

# Query Enrichment for Web-query Classification

DOU SHEN<sup>1</sup>

RONG PAN<sup>1</sup>

JIAN-TAO SUN<sup>2</sup>

JEFFREY JUNFENG PAN<sup>1</sup>

KANGHENG WU<sup>1</sup>

JIE YIN<sup>1</sup>

QIANG YANG<sup>1</sup>

<sup>1</sup>Department of Computer Science and Engineering, Hong Kong University of Science and Technology

<sup>2</sup>Microsoft Research Asia

---

Web search queries are typically short and ambiguous. To classify these queries into certain target categories is a difficult but important problem. In this paper, we present a new technique called query enrichment, which takes a short query and maps it to the intermediate objects. Based on the collected intermediate objects, the query is then mapped to the target categories. To build the necessary mapping functions, we use an ensemble of search engines to produce an enrichment of the queries. Our technique was applied to ACM Knowledge-discovery and data mining competition (ACM KDDCUP) in 2005, where we won the championship on all three evaluation metrics (precision, F1 measure, which combines precision and recall together, and creativity, which is judged by the organizers) among a total of 33 teams worldwide. In this paper, we show that, despite the difficulty in an abundance of ambiguous queries and a lack of training data, our query-enrichment technique can solve the problem satisfactorily through a two-phase classification framework. We present a detailed description of our algorithm and experimental evaluation. Our best result of F1 and precision are 42.4% and 44.4%, respectively, which are 9.6% and 24.3% higher than those from the runner-ups, respectively.

Categories and Subject Descriptors: H.2.8 [**Database Management**]: Database Applications—*Data mining*; H.3.3 [**Information Storage and Retrieval**]: Information Search and Retrieval; I.5.2 [**Pattern Recognition**]: Design Methodology— *Classifier design and evaluation*

General Terms: Algorithms, Experimentation, Performance

Additional Key Words and Phrases: Query classification, query enrichment, synonym-based classifier, ensemble learning, KDDCUP2005

---

---

Authors' addresses: Dou Shen, Rong Pan, Jeffrey Junfeng Pan, Kangheng Wu, Jie Yin and Qiang Yang, Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong, China; email: {dshen, panrong, panjf, khwu,yinjie,qyang}@cse.ust.hk; Jian-Tao Sun, Microsoft Research Asia, 49 Zhichun Road, Beijing, China; email: jtsun@microsoft.com.

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 20? ACM 0000-0000/20?/0000-0001 \$5.00

## 1. INTRODUCTION

Historically, search engine technologies and automatic text classification techniques have progressed hand in hand. Ever since the early papers by the pioneers [Jones 1971; Lewis and Gale 1994], people have recognized the possibility of conducting Web search through classification, and vice versa [Page et al. 1998; Beeferman and Berger 2000; Chekuri et al. 1997; Chen and Dumais 2000; Kang and Kim 2003]. The 2005 ACM Knowledge Discovery and Data Mining competition (KDDCUP2005 for short) made this connection even stronger. This competition series (see <http://www.acm.org/sigs/sigkdd/kddcup/>), which has a long history since 1997, is open to researchers and practitioners worldwide. Being one of the most prominent data-mining competitions, the datasets are often used later as benchmarks. The task in KDDCUP2005 is to automatically classify a subset of query log data in one month in 2005 from the MSN search engine (<http://search.msn.com>), one of the major search engines, into a set of given target categories accurately. The log data contain 800,000 queries and the target categories consist of 67 predetermined categories provided by the organizers. The dataset is available as a public domain benchmark for query classification (<http://www.acm.org/sigs/sigkdd/kdd2005/kddcup.html>). Several example queries are shown in Table I. An illustration of the hierarchical structure for the target categories is shown in Figure 1 (see Appendix A for details).

Table I. Examples of queries.

1967 shelby mustang
actress hildegard
Aldactone
alfred Hitchcock
amazon rainforest
section8rentalhouses.com
Sakpsabancnnhayat
auto price
a & r management" property management Maryland
netconfig.exe

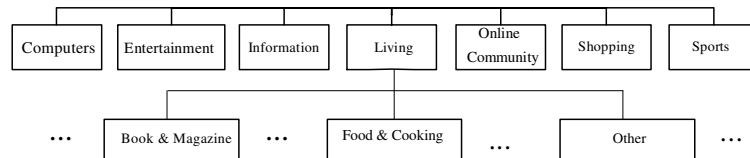


Fig. 1. Illustration of the hierarchy structure of 67 target categories.

The KDDCUP2005 task highlights the importance and difficulties of query classification, a way to understand queries [Li et al. 2005; Shen et al. 2005; Kardkovács et al. 2005; Vogel et al. 2005; Shen et al. 2006]. Query classification can be used

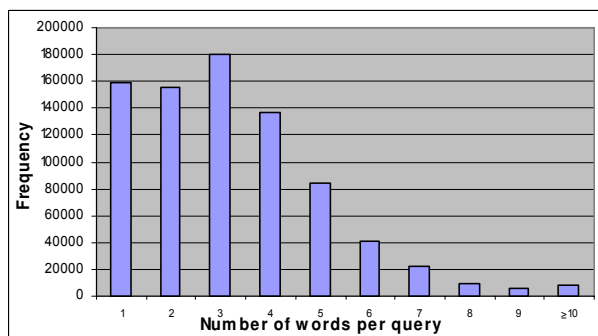


Fig. 2. Frequency of queries with different length.

to support several important tasks in information retrieval and Web search. In information retrieval, one potential area is to construct user models in order to cater to individual users' and group users' preferences. Classification of user queries is a component step both in constructing the user models and in utilizing the user models. In Web search, an important application is to organize the large number of Web pages in the search result, after the user issues a query, according to the potential categories of the results. Query classification can be used to effectively organize these results. Furthermore, in Web search, many search engine companies are interested in providing commercial services in response to the user queries, including targeted advertisement, product reviews, valued-added services such as banking and transportation, according to the categories. In these applications of Web search, query classification is very important. Instead of classifying queries directly, some previous work focuses on classifying search results as an alternative way of understanding queries [Chen and Dumais 2000].

However, there are several major difficulties which hinder the progress of query classification. First, many queries are short and query terms are noisy. As an example, in the KDDCUP2005 dataset, which is publicly available at <http://www.acm.org/kdd/kddcup>, the 800,000 queries vary a lot. They may be as simple as a single number such as "1939" or as complicated as a piece of programming code which involves more than 50 words. Figure 2 shows the statistical information about the number of words contained in each query and their frequencies in the 800,000 queries. From this figure we can see that queries containing 3 words are the most frequent (22%). Furthermore, 79% queries have no more than 4 words<sup>1</sup>. Each query is a combination of words, names of persons or locations, URLs, special acronyms, program segments and malicious codes. Some queries contain the words which are very clean while others may contain typos or meaningless strings which are totally noisy. Some words may just have their meanings as defined in a static dictionary while others may have some special meanings when used on the Internet.

A second difficulty of Web query classification is that a user query often has

<sup>1</sup>The observation is different from those given in [Silverstein et al. 1999; Jansen 2000] slightly since the 800,000 queries do not contain empty ones.

multiple meanings. For example, “Apple” can mean a kind of fruit or a computer company. In this case, query classification should provide both meanings ranked according to the likelihood of each meaning. For “Apple”, the result can be given in the form of “Computers\Hardware; Living\Food & Cooking”.

A third difficulty is that the meanings of queries may also evolve over time. For example, the word “podcast” is interpreted as a kind of Web audio blog site only recently. But such words cannot be found in a traditional static dictionary. The distribution of the meanings of this term is therefore a function of time on the Web. In order to preserve the existing meanings of words as well as finding out the new meanings, we cannot simply classify a query solely based on a static and out-of-date training set [Beitzel et al. 2005]. Instead, we should obtain the meanings of queries from the Web in an online manner. Our approach is to retrieve the most related documents for the query, and extract their semantic features. The distribution of the meanings on the Web should influence the ranking of the target categories, among other factors. Also, in order to obtain a better and unbiased understanding of each query, we should not rely on a single search engine but combine multiple results from different search engines.

In summary, an important issue is how to automatically and effectively classify a large number of queries that are inherently short, noisy and ambiguous, when there is a lack of clear definitions for these data (such as a dictionary) and a lack of additional training data (such as a labeled query log file). A major traditional approach in handling short and ambiguous queries is through query expansion using relevance feedback [Chang and Hsu 1998], whereby users provide information on which retrieved result is relevant to the query. However, this does not work for our problem, because in many Web search problems the results must be generated automatically and efficiently for the users. Another major method is query expansion by using a dictionary or thesaurus [Voorhees 1994], whereby the users’ query words are enlarged to contain additional information. However, in a Web search domain, many queries are newly created words and their intended meanings are a moving target. In some sense, generating an enough good thesaurus to handle the query expansion requires that we solve the query classification problem in the first place. According to our statistics, most queries are short; an illustration is shown for randomly selected queries in Figure 2, where the occurrence frequency of queries is shown against the number of words in each query. In addition, queries can be noisy as well, for there are many mis-spelt queries. A distribution of the meanings of queries is shown in Figure 3, where we plot the query count in percentage against the number of meanings that each query corresponds to, on the 111 randomly chosen validation samples from the KDDCUP2005 dataset. These 111 queries are labeled by human experts. As can be seen from Figure 3, many queries have more than three possible meanings. All of these characteristics indicate that we cannot solely rely on a static thesaurus to classify the queries.

Recently, some interesting works have emerged on using query logs to expand the meanings of user queries [Wen et al. 2002; Beeferman and Berger 2000]. However, such a method requires that there exists a query log for us to use, which contains a mapping from submitted queries to clicked Web pages. However, in general, such logs are not available for timely training and application of the query classification

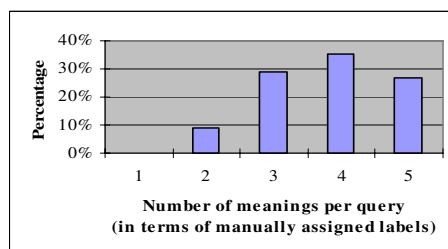


Fig. 3. Percentage of queries with different number of meanings (in terms of manually assigned labels).

models, as in the case for the KDDCUP2005 competition. Even when we can obtain such a log file, the mapping from queries to the clicked Web pages does not automatically provide the target categories either, because the pages themselves still need be mapped to different categories. In addition, such a mapping may be biased by the choice of any single search engine that collects the query logs, and thus the classification approach adopted may not be as general and objective as possible. In order to obtain an objective view on the categories that each query can be mapped to, views from multiple popular search engines should be consulted. Finally, obtaining up-to-date click-through data from multiple search engines raises serious privacy, legal and business issues, which are beyond the means of many practitioners, businesses and governments.

