is always substituted by a symbol in  $T$  not by  $\lambda$ ).

Using the OK model, one can compute the transition probability of a source query to a target query by integrating the individual probabilities of all possible paths allowed by the generative model. The explicit form of this probability is:

$$p(y|x) = \sum_{z=\max(0, m-n)}^m \frac{G(z)n!z!}{(n+z)!} \sum_{x'} \sum_{y'} \prod_{k=1}^{n+z} S(y'_k, x'_k) \quad (10)$$

where  $m = |x|$  and  $n = |y|$ . The outermost summation is over all possible values of number of insertions. The factor  $\frac{n!z!}{(n+z)!}$  is the number of different ways in which  $z$   $\xi$  characters can be inserted in  $x$  to obtain  $x'$ . Although the explicit calculation of this probability is too expensive, due to the combinatorial element, Oommen and Kashyap provide a dynamic program which runs in approximately cubic time,  $O(mn \min(m, n))$  [20].

As in the generalized edit distance models, we represent the queries as strings and the terms as characters. That is,  $q_s = w_1 w_2 \dots w_n$  is the source query and  $q_t = w_1 w_2 \dots w_m$  is the target query where  $w_i, w_j \in T$ .

### 5.3.1 Parameter estimation

In order to actually employ the OK model we need to estimate distributions  $G$  and  $S$ . The model accepts arbitrary probability distributions and estimating meaningful parameters is not trivial. We devise an estimation strategy similar to that used by Kolak and Resnik [15] who apply the OK model to optical character recognition (OCR). The idea is to generate an alignment with a simpler model and then estimate all substitutions in  $S$  directly from a large-enough

dataset. As a corpus of pairs to align we used the session data of Section 5.2.1, under the assumption that contiguous query pairs represent reasonable candidates of naturally occurring query transitions. Subsequently, we ran a generalized edit-distance model to find the least-costly alignment of each pair and then count the edit operations that make up this alignment. The application of the edit-distance model thus provides a way to reverse engineer the query transitions and obtain estimates for the term insertion, deletion and substitution probability distributions.

Summarizing, we run a generalized edit distance model on the query pairs in our session dataset<sup>4</sup>, and count the number of times each term insertion, deletion, and substitution occurred directly off the alignments. By integrating these counts over all pairs and normalizing them into probability distributions, we obtain the necessary estimates for  $G$  and  $S$ . The inspection of the outcome of this procedure revealed that the deletion probability (i.e.  $S(\lambda|x_i)$ ) is largely over-estimated (e.g. over 0.3 for some terms). This is possibly due to the noisy alignment procedure. To solve this problem we introduce one adjustable parameter called *damping factor*, denoted by  $DF$ . For each term, the deletion probability is corrected as  $S(\lambda|t) \leftarrow S(\lambda|t)/DF$ , then  $S$  is re-normalized so that  $\forall x_i \sum_{y_j \in T \cup \{\lambda\}} S(y_j|x_i) = 1$ .

## 6. EXPERIMENTS

We evaluate all models discussed so far on two datasets. As an external benchmark we use the unsupervised distributional similarity system (DistSim) of Alfonseca *et al.* [1]. DistSim implements an extension of the vector-space model of distributional similarities to query-to-query similarity, by combining different vectors using the geometric mean of the frequencies for each of the features separately. Features are  $n$ -grams collected in the context of each query in hundreds of millions of documents. The score of a query reformulation is the cosine of the vectors representing each query. DistSim generates richly lexicalized high-dimensional models which in evaluation [1] outperformed web kernel methods [24].

### 6.1 Experimental setup

The evaluation involves a query reformulation task in which several source queries are provided, each with a set of candidate reformulations scored by raters. Each model predicts a real-valued score for each source-target reformulation pair. The score represents the quality of the reformulation according to the system. While the absolute value of the score might not be meaningful in itself it is used to rank the queries in the set of possible targets for the same source. Several evaluation metrics are used to quantify the performance of a system: Spearman rank correlation, precision at  $N$  and mean average precision. We use Spearman correlation as our primal evaluation measure as it is independent of the choice of a threshold which is necessary for precision.

