

Analyzing Medical Image Search Behavior: Semantics and Prediction of Query Results

Maria De-Arteaga¹ · Ivan Eggel² · Charles E. Kahn Jr.³ · Henning Müller²

Published online: 26 March 2015
© Society for Imaging Informatics in Medicine 2015

Abstract Log files of information retrieval systems that record user behavior have been used to improve the outcomes of retrieval systems, understand user behavior, and predict events. In this article, a log file of the ARRS GoldMiner search engine containing 222,005 consecutive queries is analyzed. Time stamps are available for each query, as well as masked IP addresses, which enables to identify queries from the same person. This article describes the ways in which physicians (or Internet searchers interested in medical images) search and proposes potential improvements by suggesting query modifications. For example, many queries contain only few terms and therefore are not specific; others contain spelling mistakes or non-medical terms that likely lead to poor or empty results. One of the goals of this report is to predict the number of results a query will have since such a model allows search engines to automatically propose query modifications in order to avoid result lists that are empty or too large. This prediction is made based on characteristics of the query terms themselves. Prediction of empty results has an accuracy above 88 %, and thus can be used to automatically modify the query to avoid empty result sets for a user. The semantic analysis and data of reformulations done by users in the past can aid the development of better search systems, particularly to improve results for novice users. Therefore, this paper gives important ideas to better understand how people search and how to use

this knowledge to improve the performance of specialized medical search engines.

Keywords Image retrieval · Human-computer interaction · Machine learning · Statistic analysis · Information storage and retrieval · Medical image search · Log file analysis

Introduction

Medical imaging studies have increased significantly in both quantity and complexity over the past 30 years [1]. Images are an essential part of medical diagnosis and treatment planning, and many tools have been created to search and interpret images, as well as to give medical doctors decision support [2, 3]. Among medical specialties, radiologists are at the forefront of analyzing images, searching for specific patterns in them, and describing them in reports that form a basis for further decision making. In general, physicians increasingly use online resources to search for information. Radiologists commonly use standard search engines to look for image information for medical images [4]. Specialized radiology search engines such as ARRS GoldMiner,¹ Yottalook,² or Shambala³ allow users to search for images in the medical literature using text queries or, in some cases, image examples to search for visual similarity. Research has shown that text search, filters for imaging modality, and image and region-of-interest search are requested by radiologists [5].

In contrast to other approaches to study users' web-site usage, search log analysis is an unobtrusive method that shows significant advantages compared to surveys and

✉ Maria De-Arteaga
mdeartea@andrew.cmu.edu

¹ Carnegie Mellon University, 4800 Forbes Ave.,
Pittsburgh, PA 15213, USA

² HES-SO, Rue de TechnoPole 3, 3960 Sierre, Switzerland

³ University of Pennsylvania, Philadelphia, PA, USA

¹ <http://goldminer.rrs.org>

² <http://www.yottalook.com>

³ <http://shambala.khresmoi.eu>

laboratory studies in scale, power, scope, and location [6]. Despite limitations such as possibly imprecise user representation, less versatility, less richness, and a loose link to concepts supposed to be measured [6], search log analysis has been used in the biomedical domain to examine textual and visual retrieval systems [7].

Search logs of general search engines have been used to predict flu outbreaks and to analyze medication use [8]. They also have been used to analyze image search behavior [9, 10]. Analysis of MedLine search behavior in the medical literature was conducted based on log files [11, 12]. Closest to the presented work are the analyses of Tsikrika et al. [7] and Rubin et al. [13] that both used ARRS GoldMiner log files, but a much smaller set of queries (25,000 and 30,000, respectively, so around 10 % of the data used in this text). None of these systems performs user profiling, which would be possible with registered users of a search system. Detecting user profiles from log files was attempted in [14], but we do not try to separate queries into several user categories for ARRS GoldMiner as the technologies do not seem fully stable and our objective is to rather predict problematic queries for any user group.

Tsikrika et al. [7] analyzed 25,000 ARRS GoldMiner queries to investigate the process of query formulation and query modification in order to identify medical professionals' information needs with the aim to improve the effectiveness of the search support of such systems. This article extends the previous work using a dataset of 222,005 search queries with timestamp information. Timestamp information was not available in the previous study and was used to create user sessions with specific time limitations. Additionally, the key contribution of this paper lies in the use of machine learning algorithms to predict a query's success and the number of results for a specific query.

Similarly, Rubin et al. [13] analyzed 30,000 queries to ARRS GoldMiner and Yottalook, and implemented an algorithm for mapping search terms to RadLex,⁴ an ontology consisting of radiology terms, with the goal of determining what radiologists search for on the Web. As their research showed, giving the queries a RadLex semantic context improves the robustness of the analysis. Therefore, this paper also includes mapping to RadLex terms and axes, using an automatic text categorization system [15] that gives a robust mapping. This system does the mapping in three different ways, which allows to differentiate a query that is itself a RadLex term from one that includes several RadLex terms, among other cases.