This paper presents a new approach to classifying large quantities of search queries automatically. Our approach is called query enrichment, which takes a short query and classifies it into target categories by making use of a set of intermediate objects. In our application, these intermediate objects are the Web pages and the category taxonomies such as that from the Open Directory Project (ODP)<sup>2</sup>. Query enrichment makes the following basic assumptions:

- Given that the Web is one of the largest data sources available, although we do not know the true meanings of many user queries, their intended meanings should be reflected by the Web as a whole. In addition, although a particular search engine may not fully reflect the intended meanings of the query for a user, when we combine many different search engines and Web directories, the meanings of the query embedded in the Web as a whole can be extracted. In other words, by searching the Web for answers, users have expressed their “trust” that their answer is somewhere located on the Web. By this assumption, our approach can be to first “enrich” a query by covering its potential meanings through Web search, and then classify the search results in a subsequent step to the target categories.
- A set of objects exist that can “cover” the target categories. An example of the intermediate objects is the ODP, which contains a collection of more than 170,000 different categories and can be taken as such an extensive taxonomy.
- We use results from search engines to provide the intermediate objects. Even

<sup>2</sup><http://dmoz.com>

though a particular search engine may be biased and therefore cannot guarantee high-quality answer to all the search queries from this huge collection, each search engine might provide a different view point in interpreting the user query. The result of different view points should be combined into a coherent whole for better coverage and higher robustness of the answers.

This last point means that we can submit the queries to multiple search engines, thus the intended answers are among the returned results (i.e., resultant Web pages). In this way, the relative “weights” of multiple search engines in generating collective answers to the queries can be learned from a validation dataset, rather than being predefined.

One possibility is to exploit a meta-search engine. Meta search engines [Howe and Dreilinger 1997; Selberg and Etzioni 1995] submit queries to multiple different search engines and integrate the results into a single list to be returned to the users. Our task is similar to the use of meta-search engines, except that in our case we are interested in query classification rather than search. Therefore, we will follow a similar approach, in which we submit the query to different search engines and classify the query based on the search results from each search engine. In a final step, we integrate the classification results corresponding to each search engine together. In machine learning, this is known as an ensemble of classifiers. Similar to meta-search engines, the ensemble of classifiers combines the results of different component classifiers using a collection of weights that are learned from a validation dataset. However, different from meta-search engines, the ensemble of classifiers is aimed at classifying a query using a collection of classifiers where each one can be somewhat biased. By integrating the classification results, they compliment each other, and the results are typically more robust as compared to any single classifier. General introductions of ensembles of classifiers are given in [Hansen and Salamon 1990; Bauer and Kohavi 1999; Caruana et al. 2004; Fan et al. 1999].

We call our approach query enrichment. In general, query enrichment consists of two major steps:

- First, we replace a query by a set of objects, where the meanings of the query are embedded in the set;
- Second, we then classify the query based on the set of objects into ranked target categories.

In order to enrich an input query, in our application, we submit the query to multiple Web search engines and obtain the search results as a ranked list. This ranked list consists of two types of objects: the resultant Web pages and the intermediate taxonomies such as the ODP. For example, a query “Apple” submitted to Google (<http://www.google.com>) returns a list of resultant Web pages that can be mapped to “Computers\System; Business\Food\_and\_Related\_Products” in the ODP directory. The Web pages and the set of categories in the intermediate taxonomies serve as the basis to map to the target categories.

The rest of the paper is as follows. In the next section, we give an overview to our approach. In Sections 3 and 4, we explain the two phases of our solution, respectively. Phase I corresponds to the training phase of machine learning algorithms. Phase II corresponds to the testing phase. In phase I, two kinds of classifiers are

developed as the base classifiers. One is synonym-based and the other is statistics-based. Phase II consists of two stages. In the first stage, the queries are enriched such that for each query, its related Web pages together with their category information (if they have) are collected through the use of search engines. In the second stage, the objects in the enriched result are classified through the base classifiers trained in phase I. Based on the classification results obtained by the base classifiers, we get the last classification results through an ensemble of classifiers. In Section 5, we describe our experimental results on the KDDCUP2005 tasks. We show that using our solutions, we can achieve superior performance as compared to other competitive methods and similar performance as compared to human labelers, when we appropriately integrate the search engines and combine the query classification results.

## 2. PROBLEM DEFINITION AND OVERALL APPROACH

The query classification problem is not as well-formed as other classification problems such as text classification. The difficulties include short and ambiguous queries and the lack of training data. In this section, we give a formal definition of the query classification problem, which is inspired by the tasks of KDDCUP2005 competition.

### Query Classification:

- \* *The aim of query classification is to classify a user query  $Q_i$  into a ranked list of  $n$  categories  $L_{i1}, L_{i2}, \dots, L_{in}$ , among a set of  $N$  categories  $\{L_1, L_2, \dots, L_N\}$ . Among the output,  $L_{i1}$  is ranked higher than  $L_{i2}$ , and  $L_{i2}$  is higher than  $L_{i3}$ , and so on.*
- \* *The queries are collected from real search engines submitted by Web users. The meanings and intension of the queries are subjective.*
- \* *The target categories are a tree with each node representing a category. The semantic meanings of each category are defined by the labels along the path from the root to the corresponding node.*

In addition, the training data must be found online because in general, labeled training data for query classification are very difficult to obtain.

Table II shows several examples of queries and their categories chosen from the test data of KDDCUP2005 competition. As we can see, for each query, it may have more than one categories and the categories are ordered, according to the probabilities that the query belongs to them.

To facilitate the understanding of the definition of query classification and the formal discussion, we make the following definitions.

**DEFINITION 1 (TARGET CATEGORIES).** *Target categories are the categories that we will classify user queries into.*

For example, for the KDDCUP2005 task, these are 67 categories provided by the organizers as the final targets of query classification (see Appendix A for details).

**DEFINITION 2 (INTERMEDIATE TAXONOMIES).** *Associated with a search engine or Web directory, there is often a taxonomy of categories. We wish to distinguish*

Table II. Examples of queries and their categories

Query	Categories
Aerosols	Information\Science & Technology Living\Health & Fitness Living\Family & Kids
cross pendant	Living\Gifts & Collectables Living\Fashion & Apparel Living\Religion & Belief Shopping\Stores & Products Shopping\Buying Guides & Researching
aberdeen police department	Information\Law & Politics Information\Local & Regional

between the target categories in DEFINITION 1 and the existing taxonomies on the Web. Thus we call these the “intermediate taxonomies” in this paper.

For example, the ODP defines a taxonomy that consists of more than hundreds of thousands categories organized in a tree structure.

DEFINITION 3 (QUERY ENRICHMENT FUNCTION). *The query enrichment function  $u$  is a function that maps from a query  $Q$  to a set of intermediate objects on the Web:  $u: Q \rightarrow$  intermediate objects. An example of an Intermediate object is a Web page. Another type of object is the category labels in the intermediate taxonomy such as ODP directory.*

DEFINITION 4 (QUERY CLASSIFICATION FUNCTION). *A query classification function  $f_{Q2C}$  maps from a user query  $Q$  to one or more of the target categories:  $f_{Q2C}: Q \rightarrow$  target categories.*

DEFINITION 5 (TEXT CLASSIFICATION FUNCTION). *A Text classification function  $h_{T2C}$  maps from a body of text  $T$  to one or more of the target categories:  $h_{T2C}: T \rightarrow$  target categories.*

DEFINITION 6 (INTERMEDIATE TAXONOMY TO TARGET CATEGORY MAPPING). *The Intermediate Taxonomy to Target category mapping function  $l_{IT2C}$  is a function that maps from a category in the intermediate taxonomy to one or more target categories.*

The query classification function  $f_{Q2C}$  (DEFINITION 4) can be constructed by one of two strategies. Strategy one corresponds to using text-based statistical classifiers: for each query  $Q$ ,

- We first map the query  $Q$  to Web pages  $T$ ;
- We then apply a text classification function  $h_{T2C}$  to map  $T$  to the target categories.

Alternatively, we can build a synonym-based classification function as follows:

- We first submit the query  $Q$  to search engines and obtain category labels of some intermediate taxonomy (these categories are different from the target categories).
- We then map the intermediate categories thus obtained to the target categories.



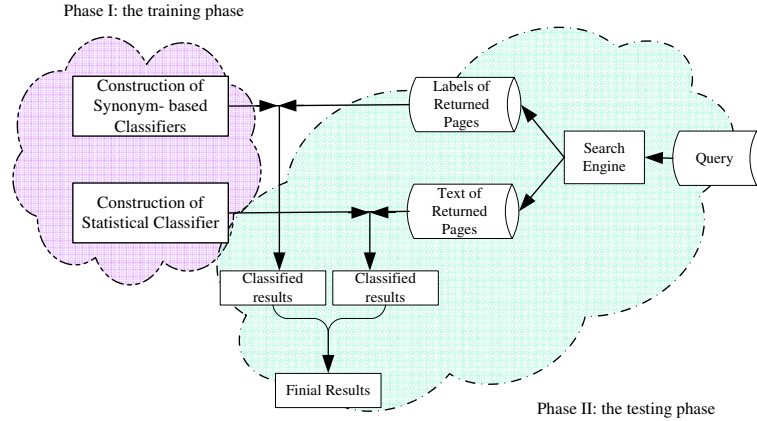


Fig. 4. The architecture of our approach.

By combining the above two strategies, we can obtain the ensemble-based approach.

As in typical machine learning applications, we also adopt two phases in our solution. In phase I, which corresponds to the training phase of machine learning algorithms, we collect data from the Web for training classifiers that can be used to classify intermediate objects to the target categories. In phase II, which corresponds to the testing phase in machine learning research, we apply the classification functions thus built to the target categories. The classifiers learned in Phase I are applied to classify the queries. The overall architecture of our approach is shown in Figure 4 and the detailed architectures of the two phases are shown in Figure 5.

In order to make our approach easy to follow, we will take the KDDCUP2005 task as an example when describing the approach.

### 3. PHASE I: CLASSIFIER TRAINING

We now discuss phase I of our approach in detail. In this phase, we train classifiers for mapping from intermediate objects into the target categories. As noted, a main problem here is the lack of training data, a difficulty which makes many previous machine learning methods unable to be applied. The objective in phase I is to collect the data from the Web that can be used to train classification functions.

#### 3.1 Synonym-based Mapping from Intermediate Taxonomies to Target Categories

We first discuss how to construct mapping functions  $l_{IT2C}$  from intermediate categories to the target categories via synonym-based mapping functions. The taxonomy from different search engines can be different from each other, and from that of the target categories. The mapping function  $l_{IT2C}$  can be built by synonym-based keyword matching. In this approach, we compare the terms used to describe the categories in the target categories with those in the intermediate taxonomies.