### 6.2 Combined models

We evaluate all the similarity measures individually as well as in combination. We experiment with one unsupervised combination method, a baseline which simply averages all signals<sup>5</sup>. We also evaluate a supervised combination,

<sup>4</sup>A sorted joint-normalized generalized edit distance model.

<sup>5</sup>All non-normalized individual signals are re-normalized before combination.

	QS1500	CC2000
Number of query pairs	1500	2000
Number of source queries	57	500
Average log-probability $q_s$	-10,61	-10,57
Average log-probability $q_t$	-9,33	-10,27
Average Number of terms $q_s$	3,40	2,08
Average Number of terms $q_t$	2,83	2,24

**Table 1.** Statistics of the evaluation datasets.

and a supervised optimization of the OK model. The OK model involves an adjustable parameter  $DF$  that should be picked empirically. Hence, we optimize the OK model separately by a supervised leave-one-out procedure. All values for  $DF = 10^i$  were evaluated at  $i = 1, 2, \dots, 10$ . As a full supervised combination we used a neural network regression model using all of the features introduced in the paper, excluding OK. This approach lets us exploit potential non-linearities in the signals. For each network model three parameters are optimized: learning rate, number of hidden units and number of iterations (epochs) over the training data. Predictions are generated in a leave-one-out scheme where in turn a source query  $q_s$  is excluded for prediction.

### 6.3 Evaluation data

The first evaluation set, QS1500 is based on the gold standard from [1]. It contains 57 source queries, each paired with up to 20 target queries. The candidate reformulations are generated from the top-20 ranked suggestions using several different systems, based on the web kernel approach [24], and distributional similarity. Two raters evaluated each pair, using the 5-Likert scale defined in [24]. The weighted Cohen’s Kappa was 0.7111 on a binary split at level 1.5, indicating substantial agreement amongst the raters for a binary decision. In the computation of the precision at  $N$  scores we use the pairs with a score of 1.5 or more [1] as positive pairs.

The second evaluation set, called CC2000, was built from scratch based on the hypothesis that two different queries are related if they lead to user clicks on the same documents. This approach is similar to the method proposed by Fitzpatrick and Dent [8]. Our technique adds click information, thus strengthening the preference for precision over recall in the extraction of related queries. For a randomly extracted set of 500 source queries, we randomly sampled 4 targets. 3 out of 4 targets are queries that have been co-clicked with at least 10 different results. The remaining one has been co-clicked only once. The latter pair acts as a control on the quality of click as a measure of relatedness. The 2,000 pairs were judged by 5 raters, with access to the search result, in blind evaluation according to a 4-point scale: Unrelated(1), Slightly Related(2), Very Related(3), Same Meaning(4). Inter-rater agreement of 5 raters on a binary classification task (class 1 = Unrelated or Slightly Related, class 2 = Very Related or Same Meaning) gave a Kappa value of 0.65. A connection between the co-click hypothesis and human ratings can be seen from computing average human scores for the automatically created distinction. This results in an average human score of 3.1 for pairs with more than 10 co-clicks, and an average human score of 2.3 for pairs with 1 co-click. This shows that the co-click hypothesis yields positively related pairs that are judged on average as Very Related by human raters, while the control set are judged only as Slightly Related. The gold standard for each pair