The first part of the paper builds on the past work to construct a detailed analysis of a larger log file of the ARRS GoldMiner search system, while also aiming to improve technical aspects of the methodology. The second part of the paper

uses machine learning to build a predictive model that is able to determine the range of the number of query results. ARRS GoldMiner retrieves all documents containing all query terms (with “AND” connection by default); additionally, if the term is in a vocabulary, the search is also done using the corresponding concept (MeSH, SNOMED, etc.). Therefore, it is possible to have queries with too many results and others with no results. Machine learning techniques, though widely used when working with search log files from search engines [8], have not been applied to analyze ARRS GoldMiner nor radiologists' image search behavior [4].

The results presented in this paper provide a better understanding of the way in which physicians search for information. Given the fact that ARRS GoldMiner was originally a search engine with only radiology articles/images and supported by the American Roentgen Ray Society, it is still very strong in this domain. The behavior of its users can be considered somewhat representative for the behavior of radiologists even though the system can be accessed by any Internet user, an assumption that is supported by the high percentage of queries that can be mapped to RadLex, especially when it is taken into account that many radiology queries cannot be mapped due to spelling mistakes or absence of terms in the RadLex terminology. It also proposes two algorithms to predict whether a query will have at least one result and in what range the number of query results will be, respectively. Both algorithms have a very high accuracy and use very simple data as input, two characteristics that make them a viable alternative to be implemented in search engines as a criterion to determine when a query modification should be suggested as the computation is extremely fast. For example, if the algorithm predicts there will be too many results, the search engine could suggest the user to narrow the search; similarly, if the prediction forecasts no results, the search engine could suggest alternative queries that return results. To propose alternative queries, the analysis of what other users have done in the past in terms of query reformulations, such as the one presented in this paper, can be extremely useful. For example, modifications that have been successful for other users in the past could work as a basis for suggestions made to new users. Such a recommendation system would potentially work better the more queries and query modifications it contains. Spelling correction can be another source for query modifications.

This paper is organized as follows: “Methods” section includes a description of the data, of the methods used to produce descriptive analysis, and of the machine learning models. “Results” section presents the descriptive analysis of user search behavior and the results of the predictive models. Finally, in “Discussion” section, results are discussed and “Conclusion” section contains the conclusions.

⁴ <http://www.radlex.org>

Methods

Data Source

The examined query log was produced by the American Roentgen Ray Society (ARRS) GoldMiner medical image search engine [16], which currently provides access to more than 485,000 selected images from peer-reviewed biomedical journals targeted mainly to clinical professionals. The images are indexed using the keywords of the caption, the imaging modality, and the age and the gender of the patient, which are all automatically extracted from the text.

The search procedure within ARRS GoldMiner always starts with a keyword search, with the possibility of filtering results at a later stage by gender, age groups, and modality. The results are returned as a set of pages, each consisting of a list of up to ten results, or a display of up to 40 image thumbnails. Each result contains the image thumbnail, the caption, the modality, and a link to the article containing the image. The acquired log file contains 222,005 consecutive queries. Each log entry includes a timestamp, a client identifier (encrypted IP address to preserve privacy), the query itself, and the number of results found for that query.

Preprocessing of the query logs was done in the same way as Tsirikas et al. [7]: all queries were converted to lowercase, various special characters were removed, and medical imaging modalities were normalized (for example, “XR,” “X-ray,” and “xray” were mapped to a single term). Consecutive identical queries in the same session and with the same number of results were considered as a single query. Such entries occur when a searcher submits a query, then views a document, and returns to the search engine. The Web server typically logs this second visit with the identical user identification and query but with a new timestamp. Also, result page navigation can cause the same logging behavior. The log also contained identical queries in the same session that yielded different result sets; these queries were kept because they could reflect the use of filters (for age group or modality, for example).

Descriptive Analysis

Understanding the user’s behavior is key to enhance information retrieval systems. The first part of this paper provides descriptive analysis of the data contained in the log files.

Log analysis at session level can provide valuable information. A session is defined as a series of queries done by a single user within a small range of time where he/she attempts to fill a single information need [17]. As commonly applied, a session cut-off time of 30 minutes was defined [18]. This means that all consecutive queries within less than 30 minutes of inactivity to the previous query will be considered a session. A query made later than the cut-off time to the previous query will be put into a new session. Query modification analysis is

conducted within session boundaries and identifies the relationship between consecutive queries with three possible outcomes: query generalization, query specification, and query reformulation.

In order to put the queries into a semantic context, a mapping from queries to RadLex terms was applied. RadLex is a reference ontology for the radiology domain that currently contains more than 30,000 terms used mainly for standardized indexing and retrieval of radiology information resources. It was developed by the Radiological Society of North America (RSNA) in order to satisfy needs of software developers, system vendors, and radiology users by adopting the best features of existing terminology systems, while producing new terms to fill critical gaps [19, 20]. Standard lexicons such as RadLex can be used to solve data-mining challenges that occur due to synonyms, negation, and inheritance⁵; for example, all synonyms are mapped to the same RadLex term. This mapping was mainly done to determine which of the RadLex axes were most often represented in the queries, as well as to count the term frequency of the mapped RadLex terms. The mapping from queries to RadLex terms was achieved by using Ruch’s system for automatic assignment of biomedical categories [15] using lexical similarity of terms. Each term that could be mapped to RadLex was classified into one of the following 15 axes of RadLex: Imaging protocol, Report, Procedure, RadLex descriptor, Property, Anatomical entity, Imaging observation, Process, Imaging modality, Non-anatomical substance, RadLex non-anatomical set, Report component, Procedure step, Object, and Clinical finding, which are the main RadLex axes.