Consider two categories,  $c_1$  and  $c_2$ , where  $c_1$  is a category label from the target categories and  $c_2$  is a label from an intermediate taxonomy. If they share some same keywords, we can map  $c_2$  to  $c_1$  directly. For example, in the KDDCUP2005 task, the target categories have two levels, such as, “Computers\Hardware”. The

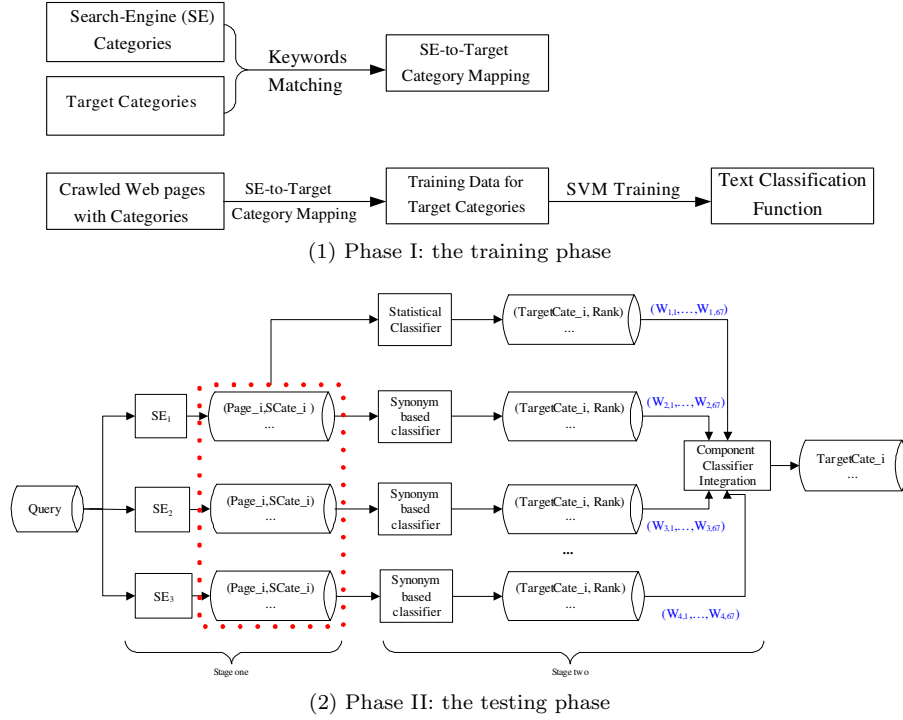


Fig. 5. The two phases: training and applying classifiers.

second level specifies a particular field within the first level. For most of the target categories, we only consider the keywords in the second level because they do not confuse with other categories. A typical example is “Computers\Internet & Intranet”. Even if we do not consider the first level “Computers”, there are no other categories which may be confused with “Internet & Intranet”. Therefore, if a category from an intermediate taxonomy contains “Internet” or “Intranet”, we will map it to “Computers\Internet & Intranet”. However, there are many categories that are more difficult to deal with. For them we require that the keywords in the first level and the second level match the categories in the intermediate taxonomy simultaneously. Otherwise, we cannot distinguish between two categories that share the same keywords only in the second level, such as “Computers\Hardware” and “Living\Tools & Hardware”. Although both of them have the keywords “Hardware” in the second level, they belong to two different domains “Computers” and “Living” as defined by the first level.

In order to improve the coverage of the mapping function  $l_{IT2C}$ , we can extend the keywords in the category names to include both the singular and plural forms in advance. For example “Living\Book & Magazine” is extended to “Living\Book & Magazine & Books & Magazines”.

After applying the above procedure, we may still miss a large number of mappings. Many categories in intermediate taxonomy do not occur in the target categories, although they share words with the same meanings. In response, we expand

the keywords in each label in the target categories through WordNet [Miller et al. 1990]. For example, the keyword “Hardware” is extended to “Hardware & Devices & Equipments” and the keyword “Movies” is extended to “Movies & Films”.

### 3.2 Text Data Collection and Statistical Classifier Training

We now discuss how to build the statistical classification function  $h_{T2C}$  to map from a body of text  $T$  to a target category. As we will discuss in section 4.2, the synonym-based classifiers have a certain drawback - they have low recall. In order to address this problem, we consider the statistical classifiers to help classify queries. A statistical classifier classifies a query based on the semantic content of a text, which can provide better recall as well as precision. Even when the synonym-based classifier fails to give any meaningful mapping for a query, the query can still be classified by a statistical classifier. Any kind of statistical text classifiers, such as Naive Bayes [McCallum and Nigam 1998], KNN [Yang 1999], Support Vector Machine (SVM) [Joachims 1998; 1999], can be applied.

To construct a statistical classifier, the first step is to collect the training data. This step is nontrivial since no training data are provided explicitly as we stated. In order to collect the training data, we use the following methods:

- First, we collect Web pages from some online manually labeled Web page directories, such as the ODP.
- By applying function  $l_{IT2C}$ , we can map the collected Web page into the target categories. Thus the mapped Web pages can be used as the training document for the target categories.
- The training data among the target categories are usually extremely unbalanced. In order to remove the impact of unbalanced distributions, as well as speed up the training step, we need to sample the training documents.

After the training examples are collected for each target category, we can follow the traditional procedure to train statistical classifiers, including some proper preprocessing steps and parameter tuning.

To make the above procedure clearer, we illustrate it through our solution to the KDDCUP2005 task. We use the SVM classifier because of its high generalization performance when used for text classification task and its easy implementation. About 1,540,000 Web pages are collected from ODP in total. After applying the mapping function between the ODP and the KDDCUP2005 categories, only 500,000 Web pages fall in the KDDCUP2005 categories. To address the unbalanced distribution of Web pages among the 67 categories, we randomly selected 15,000 Web pages from the categories with more than 15,000 Web pages and keep all the pages for the categories with less 15,000 Web pages. The Document Frequency (DF) and the Information Gain (IG) methods are used for feature selection [Yang and Pedersen 1997]. Then we use the  $SVM^{light}$  software package (<http://svmlight.joachims.org/>) to train a SVM. A linear kernel is used and the one-against-rest approach is applied for the multi-class case.

## 4. PHASE II: QUERY CLASSIFICATION

Phase I of our algorithm is designed to collect data for training the mapping functions which include the synonym-based and statistics-based classifiers. Phase II of

our algorithm is devoted to classifying a query to one or more target categories based on the classifiers.

Phase II is conducted through two stages. The first stage is to enrich the queries by searching related pages which can specify the meanings of the queries more accurately. Enrichment is necessary because the queries are rather short, and their meanings ambiguous. We perform the enrichment process by finding the relevant text from related Web pages as well as the category information of these Web pages (if they have) through Web search.

The second stage is to classify the queries based on the data collected in the first stage and the classifiers trained in phase I. In this stage, we take the two kinds of classifiers with totally different mechanisms trained in phase I as the base classifiers and develop several ensemble classifiers with different ensemble strategies to classify the queries. The experimental results show that the ensemble classifiers can improve the classification performance significantly.

#### 4.1 Query Enrichment

Query enrichment is a key step because our goal is to classify short and ambiguous queries without any additional descriptions about these queries. After this step, two kinds of information for each query are collected. One is the list of Web pages related to the target query. The other is the set of categories corresponding to the related pages. These two kinds of information will be leveraged by the two kinds of classifiers trained in phase I respectively.

In our approach, we send each query into multiple search engines which can provide options for both directory search and Web search. Directory search in a search engine refers to search algorithms that return the related pages of a query together with the pages' categories. Since these categories of Web pages are labeled by human labelers, it is appropriate to use them to classify the queries. However, not all the pages indexed by the search algorithm contain category information; in this case, Web search can return more related pages than Directory search. Based on the contents of the pages returned by Web search, we can classify the queries using a text classification algorithm.

In summary, to enrich a query through search engines, we use the following three steps:

- (1) We first try to get the related pages through "Directory Search";
- (2) If we cannot get any results from step 1, we try "Web Search";
- (3) If we still cannot get any result, the query must be too noisy or totally meaningless. Thus, we use some preprocessing approaches to clean up these queries and then resubmit them to "Directory Search" and "Web Search" in turn, as done in Steps 1 and 2. If after this step, there is still no result returned, the query will not be processed anymore, and no classification results are generated. Currently, two preprocessing approaches are employed for the cleaning-up, which will be described in detailed in the next example.

In our solution to KDDCUP2005 task, we use three search engines, including Google, Looksmart and a search engine developed by ourselves based on Lemur<sup>3</sup>.

<sup>3</sup><http://www.lemurproject.org/>

Now let us use Google as an example to illustrate the three steps in detail.

- At first, all 800,000 queries are preprocessed by removing special characters such as “, #, %” while keeping letters and digits. Then these 800,000 queries are sent into Google Directory Search. We are able to retrieve related pages for about 500,000 (63%) of all queries in this way.
- We then send the remaining 300,000 queries into Google Web Search and get results for additional 200,000 queries.
- For the remaining 100,000 queries, we conduct further preprocessing. We use the function of “Did you mean”, provided by Google to trigger the search for the most relevant queries to the original ones. For example, given the query “a cantamoeba”, neither Google “Directory Search” nor “Web Search” returned any results. However, by trying “Did you mean”, Google can return the results for the word “acanthamoeba” which is related to health, disease and medicine. In this way, we can get the results for another set of 60,000 queries from Google “Directory Search” or “Web Search”.
- However, after this step, there are still 40,000 queries left without any results. These queries are very noisy. They are usually connected words without spaces, long meaningless clobbers, or URL addresses containing parameters or even malicious codes. For these queries, we try to convert them into meaningful ones by extracting words from them through a maximum length matching method based on the dictionary WordNet. This method tries to extract as many as possible meaningful words and for the extracted words, they should be as long as possible. Take the query “wheelbearingproblemsdamage” as an example, Google cannot return any results through either “Directory Search”, “Web Search” or even “Did you mean”. Therefore, we can split the whole query into four meaningful words “wheel bearing problems damage”. After that, we can get reasonable results from Google “Directory Search” or “Web Search”. In this way, we can get the results for 30,000 out of the remaining 40,000 noisy queries.
- The remaining 10,000 queries cannot receive any pages from Google as answers. These queries are inherently noisy and meaningless, such as “dddddfdfdfdfdf”. Therefore, our classifier will not return any answers for these queries although we can assign labels to them according to the prior distribution of the target category. This potentially reduced recall (one measurement of the classification performance) for the query classifier. But because they only correspond to a tiny portion of all the queries (1.25%), they do not affect the final classification results much.

We follow the same steps when using Looksmart. Among the 800,000 queries, about 200,000 queries have “Directory Search” results. 400,000 have “Web Search” results and the remaining 200,000 have not any results.