QS1500							
#	Similarity Function	Spearman	mAP	Prec@1	Prec@3	Prec@5	Sig.
1	NN	0.500	0.806	0.836	0.741	0.637	7
2	Oommen-Kashyap	0.470	0.747	0.782	0.698	0.637	10
3	DistSim	0.438	0.744	0.768	0.679	0.628	12
4	Mean all	0.435	0.772	0.818	0.691	0.633	13
5	SortedGenEdit(S)	0.429	0.774	0.845	0.709	0.638	13
6	SortedGenEdit(G)	0.428	0.775	0.828	0.712	0.648	13
7	PMI <sub>Web</sub> (S)	0.417	0.713	0.764	0.630	0.589	13
8	PMI <sub>Web</sub> (J)	0.409	0.730	0.782	0.679	0.594	13
9	SortedGenEdit(J)	0.408	0.771	0.832	0.701	0.639	13
10	GenEdit(S)	0.382	0.743	0.796	0.695	0.619	15
11	GenEdit(G)	0.380	0.745	0.795	0.698	0.625	15
12	GenEdit(J)	0.365	0.737	0.790	0.692	0.609	15
13	SortedEdit2	0.320	0.714	0.757	0.668	0.630	18
14	SortedEdit1	0.314	0.697	0.763	0.660	0.595	18
15	PMI <sub>Web</sub> (G)	0.283	0.670	0.627	0.547	0.570	18
16	Edit2	0.270	0.649	0.715	0.618	0.571	18
17	Edit1	0.252	0.633	0.683	0.615	0.550	18
18	Length-target(Char)	0.139	0.519	0.435	0.456	0.437	20
19	Length-target(Term)	0.112	0.506	0.453	0.423	0.413	20
20	log-prob target	-0.161	0.452	0.309	0.309	0.341	-

**Table 2.** The grand table for QS1500. The column Sig. gives the index of the model with the highest Spearman correlation that the corresponding model is significantly higher than with  $p < 0.05$ . Length and log probability of target are absolute baselines.

Similarity Measure	QS1500	CC2000
Oommen-Kashyap	<b>0.470*</b>	0.391* (6)
SortedGenEdit(S)	0.429*	0.407* (4)
SortedGenEdit(G)	0.428*	0.419* (2)
SortedGenEdit(J)	0.408	0.391* (7)
GenEdit(S)	0.382	0.414* (3)
GenEdit(G)	0.380	<b>0.424* (1)</b>
GenEdit(J)	0.365	0.402* (5)
SortedEdit2	0.320	0.288 (11)
SortedEdit1	0.314	0.298 (9)
Edit2	0.270	0.292 (10)
Edit1	0.252	0.299 (8)

**Table 3.** Generalized edit distances for QS1500 and CC2000. The ranks of the features for CC2000 are given in parentheses; \* indicates a higher Spearman correlation than the highest performing edit distance baseline (SortedEdit2 for QS1500 and Edit1 for CC2000) at a significance level of 0.95.

is the average of the 5 ratings. Choosing a fixed threshold for the precision scores is not straightforward; e.g., using a threshold at three produces 132 all positive sets and 90 all negative sets, therefore we would not be able to compute a meaningful precision score for too many sets. To avoid this problem we choose in each set the positive pair as the one with the highest score. In this way we obtain 774 positive pairs and 1221 negative pairs. Thus in terms of precision we evaluate the performance of systems at identifying the best available pair. The following table summarizes some datasets statistics: Table 1 summarizes the basic properties of the datasets.

## 7. RESULTS AND DISCUSSION

In the following sections, we discuss the performance of generalized edit distance with respect to baselines (simple edit distance and distributional similarity models), comment on the effect of taxonomic normalization of PMI, and report the performance of combining different measures in supervised and unsupervised settings. In order to give a birds-eye overview, the results for all models are given in Tables 2 and 4 for QS1500 and CC2000, respectively. In these tables, we report Spearman correlation, mean average precision and

precision at various positions. Since, there are only 4 target queries per source in CC2000, we report precision values at 1, 2 and 3 for that dataset. For QS1500, the precision values at 1, 3 and 5 are reported.

### 7.1 Generalized Edit Distance

The Spearman rank correlations obtained for all edit distance models are given in Table 3 for QS1500 and CC2000. Several points are worth discussing in these results.

*Generalized edit distance is better than simple edit distance.* For both datasets, the generalized edit distance models (all variants of GenEdit and SortedGenEdit) outperform the simple edit-distance based features (Edit1 and Edit2). This observation is also supported by the significance results obtained by one-tailed t-tests reported in the same table. This result proves that our method is a powerful, yet simple, generalization of an already robust query similarity measure. To the best of our knowledge ours is the first (successful) application of such generalized algorithms to IR.