Predictive Models

A machine learning approach was applied to build a system capable of predicting the number of results a query will have. Two different tasks were defined: predicting if a query will have no results and predicting the range of the number of results (0–10 results, 10–100 results, or more than 100 results). These three classes were chosen because fewer than ten results could be considered a query with too few results and more than 100 could be considered a very broad query where no one would look at all results, whereas in between could be considered a desirable result set.

Each query was represented by 18 attributes that were used to train the machine learning algorithms. The attributes were the following:

RadLex mappings: As explained in the “Descriptive Analysis” section, queries were mapped to RadLex terms in order to place them in a semantic context. Four types of mappings were possible: *exact* (the whole query

⁵ http://www.rsna.org/RadLex_in_Your_Practice.aspx

corresponds to a term in the RadLex ontology), *all terms* (all the terms in the query can be mapped to a RadLex concept), *partial* (at least one, but not all, the terms in the query are mapped to RadLex), and *none* (no term in the query can be mapped to RadLex). The first RadLex-related attribute is the type of mapping done. Given there are multiple types of mappings, each query can have between 0 and N RadLex mappings, N being the number of terms in the query. Therefore, 13 attributes were created, one for every RadLex axis present in the log files. These are binary attributes; every query is assigned a 0 or 1 in each of these variables, depending on whether the query was mapped to the axis or not.

Number of tokens in query: Two attributes were created based on the number of tokens in the query: total number of tokens and number of tokens without stopwords. The query “tumor in lung”, for example, has three tokens and two non-stopword tokens.

Appearances of terms in log files: A dictionary with all the words in the queries was created, and for each of them the total number of queries in which it appears was counted. Later, this information was used to build two attributes of the vector representation of each query: *min logfile appearances* and *max logfile appearances*. In the previous example, “tumor in lung”, let us assume “tumor” appears 108 times, “in” appears 2000 times, and “lung” appears 520 times. Then, for this query, *min appearances*=108 and *max appearances*=2000.

To prevent deceitful results due to unbalanced classes, the Synthetic Minority Over-Sampling Technique (SMOTE) [21] was used to balance the classes. Once this was done, a set of machine learning algorithms was selected based on the state-of-the-art tools used in the field, and experiments were conducted on these methods to determine which has the best performance in this specific case. The experiments were done using 10-fold cross-validation on all the data. The methods considered were support vector machines [22], logistic regression [23], random forests [24], and other decision trees. The criteria used to compare them were based on correctly classified instances, kappa statistic [25], F-measure [26], and the area under the receiver operating characteristic (ROC) curve [26].

Finally, in order to analyze the impact of each attribute in the predictive model, providing understanding on which elements are relevant for prediction and which are not, an information gain attribute ranking [27] was applied to determine the importance of each attribute.

Results

This section describes the main outcomes of this article. In the first part, the descriptive analysis is presented. Then, the predictive models, their accuracy, and other metrics are exposed.

Descriptive Analysis

Terms and Queries A query corresponds to the exact text a user types into the search engine, whereas terms are extracted from the queries and might constitute the whole or part of a query. The total number of queries was reduced from 222,005 to 200,361 after preprocessing, with 92,909 queries (46 %) being distinct and 75,118 queries (37.4 %) appearing only a single time. In comparison to these results, the study in [7], working with 25,000 records, 63 % of the queries appeared a single time; the difference between these two numbers shows there is a gain in information when working with a larger dataset.

Each query was repeated on average twice, and 17,791 of the 200,361 queries (8.9 %) occurred more than once. This shows that relatively few queries are repeated. The high average can be explained by the fact that the ten most frequently occurring queries represented approximately 2 % of all queries. Queries that occurred only once were extremely specific terms, minor spelling mistakes that did not occur frequently, or totally off-topic queries.

Regarding the most frequently occurring terms, 33,903 (17 %) of the queries contained at least one of the ten most frequently occurring terms, and 91,589 (46 %) contained one of the top 100 terms, with “cyst” being the most frequent. Figure 1 shows the proportion of queries containing the most frequently occurring terms. Tables 1 and 2 show the most frequently occurring queries and terms, respectively. Results are very similar to [7] with seven of the most frequent queries and nine of the most frequent terms occurring in both albeit with a slightly changing order and very different absolute numbers.