The third engine we used is developed by ourselves based on Lemur. We first crawled more than one million Web pages with the category information from ODP. Then we indexed the collection of these pages with Lemur. Given a query, Lemur can retrieve a number of pages which are most relevant to the query together with their corresponding categories. Therefore the function of this search engine based on Lemur and the ODP data is similar to the “Directory Search” provided by

Google and Looksmart. Using this search engine, we can retrieve the related pages for most of the 800,000 queries except 35,000 queries.

In summary, after enriching the queries through a search engine, we can obtain two lists for each query. One is the returned Web pages list from a search engine. The other is a category list attached to the Web pages in the Web page list. Note that some Web pages have no category information. These two lists will be leveraged by different kinds of classifiers individually.

## 4.2 Query Classification using Synonym-based Mapping Function

Using the synonym-based category mapping functions discussed in Section 3.1, we can now build a query-to-target-category mapping function  $f_{Q2C}$ . In particular, for each query  $Q$ , through function  $u$  constructed in Section 4.1, we can obtain a set of intermediate categories  $S$  of the related Web pages returned by a certain search engine. Then we apply the category-mapping function  $l_{IT2C}$  to  $S$ , which returns an ordered list of target categories.

Using different search engines, we might obtain different intermediate categories. For each category of a certain search engine's taxonomy, we can construct a mapping function  $l_{IT2C}$  between it and target categories as shown in Section 3.1. After obtaining these mapping functions, we can perform query classification based on the intermediate categories  $S$  returned by a search engine. We then map the intermediate categories to the target categories using the constructed mapping functions. We keep track of the number of times for each target category being mapped onto. Then we can obtain the target categories together with their occurrence frequency. By ranking the categories in terms of the occurrence frequency, we get a ranked list of final target categories into which the query can be classified. Based on different intermediate category lists and their corresponding mapping function, we obtain different classification results. The classification functions thus obtained are known as synonym-based classifiers.

Based on the above mentioned approach, the synonym-based classifiers tend to produce results with high precision but low recall. They produce high precision because the synonym-based classifiers are based on the manually annotated Web pages and can utilize the classification knowledge of the human editors. Therefore, once a mapping function is constructed, the classification result is highly probable to be correct. For example, we have shown that the categories in the intermediate taxonomy such as "... \Computers\... \ Hardware\..." are mapped to the target category "Computers\Hardware". Therefore, once the Web pages associated with a query fall into the category "Computers\Hardware\Storage\Hard Drives", we can assign "Computers\Hardware" to the query with high confidence. They produce low recall because it is hard to find out all the mappings from the intermediate taxonomy to the target categories by keyword mapping. For example, about 80,000 out of 354,000 categories in Google's taxonomy are not mapped onto target categories. Therefore, we cannot map all the intermediate categories in the category list for a query to the 67 target categories and may miss some correct categories for the query.

Another reason for the low-recall problem is that a search engine may return only a few or even no Web pages with categories. The synonym-based classifier may fail because of the search results shortage problem. Therefore, we need to construct

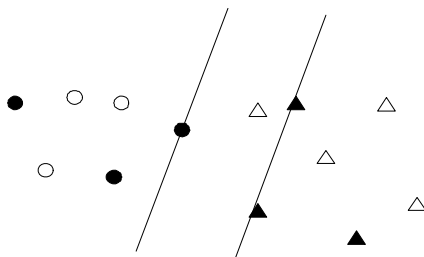


Fig. 6. Illustration of the advantage of statistical classifiers.

other classifiers to help handle the low-recall problem which will be described in the next section.

### 4.3 Query Classification Using Statistical Classifiers

In this section we discuss the use of statistical classifiers to classify queries. As shown in Section 4.1, after submitting the query  $Q$  to a search engine, we get a list of Web pages. Then, we can extract a body of text from the Web pages to capture the occurrence context of the issued query which can help determine its potential meanings. To accomplish this, we keep the top  $N$  results in the returned list (the parameter  $N$  will be studied in Section 5) and use the aggregate terms of the corresponding snippets, titles and URLs terms to represent the query. The query's bag of terms will be processed by stop-word removing, stemming and feature selection. The resultant term vector can be used as input of the statistical classifiers and a ranked list of categories for each query will be produced.

Statistical classifiers are expected to achieve higher recall compared to the synonym-based classifiers. The reason is illustrated in Figure 6. For simplicity, only two categories are considered. The circles and triangles represent the samples belonging to two categories respectively. The black symbols represent the training data we collect for the two categories by the way shown in Section 3.2. The white symbols represent the Web pages that we crawled from the ODP which should be mapped to the target categories but they are not, because of the drawbacks of the mapping function  $l_{IT2C}$  from the intermediate taxonomy to the target categories. Given a query  $Q$ , after the query enrichment step, we may happen to get only the white circles or triangles to represent the query. Thus it is clear that we cannot judge the labels of the query by the synonym-based classifiers since there is no mapping relationship between the white circles and triangles to the target categories. However, the statistical classifiers can still obtain the separating hyperplane based on the black circles and triangles which can be used to classify the white circles and triangles, thus can further be used for classifying the query  $Q$ .

### 4.4 Ensemble of Classifiers

From the previous two approaches, we can build various classifiers independently which can classify the input queries into the target categories. These approaches are based on different mechanisms and can be complementary to each other. Previous

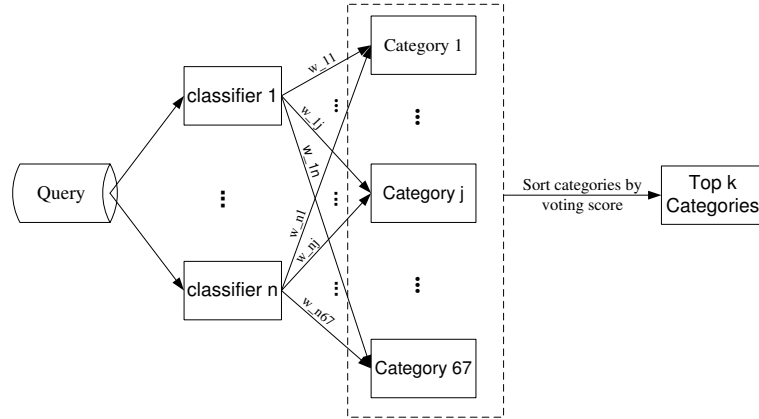


Fig. 7. Illustration of the ensemble classifiers.

works have shown that a proper combination of different base classifiers can improve the final classification performance [Hansen and Salamon 1990; Kittler et al. 1998; Bauer and Kohavi 1999; Cann 2003; Fan et al. 1999]. In this section, we consider how to combine them together.

[Dietterich 2000; Alpaydin 2004] have categorized different ensemble based methods. Of these, voting, bagging and boosting are the major ensemble methods. In voting, which defines a linear combination of existing classifiers, the main task is to decide the weights used in the combination. In bagging, the base classifiers are generated to be different by training them over slightly different training data that are sampled from the given training set by bootstrap. In boosting, the base classifiers are generated sequentially where each one is trained on the mistakes of the previous classifiers. Both bagging and boosting require that there are sufficient labeled data for generating base classifiers [Meyer and Brown 1998; Dietterich 2000]. Because in query classification, labeled data are very scarce (in our case, there are only 111 of them, which correspond to 0.014% of all the data), voting becomes the most natural choice. We thus focus on how to set the weights of different classifiers.

Figure 7 illustrates our ensemble-based method in query classification. As we can see from the figure, the results of the classifiers are linearly combined to generate the final category ranking. In our approach, we consider two ways to assign the combination weights for different classifiers. The first is to make use of the validation dataset involving a small number of queries labeled by human labelers. Considering that the validation dataset is too small and it is easy to be overfitting, a second way is to ignore this validation dataset and give each base classifier equal weight. More details about the first way are shown below.

As we can imagine, different classifiers introduced in the above sections have different performance. Some may work better than others on certain categories. For example, a classifier may achieve high precision on one category while having high recall on the other category. This indicates that it is not proper to assign a single weight to a classifier. Instead, we should differentiate the weight of a classifier on different categories according to its performance. To determine the



weights, we validate each base classifier on the validation samples. The higher precision a classifier achieves on a category, the higher weight is assigned to the classifier on that category. As a result, each classifier may obtain a weight value  $W_{ij}$  on a target category  $j$ .  $W_{ij}$  is defined by:

$$W_{ij} = \frac{p_{ij}}{\sum_{k=1..n} p_{kj}}$$

where  $p_{ij}$  is the precision of classifier  $i$  on categories  $j$ . The definition of precision will be given in Section 5.2.

Three more ways to determine the combination weights similar to the above will also be discussed and tested in Section 5.

## 5. EXPERIMENTS AND DISCUSSIONS

To test our proposed approach, we conduct extensive experiments on the KDD-CUP2005 datasets. The experimental results validate the effectiveness of our approach. In addition, we give the analysis of the consistency of the three labelers on their judgments of the performance of the classifiers in this section.

As introduced in the previous section, we have six classifiers in total up to now, three synonym-based classifiers, one statistical classifier that is SVM and two ensemble classifiers. For simplicity, we refer to the three synonym-based classifiers that rely on Google, Looksmart, and Lemur as S1, S2 and S3 respectively. We denote the ensemble classifier that relies on the validation dataset as EDP (since it assigns different weights for each base classifier on Different categories according to its Precision on the category) and the one that does not use validation data as EN.

### 5.1 Datasets

One of the KDDCUP2005 datasets contains 111 sample queries together with the human labeled categories. These samples help exemplify the format of the queries and provide the semantics of a tiny number of queries. In fact, since the category information of these queries is truthful, they can serve as the validation data for our proposed classifiers. Another dataset provided by the organizers contains 800,000 queries in total which are selected from the MSN search logs to test the submitted solutions. Since manually labeling all the 800,000 queries is too expensive and time consuming, the organizers randomly selected 800 queries at last and invited three human labelers to label them. We denote the three labelers (and sometimes the datasets labeled by them if no confusion is caused) as L1, L2 and L3, respectively. We refer to the former as *sample dataset* and the latter as *testing dataset* in the following sections. Both of these datasets can be used to evaluate the different classification approaches. The *sample dataset* can be used to determine the ensemble weights as well in this paper.