*Sorting has an effect on results.* For QS1500, sorted edit-distance based features (SortedGenEdit) outperformed their unsorted counterparts (GenEdit) by margins of more than 4 percent (though, we were unable to confirm a significant difference between them). The pattern is different in CC2000 where the unsorted features outperform their sorted counterparts albeit with smaller margins. We hypothesize that this effect might be related to the query length, greater in QS1500 both for  $q_s$  and  $q_t$ . It is possible that as longer queries are more subject to permutations, sorted distance measures emerge as more robust.

*Generative models are promising.* Especially for QS1500, we see that the increased complexity of the OK model pays off in terms of performance. Since the OK model uses the substitution probability matrices computed by the alignments obtained by SortedGenEdit(J) model, the difference between the SortedGenEdit(J) and the OK model becomes even more encouraging. Although we picked the damping factor value for each set by supervision (i.e., leave-one-out), we should note that the OK model is robust with respect to varying values of  $DF$ . In the optimization experiments we

CC2000							
#	Similarity Function	Spearman	mAP	Prec@1	Prec@2	Prec@3	Sig.
1	NN	0.432	0.739	0.583	0.511	0.451	10
2	GenEdit(G)	0.424	0.716	0.545	0.485	0.447	11
3	SortedGenEdit(G)	0.419	0.712	0.550	0.476	0.446	11
4	GenEdit(S)	0.414	0.714	0.543	0.488	0.447	11
5	SortedGenEdit(S)	0.407	0.710	0.546	0.477	0.445	11
6	GenEdit(J)	0.402	0.714	0.541	0.482	0.448	11
7	Oommen-Kashyap	0.391	0.704	0.515	0.484	0.451	11
8	SortedGenEdit(J)	0.391	0.706	0.538	0.474	0.445	11
9	Mean all	0.386	0.711	0.531	0.485	0.448	12
10	PMI <sub>Web</sub> (G)	0.369	0.698	0.506	0.473	0.449	13
11	PMI <sub>Web</sub> (J)	0.330	0.692	0.485	0.474	0.444	17
12	DistSim	0.322	0.707	0.532	0.492	0.427	17
13	Edit1	0.299	0.656	0.441	0.438	0.418	17
14	SortedEdit1	0.298	0.662	0.456	0.444	0.420	17
15	Edit2	0.292	0.676	0.478	0.458	0.432	17
16	SortedEdit2	0.288	0.685	0.487	0.461	0.432	17
17	PMI <sub>Web</sub> (S)	0.264	0.681	0.477	0.461	0.437	17
18	log-prob target	0.114	0.626	0.384	0.416	0.404	19
19	Length-target(Char)	-0.036	0.603	0.362	0.390	0.392	-
20	Length-target(Term)	-0.077	0.603	0.358	0.387	0.388	-

**Table 4.** The grand table for CC2000. The column Sig. gives the index of the model with the highest Spearman correlation that the corresponding model is significantly higher than with  $p < 0.05$ . Length and log probability of target are absolute baselines.

observed that the Spearman scores of the OK models remain at levels either comparable to or superior than other generalized models’ in a range of almost 7 orders of magnitude for *DF*. However, the simpler generalized models are even more robust. Better formulations of the generative model constitute an interesting direction for future research.

## 7.2 Taxonomic normalization

*Type of normalization is important.* Especially for the query-level similarity measures (variants of PMI<sub>Web</sub>), the type of normalization has a significant effect on performance. E.g., as it can be seen in Table 2, PMI<sub>Web</sub>(G) performs badly in QS1500 (Spearman 0.283, rank 15), but PMI<sub>Web</sub>(S) is more competitive (Spearman 0.417, rank 7) even though both measures are based on the same PMI values and only differ by their type of normalization. Interestingly, a similar but reversed pattern is observed for CC2000 in Table 4. For this dataset, PMI<sub>Web</sub>(G) is the best measure among the query-level measures with a Spearman correlation of 0.369 and overall rank of 10. PMI<sub>Web</sub>(S), on the other hand, can achieve a correlation of 0.264 and is placed quite low in the overall ranking. The difference between PMI<sub>Web</sub>(S) and PMI<sub>Web</sub>(G) is significant ( $p < 0.05$ ) for both datasets in opposite directions. This evidence alone suggests that different normalizations can capture different properties of different datasets. A similar pattern also emerges from the generalized edit distance models, specialization works best for QS1500 and generalization for CC2000.