The majority of the queries consisted of two terms, followed by queries with one term, and then by those with three terms. The mean number of terms per query was 2.21; the median was 2. Among all queries, 182,004 (90.8 %) consisted of three or fewer terms. In contrast, PubMed averages 3.54 terms per query [12], with a median of 3 terms per query; 80 %

Table 1 The most frequent queries in the logfile

	Query	Frequency
1	mega cisterna magna	820
2	bastrup disease	798
3	limbus vertebra	462
4	toxic	428
5	cystitis cystica	405
6	buford complex	274
7	thornwaldt cyst	274
8	splenic hemangioma	254
9	double duct sign	249
10	cystitis glandularis	245

Table 2 The most common terms occurring in the queries

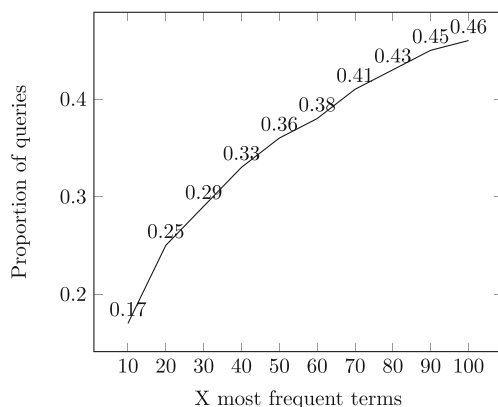
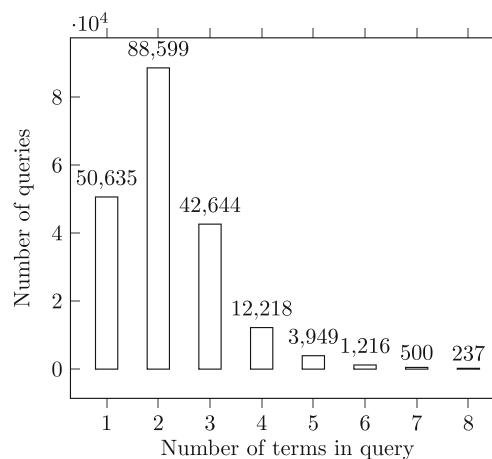
	Term	Frequency
1	cyst	6346
2	mri	3536
3	disease	3536
4	ct	3504
5	fracture	3366
6	tumor	3233
7	syndrome	2994
8	liver	2486
9	pulmonary	2424
10	sign	2293

of all queries have no more than 4 terms. Figure 2 shows the number of queries given the number of terms in it. Again, these results are very similar to results in [7].

Sessions In the log files, 103,029 user sessions were identified. Among these, 100,761 (97.8 %) contain less than seven queries, 64,679 (62.7 %) contain only one query, 17,379 (16.9 %) have two queries, and 8453 (8 %) have three queries. The longest session has 126 queries.

Studying 97,315 query pairs of consecutive queries in sessions showed that, out of these, 36,056 (37.1 %) do not share any common terms and only 741 (0.76 %) are identical (this is influenced by result filtering), making 61,259 (62.9 %) of consecutive queries in a session share at least one common term.

When analyzing the modifications done by a user in a session, 30,622 (31.4 %) query pairs represent a query reformulation, followed by query generalization 16,757 (17.2 %) and query specification 13,139 (13.5 %). This confirms results obtained by Tsikrika et al. [7] and thus opposes the large majority of studies analyzing Web search engines logs, where reformulation is also the mostly frequently observed query modification type, but it is followed by specification and generalization [28]. Unlike Tsikrika et al. [7], available query time

**Fig. 1** Proportion of the queries containing the most frequently occurring terms**Fig. 2** The number of queries with a specific number of terms in the query

information allowed this study to limit the analysis to consecutive queries inside a search session, instead of all consecutive queries by the same client IP, leading to a much smaller number of query pairs relative to the search log size. According to our analysis, among the 91,375 subsequent queries in a session, the vast majority of queries 66,819 (73.1 %) have a time span of less than 1 min between two queries.

RadLex Mapping From the 200,361 queries left after pre-processing, 124,719 (62.2 %) queries could be mapped to RadLex with one of the three techniques used: 36,372 (18.2 %) queries where an exact match to a RadLex concept, while 76,928 (38.4 %) could be partially mapped, and 11,419 (5.7 %) had every term mapped to a concept in the ontology. The remaining 75,642 (37.8 %) queries could not be mapped to RadLex at all. The terms include non-medical terms, many spelling mistakes, and terms that are too specific and not part of RadLex. In [13], 52 % of the terms could be mapped to a smaller and older version of RadLex.

The most common RadLex axis is *clinical finding*, with 79,721 queries being or containing a term that could be mapped to it, which represents 40 % of all queries. The second most common axis is *anatomical entity* with 38,791 (19.3 %) queries, having a huge gap with the third most common axis, *RadLex descriptor*, which is only present in 22,321 (1.1 %) queries (for analyzing these percentages, it is very important to remember every query can be mapped to more than one or to none RadLex terms). Figure 3 shows the relationship between number of queries and RadLex axes. In [13], the most frequent axis was anatomic location (52.3 %), but RadLex was much smaller at the time and it is possible that this is responsible for part of these differences with findings only covering 10.7 % of the queries in this older analysis.