### 5.2 Evaluation Criteria

The evaluation criteria adopted by the KDDCUP2005 organizers is the standard measures for evaluating the performance of classification in Information Retrieval (IR), including precision, recall and F1-measure [van 1979]. The definitions of

precision, recall and F1 in the query classification context are given below:

$$\begin{aligned}
 A &: \sum_i \# \text{ of queries are correctly tagged as } c_i \\
 B &: \sum_i \# \text{ of queries are tagged as } c_i \\
 C &: \sum_i \# \text{ of queries whose category is } c_i \\
 \text{Precision} &= \frac{A}{B} \\
 \text{Recall} &= \frac{A}{C} \\
 \text{F1} &= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}
 \end{aligned}$$

For the sample dataset, we report the precision, recall and F1 evaluation results for each classifier. For the testing dataset, since three labelers are asked to label the queries, the results reported on this dataset are the average values [Li et al. 2005]. Take the calculation of F1 as an example:

$$\text{Overall F1} = \frac{1}{3} \sum_{i=1}^3 (\text{F1 against human labeler } i)$$

For some of the ensemble classifiers, we need to know the performance of a classifier on a certain category. It is easy to define such criteria according to the above definition. For example, the precision of classifier  $i$  on category  $j$  could be defined as:

$$P_{ij} = \frac{\# \text{ of queries are correctly tagged as } c_j \text{ by classifier } i}{\# \text{ of queries are tagged as } c_j \text{ by classifier } i}$$

Because there are 67 target categories, a random classification algorithm would give a precision of  $1/67 = 1.5\%$  if one category is to be returned, and give a precision of  $5/67 = 7.5\%$  if five categories are returned. Therefore, we are comparing against a baseline of 7.5% for the KDDCUP2005 task, whose requirement is to return at most five results.

### 5.3 Quality of the Testing Dataset

Because the testing dataset is provided by three human labelers, we wish to evaluate the quality of this dataset. In particular, in this section, we analyze the consistency between the three labelers in the testing dataset.

Table III gives the distribution of the number of categories assigned by the human labelers to each query, that is, how many queries are assigned  $N$  categories where  $N$  changes from 1 to 5. The ‘‘Ave’’ row shows the average number of categories for each query. From the table, it seems that the three labelers disagree with each other very much. However, from Figure 8 which shows the distribution of the 67 categories assigned by the three labelers to the 800 testing queries, we can see that the distributions for the three human labelers are very similar. From Table III and Figure 8 we can conclude that the general distributions of categories are very similar. Though the labelers, L1 and L3 tend to assign more categories for each query than L2.

Figure 9 shows the F1 values of the six classifiers on the testing data labeled by the three labelers. From Figure 9, we can see that the three labelers have a high

Table III. The number of queries with  $N$  labels.

$N$	L1	L2	L3
1	14	118	16
2	138	365	79
3	186	225	212
4	222	71	199
5	240	21	294
Ave	3.67	2.39	3.845

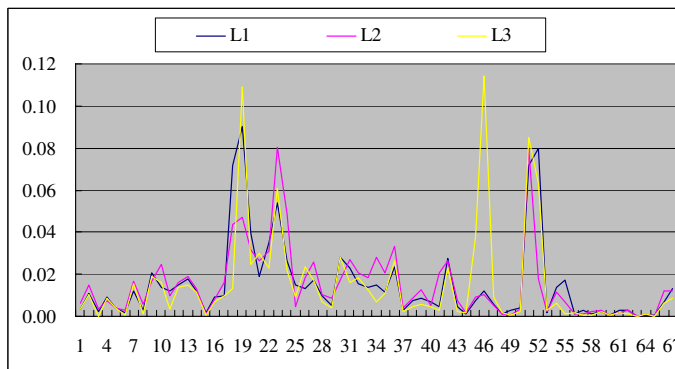


Fig. 8. The distribution of the labels assigned by the three labelers.

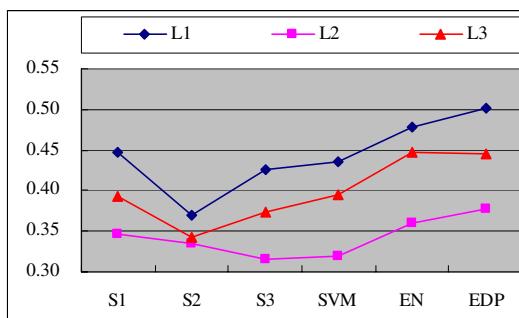


Fig. 9. Performance of different classifiers evaluated by different Labels.

correlation with respect to the relative performance of the classifiers, especially labeler 1 and labeler 3. After ranking the six classifiers according to each labeler, we calculate the Spearman rank-order correlation coefficient between each pair of the labelers [Hoel 1966]. Spearman correlation is a nonparametric approach to calculating the relationship between two variables based on ranking the two variables and no assumption about the distribution of the values is made. The results are shown in Table IV.

In general terms, correlation coefficients over 0.67 indicate strong relationships [Cann 2003]. So we can conclude that the three labelers are highly correlated

Table IV. Spearman correlation between each pair of labelers.

L1.vs.L2	L1.vs.L3	L2.vs.L3
0.829	0.943	0.771

Table V. Performance of each labeler against another labeler.

## (1) Precision

	L1	L2	L3
L1	1.000	0.637	0.561
L2	0.415	1.000	0.367
L3	0.588	0.590	1.000

## (2) Recall

	L1	L2	L3
L1	1.000	0.415	0.588
L2	0.637	1.000	0.590
L3	0.561	0.367	1.000

## (3) F1

	L1	L2	L3
L1	1.000	0.502	0.574
L2	0.502	1.000	0.452
L3	0.574	0.452	1.000

when they determine the performance of classifiers.

Besides the above correlation analysis, we also investigate the performance of each labeler when we take the other two labelers as the ground truth. The results are shown in Table V. The labelers in the columns are tested against the labelers in the rows which are taken as the ground truth. The average F1 value of the labelers is 0.509. Therefore, we can conclude that the query classification problem is not easy and the performance of our proposed ensemble classifier (with F1 equal to 0.444) is close to the average performance of the three human labelers.

#### 5.4 Experimental Results and Explanation

In the following part, we first investigate the two main parameters which affect the performance of our proposed classifiers on the two datasets and then we compare the performance of these classifiers. After that, we investigate several other ensemble strategies besides the two we introduced in section 4.4. Finally, we compare our approaches with other participants as well. Table VI shows a summary of the classifiers under comparison including the six classifiers we have introduced and three more ensemble classifiers we will introduce in section 5.4.3.

**5.4.1 Effect of Parameter Tuning.** There are two main parameters that significantly impact the performance of our proposed classifiers. The first parameter is the number of result pages returned by search engines, that we should use for enriching each query. If we use too few pages, we may fail to cover the diverse topics of a query. However if we use too many pages, we may introduce much noise. Figure 10 shows the performance of different classifiers with respect to an increasing

Table VI. A summary of the classifiers.

Symbols of classifiers	Meaning of the classifier
S1	The synonym-based classifier based on Google
S2	The synonym-based classifier based on Looksmart
S3	The synonym-based classifier based on Lemur
EDP	The <i>Ensemble</i> classifier in which the weights for each base classifier on <i>Different</i> categories are set in proportion to its <i>Precision</i> on the category when tested on the validation dataset
EDF	The <i>Ensemble</i> classifier in which the weights for eachbase classifier on <i>Different</i> categories are set in proportion to its <i>F1</i> on the category when tested on the validation dataset
EP	The <i>Ensemble</i> classifier in which each base classifier is assigned a single weight in proportion to its overall <i>Precision</i> when tested on the validation dataset
EF	The <i>Ensemble</i> classifier in which each base classifier is assigned a single weight in proportion to its overall <i>F1</i> when tested on the validation dataset
EN	The <i>Ensemble</i> classifier in which each base classifier is equally weighted and do <i>Not</i> rely on the validation dataset

number of related pages used for classification. Here, the related pages are taken into consideration in the order they are returned by the search engines; That is in the order of the degree of relevance with the query. The results shown in Figure 10 verify our conjecture. As the number of related pages increases, the precision increases at first and then tends to decrease. The reason is that we need a certain amount of pages to get the meanings of the query. However, if we include too many pages, noise may be introduced which can reduce the precision. From Figure 10, we can see that the critical point for most classifiers is 40 pages. Before that point, the performance of the classifiers increases when we use more pages while after that point, the performance begins to decrease. Therefore in the following experiments we just keep the top 40 pages for query classification.

Another parameter is the number of labels we should assign to each query. The KDDCUP2005 competition rule allows us to assign at most 5 labels for each query. However, in order to achieve higher precision, we should assign few but accurate labels. If we hope to achieve higher recall, we need to assign more possible labels. Figure 11 shows the performance of different classifiers by varying the number of classified categories. As we expected, when the number of classified categories increases, the precision of all the classifiers decreases while the recall increases significantly. In contrast, the value of F1 increases at the beginning and then decreases. For most of the classifiers, the maximum values of F1 are achieved when four categories are generated for each query. Although the F1 values are close to those obtained when five categories are assigned, the precision values are much lower. We also adopted some heuristic rules to decide the number of classified categories. As shown before, for each query, our classifiers can return a ranked

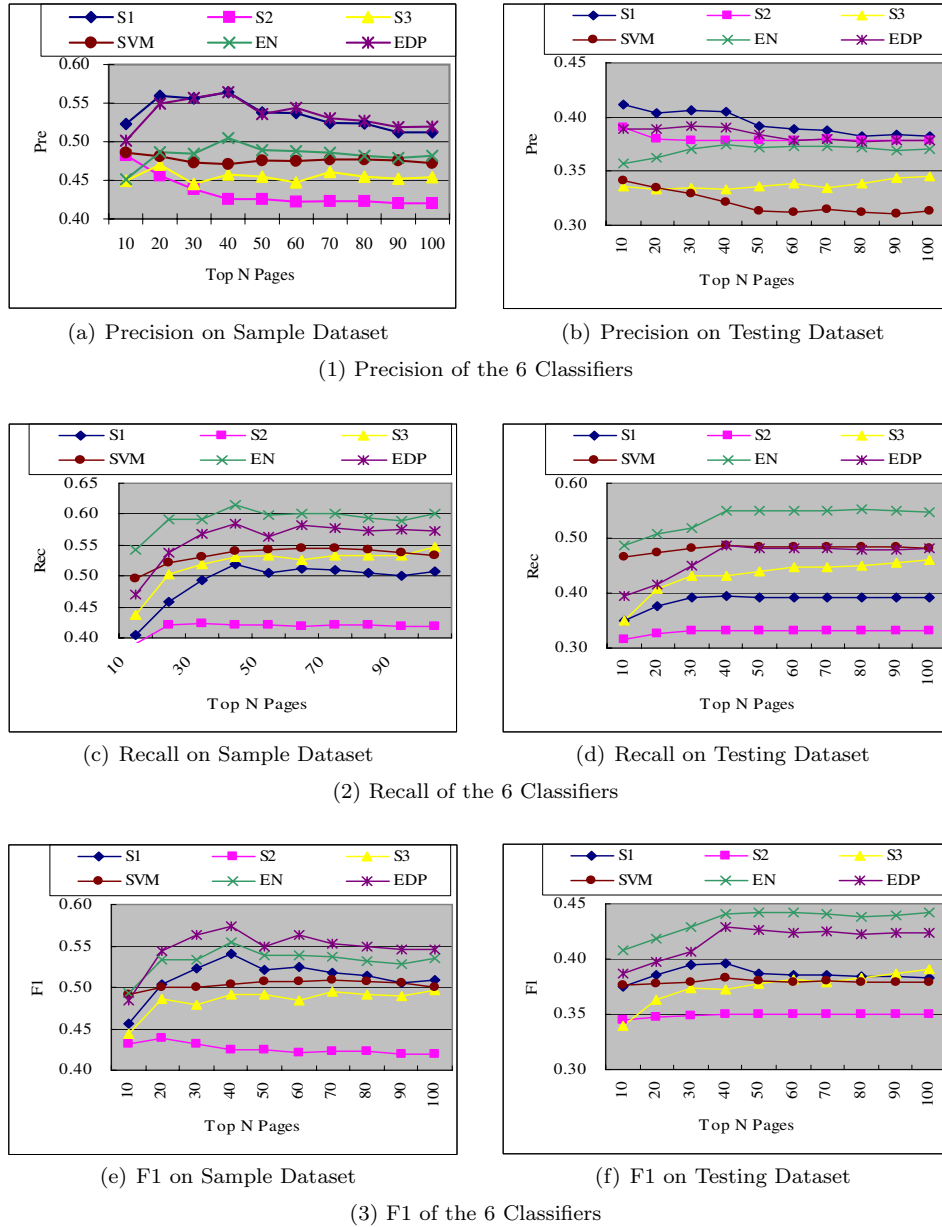
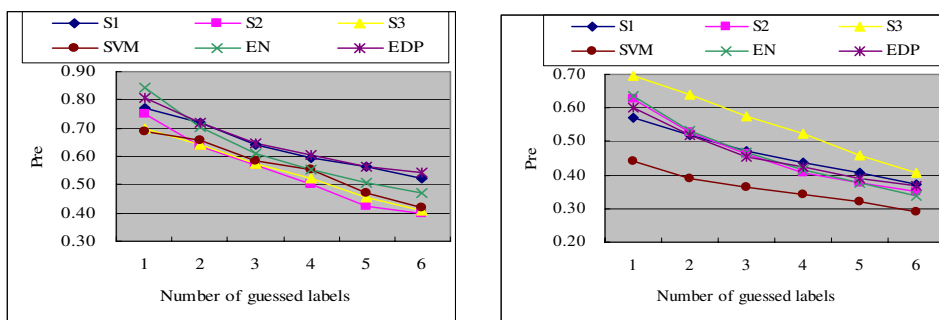