One intriguing explanation for this pattern involves the dominant directionality in the datasets. We know that in QS1500 target queries are shorter than source queries on average, suggesting that transitions are more likely in the generalization direction. It is possible that a feature which favors transitions in the generalization direction loses its discriminative power and do not correlate well with human judgments because of the bias in the dataset. A similar effect in the reversed direction is compatible with results on CC2000. Further investigations of such information-theoretic approximations of taxonomic notions related to

generality and inclusion seem also particularly promising as they don’t rely on pre-computed resources.

## 7.3 Combined models

*Supervision works.* In Tables 2 and 4, we see that the neural network model constructed on all features (NN) outperformed all other methods in both datasets. This was rather expected but it is important to see that there is room for further improvements, and that the features we propose provide complementary information. Experiments with linear regression were less successful which suggests there might be non-linear interactions between features that can be captured by the neural network. It is also interesting to notice that supervised combination improves but not by a large margin indicating that our single generalized features have good discriminative power in absolute terms. The comparison with the high-dimensional distributional similarity model (DistSim) is also positive, DistSim performs only marginally better than the GenEdit models on one dataset (QS1500) – although, in terms of precision, GenEdit models are still better – and worse on CC2000. Naive unsupervised combination yields mixed results.

## 8. CONCLUSION

In this paper we proposed an approach to query reformulation aiming at the combination of string similarity measures and corpus-based semantic association measures. Generalized Levenshtein distance algorithms provide a principled framework for this combination. By manipulating the edit distance cost function our models can incorporate naturally useful statistical association measures, including variants of pointwise mutual information which, to some extent, capture directly taxonomic relations between terms. The models we proposed are mostly unsupervised, compact and efficient, and we provided empirical evidence of their effectiveness. We also explored a generative query reformulation model which provides further improvements at some additional computational cost and estimation complexity. Finally, we evaluated supervised combinations proving that the features capture complementary aspects of the data.

This framework offers several opportunities for further research. In a related work [7] we investigate supervised models based on our features trained on noisy data within a learning to rank framework. Another interesting topic involves, as in bioinformatics, controlling the costs of all edit operations, including insert and delete, by applying algorithms such as Needleman-Wunsch [19]. Another interesting and promising topic involves moving beyond context-free reformulation methods in the generalized framework. In our approach substitution costs involve pairs of terms independent of the surrounding context while it seems likely that dependencies between terms should be taken into account in the computation of the best reformulation.

## 9. ADDITIONAL AUTHORS

Additional authors: Enrique Alfonseca (Google, Zürich, Switzerland, email: ealfonseca@google.com).