Among the queries, 99,060 (49.4 %) are mapped to one single RadLex axis, while 23,477 (11.7 %) were mapped to two axes; 2130 (1.1 %) contained terms belonging to three

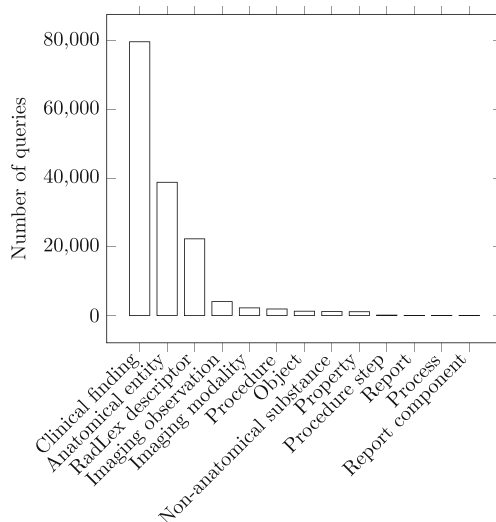


Fig. 3 Number of queries mapped to each RadLex axis

different axes and 52 (0.03 %) to four different axes. No query was mapped to more than four axes. A similar analysis was not done in the prior work of [13].

At this point, an important question is: what axes do users tend to combine for formulating their information needs? To answer this question, the matrices in Tables 3 and 4⁶ show the number of times each pair of axes co-occurs. As expected, *clinical findings* and *anatomical entities*, being the most frequent axes, co-occur with others frequently. For example, the two of them co-occur in 11,787 queries, which correspond to 20 % of the queries mapped to *anatomical entity*. Among the queries mapped to *RadLex descriptor*, 8272 were also mapped to *clinical findings*, which corresponds to 22 %. The distribution of co-occurrences, however, is not only due to the frequency with which each axis appears; for *imaging observation*, for example, *clinical findings* is only present in 1.9 % of the queries containing it, while *anatomical entity* co-occurs with it on 9.4 % of its queries.

Predictive Models

Machine learning algorithms were used to perform two tasks: predicting the range in which the number of results will be and predicting whether a query will or will not have results. This is a classification task, for which we aim to obtain the highest possible accuracy. Several experiments were conducted to determine which algorithm to use. In a first set of experiments, logistic regression, support vector machines (sequential minimal optimization), and random forests were tested. A model to predict the number of query results using the features based on

Table 3 Co-occurrence of RadLex axes in the queries (first part containing CF, O, AE, NS, RD, PP)

	CF	O	AE	NS	RD	PP
CF	79,721	175	11,787	150	8272	225
O	175	1243	229	4	89	7
AE	11,787	229	38,791	116	5217	166
NS	150	4	116	1161	55	7
RD	8272	89	5217	55	22,321	18
PP	225	7	166	7	189	109
P	280	18	357	4	163	16
PS	0	1	12	0	1	0
IO	97	6	488	2	543	16
IM	552	25	580	2	249	9
RC	2	0	5	0	3	0
R	4	0	1	0	0	0
PC	1	1	5	0	0	0

appearances of terms in log files and *number of terms in query* gave an accuracy of 50.19 % for logistic regression, 49.99 % for support vector machines, and 81.32 % for random forests. This accuracy is obtained using a 10-fold cross-validation using the entire dataset, which is the evaluation technique used in all the experiments mentioned here. Note that the accuracy of random forests is lower than the accuracy finally reported since these experiments were conducted in the first phase of the project, without taking into account the features based on RadLex mapping. Nonetheless, after finding random forests to perform radically better than the other techniques, which do not even outperform the baseline (49.99 % if every query is assigned to the majority class), random forests were chosen as the preferred method for the task. The default Weka⁷ parameters for random forests allow the model to choose how deep each tree will be and sets the number of trees to 10. Once the model had been trained using the whole set of features, experiments were conducted to determine if increasing the number of trees would improve the results. However, increasing the number of trees to 15 had a barely null impact on the accuracy (in the order of 10^{-3}), and therefore the final choice of algorithm uses 10 trees.

The dataset used for building the model is unbalanced, which means it is not divided evenly among the classes. Therefore, after representing each query as a vector in \mathbb{R}^{18} , the data were preprocessed with SMOTE, in order to prevent unbalanced classes in the training data from altering the results, and used to train a predictive model. To assess the performance of the algorithm, a 10-fold cross-validation was used. Promising results were obtained: an accuracy of 85.19 %, with an average ROC area of 0.95 and a Kappa

⁶ CF: clinical findings, O: object, AE: anatomical entity, NS: non-anatomical substance, RD: RadLex descriptor, PP: property, P: procedure, PS: procedure step, IO: imaging observation, IM: imaging modality, RC: report component, R: report, PC: process.