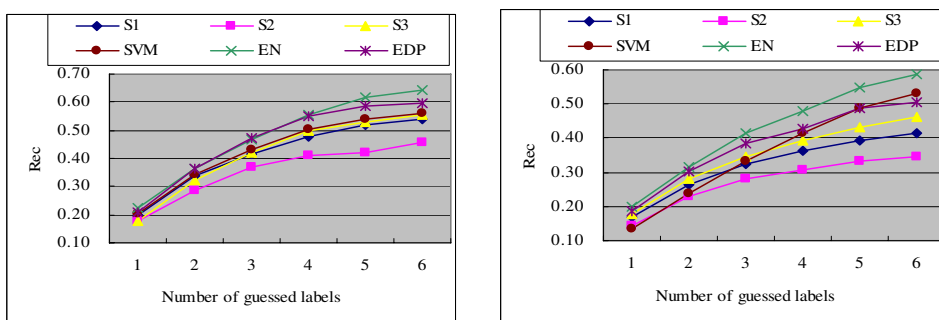
Fig. 10. Performances of different classifiers vary with the number of used related pages on the two datasets.

category list according to a certain criterion which is referred to as confidence. Among the top 5 categories, if the confidences of two neighboring categories  $c_i$  and  $c_{i+1}$  vary too much, we stop at  $c_i$  and discard the categories after  $c_i$ . This heuristic rule can help us to some extent, but not very much.



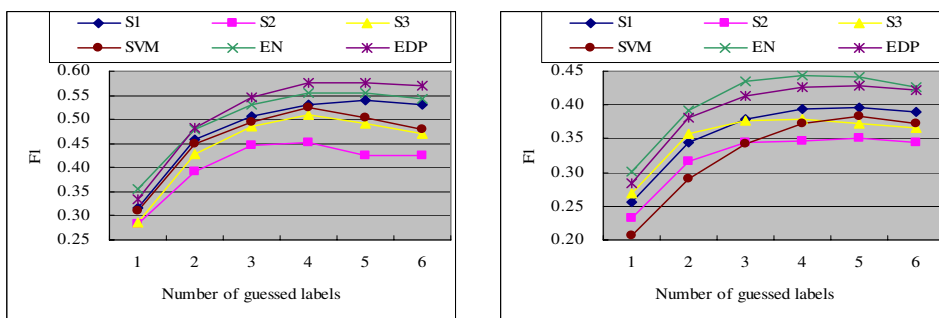
(a) Precision on Sample Dataset (b) Precision on Testing Dataset

(1) Precision of the 6 Classifiers



(c) Recall on Sample Dataset (d) Recall on Testing Dataset

(2) Recall of the 6 Classifiers



(e) F1 on Sample Dataset (f) F1 on Testing Dataset

(3) F1 of the 6 Classifiers

Fig. 11. Performances of different classifiers vary with the number of classified categories on the two datasets.

5.4.2 *Comparison between Classifiers.* From the above experimental results on both the sample dataset and the testing dataset, we can see that among the four base classifiers, S1 works best while S2 works worst. The lack of overlap among the search results from different search engines explains the performance difference among S1, S2 and S3. The reason why S2 does not work well is that for many queries, we cannot obtain enough related pages through the Looksmart search engine. For SVM, we expect that it can solve the low-recall problem caused by the three synonym-based classifiers, as discussed in section 4. Figure 10 and Figure 11 show that SVM does obtain the highest recall in most cases as compared to the synonym-based classifiers. We also notice that the two ensemble classifiers can achieve better performance in terms of F1 than any other base classifier. For the peak F1 values of the three best classifiers on the testing dataset (EN, EDP and S1), we can see that, compared with S1, EN and EDP improve the F1 by 12.1% and 8.3% relatively. In fact, when we design these two ensemble classifiers, EDP is expected to achieve higher precision because each component classifier is highly weighted on the categories where it achieves high precision and EN is expected to achieve higher F1 performance since the recall is relatively high. According to our two submitted results, the F1 value of EN (0.444) achieves 4.2% improvement relatively compared with the F1 value of EDP (0.426). While the precision value of EDP (0.424) improves by 2.3% compared with that of EN (0.414).

5.4.3 *Study of Different Ensemble strategies and the Effect of Validation Dataset.* As shown above, we employed two strategies to decide the weights for ensemble classifiers. In this section, we will study three more strategies. Besides that, we also test the effect of the validation data. The strategies we will study are shown below. For clarity, the strategy EDP introduced in section 4.4 is repeated here.

- The weight for classifier  $i$  on category  $j$  is in proportion to  $P_{ij}$ .  $P_{ij}$  refers to the precision of classifier  $i$  on category  $j$ . We denote this strategy as EDP.
- The weight for classifier  $i$  on category  $j$  is in proportion to  $F1_{ij}$ .  $F1_{ij}$  refers to the F1 value of classifier  $i$  on category  $j$  where  $F1_{ij}$  can be defined in the similar way as  $P_{ij}$ . We denote it as EDF.
- Each classifier is assigned a single weight in proportion to its overall precision. We denote it as EP.
- Each classifier is assigned a single weight in proportion to its overall F1. We denote it as EF.

EF and EP are different from EDP and EDF in that we assign a unique weight for a base classifier across all categories in EP and EF while we assign each base classifier different weights on different categories in EDP and EDF.

For the four strategies to determine the weights, we rely on the validation data (the 111 samples). The size of the validation data may affect the performance of different strategies. To test the effect, we construct five subsets by randomly picking out 20%, 40%, 60%, 80% and 100% samples. Then we construct the ensemble classifiers based on these subsets respectively. In order to remove the variance, we conduct the experiments three times. The reported values are the averaged values of the three runs.



Table VII. Comparison between different ensemble strategies.

(1) Precision

	EDP	EDF	EP	EF
20%	0.375	0.364	0.389	0.385EF
40%	0.381	0.372	0.388	0.387EF
60%	0.386	0.376	0.387	0.386EF
80%	0.388	0.378	0.387	0.386EF
100%	0.390	0.380	0.387	0.386EF

(2) Recall

	EDP	EDF	EP	EF
20%	0.419	0.422	0.503	0.500
40%	0.453	0.456	0.500	0.500
60%	0.478	0.482	0.500	0.500
80%	0.483	0.488	0.500	0.500
100%	0.488	0.492	0.500	0.500

(3) F1

	EDP	EDF	EP	EF
20%	0.393	0.388	0.435	0.431
40%	0.410	0.406	0.433	0.432
60%	0.424	0.419	0.433	0.432
80%	0.425	0.422	0.433	0.432
100%	0.429	0.425	0.433	0.432

From Table VII we can see that, with the increment of the size of the validation dataset, the performance of EDP and EDF increases steadily. When the validation dataset is small, the performance of a classifier on a certain category may not reflect the real performance of the classifier. Therefore, the weights generated tend to be unreliable and cause the bad performance of the ensemble classifiers. However, by increasing the size of the validation dataset, the performance of the base classifiers becomes more and more reliable, so do the weights we obtain. Consequently, the performance of the ensemble classifiers becomes better. Note that the performance of EP and EF does not change much with the change of the size of the validation dataset. By considering the performance of EN which assigns equal weights to each base classifier, we can obtain a reasonable explanation. As we can see, the performance among different base classifiers does not vary too much on a certain dataset in terms of precision and F1. Therefore, the weights for each base classifier generated according to EP and EF are also similar. That is, EP and EF both perform like EN no matter which validation dataset is used. Hence, the performance of EP and EF does not change much with the change of the validation dataset.

From Table VII, we can also find that EDP (EP) always achieves better precision than EDF (EF) across different validation datasets without sacrificing the value of F1. We can conclude that, in order to achieve higher precision from ensemble classifiers, we should determine the weights of the base classifiers based on their precision values. We can see that the precision achieved by EDP increases steadily while that achieved by EP is relatively stable with the increase of the validation dataset. When 80% or more of the 111 samples are used as the validation dataset,

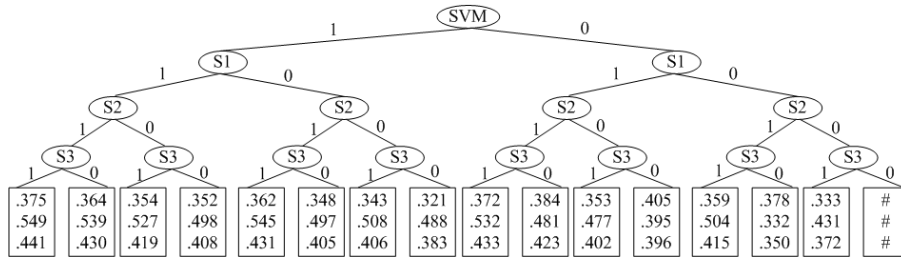


Fig. 12. Performance of all ensemble classifiers constructed from the four base classifiers.