## 10. REFERENCES

- [1] E. Alfonseca, K. Hall, and S. Hartmann. Large-scale computation of distributional similarities for queries. In *Proceedings of NAACL-HLT*, pages 29–32. Association for Computational Linguistics, 2009.
- [2] J. Allan. Relevance feedback with too much data. In *Proceedings of SIGIR*, pages 337–343. ACM, 1995.
- [3] E. Aygün, B. Oommen, and Z. Cataltepe. On utilizing optimal and information theoretic syntactic modeling for peptide classification. In *Pattern Recognition in Bioinformatics*, volume 5780 of *Lecture Notes in Computer Science*, pages 24–35. Springer Berlin, 2009.
- [4] P. Boldi, F. Bonchi, C. Castillo, and S. Vigna. From ‘Dango’ to ‘Japanese cakes’: Query reformulation models and patterns. In *Proceedings of Web Intelligence*. IEEE Cs Press, 2009.
- [5] K. Church, W. Gale, P. Hanks, and D. Hindle. Using statistics in lexical analysis. In *Lexical Acquisition: Exploiting On-Line Resources to Build a Lexicon*, pages 115–164. Erlbaum, 1991.
- [6] B. Cucerzan and E. Brill. Spelling correction as an iterative process that exploits the collective knowledge of web users. In *Proceedings of EMNLP*, pages 293–300. Association for Computational Linguistics, 2004.
- [7] F. De Bona, S. Riezler, K. Hall, M. Ciaramita, A. Herdağdelen, and M. Holmqvist. Learning dense models of query similarity from user click logs. In *Proceedings of NAACL-HLT*. Association for Computational Linguistics, 2010.
- [8] L. Fitzpatrick and M. Dent. Automatic feedback using past queries: Social searching? In *Proceedings of SIGIR*, pages 306–313. ACM, 1997.
- [9] D. He, A. Göker, and D. Harper. Combining evidence for automatic web session identification. *Information Processing and Management*, 38(5):727–742, 2002.
- [10] M. Hearst. *Search user interfaces*. Cambridge University Press, 2009.
- [11] J. Huang and E. Efthimiadis. Analyzing and evaluating query reformulation strategies in web search logs. In *Proceedings of CIKM*, pages 77–86. ACM, 2009.
- [12] B. Jansen, A. Spink, and S. Koshman. Web searcher interaction with the dogpile.com metasearch engine. *Journal Of The American Society For Information Science And Technology*, 58(5):744–755, 2007.
- [13] R. Jones and K. Klinkner. Beyond the session timeout: automatic hierarchical segmentation of search topics in query logs. In *Proceedings of CIKM*, pages 699–708. ACM, 2008.
- [14] R. Jones, B. Rey, O. Madani, and W. Greiner. Generating query substitutions. In *Proceedings of WWW*, pages 387–396. ACM, 2006.
- [15] O. Kolak and P. Resnik. OCR error correction using a noisy channel model. In *Proceedings of HLT*, pages 29–32. Association for Computational Linguistics, 2002.
- [16] T. Lau and E. Horvitz. Patterns of search: analyzing and modeling web query refinement. In *Proceedings of the seventh international conference on User modeling*, pages 119–128. Springer-Verlag New York, Inc., 1999.
- [17] V. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10(8):707–710, 1966.
- [18] M. Mitra, A. Singhal, and C. Buckley. Improving automatic query expansion. In *Proceedings of SIGIR*, pages 206–214. ACM, 1998.
- [19] S. Needleman and C. Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48(3):443–453, 1970.
- [20] B. Oommen and R. Kashyap. A formal theory for optimal and information theoretic syntactic pattern recognition. *Pattern Recognition*, 31(8):1159–1177, 1998.
- [21] G. Recchia and M. Jones. More data trumps smarter algorithms: comparing pointwise mutual information with latent semantic analysis. *Behavioral Research Methods*, 41(3):647–656, 2009.
- [22] S. Rieh and H. Xie. Analysis of multiple query reformulations on the web: the interactive information retrieval context. *Inf. Process. Manage.*, 42(3):751–768, 2006.
- [23] S. Riezler, Y. Liu, and A. Vasserman. Translating queries into snippets for improved query expansion. In *Proceedings of Coling*, pages 737–744, 2008.
- [24] M. Sahami and T. Heilman. A web-based kernel function for measuring the similarity of short text snippets. In *Proceedings of WWW*, pages 377–386. ACM, 2006.
- [25] P. Turney. Mining the web for synonyms: PMI-IR versus LSA on TOEFL. *Lecture Notes in Computer Science*, 2167:491–503, 2001.
- [26] J. Wen, J. Nie, and H. Zhang. Clustering user queries of a search engine. In *Proceedings of WWW*, pages 162–168. ACM, 2001.
- [27] J. Wen, J. Nie, and H. Zhang. Query clustering using user logs. *ACM Trans. Inf. Syst.*, 20(1):59–81, 2002.
- [28] J. Xu and B. Croft. Query expansion using local and global document analysis. In *Proceedings of SIGIR*, pages 4–11. ACM, 1996.
- [29] Z. Zhang and O. Nasraoui. Mining search engine query logs for query recommendation. In *Proceedings of WWW*, pages 1039–1040. ACM, 2006.