⁷ <http://www.cs.waikato.ac.nz>

Table 4 Co-occurrence of RadLex axes in the queries (second part containing P, PS, IO, IM, RC, R, PC)

	P	PS	IO	IM	RC	R	PC
CF	280	0	97	552	2	4	1
O	18	1	6	25	0	0	1
AE	357	12	488	580	5	1	5
NS	4	0	2	2	0	0	0
RD	163	1	543	249	3	0	0
PP	16	0	16	9	0	0	0
P	1889	1	11	23	0	1	0
PS	1	101	0	0	0	0	0
IO	11	0	4044	12	0	0	0
IM	23	0	12	2211	0	0	0
RC	0	0	0	0	10	0	0
R	1	0	0	0	0	16	0
PC	0	0	0	0	0	0	12

statistic of 0.77. More detailed information is included in Table 5.

For predicting whether a query would have results, SMOTE was also used to balance the classes in the training data and the algorithm with the best performance was also random forests. Once again, increasing the number of trees gave almost null variation in accuracy. While the first one was a classification task between two classes, the second one classifies into three classes: 0–10 results, 10–100 results, and more than 100 results. The evaluation was also done using 10-fold cross-validation and the performance is also remarkable: an accuracy of 88.29 %, with a ROC area of 0.95 and a Kappa statistic of 0.76. More details about the performance can be seen in Table 6.

The downside of several machine learning algorithms, such as random forests, is the low interpretability; it is hard to understand which variables are important and which are not. In order to gain insight into the role variables play in the prediction, Information Gain Attribute Ranking was used. For a class C and an attribute A , Ent being the entropy, the information gain, I , is measured by

$$I(C, A) = Ent(C) - Ent(C|A)$$

Table 5 Results of Random Forests for predicting the range of the number of query results—R1 has less than ten results (including no results), R2 has between 10 and 100 results, and R3 has more than 100 results

	R1	R2	R3	Weighted average
Precision	0.842	0.819	0.874	0.85
Recall	0.876	0.688	0.92	0.851
F-measure	0.899	0.748	0.897	0.849
ROC area	0.955	0.92	0.971	0.953

Table 6 Performance of Random Forests for predicting if a query will have results or not

	# Res>0	# Res=0	Weighted average
Precision	0.899	0.865	0.884
Recall	0.899	0.864	0.884
F-measure	0.899	0.865	0.884
ROC area	0.951	0.951	0.951

Table 7 and 8 show the attributes' information gain for both tasks.

Given the information gain is the difference between two entropies and for each task the entropy of the class is different, the numbers cannot be directly compared (for example, the fact that in both cases *min logfile appearances* is around 35 does not mean anything). However, conclusions can be drawn from the distribution of the values, as well as for values close to zero, since these ones mean the entropy of the class and the entropy of the class given the attribute is almost the same, meaning there is no information gain from this attribute.

In both cases, *min logfile appearances* is by far the most relevant attribute. The type of RadLex mapping done to the query, the number of tokens (both with and without stopwords), and the *max logfile appearances* are important in both cases, although this last one is more relevant in the

Table 7 Relative influence of variables for predicting if a query will have no results, according to Info Gain Evaluation

Variable	Info Gain
min logfile appearances	0.35278316
Type of RadLex mapping	0.10706495
max logfile appearances	0.09828245
Number of tokens	0.07604782
Number of non-stopword tokens	0.07554565
RadLex: clinical finding	0.02718913
RadLex: non-anatomical substance	0.00130726
RadLex: imaging observation	0.00129999
RadLex: anatomical entity	0.00082734
RadLex: procedure	0.00047458
RadLex: property	0.00042359
RadLex: RadLex descriptor	0.00035407
RadLex: imaging modality	0.00033401
RadLex: object	0.00026038
RadLex: procedure step	0.00016858
RadLex: process	0.00001056
RadLex: report component	0.00000342
RadLex: report	0.00000335

Table 8 Relative influence of variables for predicting the range of the number of query results, according to Info Gain Evaluation

Variable	Info Gain
min logfile appearances	0.3625514
max logfile appearances	0.1735592
Number of non-stopword tokens	0.1498272
Number of tokens	0.1497191
Type of RadLex mapping	0.1130494
RadLex: clinical finding	0.0122519
RadLex: RadLex descriptor	0.0091736
RadLex: imaging observation	0.0018093
RadLex: property	0.0016000
RadLex: non-anatomical substance	0.0013986
RadLex: anatomical entity	0.0013594
RadLex: imaging modality	0.0009119
RadLex: object	0.0006390
RadLex: procedure	0.0001619
RadLex: procedure step	0.0001126
RadLex: report	0.0000384
RadLex: process	0.0000363
RadLex: report component	0.0000165

second task, which could be expected since this task also aims to predict when a query will have too many results. In both cases, RadLex axes do not provide much additional information.

Discussion

In this paper, image search behavior of physicians and other web searchers for medical image information is analyzed based on the usage of log files and predictive models to determine how many results a query will have are presented. The high accuracy of the predictive models, combined with the strong patterns identified in the descriptive analysis of user behavior, can be used to improve medical image search engines. The process of suggesting query modifications to users can be divided into two questions: when to suggest a modification and what to suggest. The findings of this paper can provide answers to both questions.