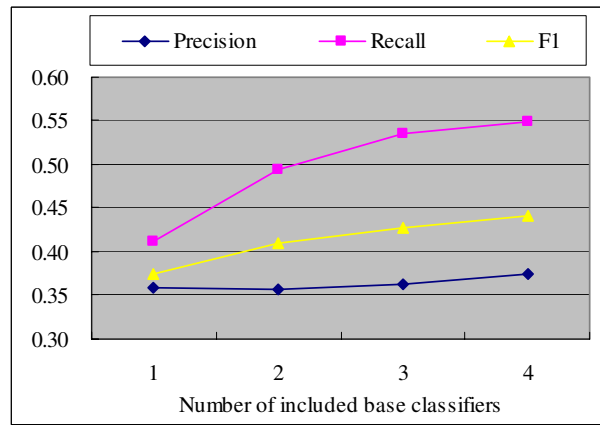


Fig. 13. Averaged performance of different kinds of ensemble classifiers categorized based on the number of base classifiers.

EDP outperforms EP in terms of precision, although the advantage of EDP is not obvious. Therefore, in order to reach higher precision, it is necessary to try EDP if we have a large validation dataset.

5.4.4 *Effect of Base Classifiers.* In the above section, we studied the effect of the different ensemble strategies as well as the effect of the validation dataset. In this section, we will study the effect of the number of base classifiers. For simplicity, we assign each base classifier equal weight when constructing an ensemble classifier as EN does. Given four base classifiers, we can obtain 15 ensemble classifiers in total.

Figure 12 shows the performance of all the ensemble classifiers on the testing dataset. Each leaf in the tree represents an ensemble classifier. On the path from a leaf to the root, “1” indicates that the corresponding base classifier is included to construct the target ensemble classifier; “0” indicates the opposite case. The three numbers in each leaf reflect the performance of the corresponding ensemble classifier in terms of precision, recall and F-measure. From Figure 12, we can conclude that the more base classifiers included, the better the performance. Figure 13 makes the

Table VIII. Average running time of each step for testing a query (seconds).

Retrieve Pages from Search Engines			Process returned Pages	Classify through Classifiers	
Google	Looksmart	Lemur + ODP		Synonym-Based	SVM
0.98s	0.39s	0.11s	0.002s	0.0006s	0.0035s

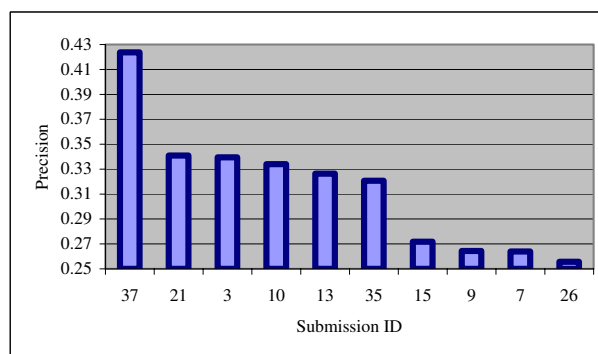
conclusion clearer. In fact, by categorizing the 15 ensemble classifiers according to the number of base classifiers included, we have four kinds of ensemble classifiers. Figure 13 is obtained by averaging the different ensemble classifiers in each kind. From Figure 12, we can also find that whenever the base classifier SVM is included, the recall will be improved obviously while the precision does not change much, and thus F1 is improved. This observation validates our idea once again that the two kinds of base classifiers complement to each other; then the combination can improve the classification results.

*5.4.5 Running Time Analysis.* In this section, we analyze the time complexity of our approaches. Since the training stage of our approaches can be completed offline, we only consider the test stage. As shown before, to test a query, we need four steps: 1) submit the query to search engines and fetch the related pages; 2) process the returned pages to obtain the intermediate category information for each page and the textual representation of the query; 3) apply the two kinds of base classifiers on the enriched representation of the query; 4) combine the results of the base classifiers. Table VIII shows the running time of each step for testing a query (averaged over the 800 queries in the testing dataset) on a PC with 512M memory and a 3.2Ghz CPU. We do not show the time for step 4 since its running time is neglectable as compared to other steps. For example, when combing the results from any two classifiers, the time for each query is about  $1.8 \times 10^{-5}$  seconds.

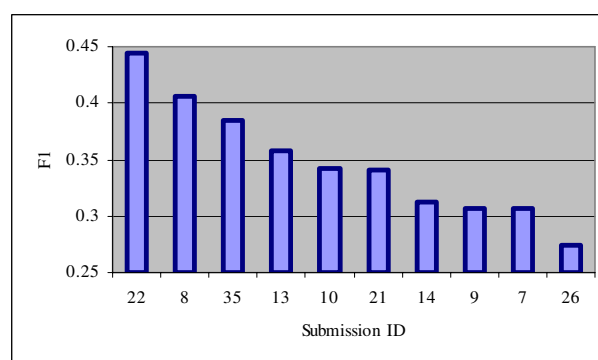
From Table VIII we can observe that the bottleneck of the test stage is to retrieve pages from search engines which includes the time for finding the related pages taken by the search engine and the time for crawling the related pages. For the search engine developed by us based on Lemur and the crawled ODP pages (denoted by Lemur+ODP), the time for retrieving pages is far less than that for Google and Looksmart since the pages are indexed locally and we do not have to crawl them. In fact, the difference of the time depends on the speed of the network. Therefore, when we run our approach on the sever of a search engine, a query can be classified in real time ( in the order of  $10^{-1}$ ).

*5.4.6 Comparison with other Participants.* To further validate our algorithm, we compare it with other KDDCUP2005 participants' systems. First, we briefly introduce two systems of the runner-ups of KDDCUP2005 [Kardkovács et al. 2005; Vogel et al. 2005]. The runner-up of "F1" criterion is the team from the MEDai/AI Insight/Humboldt University. The runner-up of "Precision" criterion is from Budapest University of Technology.

The MEDai/AI Insight/ Humboldt University team built a model containing 67 bi-class classifiers each of which is corresponding to a target category. This model can predict 67 probabilities for each query and the categories with the highest probabilities are chosen for each query. To build the model, they need to train 67 bi-class classifiers. The training dataset for each classifier consists of the 111 queries



(1) Precision of the top 10 solutions



(2) F1 of the top 10 solutions

Fig. 14. Top 10 solutions in KDDCUP2005 in terms of Precision and F1.

given by the KDDCUP2005 organizers. For each of these queries, if it belongs to category  $i$ , it is a positive training sample of the classifier of category  $i$ ; otherwise, it is a negative training sample. For each query, they send it to the Google search engine and get some relevant pages. They use these retrieved pages to represent the query. With these training datasets, they train 67 bi-class SVM classifiers. Given a test query, they input it to the classifiers, and then they can get the probability of each category. They choose the top categories as the final submitted result.

In the KDDCUP2005 Competition, the Budapest University of Technology team proposed an algorithm called the Ferrety, which is based on Internet search engines and a document categorizer [HITEC]. As we observed, they also find that both query words and KDDCUP2005 categories lack semantics while no training data are available. Thus, they seek the meanings of words by asking the Internet which acts as open dictionaries and knowledge base. They send each query into search engines and retrieve possible categories defined by these search engines. Since the obtained categories are different from the target categories, proper category mapping algorithms are needed. By observing that formal definitions of categories are available for a lot of search engines, they do the mapping in this way: first, relevant word collections are expanded by WordNet synonyms of target category names;

secondly, attach probable search engine categories for target ones by matching similarity in the word collections and category definitions. By calculating TF-IDF, more relevant words can be added to the collections. Repeat this process until the word collections become stable. Every search engine category is then mapped to a proper category defined by KDDCUP2005. Half of the 800,000 queries have results from search engines. The left are sent to their own document categorizer, which is a supervised algorithm for classifying text. They get another 320,000 answers with their categorizer. The precision of their algorithm is 0.34088.

Figure 14 contains the evaluation results for the submitted solutions ranked top 10 by organizers. For “F1” criterion, our value is 0.444 (Submission ID: 22), the F1 of the MEDai/AI Insight/Humboldt University team is 0.405 (Submission ID: 8). Our F1 value is higher than the runner-cup team by 9.6% and is higher than the mean of the other 9 teams among the top 10 by 32%. For the “Precision” criterion, our value is 0.424 (Submission ID: 37), the Precision of the Budapest University of Technology team is 0.341 (Submission ID: 21). Our Precision value is higher than that of the runner-cup team by 24.3% and is higher than the mean of the other 9 by 40.3%. Besides that, we also compared the values of precision and F1 of our solutions to the mean values of all other participants. Our precision and F1 are 98.5% and 73.5% higher than the averaged precision and F1 of the other participants.

### 5.5 Some Failed Methods

In fact, we have tried several other approaches which did not work well. Here is an example. The main idea is to build a bridge between a query and the 67 KDDCUP2005 categories by counting the number of pages related to both of them. We submitted a query into search engines and got its related pages. This set of pages is denoted by  $P_q$ . Similarly, we can get the related pages of a category, using the category name as a query directly. This set of pages is denoted by  $P_c$ . In practice, each  $P_q$  includes 100 pages for each query, and each  $P_c$  includes 10,000 pages. Then we can define the similarity between the target query and category as  $\|P_q \cap P_c\|$ , where  $\|\cdot\|$  is the size of a set. For each query, we can return the most related categories according to the similarity. However, this method does not seem to work. One possible reason is that there is correlation between some pairs of categories. For example, we found that the overlap of the 10,000 retrieved pages of the categories “Computer\Hardware” and “Living\Tools & Hardware” are about 1000 pages. Therefore we removed the overlap pages between each pair of categories. We repeated the above approach. However, the final result is still not satisfactory. One reason for the failure of this method is that we cannot judge the similarity between a query and a category simply by the size of the intersected set of the retrieved pages. The essence of the problem is how to automatically find the exact collection of pages that can well represent the semantics of a query and a category and how to determine the similarity between them based on the pages. This is a problem that requires further study.

We also tried another method based on a dictionary which does not work well either. We submitted a query  $q$  into a dictionary software, e.g., WordNet, to get the related words of a query. The result is denoted as  $R_q$ . The result includes the synonyms or antonyms of the given query, which can be a verb, a noun, or

an adjective. Take the query “car” as an example, the result contains: car, auto, automobile, machine, motorcar, railcar, railway car, railroad car and cable car, etc. Similarly, we can get the result of every category  $c$  in the same way which is denoted as  $R_c$ . A similarity between the target query and category is defined as  $\|R_q \cap R_c\|$  as shown before. We can classify a query into the top categories ranked according to similarity. When we tested this method on the validation dataset, the F1 is only about 20%. The main reason for the bad performance is that this method cannot obtain proper results for many queries through WordNet. If more dictionaries like Wikipedia are leveraged, the performance of this method may be improved.