Predicting the range of the number of query results, or predicting whether a query will have results or not (depending on the desired complexity), can be used as a criterion to determine when the engine should suggest to the user a query modification. The good performance of both classifiers makes them suitable candidates for being used by search engines. As these parameters are extremely simple when removing the RadLex categories, they are also extremely fast to execute, much faster than executing a query; without optimization, much less than half a second could be obtained. Adding this

time to a query is invisible for the user and the user can then be informed on the modifications done and the reasons for it, allowing potentially to reuse the initial query.

Once the system predicts that the query will probably not give a suitable number of results, it can make a suggestion. The information obtained from session analysis can be useful for this. Successful reformulations made by other users in the past can be used as suggestions for new users. This could be an appropriate approach whenever the query was made by another user in the past; however, as previously shown, less than 10 % of the queries occur more than once, so many queries would not have a candidate for suggestion unless the log file grows massively and is available over a long period of time. Therefore, complementary methods have to be developed. The first element that can help in improving a search engine is applying orthographic correction. This can reduce the number of queries with no results. As a second step, considering many searches give no results because they are too specific and others give too many results because they are too broad, it would be desirable to suggest a less or a more specific query, respectively. For the first case, a query in the log files which is contained in the current query and has obtained results could be a good candidate for a suggestion. For example, *aortitis retroperitoneal fibrosis* gives no results, so the search engine could propose the user to look for *retroperitoneal fibrosis*, which does have results. In the second case, the most common queries which contain the current query could be suggested as possible modifications. For example, if the initial input is *fibrosis*, the search engine could suggest a set of more specific queries for the user to choose from, such as *cystic fibrosis*, *interstitial pulmonary fibrosis*, *retroperitoneal fibrosis*. In this case, initial results can be shown in addition to the recommended reformulations.

To further improve the results, an interesting task would be to identify off-topic queries, such as “happy new year” and “San Valentine’s” that occurred in the log files. For these cases, there would be no suitable suggestion that improves the results, so the search engine could warn the user about this.

As described, the main contribution of this paper on user search log file analysis is to propose a model for medical image search engines to suggest query modifications to the users based on automatic predictions based on single queries. However, the results can also be useful for other purposes. The frequency with which certain RadLex axes appear in searches and the way in which they are combined answers the question “what are physicians looking for?”. This gives valuable information to those proposing medical image retrieval tasks as benchmarks, as it is the case of CLEF eHealth [29] or ImageCLEF [30]. Knowing what radiologists or physicians in general search for is key to establishing useful tasks.

In the machine learning portion of the research, the information gain measure provides valuable insight. The fact that the most relevant attribute is *min logfile appearances* suggests

there is an “offer-demand” relation since the number of times a query has been done is useful for predicting the number of results it will get. The same happens with *max logfile appearances*.

The fact that RadLex axes are not useful for prediction is an unexpected result since according to the hypothesis it was expected this would have impact on the number of results. However, RadLex mapping is still useful since the type of mapping has a high information gain. This classification between the three types of mapping, part of Ruch’s method [15], is particularly useful for this analysis.

Conclusion

This paper focuses on understanding how medical image search is performed and using this knowledge to improve specialized search engines. Data mining and machine learning techniques are applied to layout solid bases for a model of query modification suggestions. Two accurate predictive models are presented; the first one to determine when a query will have no results and the second one to determine the range of the number of query results. In a search engine, giving no results is always a bad performance. Suggestions and modifications should be used to prevent this, and therefore predicting when it will happen is key to improving the system. The findings are promising, proving search log files can be used to train a system able to predict the level of success a search will have based on the query terms. Furthermore, a viable model that can be used by medical search engines for identifying problematic queries and modifying them to get better results is presented.

Larger log files can even improve results since this can help to create self-learning systems. Past session information can be a valuable asset for modification suggestions to users, a field in which medical search engines still have some road ahead. In standard search engines such as Google or Bing, already queries are auto-completed while typing based on past queries and their frequencies. A similar possibility exists for medical image search if sufficiently large log files are available. Even dictionaries with standard spelling mistakes can be build based on such log files. Mapping of queries to RadLex is reliable and also allows to avoid problems with synonyms as they are all mapped to a single term. Like this, more can be found out on user intentions when querying, which can again be used to deliver better results than simply using key words.

Within log files, there is potentially more information that could be used to good advantage, such as click information and time spent visiting links. For GoldMiner, we unfortunately did not have this information available, but it is again a technique frequently used in web search log files that could be transferred to medical search. The strong patterns identified in user behavior corroborate this is a subject that should be

studied further, aiming to improve image retrieval and search engines performance for medical search. Already the described analyses potentially allow to adapt the GoldMiner system much better to the user needs by only small modifications in its functionality.