## 6. CONCLUSION AND FUTURE WORK

In this paper, we presented our approach to solving the query classification problem with its application on the task provided by KDDCUP2005. Query classification is an important as well as a difficult problem in the field of information retrieval. Once the category information for a query is known, a search engine can be more effective and can return more representative Web pages to the users. However, since the queries usually contain too few words, it is hard to determine their meanings. Another practical challenge, as shown in KDDCUP2005, is that no training data are provided for the classification task explicitly.

To solve the query classification problem, we designed an approach based on query enrichment which can map the queries to some intermediate objects. With the intermediate objects, we designed several ensemble classifiers based on two kinds of different base classifiers, a statistics-based classifier and a synonym-based one. The experimental results on the two datasets provided by the KDDCUP2005 organizers validate the effectiveness of the base classifiers as well as the ensemble strategies. We have designed a demonstration system called  $Q^2C@UST$ , with a dedicated Web site at <http://q2c.cs.ust.hk/q2c/>.

The success of our approach in the KDDCUP2005 competition can be attributed to two factors: one is the way to enrich the queries and the other is the way to combine the base classifiers. Therefore, in the future, we will conduct more research following these two directions: 1) we will try to find more valuable extra information for the queries based on which we can build the base classifiers, including those from some search engines with special features (such as <http://vivisimo.com/> which clusters the returned pages); 2) we will conduct some further research to find more effective strategies to generate ensemble classifiers.

## APPENDIX A: TARGET CATEGORIES FROM THE KDDCUP2005 TASK

Computers\Hardware	Entertainment\Humor & Fun
Computers\Internet & Intranet	Entertainment\Movies
Computers\Mobile Computing	Entertainment\Music
Computers\Multimedia	Entertainment\Pictures & Photos
Computers\Networks & Telecommunication	Entertainment\Radio
Computers\Security	Entertainment\TV
Computers\Software	Entertainment\Other
Computers\Other	Information\Arts& Humanities
Entertainment\Celebrities	Information\Companies & Industries
Entertainment\Games & Toys	Information\Science & Technology

Information\Education	Online Community\Chat & Instant Messaging
Information\Law & Politics	Online Community\Forums & Groups
Information\Local & Regional	Online Community\Homepages
Information\References & Libraries	Online Community\People Search
Information\Other	Online Community\Personal Services
Living\Book & Magazine	Online Community\Other
Living\Car & Garage	Shopping\Auctions & Bids
Living\Career & Jobs	Shopping\Stores & Products
Living\Dating & Relationships	Shopping\Buying Guides & Researching
Living\Family & Kids	Shopping\Lease & Rent
Living\Fashion & Apparel	Shopping\Bargains & Discounts
Living\Finance & Investment	Shopping\Other
Living\Food & Cooking	Sports\American Football
Living\Furnishing & Houseware	Sports\Auto Racing
Living\Gifts & Collectables	Sports\Baseball
Living\Health & Fitness	Sports\Basketball
Living\Landscaping & Gardening	Sports\Hockey
Living\Pets & Animals	Sports\News & Scores
Living\Real Estate	Sports\Schedules & Tickets
Living\Religion & Belief	Sports\Soccer
Living\Tools & Hardware	Sports\Tennis
Living\Travel & Vacation	Sports\Olympic Games
Living\Other	Sports\Outdoor Recreations
Online Community	Sports\Other

#### ACKNOWLEDGMENTS

We thank a grant from NEC China Lab (NECLC05/06.EG01), Hong Kong RGC Grant HKUST 6187/04E, Hong Kong University of Science and Technology, Department of Computer Science and Engineering and Computer Systems Group for their kind support.

#### REFERENCES

- ALPAYDIN, E. 2004. *Introduction to Machine Learning*. The MIT Press.
- BAUER, E. AND KOHAVI, R. 1999. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning* 36, 1-2, 105–139.
- BEEFERMAN, D. AND BERGER, A. 2000. Agglomerative clustering of a search engine query log. In *KDD '00: Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM Press, New York, NY, USA, 407–416.
- BEITZEL, S. M., JENSEN, E. C., FRIEDER, O., GROSSMAN, D., LEWIS, D. D., CHOWDHURY, A., AND KOLCZ, A. 2005. Automatic web query classification using labeled and unlabeled training data. In *SIGIR '05: Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM Press, New York, NY, USA, 581–582.
- CANN, A. J. 2003. *Maths from Scratch for Biologists*. John Wiley & Sons.
- CARUANA, R., NICULESCU-MIZIL, A., CREW, G., AND KSIKES, A. 2004. Ensemble selection from libraries of models. In *ICML '04: Proceedings of the twenty-first international conference on Machine learning*. ACM Press, New York, NY, USA, 18.
- CHANG, C.-H. AND HSU, C.-C. 1998. Integrating query expansion and conceptual relevance feedback for personalized web information retrieval. In *WWW7: Proceedings of the seventh international conference on World Wide Web 7*. Elsevier Science Publishers B. V., Amsterdam, The Netherlands, The Netherlands, 621–623.
- CHEKURI, C., GOLDWASSER, M., RAGHAVAN, P., AND UPFAL, E. 1997. Web search using automated classification. Sixth International World Wide Web Conference (WWW6). Poster presentation.

- CHEN, H. AND DUMAIS, S. 2000. Bringing order to the web: automatically categorizing search results. In *CHI '00: Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM Press, New York, NY, USA, 145–152.
- DIETTERICH, T. G. 2000. Ensemble methods in machine learning. In *Multiple Classifier Systems*. 1–15.
- FAN, W., STOLFO, S. J., AND ZHANG, J. 1999. The application of adaboost for distributed, scalable and on-line learning. In *KDD '99: Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM Press, New York, NY, USA, 362–366.
- HANSEN, L. K. AND SALAMON, P. 1990. Neural network ensembles. *IEEE Trans. Pattern Anal. Mach. Intell.* 12, 10, 993–1001.
- HITEC. <http://categorizer.tmit.bmr.hu>.
- HOEL, P. G. 1966. *Elementary Statistics*, Second Edition ed. Wiley.
- HOWE, A. E. AND DREILINGER, D. 1997. SAVVYSEARCH: A metasearch engine that learns which search engines to query. *AI Magazine* 18, 2, 19–25.
- JANSEN, B. J. 2000. The effect of query complexity on web searching results. *Inf. Res.* 6, 1.
- JOACHIMS, T. 1998. Text categorization with support vector machines: Learning with many relevant features. In *ECML*. 137–142.
- JOACHIMS, T. 1999. Transductive inference for text classification using support vector machines. In *ICML '99: Proceedings of the Sixteenth International Conference on Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 200–209.
- JONES, K. S. 1971. *Automatic Keyword Classifications for Information Retrieval*. Butterworth, London.
- KANG, I.-H. AND KIM, G. 2003. Query type classification for web document retrieval. In *SIGIR '03: Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM Press, New York, NY, USA, 64–71.
- KARDKOVÁCS, Z. T., TIKK, D., AND BÁNSÁGHI, Z. 2005. The ferrety algorithm for the kdd cup 2005 problem. *SIGKDD Explor. Newsl.* 7, 2, 111–116.
- KITTLER, J., HATEF, M., DUIN, R. P. W., AND MATAS, J. 1998. On combining classifiers. *IEEE Trans. Pattern Anal. Mach. Intell.* 20, 3, 226–239.
- LEWIS, D. D. AND GALE, W. A. 1994. A sequential algorithm for training text classifiers. In *SIGIR*, W. B. Croft and C. J. van Rijsbergen, Eds. ACM/Springer, 3–12.
- LI, Y., ZHENG, Z., AND DAI, H. K. 2005. Kdd cup-2005 report: facing a great challenge. *SIGKDD Explor. Newsl.* 7, 2, 91–99.
- MCCALLUM, A. AND NIGAM, K. 1998. A comparison of event models for naive bayes text classification. AAAI-98 Workshop on Learning for Text Categorization.
- MEYER, D. A. AND BROWN, T. A. 1998. Statistical mechanics of voting. *Physics Review Letters* 81, 8, 1718–1721.
- MILLER, G., BECKWITH, R., FELLBAUM, C., GROSS, D., AND MILLER, K. 1990. Introduction to wordnet: an on-line lexical database. *International Journal of Lexicography* 3, 4, 23–244.
- PAGE, L., BRIN, S., MOTWANI, R., AND WINOGRAD, T. 1998. The pagerank citation ranking: Bringing order to the web. Tech. rep., Stanford Digital Library Technologies Project.
- SELBERG, E. AND ETZIONI, O. 1995. Multi-service search and comparison using the MetaCrawler. In *Proceedings of the 4th International World-Wide Web Conference*. Darmstadt, Germany.
- SHEN, D., PAN, R., SUN, J.-T., PAN, J. J., WU, K., YIN, J., AND YANG, Q. 2005. Q2c@ust: our winning solution to query classification in kddcup 2005. *SIGKDD Explor. Newsl.* 7, 2, 100–110.
- SHEN, D., SUN, J.-T., YANG, Q., AND CHEN, Z. 2006. Building bridges for web query classification. In *SIGIR '06: Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*.
- SILVERSTEIN, C., MARAIS, H., HENZINGER, M., AND MORICZ, M. 1999. Analysis of a very large web search engine query log. *SIGIR Forum* 33, 1, 6–12.
- VAN, R. C. 1979. *Information Retrieval*, Second Edition ed. Butterworths, London.
- ACM Journal Name, Vol. ?, No. ?, ? 20?.



- VOGEL, D., BICKEL, S., HAIDER, P., SCHIMPFKY, R., SIEMEN, P., BRIDGES, S., AND SCHEFFER, T. 2005. Classifying search engine queries using the web as background knowledge. *SIGKDD Explor. Newsl.* 7, 2, 117–122.
- VOORHEES, E. M. 1994. Query expansion using lexical-semantic relations. In *SIGIR '94: Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval*. Springer-Verlag New York, Inc., New York, NY, USA, 61–69.
- WEN, J.-R., NIE, J.-Y., AND ZHANG, H.-J. 2002. Query clustering using user logs. *ACM Transactions on Information Systems (TOIS)* 20, 1, 59–81.
- YANG, Y. 1999. An evaluation of statistical approaches to text categorization. *Inf. Retr.* 1, 1-2, 69–90.
- YANG, Y. AND PEDERSEN, J. O. 1997. A comparative study on feature selection in text categorization. In *ICML '97: Proceedings of the Fourteenth International Conference on Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 412–420.

Received Month Year; revised Month Year; accepted Month Year