References

1. High-level Expert Group on Scientific Data. Riding the wave: How Europe can gain from the rising tide of scientific data. Submission to the European Commission, available online at <http://cordis.europa.eu/fp7/ict/e-infrastructure/docs/hlg-sdi-report.pdf>, 2010
2. Doi K: Computer-aided diagnosis in medical imaging: historical review, current status and future potential. *Comput Med Imaging Graph* 31:198–211, 2007
3. Müller H, Michoux N, Bandon D, Geissbuhler A: A review of content-based image retrieval systems in medicine—clinical benefits and future directions. *Int J Med Inform* 73:1–23, 2004
4. Markonis D, Holzer M, Dungs S, Vargas A, Langs G, Kriewel S, et al: A survey on visual information search behavior and requirements of radiologists. *Methods Inf Med* 51:539–548, 2012
5. Markonis D, Baroz F, de Castaneda RL R, Boyer C, Müller H: User tests for assessing a medical image retrieval system: a pilot study. *Stud Health Technol Inf* 192:224–228, 2013
6. Jansen BJ, Spink A, Taksai I. Handbook of research on web log analysis. IGI Global, 2009
7. Tsikrika T, Müller H, Kahn Jr, CE: Log analysis to understand medical professionals’ image searching behaviour. *Stud Health Technol Inf* 180:1020–1024, 2012
8. Yom-Tov E, White RW, Horvitz E: Seeking insights about cycling mood disorders via anonymized search logs. *J Med Internet Res* 16:e65, 2014
9. Müller H, Boyer C, Gaudinat A, Hersch W, Geissbuhler A: Analyzing web log files of the health on the net HONmedia search engine to define typical image search tasks for image retrieval evaluation. *Stud Health Technol Inf* 129(Pt 2):1319–1323, 2007
10. Müller H, Kalpathy-Cramer J, Hersch W, Geissbuhler A: Using Medline queries to generate image retrieval tasks for benchmarking. *Stud Health Technol Inf* 136:523–528, 2008
11. Herskovic JR, Tanaka LY, Hersch W, Bernstam EV: A day in the life of PubMed: analysis of a typical day’s query log. *J Am Med Inform Assoc* 14:212–220, 2007
12. Islamaj Dogan RI, Murray GC, Névél A, Lu Z. Understanding PubMed user search behavior through log analysis. *Database (Oxford)* 2009:bap018, 2009
13. Rubin DL, Flanders A, Kim W, Siddiqui KM, Kahn Jr, CE: Ontology-assisted analysis of web queries to determine the knowledge radiologists seek. *J Digit Imaging* 24:160–164, 2011
14. Palotti J, Hanbury A, Müller H. Exploiting Health Related Features to Infer User Expertise in the Medical Domain. Web Search Click Data workshop at WSCM, New York City, NY, USA, 2014.
15. Ruch P: Automatic assignment of biomedical categories: toward a generic approach. *Bioinformatics* 22:658–664, 2006
16. Kahn Jr, CE, Thao C: GoldMiner: a radiology image search engine. *AJR Am J Roentgenol* 188:1475–1478, 2008
17. Silverstein C, Marais H, Henzinger M, Moricz M: Analysis of a very large web search engine query log. *SIGIR Forum* 33(1):6–12, 1999
18. Jones R, Klinkner KL. Beyond the session timeout: automatic hierarchical segmentation of search topics in query logs. In:

- Proceedings of the 17th ACM Conference on Information and Knowledge Management. ACM, 2008. p. 699–708
19. Langlotz CP: RadLex: a new method for indexing online educational materials. *RadioGraphics* 26:1595–1597, 2006
20. Rubin DL: Creating and curating a terminology for radiology: ontology modeling and analysis. *J Digit Imaging* 21:355–362, 2008
21. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP: SMOTE: Synthetic Minority Over-sampling Technique. *J Artif Intell Res* 16:321–357, 2002
22. Chang CC, Lin CJ. LIBSVM: a library for support vector machines, 2001
23. Le Cessie S, Van Houwelingen J. Ridge estimators in logistic regression. *Applied Statistics*. 1992; p. 191–201.
24. Breiman L: Random forests. *Mach Learn* 45:5–32, 2001
25. Viera AJ, Garrett JM: Understanding interobserver agreement: the kappa statistic. *Fam Med* 37:360–363, 2005
26. Manning CD, Raghavan P, Schütze H. *Introduction to Information Retrieval*. Cambridge University Press, 2008
27. Hall MA, Holmes G: Benchmarking attribute selection techniques for discrete class data mining. *IEEE Trans Knowl Data Eng* 15: 1437–1447, 2003
28. Hollink V, Tsikrika T, de Vries AP: Semantic search log analysis: a method and a study on professional image search. *J Am Soc Inf Sci Technol* 62:691–713, 2011
29. Goeuriot L, Kelly L, Li W, Palotti J, Pecina P, Zuccon G, et al. ShARE/CLEF eHealth Evaluation Lab 2014, Task 3: User-centred health information retrieval CLEF eHealth overview. In: *CLEF Proceedings*. Springer LNCS, 2014
30. Seco de Herrera AG, Kalpathy-Cramer J, Demner Fushman D, Antani S, Müller H. Overview of the ImageCLEF 2013 medical tasks, CLEF working notes 2013, Valencia, Spain, 2013